

Analysis of Utterance Embeddings and Clustering Methods Related to Intent Induction for Task-Oriented Dialogue

Jeiyoon Park^{1,3}, Yoonna Jang¹, Chanhee Lee², Heuseok Lim¹

¹Department of Computer Science and Engineering, Korea University

²Naver Corporation

³LLSOLLU

{k4ke, morelychee, limhseok}@korea.ac.kr, chanhee.lee@navercorp.com

Abstract

The focus of this work is to investigate unsupervised approaches to overcome quintessential challenges in designing task-oriented dialog schema: assigning intent labels to each dialog turn (*intent clustering*) and generating a set of intents based on the intent clustering methods (*intent induction*). We postulate there are two salient factors for automatic induction of intents: (1) clustering algorithm for intent labeling and (2) user utterance embedding space. We compare existing off-the-shelf clustering models and embeddings based on DSTC11 evaluation. Our extensive experiments demonstrate that the combined selection of utterance embedding and clustering method in the intent induction task should be carefully considered. We also present that pretrained MiniLM with Agglomerative clustering shows significant improvement in NMI, ARI, F1, accuracy and example coverage in intent induction tasks. The source codes are available at <https://github.com/Jeiyoon/dstc11-track2>.

1 Introduction

Why Intent Induction? We humans are generalists. During a conversation, we listen to the other person’s utterance and naturally grasp which intent of the utterance is. With the skyrocketing demand for conversational AI, however, the more user utterances a dialogue system encounters, the more unknown intents it does. Predefining user intent is expensive and it is impossible to annotate all the user intents.

Since provided user utterances are unlabeled, intent induction (Haponchyk et al., 2018; Perkins and Yang, 2019; Chatterjee and Sengupta, 2020; Zeng et al., 2021; Zhang et al., 2022) focuses on discovering user intents from user utterances. However, previous studies did not conduct an in-depth analysis of the application of existing models to intent induction that might cause performance degradation problems (Zeng et al., 2021).

Agent: Hello, you are currently speaking with Rivetown Insurance customer service. My name is Julian, How may I be of service to you?

Customer: The services I have been receiving from your company has been encouraging and my intent is to increase or enroll for more plans with you.

Agent: Whoa, that is such an encouraging word coming from you. So, which of the plans do you wish to register for?

Customer: I would like to enroll for the life insurance policy.

`dialog act: ['InformIntent'], intent: ['EnrollPlan']`

Agent: That won't be an issue, I can help you get registered right away.

Customer: Ok then, but first can I get my current policy number.

`dialog act: ['InformIntent'], intent: ['GetPolicyNumber']`

Agent: Of course, I can help you get that.

Figure 1: A sample segment of conversation transcript.

Dataset	Domain	#Intents	#Utterances
DSTC11 _{dev}	insurance	22	66,875
DSTC11 _{test}	insurance	22	913

Table 1: Statistics of development dataset.

Our Approach. Intuitively, for good intent induction, user utterances must be well represented in the embedding space, and good clustering algorithms must be employed to capture latent barycenters of user intent clusters to handle both predefined and unseen intents well.

In this paper, We postulate there are two salient factors for automatic induction of intents: (1) clustering algorithm for intent labeling (Cheung and Li, 2012; Hakkani-Tür et al., 2015; Padmasundari, 2018) and (2) user utterance embedding space (Wang et al., 2020; Song et al., 2020; Gao et al., 2021; Chuang et al., 2022; Nishikawa et al., 2022).

We analyze how the two key factors affect to user intent clustering and intent induction. Our extensive experiments with existing models demonstrate

that pretrained MiniLM with Agglomerative clustering shows significant improvement in NMI, ARI, F1, accuracy and example coverage.

2 Task 1: Intent Clustering

2.1 Task Description.

A set of conversation transcripts are given as Figure 1, intent clustering model aims to (i) generate intent labels and (ii) align each utterance annotated with dialog act (i.e., "InformIntent").

Dataset. We conduct experiments on DSTC11 development dataset. It consists of 948 human-to-human conversation transcript with speaker role, utterance, dialog act, and intent. Testset is composed of 913 customer utterances and corresponding user intents. The dataset statistics are summarized in Table 1.

Metrics. We follow the same experimental metrics employed in the DSTC11 task proposal: Normalized mutual information (NMI), Adjusted rand index (ARI), Accuracy (ACC), Precision, Recall, F1 score, Intent example coverage, and the number of clusters (#K).

NMI is used for measuring dependency between two different distributions:

$$NMI(X, Y) = \frac{\mathbb{I}(X; Y)}{\min(\mathbb{H}(X), \mathbb{H}(Y))} \quad (1)$$

, where $X = [X_1, \dots, X_r]$ denote clustered labels, $Y = [Y_1, \dots, Y_s]$ are reference labels, \mathbb{I} stands for mutual information, and \mathbb{H} is entropy.

ARI is a measure for computing similarities between clustered results and reference labels:

$$ARI(X, Y) = \frac{\sum_{ij} \binom{n_{ij}}{2} - [AB]/\binom{n}{2}}{\frac{1}{2}[A + B] - [AB]/\binom{n}{2}} \quad (2)$$

, where $A = \sum_i \binom{a_i}{2}$, $B = \sum_j \binom{b_j}{2}$, $n_{ij} = |X_i \cap Y_j|$, a_i is the number of clustered label X_i and b_j is the number of reference label Y_j .

Both Precision and Recall measure many-to-one alignments from clustered labels to reference labels. F1 score is a harmonic mean between precision and recall. After performing a many-to-one alignment, Intent example coverage is determined as percent of examples whose reference has a corresponding predicted result.

2.2 Methods

Clustering Algorithm. Intent clustering focuses on assigning dialog intents to each dialog. K -means clustering (Lloyd, 1982), for example, regards a set of conversation transcript $\mathcal{T}_1, \dots, \mathcal{T}_m$.

Each transcript \mathcal{T} consists of turn-level dialog acts, speaker's role, whether it is Agent or Customer, and dialog utterances $X_1, \dots, X_n \in \mathbb{R}^p$. K -means minimizes the following equation in the Euclidean embedding space:

$$\min_{\beta_1, \dots, \beta_K \in \mathbb{R}^d} \sum_{i=1}^n \min_{k \in [K]} \mu_{ik} \|X_i - \beta_k\|_2^2 \quad (3)$$

, where $\{\beta_k\}_{k=1}^K$ denotes centroid of each intent cluster, $[K] = 1, 2, \dots, K$, μ_{ik} is alignment factor. Equation 1 can be expressed as:

$$\min_{V_1, \dots, V_K} \left\{ \sum_{k=1}^K \sum_{i \in V_k} \mu_{ik} \|X_i - \beta_k\|_2^2 : \bigsqcup_{k=1}^K V_k = [n] \right\} \quad (4)$$

, where \bigsqcup stands for disjoint union, and $\{V_k\}_{k=1}^K$ is intent cluster, determined by Voronoi diagram in utterance embedding space.

Given initial intent centroids $\{\beta_k^1\}_{k=1}^K$, we assign each utterance embedding to its nearest centroid (a.k.a., Expectation step):

$$V_k^t = \{i \in [n] : \|X_i - \beta_k\|_2^2 \leq \|X_i - \beta_j\|_2^2, \} \quad (5)$$

, where $\forall j \in [K]$. Then, we update the location of centroids (a.k.a., Maximization step):

$$\beta_k^{t+1} = \frac{1}{|V_k^t|} \sum_{i \in V_k^t} X_i \quad (6)$$

We iterate equation 5 and equation 6 alternatively until equation 4 converges.

In this paper, we conduct experiments with K -means (Lloyd, 1982), BIRCH (Zhang et al., 1996), Agglomerative clustering (Steinbach et al., 2000), Spectral clustering (Yu and Shi, 2003), Bisecting K -means (Di and Gou, 2018), and Variational optimal transportation (VOT) (Mi et al., 2018).

Embeddings. User utterance should be represented in the embedding space which is able to capture universal and rich semantic information. MiniLM (Wang et al., 2020) is an effective task-agnostic distillation to compress transformer-based language models. MPNet (Song et al., 2020) proposes permuted language model for dependency among predicted tokens and makes the model to see a full sentence and auxiliary position information. SimCSE (Gao et al., 2021) leverages a simple contrastive learning framework. DiffCSE (Chuang

Clustering Method	K-means Clustering								
	Metric	NMI	ARI	ACC	Precision	Recall	F1	Example Coverage	#K
EASE ^m _{ROBERTA}	33.4	14.2	28.0	28.0	69.0	39.9	43.9	43.9	5
EASE _{BERT}	36.1	18.1	35.3	36.8	58.9	45.3	53.5	53.5	8
EASE ^m _{BERT}	38.4	10.9	23.6	42.7	24.2	30.9	87.6	87.6	44
EASE _{ROBERTA}	45.8	25.4	40.0	43.7	60.7	50.9	65.6	65.6	12
DiffCSE ^{trans} _{BERT}	43.8	15.5	29.4	49.5	30.6	37.8	90.1	90.1	40
DiffCSE ^{sts} _{BERT}	46.6	16.9	29.9	50.4	30.8	38.2	89.9	89.9	43
DiffCSE ^{sts} _{ROBERTA}	53.9	29.0	45.1	57.2	46.8	51.5	91.1	91.1	30
DiffCSE ^{trans} _{ROBERTA}	55.0	23.1	35.7	60.6	35.7	44.9	96.7	96.7	50
SimCSE ^u _{BERT}	31.7	13.3	27.9	27.9	<u>65.0</u>	39.0	43.9	43.9	5
SimCSE ^u _{BERT+}	47.0	25.7	38.4	46.9	46.9	46.9	76.8	76.8	19
SimCSE ^u _{ROBERTA}	51.2	29.0	44.4	49.7	61.6	55.0	70.5	70.5	14
SimCSE _{BERT+}	53.0	27.8	39.8	55.0	42.2	47.8	91.0	91.0	31
SimCSE _{BERT}	53.1	24.4	39.0	60.5	39.5	47.8	96.3	96.3	44
SimCSE ^u _{ROBERTA+}	53.2	25.9	42.2	56.8	43.8	49.5	91.7	91.7	32
SimCSE _{ROBERTA}	56.6	28.8	42.7	60.8	43.3	50.6	91.3	91.3	36
SimCSE _{ROBERTA+}	56.8	28.9	41.3	62.5	41.6	49.9	98.4	98.4	42
Glove ^{avg}	30.5	7.0	20.6	34.6	22.2	27.0	92.2	92.2	50
MPNet	59.3	32.3	46.1	66.0	47.1	54.9	96.5	96.5	42
MiniLM _{L6}	59.3	35.7	52.6	62.2	54.9	58.4	92.4	92.4	28
MiniLM _{MULTIQA}	<u>61.7</u>	<u>38.2</u>	55.1	<u>66.6</u>	55.4	<u>60.5</u>	<u>98.8</u>	<u>98.8</u>	30
MiniLM _{L12}	63.1	38.9	<u>54.9</u>	68.0	54.9	60.8	100.0	100.0	31

Table 2: Clustering results on DSTC11 dataset. We employ K-means clustering algorithm to all utterance embeddings. *m* denotes multilingual model, *u* stands for unsupervised model, and + means large model.

Method	DiffCSE	SimCSE	MPNet	MiniLM
# Param	250M	125M	110M	21.3M

Table 3: The number of parameters. Both DiffCSE and SimCSE denote RoBERTa_{base} model.

et al., 2022) learns the difference between the original sentence and a stochastically masked sentence. EASE (Nishikawa et al., 2022) exploits sentence embedding via contrastive learning between the original sentence and its related entities.

2.3 Result Analysis

To analyze the effects of both embedding and clustering algorithm, We first heuristically fix the clustering algorithm and find the most meaningful embedding. Then, we opt for the most suitable clustering method based on the embedding.

Analysis of Embeddings. We show the experimental results for analyzing the effect of embedding space in Table 2. We observe that EASE records poor performance in all metrics, followed by averaged Glove embedding which shows the worst result, though EASE is a large-scale model. Note that both entity-aware contrastive learning

and multilingual setting cause performance degradation in dialog intent clustering task.

The results for DiffCSE also show that unsupervised contrastive learning between original utterance and edited utterance exacerbates model performance in both STS and Trans tasks. Despite the model size being twice smaller than DiffCSE as shown in Table 3, SimCSE gets comparable scores to MPNet, with a similar model size. Note that MiniLM achieves remarkable performances in both *L12* setting¹ and *MULTIQA* setting² in all metrics. These results demonstrate that (i) performance increases as the number of parameter decreases, which means excessively large embedding model leads to performance degradation, and (ii) the use of a student network which is trained by the teacher’s self-attention distributions and guiding layer improves intent clustering performance.

Visualization. Figure 2 gives UMAP (McInnes et al., 2018) visualization of clustering results. Note that we employ UMAP, instead of t-SNE (van der Maaten and Hinton, 2008), because

¹<https://huggingface.co/sentence-transformers/all-MiniLM-L12-v2>

²<https://huggingface.co/sentence-transformers/multi-qa-MiniLM-L6-cos-v1>

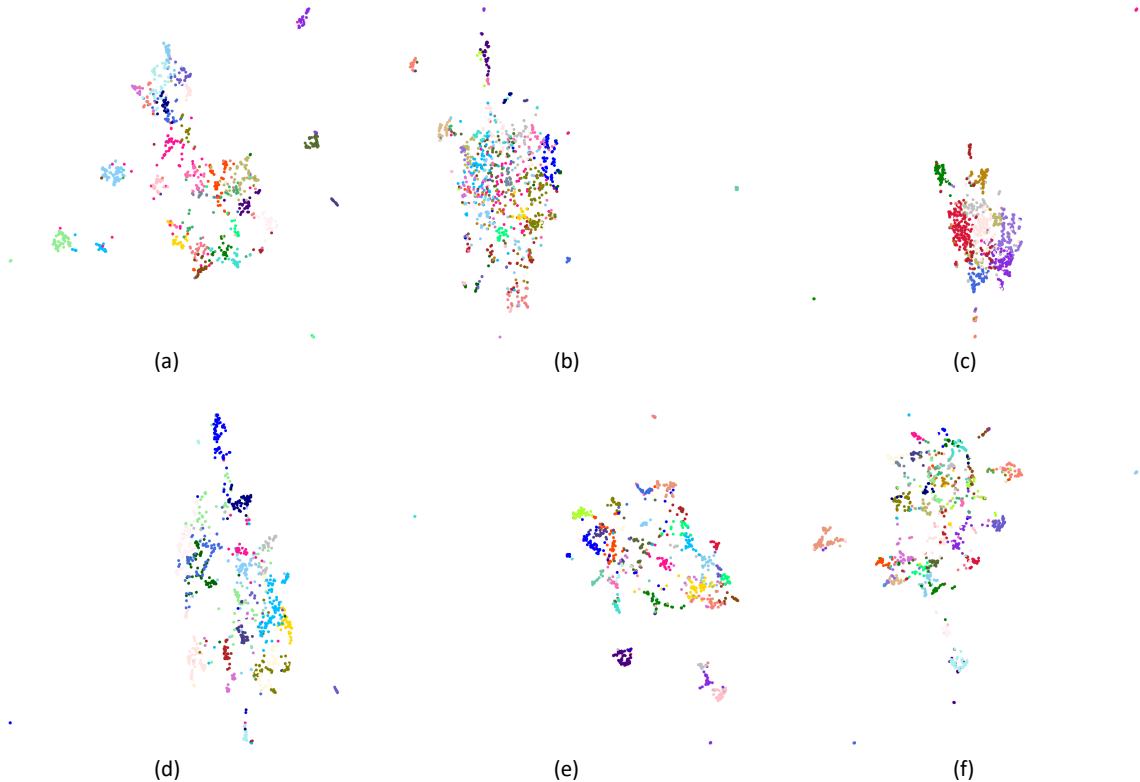


Figure 2: UMAP visualization of intent clustering results with different embeddings based on DSTC11 development dataset: (a) MiniLM, (b) Glove, (c) Ease-roberta, (d) DiffCSE-roberta, (e) SimCSE-roberta, and (f) MPNet. We apply K -means clustering algorithm to (a)-(e) experiments.

Utterance Embedding	MiniLM _{L12}							
	Metric	NMI	ARI	ACC	Precision	Recall	F1	Example Coverage
Bisect K -means	32.2	13.9	23.8	26.0	67.4	37.5	37.0	5
VOT	53.8	31.8	48.5	54.1	53.5	53.8	84.8	20
Spectral	57.6	<u>34.9</u>	51.3	58.8	<u>56.1</u>	<u>57.4</u>	82.9	24
Agglomerative	57.9	34.2	<u>51.8</u>	58.0	55.4	56.7	<u>88.6</u>	23
BIRCH	<u>59.9</u>	32.9	46.3	<u>64.6</u>	47.0	54.4	100.0	47
K -means	63.1	38.9	54.9	68.0	54.9	60.8	100.0	31

Table 4: Clustering results on DSTC11 dataset. We apply MiniLM_{L12} utterance embedding to all clustering methods.

UMAP preserves more global structure than t-SNE. We observe that (b) - (d) vertically degenerated into each embedding space. We also find that (a) presents the most well-clustered result and covers all intents (i.e., Example Coverage is 100.0) while (e) and (f) embeddings suffer from relatively ill-clustered result and outliers.

Analysis of Clustering Methods. In Table 4, we observe that K -means method outperforms the other models, followed by BIRCH. Note that VOT records poor performance which means optimal transportation with the variational principle deteriorates the result in intent clustering task. Bisect K -means clustering algorithm calculates the point density and average density of all points to initial-

ize cluster barycenters. However, K -means with K -means++ (Arthur and Vassilvitskii, 2007) initializer shows much better performance, contributing to the overall cluster inertia.

Voronoi Diagram. Figure 3 gives UMAP visualization of clustering results with Voronoi diagram. The number of barycenters of Bisect K -means is five which means this ill-clustering model is not able to cover all latent user intents. Though BIRCH records compatible results to K -means, the number of predicted barycenters is excessively large compared to the number of reference K . It causes outlier barycenter problem that barycenter contains a few utterances and disrupts universal clustering representation. It also demonstrates that hierarchi-

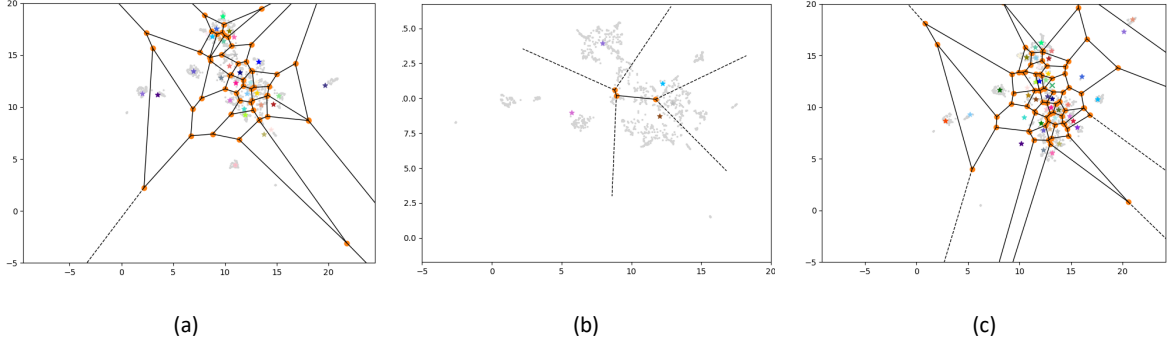


Figure 3: Voronoi diagram visualization of intent clustering results with different clustering algorithms based on DSTC11 development dataset: (a) K -means, (b) Bisect K -means, and (c) BIRCH. We apply MiniLM_{L12} user utterance embedding to (a)-(c) experiments. * stands for barycenter of each cluster and x represents outlier barycenter with fewer than five utterances.

Clustering Method	K -means Clustering							
	Metric	NMI	ARI	ACC	Precision	Recall	F1	Example Coverage
EASE _{BERT} ^m	27.0	12.3	23.9	24.1	72.9	36.2	26.4	5
EASE _{BERT}	53.1	25.6	42.8	50.2	57.4	53.5	80.0	26
EASE _{ROBERTA} ^m	59.6	41.8	53.7	57.8	68.8	62.8	86.0	34
EASE _{ROBERTA}	60.5	50.7	52.7	55.8	70.0	62.1	83.1	28
DiffCSE _{BERT} ^{trans}	21.0	11.5	23.3	23.7	<u>77.5</u>	36.3	32.5	6
DiffCSE _{BERT} ^{sts}	49.6	22.3	40.2	46.1	58.4	51.5	93.2	31
DiffCSE _{ROBERTA} ^{sts}	65.4	42.3	53.7	59.4	65.9	62.5	89.8	31
DiffCSE _{ROBERTA} ^{trans}	65.4	53.3	57.6	63.1	71.1	66.8	<u>90.0</u>	28
SimCSE _{BERT+} ^u	33.1	19.1	28.0	28.0	<u>77.5</u>	41.2	33.4	5
SimCSE _{BERT} ^u	58.3	29.6	46.4	58.5	59.1	58.8	83.1	32
SimCSE _{BERT+}	59.4	29.1	47.0	54.2	60.5	57.2	86.2	30
SimCSE _{ROBERTA}	61.4	34.3	47.1	54.3	62.0	57.9	80.3	29
SimCSE _{BERT}	62.7	37.3	52.6	63.4	59.8	61.6	89.6	34
SimCSE _{ROBERTA+} ^u	67.8	51.1	57.7	64.2	70.2	67.1	86.4	32
SimCSE _{ROBERTA} ^u	68.6	48.0	54.7	65.6	71.2	68.3	86.1	31
SimCSE _{ROBERTA+}	69.2	38.7	52.1	62.4	70.0	66.0	93.8	32
Glove ^{avg}	35.0	18.8	29.1	36.3	51.2	42.4	77.0	29
MPNet	72.8	41.6	60.2	65.6	76.3	70.6	86.5	26
MiniLM _{L12}	73.2	47.1	57.1	<u>66.0</u>	72.2	69.0	83.1	25
MiniLM _{L6}	74.7	<u>52.8</u>	<u>61.2</u>	70.5	74.9	<u>72.7</u>	83.4	25
MiniLM _{MULTIQA}	77.4	54.5	63.2	70.5	79.1	74.6	80.0	24

Table 5: Intent induction results on DSTC11 dataset. We employ K -means clustering algorithm to all utterance embeddings. m denotes multilingual model, u stands for unsupervised model, and $+$ means large model.

cal process of BIRCH including removing outliers and cluster refining does not have a substantial impact.

3 Task 2: Intent Induction

3.1 Task Description

Unlike Task 1, intent induction takes different transcripts which annotate only the speaker’s role and intent label. The goal of intent induction is to create a set of intents and match them to each utterance

without access to the ground-truth dialog acts. In this paper, we use DSTC11 dataset as shown in Table 1, exploit the same clustering algorithms and utterance embeddings, and evaluate methods using NMI, ARI, Accuracy, F1, and Example coverage. We employ the provided automatic dialog act predictions.

3.2 Result Analysis

Analysis of Embeddings. Table 5 demonstrates that MiniLM-based utterance embedding can im-

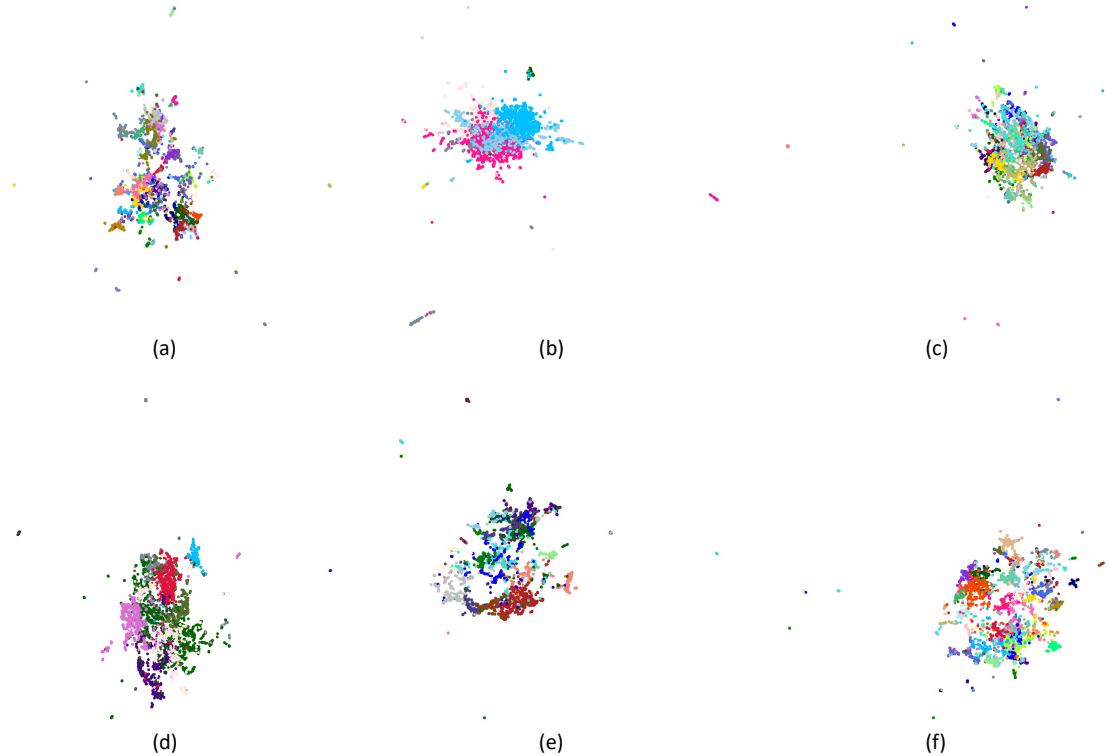


Figure 4: UMAP visualization of intent induction results with different embeddings based on DSTC11 development dataset: (a) MiniLM, (b) Glove, (c) Ease-roberta, (d) DiffCSE-roberta, (e) SimCSE-roberta, and (f) MPNet. We apply Agglomerative clustering algorithm to (a)-(e) experiments.

Utterance Embedding	MiniLM _{MULTIQA}							
Metric	NMI	ARI	ACC	Precision	Recall	F1	Example Coverage	#K
Bisect K -means	72.3	55.7	57.4	63.0	<u>81.2</u>	70.9	76.8	27
VOT	75.3	50.6	60.0	71.9	72.3	72.1	93.4	33
Spectral	75.0	51.7	59.0	67.7	75.9	71.6	86.3	23
K -means	77.4	54.5	63.2	70.5	79.1	74.6	80.0	24
BIRCH	<u>79.5</u>	<u>62.8</u>	68.1	<u>71.9</u>	84.6	<u>77.7</u>	<u>89.9</u>	24
Agglomerative	81.0	64.2	<u>66.5</u>	75.5	80.6	78.0	86.2	25

Table 6: Intent induction results on DSTC11 dataset. We apply MiniLM_{MULTIQA} utterance embedding to all clustering methods.

prove induction performance for user intents. Unlike Task 1, MiniLM_{MULTIQA} outperforms the other methods, followed by MiniLM_{L6}³. It represents that MiniLM with 6 layers and 384 hidden size can learn universal utterance representation and capture meaningful information. Though the performance tends to increase as the size of each embedding model increases, we observe that larger models do not always perform better. Note that DiffCSE_{BERT}^{trans} and EASE_{BERT}^m record high recall because the number of clusters K is disastrously low. We observe that it dampens overall model performance. As shown in both Task 1 and Task

³<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

2 results, (i) entity-aware contrastive learning, (ii) multilingual setting, (iii) contrastive learning between the original utterance and edited utterance are ultimately not the optimal choice to improve intent induction performance.

Visualization. We present UMAP visualization of intent induction results with different utterance embeddings (Figure 4). It corroborates that MiniLM-based intent induction can provide wider and well-separated results preserving its meaningful embedding space.

Analysis of Clustering Methods. In Table 6, we show the results of intent induction comparing models that clustering algorithms. Unlike the results in Task 1, Agglomerative clustering and BIRCH out-

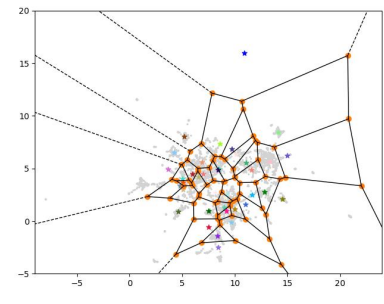
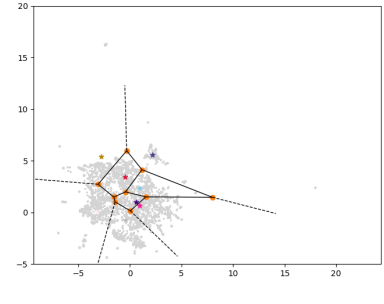
Method	Voronoi Diagram	Examples	Label
Agglomerative + MiniLM		<i>Would it be on like any of my other bills?</i> <i>Oh okay. Where at on the bill. I'm looking now.</i> <i>Okay. When will I see that in my bill?</i> <i>Hi I am calling because I received a wrong bill.</i> <i>Thank you. How much will this make my bill go up?</i> <i>Yeah, so exactly how does this work? Do I just send a bill or what?</i> <i>crap! Well what can I do now? Can I make the payment and have the same policy that I had before?</i>	34
Agglomerative + DiffCSE		<i>Would it be on like any of my other bills?</i> <i>Oh okay. Where at on the bill. I'm looking now.</i> <i>Okay. When will I see that in my bill?</i> <i>Hi I am calling because I received a wrong bill.</i> <i>Thank you. How much will this make my bill go up?</i> <i>Yeah, so exactly how does this work? Do I just send a bill or what?</i> <i>crap! Well what can I do now? Can I make the payment and have the same policy that I had before?</i>	5 2 7 6 0 1 4

Figure 5: Qualitative result on DSTC11 dataset comparing MiniLM-based and DiffCSE-based models. We fix the clustering model as Agglomerative clustering in this experiment.

perform the other baselines. Agglomerative clustering starts with the utterances as individual clusters and merges them if they have similar intents. The experiment shows that the approach exploiting MiniLM_{MULTIQA} and Agglomerative clustering achieves state-of-the-art results among the other methods in the intent induction task.

3.3 Qualitative Results

Figure 5 shows the qualitative result for intent induction on DSTC11 dataset. We generate user intent labels using two different models. We observe that Agglomerative clustering on MiniLM utterance embedding can classify bill-related utterances well. However, Agglomerative clustering on DiffCSE utterance embedding is not able to discern user intent well and amalgamate completely different user utterances into a single intent cluster. Note that all clusters have bill-related user utterances which means it is an ill-clustered result.

We conclude that user utterance embedding is one of the most important factors affecting performance, and we should add a huge caveat that selection of utterance embedding in intent induction task should be very careful. Further analyzing how to leverage an embedding model is an interesting direction for future work.

Our extensive experiments demonstrate that the

combined selection of utterance embedding and clustering method in the intent induction task should be carefully considered.

3.4 Quantitative Results

We conduct quantitative experiments to analyze (i) the correlation between the number of induced clusters and performance, and (ii) the correlation between the number of utterances per intent and performance (Figure 6). Note that the maximum number of sample utterances aligned to intent is limited to 50. First, We find that NMI and ARI increase as the number of induced clusters approaches the reference K. We also observe that Example coverage decreases as the number of induced clusters go toward the reference K. We demonstrate that there is a trade-off relationship between NMI, ARI and Example coverage. Second, the correlation between the number of utterances per intent and performance shows the same pattern as the correlation between the number of induced clusters and performance.

4 Conclusion

The conclusion of this paper is threefold:

- We empirically demonstrate that the combined selection of utterance embedding and clustering method in the intent induction task should

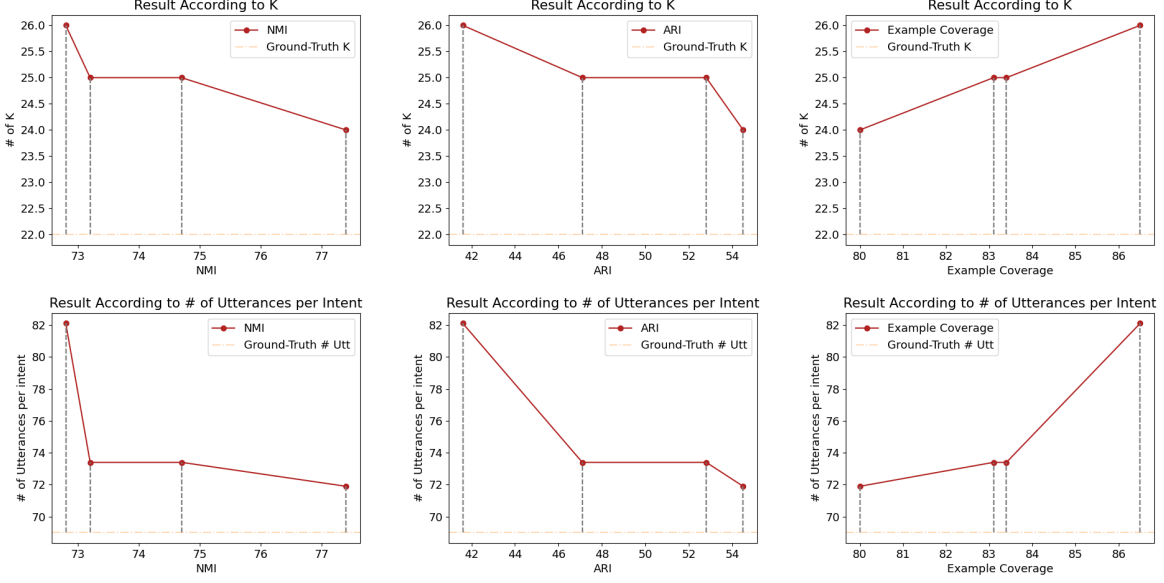


Figure 6: Quantitative results on DSTC11 dataset related to NMI, ARI, F1 score, and Example coverage. We analyze using four results with the best performance introduced in Table 5.

be carefully considered.

- We also present that pretrained MiniLM with Agglomerative clustering shows significant improvement in NMI, ARI, F1, accuracy and example coverage in intent induction tasks.
- We find that there is a trade-off relationship between NMI, ARI and Example coverage.

5 Discussion and Future Work

In clustering, broadly, there are three categories of methods: (i) Barycenter-based formulation, (ii) Density-based formulation, and (iii) Distance-based formulation. In this paper, we mainly dealt with barycenter-based methods. Indeed, K -means clustering method, for instance, theoretically produce the same result as barycenter-based formulation and distance-based formulation in Euclidean embedding space:

$$\sum_{i,j=1}^n \|X_i - X_j\|_2^2 = 2n \sum_{i=1}^n \|X_i - \beta\|_2^2 \quad (7)$$

, where $\beta = \frac{1}{n} \sum_{i=1}^n X_i$. Besides, barycenter-based K -means in Euclidean space can circumvent the problem frequently caused by the location of barycenter in different measure space (e.g., Wasserstein space (Zhuang et al., 2022)):

$$\begin{aligned} \sum_{i=1}^n \|X - X_i\|_2^2 &= n\|X - \beta\|_2^2 + \sum_{i=1}^n \|X - X_i\|_2^2 \\ &\geq n\|X - \beta\|_2^2 \end{aligned} \quad (8)$$

However, K -means clustering in Euclidean space often loses salient geometric information of dataset, due to its formulation. In the same vein, though we didn't add the experimental result of Density-based formulation to the table, the performance was disastrous. DBSCAN (Ester et al., 1996), for example, recorded 7.7 NMI, 0.0 ARI, 16.1 Accuracy, 27.6 F1 score, and 41.0 Example coverage in intent clustering task. Therefore, both clustering in different measure spaces and clustering using embedding density should be investigated. Further analyzing how to leverage an embedding model is also an interesting direction for future work.

Acknowledgements

This work was supported by Institute of Information & communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No. 2020-0-00368, A Neural-Symbolic Model for Knowledge Acquisition and Inference Techniques). This research was supported by the MSIT(Ministry of Science and ICT), Korea, under the ITRC(Information Technology Research Center) support program(IITP-2022-2018-0-01405) supervised by the IITP(Institute for Information & Communications Technology Planning & Evaluation).

References

David Arthur and Sergei Vassilvskii. 2007. K -means++: The advantages of careful seeding. In

- Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '07*, page 1027–1035, USA. Society for Industrial and Applied Mathematics.
- Ajay Chatterjee and Shubhashis Sengupta. 2020. [Intent mining from past conversations for conversational agent](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4140–4152, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Jackie Chi Kit Cheung and Xiao Li. 2012. [Sequence clustering and labeling for unsupervised query intent discovery](#). In *Proceedings of the Fifth International Conference on Web Search and Web Data Mining, WSDM 2012, Seattle, WA, USA, February 8-12, 2012*, pages 383–392. ACM.
- Yung-Sung Chuang, Rumen Dangovski, Hongyin Luo, Yang Zhang, Shiyu Chang, Marin Soljagic, Shang-Wen Li, Scott Yih, Yoon Kim, and James Glass. 2022. [DiffCSE: Difference-based contrastive learning for sentence embeddings](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4207–4218, Seattle, United States. Association for Computational Linguistics.
- Jian Di and Xinyue Gou. 2018. Bisecting k-means algorithm based on k-valued selfdetermining and clustering center optimization. *J. Comput.*, 13(6):588–595.
- Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, KDD'96*, page 226–231. AAAI Press.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. [SimCSE: Simple contrastive learning of sentence embeddings](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Dilek Hakkani-Tür, Yun-Cheng Ju, Geoffrey Zweig, and Gokhan Tur. 2015. Clustering novel intents in a conversational interaction system with semantic parsing. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- Iryna Haponchyk, Antonio Uva, Seunghak Yu, Olga Uryupina, and Alessandro Moschitti. 2018. Supervised clustering of questions into intents for dialog system applications. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2310–2321.
- Stuart P. Lloyd. 1982. [Least squares quantization in pcm](#). *IEEE Trans. Inf. Theory*, 28(2):129–136.
- L. McInnes, J. Healy, and J. Melville. 2018. [UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction](#). *ArXiv e-prints*.
- Liang Mi, Wen Zhang, Xianfeng Gu, and Yalin Wang. 2018. Variational Wasserstein clustering. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 322–337.
- Sosuke Nishikawa, Ryokan Ri, Ikuya Yamada, Yoshimasa Tsuruoka, and Isao Echizen. 2022. [EASE: Entity-aware contrastive learning of sentence embedding](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3870–3885, Seattle, United States. Association for Computational Linguistics.
- Srinivas Bangalore Padmasundari. 2018. Intent discovery through unsupervised semantic text clustering. In *Proc. Interspeech*, volume 2018, pages 606–610.
- Hugh Perkins and Yi Yang. 2019. Dialog intent induction with deep multi-view clustering. In *EMNLP*.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tiejun Liu. 2020. [MpNet: Masked and permuted pre-training for language understanding](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 16857–16867. Curran Associates, Inc.
- M. Steinbach, G. Karypis, and V. Kumar. 2000. [A comparison of document clustering techniques](#). In *KDD Workshop on Text Mining*.
- Laurens van der Maaten and Geoffrey E. Hinton. 2008. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9:2579–2605.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. [Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 5776–5788. Curran Associates, Inc.
- Stella X. Yu and Jianbo Shi. 2003. [Multiclass spectral clustering](#). In *ICCV*, pages 313–319. IEEE Computer Society.
- Zengfeng Zeng, Dan Ma, Haiqin Yang, Zhen Gou, and Jianping Shen. 2021. [Automatic intent-slot induction for dialogue systems](#). In *Proceedings of the Web Conference 2021, WWW '21*, page 2578–2589, New York, NY, USA. Association for Computing Machinery.
- Tian Zhang, Raghu Ramakrishnan, and Miron Livny. 1996. [Birch: an efficient data clustering method for very large databases](#). In *SIGMOD '96: Proceedings of the 1996 ACM SIGMOD international conference on Management of data*, pages 103–114, New York, NY, USA. ACM.

Yuwei Zhang, Haode Zhang, Li-Ming Zhan, Xiao-Ming Wu, and Albert Lam. 2022. [New intent discovery with pre-training and contrastive learning](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 256–269, Dublin, Ireland. Association for Computational Linguistics.

Yubo Zhuang, Xiaohui Chen, and Yun Yang. 2022. Wasserstein k -means for clustering probability distributions. *Thirty-sixth Conference on Neural Information Processing Systems*.