# Similarity-Based Intent Detection Using an Enhanced Siamese Network

Anonymous Full Paper Submission 12

#### Abstract 001

In Natural Language Understanding (NLU), intent 002 detection is crucial for improving human-computer 003 interaction. However, traditional supervised learn-004 ing models rely heavily on large annotated datasets, 005 006 limiting their effectiveness in low-resource scenarios with limited labeled data. Siamese networks, which 007 are effective at learning similarity-based representa-008 tions, provide a promising alternative by enabling 009 few-shot learning. However, Siamese networks typi-010 cally rely on contrastive loss or triplet loss, both of 011 which introduce challenges. This study introduces 012 a similarity-based intent detection model using an 013 enhanced Siamese network to address these limita-014 tions. Our model employs Manhattan, Euclidean, 015 and Cosine similarity metrics combined with a fu-016 sion layer to improve intent classification accuracy. 017 We evaluated the model on the Airline Travel In-018 formation System (ATIS) and SNIPS datasets and 019 demonstrated its superiority over state-of-the-art 020 methods, particularly in low-resource and few-shot 021 learning scenarios. The results highlight significant 022 accuracy gains while maintaining computational ef-023 ficiency, making it a robust solution for real-world 024 025 dialog systems.

#### Introduction 1 026

In today's interconnected world, dialog systems such 027 as chatbots play a crucial role in facilitating human-028 machine interaction across applications such as cus-029 tomer service and digital assistants. Central to these 030 systems is intent detection, which classifies users' ut-031 terances into predefined classes. A key challenge in 032 intent detection is its reliance on labeled data, mak-033 ing data acquisition and annotation labor-intensive, 034 time-consuming, and costly. Moreover, models of-035 ten struggle with out-of-domain or unseen intents, 036 reducing their real-world adaptability. 037

Traditional approaches to intent detection, such 038 as rule-based methods and shallow machine learning, 039 struggle with the complexities of natural language. 040 The rise of deep learning models has significantly 041 improved performance under the supervised learning 042 paradigm. However, these approaches require large 043 annotated datasets for each intent class, which limits 044 their adaptability to new labels. 045

Siamese networks are designed to compute the 046

similarity or dissimilarity between pairs of inputs, 047 making them well suited for tasks where a model 048 must identify novel classes based on minimal labeled 049 examples [1]. These architectures have shown ef-050 fectiveness in various domains, including computer 051 vision [2], text similarity [1], and domain represen-052 tation [3]. However, Siamese networks are often 053 trained with contrastive loss [4], or triplet loss [5]. 054 Both approaches aim to minimize the distance be-055 tween similar pairs while maximizing the distance 056 between dissimilar pairs. Triplet loss, in particular, 057 requires the careful selection of triplets, consisting 058 of an anchor, a positive example (similar to the 059 anchor), and a negative example (dissimilar to the 060 anchor). This selection process can be computa-061 tionally expensive and is sensitive to the choice of 062 margin parameter. 063

To address these challenges, this study pro-064 poses an enhanced intent detection model using 065 a similarity-based Siamese network with multiple 066 distance layers and a fusion layer. The fusion layer 067 improves the similarity measures between sentences 068 with the same intent and helps to better separate 069 dissimilar intents using a binary classification ap-070 proach. This eliminates the need for triplet loss, 071 reducing model complexity while maintaining the 072 model's ability to differentiate between similar and 073 dissimilar intents. 074 075

Our main contributions are as follows:

- We learn sentence representations using an en-076 coder in a Siamese network which serves as a 077 feature extractor. 078
- We demonstrate the efficacy of using multiple 079 distance layers combined with a fusion layer 080 through ablation studies. 081
- In comparison to benchmark approaches, the 082 proposed model has shown state-of-the-art 083 (SOTA) performance, with accuracy well above 084 99%. 085

086

#### 2 **Related Work**

Early intent detection relied on rule-based systems 087 and statistical models, such as HMM, SVM, and 088 Naïve Bayes, which struggled with scalability, seman-089 tic nuances, and context. [6] linked seen and unseen 090 intents using manual attributes such as "action," 091 NLDL

#12

"object," "location," and "time." Meanwhile, [7] and
improved intent detection using prosodic cues
and n-grams, although extensive feature engineering
was still required [9].

The introduction of deep learning revolutionized 096 intent detection by enabling neural networks to au-097 tomatically learn data patterns and representations. 098 [10] applied a modified RNN with pre-trained embed-099 dings for dialog act classification. CNNs, introduced 100 by [11] for encoding, struggled with long-range de-101 pendencies, which [12] addressed through dual fea-102 ture fusion with capsule networks. Similarly, [13] 103 proposed a Bi-model based RNN for joint intent 104 detection and slot filling, achieving state-of-the-art 105 results on benchmark datasets. 106

Few-shot learning approaches have further ad-107 vanced the field by allowing rapid adaptation to 108 unseen intents with minimal labeled data. For in-109 stance, [14] leveraged pretrained models for few-shot 110 intent detection, overcoming large-scale data depen-111 dency challenges. Building on this, [15] introduced 112 self-supervised pretraining with prototype-aware at-113 tention for few-shot intent detection, which improved 114 performance in scenarios with limited labeled data. 115 Siamese networks, widely successful in computer 116 vision [2], have been adapted for few-shot learning in 117 text classification. [1] applied Siamese networks to 118 measure text similarity in ambiguous domains, while 119 [3] extended them to learn semantic relationships 120 with triplet loss in domain-specific contexts. Despite 121 their success, Siamese networks with triplet loss are 122 computationally expensive and sensitive to margin 123 selection [16, 17]. 124

While these approaches have significantly ad-125 vanced intent detection, they often require extensive 126 preprocessing and computational resources. Our ap-127 proach builds on these advancements by integrating 128 combined distance metrics for similarity checking, 129 balancing efficiency and accuracy. This method by-130 passes large-scale triplet selection and margin tuning, 131 making it more robust to noisy data and better at 132 generalizing across domains. 133

## 134 3 Methodology

Figure 1 shows the architecture of the proposed similarity-based intent detection model using a Siamese Neural Network. The model consists of two identical subnetworks, a distance layer, a metric transformation layer, a fusion layer, and dense layers. The details of the model are presented in the following subsections.

#### 142 3.1 Data Preprocessing

In the proposed model, the dataset was preprocessed
by first creating pairs of texts and their corresponding intents. For each pair, a label of 1 was assigned if
the intent matched the actual intent of the text, and



Figure 1. Proposed similarity based Siamese architecture

0 otherwise. Both positive pairs(matching intents) 147 and negative pairs (non-matching intents) were used 148 to help the model learn to distinguish between simi-149 lar and dissimilar intents. The text and intent were 150 then tokenized using the Keras Tokenizer, which 151 converted each word into a unique integer based on 152 its frequency in the entire dataset. This process 153 involved fitting the tokenizer on both the text and 154 intent columns to build a shared vocabulary. 155

After tokenization, each text and intent pair was transformed into numerical sequences. To ensure uniform input dimensions for the model, the sequences were padded to a fixed length. Specifically, each sequence was padded to a maximum length of 45 tokens for the ATIS dataset and 36 tokens for the SNIPS dataset for efficient batch processing. 162

#### 3.2 Embedding Layer

An embedding layer is created to convert the in-164 put sequences into dense vector representations. 165 The embedding weights were initialized with pre-166 trained Word2Vec embeddings trained on 100 bil-167 lion words from Google News [18]. The pretrained 168 embeddings help capture the semantic relationship 169 between words, where words with similar mean-170 ings have similar vector representations, aiding in 171 better text processing. The embedding layer is 172 shared between both input sequences to ensure iden-173 tical embeddings for the same words, regardless of 174 their position in the pair. Each word was repre-175 sented by 300-dimensional vectors, with padding 176 applied to maintain a uniform sentence length. 177 Formally, for an utterance of length T, the  $i^{th}$ 178 word is mapped to *d*-dimensional embedding. Let 179  $X_1 = [x_{1,1}, x_{1,2}, \dots, x_{1,T}]$  be the first input sequence 180 of length T and  $X_2 = [x_{2,1}, x_{2,2}, \dots, x_{2,T}]$  be the sec-181 ond input sequence of length T, both mapped to 182 *d*-dimensional embeddings. The embedding matrices 183  $E_1$  and  $E_2$  can be expressed as: 184

$$E_1 = [E(x_{1,1}), E(x_{1,2}), \dots, E(x_{1,T})]$$
 (1) 185

$$E_2 = [E(x_{2,1}), E(x_{2,2}), \dots, E(x_{2,T})]$$
(2) 186

163

(3)

187 Where  $E_1, E_2 \in \mathbb{R}^{T \times d}$ .

#### 188 3.3 Siamese Layer

The proposed model employs Bidirectional long 189 short-term memory (BiLSTM) as the subnetwork 190 of the Siamese layer. BiLSTMs are particularly ef-191 fective for processing sequences of data [19]. Unlike 192 traditional feedforward neural networks, BiLSTMs 193 have connections that form directed cycles, enabling 194 them to maintain a hidden state that captures in-195 formation from the previous steps in the sequence. 196 This makes BiLSTMs particularly effective for tasks 197 in which the order and context of input data are 198 crucial, such as language modeling [20], time series 199 prediction [21], and sequence-to-sequence tasks [22]. 200 In the proposed model, the embedded sequences 201  $E_1$  and  $E_2$  are passed through a BiLSTM, which 202 processes the sequence step by step and generates a 203 hidden state at each step. The recurrence relations 204 for the hidden states at time step t for the first and 205 second sequences can be expressed as: 206

207 
$$h_t^{(1)} = \sigma(W_h \cdot h_{t-1}^{(1)} + U \cdot E(x_{1,t}) + b_h)$$

208 
$$h_t^{(2)} = \sigma(W_h \cdot h_{t-1}^{(2)} + U \cdot E(x_{2,t}) + b_h)$$
 (4)

where  $h_t^{(1)}$  and  $h_t^{(2)}$  are the hidden states of BiL-STM at time step t for the first and second sequences, respectively.  $\sigma$  represents the activation function,  $W_h$  and U are weight matrices, and  $b_h$  is the bias term. The final hidden states,  $h_T^{(1)}$  and  $h_T^{(2)}$ , represent the encoded information for the complete sequences  $X_1$  and  $X_2$ .

### 216 3.4 Distance Layer

To compute the similarity between encoded inputs 217  $h_T^{(1)}$  and  $h_T^{(2)}$ , a distance layer with three different 218 metrics: Euclidean distance, Cosine similarity, and 219 Manhattan distance. These metrics were chosen 220 carefully for their complementary strengths. Eu-221 clidean distance provides overall similarities, while 222 Cosine similarity focuses on the orientation of vec-223 tors [23]. Manhattan distance captures absolute 224 differences and is efficient for high dimensional data 225 [24]. This combination enhances the model's ability 226 to generalize across varied intents and domains. 227

228 
$$D_{\text{Euclidean}} = \sqrt{\sum_{i=1}^{d} \left(h_T^{(1)}[i] - h_T^{(2)}[i]\right)^2} \quad (5)$$

229 
$$D_{\text{Manhattan}} = \sum_{i=1}^{d} \left| h_T^{(1)}[i] - h_T^{(2)}[i] \right|$$
(6)

$$D_{\text{Cosine}} = 1 - \frac{h_T^{(1)} \cdot h_T^{(2)}}{\left\| h_T^{(1)} \right\| \left\| h_T^{(2)} \right\|}$$
(7) 230

Where  $D_{\text{Euclidean}}$ ,  $D_{\text{Manhattan}}$ , and  $D_{\text{Cosine}}$  represent the Euclidean, Manhattan, and cosine similarity distance metrics, respectively. 233

A logarithmic function was used to scale the distances for better learning and generalization [25]. 235 The transformed distance values are then concatenated into a feature vector that captures multiple types of similarities between the two input sequences. 239

#### 3.5 Dense Layer 240

The feature vector is passed through dense layers 241 to refine the similarity score. The first dense layer 242 was activated with ReLU function to learn complex 243 patterns from the concatenated distance metrics 244 This layer ensures that the different distance metrics 245 are jointly processed, enhancing the model's ability 246 to capture nuanced relationships. 247

Another dense layer with a sigmoid activation 248 function was added to output the probability scores. 249

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$
 (8) 250

261

269

270

where z is the output from the final layer before 251 the activation function. 252

This score represents the likelihood that the input 253 sequences corresponds to same intent. 254

For training, we used binary cross-entropy loss 255 as the objective function, given that the task is 256 framed as a binary classification problem: determining whether two input sequences correspond to the 259 same intent. The binary cross-entropy loss is defined 259 as: 260

Loss = 
$$-\frac{1}{N} \sum_{i=1}^{N} (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$
(9)

where N is the total number of samples,  $y_i$  is the 262 ground truth label (0 or 1), and  $\hat{y}_i$  is the predicted 263 probability for sample *i*. 264

This loss function ensures that the model effectively learns to distinguish between similar and dissimilar intent pairs by minimizing the error in predicting similarity probabilities.

### 4 Experimental Study

### 4.1 Dataset

To verify the validity of our proposed model, we 271 conducted experiments on two widely used NLU 272 datasets, ATIS [26] and SNIPS [27], chosen for their 273 complementary characteristics: ATIS, with 4978 274

327

training samples, 21 intents, and 128 slots [28], offers 275 a domain-specific challenge focused on flight-related 276 information, while SNIPS, with 15884 utterances, 277 7 intents, and 72 slots, provides a broader, more 278 diverse set of challenges across various domains such 279 as weather and entertainment, offering a comprehen-280 sive test bed for evaluating our model's effectiveness 281 and generalizability. 282

#### 283 4.2 Experimental Setup

In the proposed model, We set the BiLSTM encoders 284 to 128 hidden units and used the ReLU activation 285 function. To further improve training, we apply a 286 recurrent dropout and regular dropout rate of 0.5 287 to randomly drop some units. The batch size was 288 289 set to 32, and only ten epochs were used to train the model due to computational cost constraints. 290 A learning rate of 0.001 was used with the Adam 291 optimizer [29]. For the dense layers, we set the unit 292 to 32 and 1, respectively. 293

For both the ATIS and SNIPS datasets, the 294 train\_test\_split function from scikit-learn library 295 was used to split the data into an 80% training set 296 and a 20% test set, following the practice used in 297 previous studies [28]. To ensure robust result, the ex-298 periments were conducted with five different random 299 seeds, and the average performance across these runs 300 was reported. Accuracy was used as the evaluation 301 metric, as it is the most widely adopted metric in 302 existing models [30], [5], [31]. 303

The model was implemented using the TensorFlow framework and trained on a machine equipped with an Intel Core i7 processor and 16.0 GB of RAM. Due to these hardware limitations, we focused on optimizing model performance within the constraints of the available computational resources.

#### 310 4.3 Comparative Method

To further demonstrate the efficiency of the proposed model, we identified the best-performing settings of our model and subsequently compared it with the following baseline models:

- C2A-SLU [32]: This uses a contrastive attention mechanism to compare input sets and extract features for intent detection.
- LIDSNet [30]: A Siamese model with triplet loss was to reduce the distance between anchor and positive examples relative to negative examples.
- **BERT+PSN** [33]: Proposes a pseudo Siamese Network for intent detection using BERT encoders.
- SN-TripletLoss [5]: Proposes a Siamese network with a triplet training framework.

## 5 Results and Discussion

The performance results of the proposed Siamese 328 network across the ATIS and SNIPS datasets, as 329 presented in Table 1, demonstrate the significant im-330 pact of distance metric selection on intent detection 331 task. The choice of distance metrics plays a role in 332 determining the model's ability to generalize across 333 datasets, particularly when the datasets vary in lin-334 guistic diversity and domain specificity. A detailed 335 analysis of these results reveals key patterns regard-336 ing how individual and combined metrics influence 337 model performance, providing valuable insights into 338 the effectiveness of the Siamese network for intent 339 detection. 340

The models employing individual distance met-341 rics: Manhattan\_distance, Cosine\_similarity, 342 and Euclidean\_distance, display notable varia-343 tions in performance across the two datasets. On 344 ATIS dataset, Manhattan and Euclidean metrics 345 achieve relatively high accuracy with score of 95.41% 346 and 95.42%, respectively. This can be attributed to 347 the structured and domain-specific nature of ATIS 348 queries, which often involve repetitive patterns and 349 similar syntax. Metrics like Manhattan and Eu-350 clidean are particularly well-suited for this scenario 351 because they measure numerical differences and ge-352 ometric distances between vector representations. 353 These metrics help the model capture small differ-354 ences in query formulations, which is essential for 355 distinguishing intents that are syntactically close but 356 semantically distinct, such as flight queries differing 357 by destination or departure time. 358

However, the same metrics do not perform as well 359 on the SNIPS dataset, with accuracies dropping 360 to 86.10% (Manhattan) and 85.23% (Euclidean). 361 SNIPS contains more diverse and varied language in-362 puts, covering multiple domains and informal phras-363 ing patterns. This variability makes Manhattan and 364 Euclidean metrics less effective, as they struggle to 365 handle the semantic richness and lexical variability 366 present in the dataset. In this context, the reliance 367 on numerical distance alone limits the model's ability 368 to capture the underlying meaning of the queries. 369

On SNIPS, Cosine similarity performs slightly bet-370 ter than Manhattan and Euclidean metrics, achiev-371 ing an accuracy of 85.75%. Cosine similarity fo-372 cuses on the angular relationship between vectorized 373 inputs, disregarding magnitude differences. This 374 makes makes it more suitable for datasets like SNIPS, 375 where different word choices and query lengths might 376 convey the same intent. Cosine similarity helps the 377 model recognize semantic alignment between differ-378 ent formulations of the same query, even when the 379 exact phrasing differs significantly. 380

The result show a substantial improvement <sup>381</sup> in performance when multiple metrics are combined, highlighting the benefits of metric fusion. <sup>383</sup>

Model4 (Euclidean + Cosine) achieves the high-384 est accuracy ob both datasets, with 99.81% on ATIS 385 and 99.67% on SNIPS. The success of this combi-386 nation suggests that Euclidean distance and Cosine 387 similarity play complementary roles in capturing 388 intent. Euclidean distance provides a measure of 389 positional and magnitude-based differences, which is 390 useful for distinguishing between queries with over-391 lapping terms but different meanings. On the other 392 hand, Cosine similarity captures the semantic re-393 lationships between queries, allowing the model to 394 generalize across variations in language use. 395

Similarly, Model6 (Cosine + Manhattan) demon-396 strates strong performance, especially on SNIPS, 397 with an accuracy of 99.36%. This combination lever-398 ages Manhattan distance's ability to measure po-399 sitional differences along with Cosine similarity's 400 strength in detecting semantic alignment. The im-401 proved performance on SNIPS reflects the impor-402 tance of accounting for both the semantic meaning 403 and the positional variations in user queries, espe-404 cially when dealing with multi-domain data. 405

However, adding a third metric does not always 406 lead to further improvements. Model7 (Cosine + 407 Manhattan + Euclidean), which combines all three 408 metrics, achieves slightly lower accuracy than sim-409 pler combinations, with 99.65% on ATIS and 98.91% 410 on SNIPS. This drop in performance can be at-411 tributed to redundancy between the metrics, as well 412 as the increased complexity of balancing the influ-413 ence of multiple metrics during training. In some 414 cases, adding more metrics introduces noise and 415 makes it harder for the model to learn effectively, 416 leading to overfitting or diminishing returns in ac-417 curacy. These results highlight the importance of 418 careful metric selection, as simpler combinations 419 may often be more effective than using all available 420 metrics. 421

The ability of the fused models to perform con-422 sistently across both datasets indicates that metric 423 fusion improves the generalization of the Siamese 424 network. While individual metrics perform well 425 only on specific datasets, such as Manhattan and 426 Euclidean on ATIS or Cosine similarity on SNIPS. 427 their combination allows the model to leverage the 428 strengths of each metric. This enables the model 429 to handle both domain-specific queries (ATIS) and 430 multi-domain queries (SNIPS) effectively. For in-431 stance, the combination of Euclidean and Cosine 432 433 metrics captures both the magnitude-based distinctions needed for structured queries and the semantic 434 alignment needed for diverse inputs. 435

The logarithmic transformation applied to each
metric further enhances the model's generalization.
By normalizing the values of the metrics, the transformation smooths out large differences and prevents
any single metric from dominating the fusion layer.
This ensures that the model benefits equally from

the complementary strengths of multiple metrics, 442 leading to improved convergence during training. 443

The results suggest that metric selection should 444 align with the nature of the dataset. For datasets 445 like ATIS, where queries are relatively uniform and 446 domain-specific, metrics that measure numerical dif-447 ferences or geometric distances are more effective. 448 In contrast, for datasets like SNIPS, where semantic 449 richness and diversity are key characteristics, Cosine 450 similarity or combinations of metrics that capture 451 both semantic and positional differences yield better 452 performance. 453

Additionally, the performance decline observed in 454 Model7 highlights the need to balance model com-455 plexity with performance gains. While combining 456 metrics can enhance generalization, using too many 457 metrics may lead to redundancy and hinder perfor-458 mance. These findings emphasize the importance of 459 selecting complementary metrics that align with the 460 specific requirements of the task and dataset. 461

**Table 1.** Performance of Different Versions of theProposed Model

Models	Distance Metrics	Accuracy (%)	
		ATIS	SNIPS
Model1	Manhattan	95.41	86.10
Model2	Cosine Similarity	95.32	85.75
Model3	Euclidean	95.42	85.23
Model4	Euclidean + Cosine	99.81	99.67
Model5	Euclidean + Manhattan	99.80	86.13
Model6	Cosine + Manhattan	99.74	99.36
Model7	Cosine + Manhattan + Euclidean	99.65	98.91

# 6 Comparison with State-of- 462 the-Art Models 463

To evaluate our proposed model, we compared the 464 best performing setting, Model4 against state-of-theart models. The results in Table 2 demonstrate the 466 superiority of our model on both ATIS and SNIPS 467 datasets.

**Table 2.** Comparison with Published Results on ATISand SNIPS Datasets

Model	ATIS (%)	SNIPS (%)
C2A-SLU [32]	96.84	_
LIDSNet [30]	95.97	98.00
BERT+PSN [33]	-	92.89
SN-TripletLoss [5]	99.56	99.31
Ours	<b>99.81</b>	99.67

468

Our model outperformed the C2A-SLU model by 469 3.06% on the ATIS dataset. This improvement can 470 be attributed to the fact that contrastive learning 471 primarily focuses on representation learning, which 472

NLDL

#12

528

may not be as directly optimized for task-specificobjectives as in the proposed approach.

Compared to LIDSNet, which uses triplet loss for 475 training, our model achieved a 4% higher accuracy 476 on the ATIS dataset and a 1.70% improvement on 477 the SNIPS dataset. Compared with SN-TripletLoss, 478 which also uses a triplet loss framework, our model 479 showed an improvement in accuracy of 0.25% and 480 0.36% on the ATIS and SNIPS datasets, respectively. 481 The more efficient performance of our model can 482 be attributed to its reliance on distance metrics, 483 which involve fewer hyperparameters and are less 484 prone to the challenges of tuning margin parameters, 485 learning rates, and triplet mining strategies that 486 often lead to suboptimal performance in LIDSNet 487 and SN-TripletLoss. 488

In contrast to BERT+PSN, which uses a pseudoSiamese network for few-shot intent detection, our
model demonstrated a notable 7.3% improvement
in accuracy on the SNIPS dataset. This significant
margin underscores the robustness of our approach,
particularly in scenarios with limited labeled data
for intent detection.

# 496 7 Conclusion

This study presented a novel intent detection ap-497 proach using an enhanced Siamese network that 498 integrates multiple distance metrics with a fusion 499 layer. The proposed model demonstrated superior 500 performance on the ATIS and SNIPS datasets, out-501 performing state-of-the-art methods. The combi-502 nation of Manhattan, Euclidean, and Cosine sim-503 ilarity metrics proved crucial in handling diverse 504 and domain-specific tasks, improving generalization 505 506 and reducing dependence on annotated datasets. By simplifying the architecture and minimizing hyperpa-507 rameter tuning, our model offers an efficient, scalable 508 solution, particularly in low-resource environments. 509

Despite the promising results, this study has some 510 limitations. A key limitation lies in the use of ATIS 511 and SNIPS datasets, which, although widely adopted 512 benchmarks, have become overused in recent re-513 search. As a result, the performance gains observed 514 on these datasets may not translate directly to real-515 world applications with more complex and evolving 516 intent structures. Additionally, while the fusion of 517 multiple metrics improved accuracy, the individual 518 metrics produced only marginal improvements. This 519 suggests that the impact of individual metrics might 520 be limited when dealing with datasets that are not 521 as saturated with patterns as ATIS. Another limita-522 tion is the absence of a detailed investigation into 523 the learned representations and the specific contri-524 butions of each metric. This leaves room for further 525 exploration into how the metrics complement each 526 527 other.

## References

- K. Amin, G. Lancaster, S. Kapetanakis, K.-D. 529 |1|Althoff, A. Dengel, and M. Petridis. "Ad-530 vanced similarity measures using word embed-531 dings and siamese networks in CBR". In: Intel-532 ligent Systems and Applications: Proceedings 533 of the 2019 Intelligent Systems Conference (In-534 telliSys) Volume 2. Springer. 2020, pp. 449-535 462. DOI: 10.1007/978-3-030-29513-4\_32. 536
- [2] N. Fatima, Q. M. Areeb, I. M. Khan, and 537
  M. M. Khan. "Siamese network-based computer vision approach to detect papaya seed 539
  adulteration in black peppercorns". In: Jour-540
  nal of Food Processing and Preservation 46.9
  (2022), e16043. DOI: 10.1111/jfpp.16043.
- [3] M. P. Dhaliwal, H. Tiwari, and V. Vala. "Automatic creation of a domain specific thesaurus sum sign siamese networks". In: 2021 IEEE 15th sum sum sign sign conference on Semantic Computing (ICSC). IEEE. 2021, pp. 355–361. DOI: 547 10.1109/ICSC50631.2021.00066. 548
- [4] R. Hadsell, S. Chopra, and Y. LeCun. "Dimensionality reduction by learning an invariant mapping". In: 2006 IEEE computer society conference on computer vision and pattern 552 recognition (CVPR'06). Vol. 2. IEEE. 2006, 553 pp. 1735–1742. 554
- [5] F. Ren and S. Xue. "Intention detection based 555 on siamese neural network with triplet loss". 556 In: *IEEE Access* 8 (2020), pp. 82242–82254. 557 DOI: 10.1109/ACCESS.2020.2991484. 558
- [6] E. Ferreira, B. Jabaian, and F. Lefevre. "Online adaptative zero-shot learning spoken language understanding using word-embedding". 561 In: 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE. 2015, pp. 5321–5325. DOI: 564 10.1109/ICASSP.2015.7178987. 565
- [7] A. Stolcke, K. Ries, N. Coccaro, E. Shriberg, 566
  R. Bates, D. Jurafsky, P. Taylor, R. Martin, 567
  C. V. Ess-Dykema, and M. Meteer. "Dialogue 568
  act modeling for automatic tagging and recognition of conversational speech". In: Computational linguistics 26.3 (2000), pp. 339–373. 571
  DOI: 10.1162/089120100561737. 572
- S. Grau, E. Sanchis, M. J. Castro, and D. Vilar. 573
   "Dialogue act classification using a Bayesian approach". In: 9th Conference Speech and Computer. Citeseer. 2004. DOI: 10.21437/SPECOM. 576
   2004–87. 577
- [9] S. L. Kivisaari, A. Hultén, M. van Vliet, T. 578 Lindh-Knuutila, and R. Salmelin. "Seman-579 tic feature norms: a cross-method and cross-580 language comparison". In: *Behavior Research*581

 582
 Methods (2023), pp. 1–10. DOI: 10.3758 /

 583
 \$13428-023-02311-1.

- H. Khanpour, N. Guntakandla, and R. [10]584 Nielsen. "Dialogue act classification in domain-585 independent conversations using a deep recur-586 rent neural network". In: Proceedings of col-587 ing 2016, the 26th international conference 588 on computational linguistics: Technical papers. 589 2016, pp. 2012-2021. DOI: 10.18653/v1/C16-590 1189. 591
- [11] Y. Kim. "Convolutional Neural Networks for 592 Sentence Classification". In: Proceedings of 593 the 2014 Conference on Empirical Methods 594 in Natural Language Processing (EMNLP). 595 Ed. by A. Moschitti, B. Pang, and W. Daele-596 mans. Doha, Qatar: Association for Computa-597 598 tional Linguistics, Oct. 2014, pp. 1746–1751. DOI: 10.3115/v1/D14-1181. URL: https: 599 //aclanthology.org/D14-1181. 600
- L. Wang, H. Yang, F. Li, W. Yang, and Z. Zou.
  "Intent detection model based on dual-channel feature fusion". In: 2022 IEEE 6th Information Technology and Mechatronics Engineering Conference (ITOEC). Vol. 6. IEEE. 2022, pp. 1862–1867. DOI: 10.1109/ITOEC53115.
  2022.9734626.
- [13] Y. Wang, Y. Shen, and H. Jin. "A bi-model
  based RNN semantic frame parsing model for
  intent detection and slot filling". In: arXiv
  preprint arXiv:1812.10235 (2018).
- [14] J. A. Rodriguez, N. Botzer, D. Vazquez, C.
  Pal, M. Pedersoli, and I. H. Laradji. IntentGPT: Few-Shot Intent Discovery with
  Large Language Models. 2024. URL: https:
  //openreview.net/forum?id=2kvDzdC5rh.
- 617 [15] S. Yang, Y. Du, X. Zheng, X. Li, X. Chen,
  618 Y. Li, and C. Xie. "Few-shot intent detection
  619 with self-supervised pretraining and prototype620 aware attention". In: *Pattern Recognition*621 (2024), p. 110641.
- [16] X. Dong and J. Shen. "Triplet Loss in Siamese
  Network for Object Tracking". In: Computer Vision-ECCV 2018: 15th European Conference, Munich, Germany, September 8-14,
  2018, Proceedings, Part XIII 15. Springer.
  2018, pp. 472–488. DOI: 10.1007/978-3030-01261-8\_28.
- [17] A. E. Khan, M. J. Hasan, H. Anjum, and N.
  Mohammed. "Shadow: A Novel Loss Function
  for Efficient Training in Siamese Networks".
  In: arXiv preprint arXiv:2311.14012 (2023).
- [18] T. Mikolov, W.-t. Yih, and G. Zweig. "Linguistic regularities in continuous space word
  representations". In: *Proceedings of the 2013 conference of the north american chapter of*the association for computational linguistics:

Human language technologies. 2013, pp. 746–638 751. DOI: 10.3115/v1/N13-1090. 639

- M. Kamyab, G. Liu, A. Rasool, and M. Adjeisah. "ACR-SA: attention-based deep model 641 through two-channel CNN and Bi-RNN for 642 sentiment analysis". In: *PeerJ Computer Sci-*643 *ence* 8 (2022), e877. DOI: is10.7717/peerj-644 cs.877. 645
- [20] N. A. Tu, D. X. Hieu, T. M. Phuong, and 646
  N. X. Bach. "A bidirectional joint model for 647 spoken language understanding". In: *ICASSP* 648 2023-2023 IEEE International Conference 649 on Acoustics, Speech and Signal Processing 650 (ICASSP). IEEE. 2023, pp. 1–5. DOI: 10.651 1109/ICASSP49357.2023.10096195. 652
- [21] J. Wang, X. Li, J. Li, Q. Sun, and H. Wang. 653
   "NGCU: A new RNN model for time-series 654
   data prediction". In: *Big Data Research* 27 655
   (2022), p. 100296. DOI: 10.1016/j.bdr.2021. 656
   100296. 657
- J. C.-W. Lin, Y. Shao, Y. Djenouri, and U. 658
  Yun. "ASRNN: A recurrent neural network 659
  with an attention model for sequence label- 660
  ing". In: *Knowledge-Based Systems* 212 (2021), 661
  p. 106548. DOI: 10.1016/j.knosys.2020. 662
  106548. 663
- [23] S. Mukherjee and R. Sonal. "A reconciliation 664 between cosine similarity and Euclidean distance in individual decision-making problems". 666 In: Indian Economic Review 58.2 (2023), 667 pp. 427–431. DOI: 10.1007/s41775-023-668 00206-8. 669
- [24] C. C. Aggarwal, A. Hinneburg, and D. A. 670 Keim. "On the surprising behavior of distance metrics in high dimensional space". In: 672 Database theory—ICDT 2001: 8th international conference London, UK, January 4-6, 674 2001 proceedings 8. Springer. 2001, pp. 420-434. DOI: 10.1007/3-540-44503-X\_27. 676
- [25] R. M. West. "Best practice in statistics: The G77 use of log transformation". In: Annals of Clini-G78 cal Biochemistry 59.3 (2022), pp. 162–165. DOI: G79 10.1177/00045632211050531. G80
- [26] C. T. Hemphill, J. J. Godfrey, and G. R. Doddington. "The ATIS spoken language systems 632 pilot corpus". In: Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990. 635 1990. DOI: Doddingtonis10.3115/116580. 636 116613. 637
- [27] A. Coucke, A. Saade, A. Ball, T. Bluche, 688
  A. Caulier, D. Leroy, C. Doumouro, T. Gis- 689
  selbrecht, F. Caltagirone, T. Lavril, et al. 690
  "Snips voice platform: An embedded spoken 691
  language understanding system for private- 692
  by-design voice interfaces. arXiv 2018". In: 693

NLDL

#12

arXiv preprint arXiv:1805.10190 (2018). DOI:
 10.48550/arXiv.1805.10190.

- E28 H. Weld, X. Huang, S. Long, J. Poon, and
  S. C. Han. "A survey of joint intent detection
  and slot filling models in natural language
  understanding". In: ACM Computing Surveys
  55.8 (2022), pp. 1–38. DOI: 10.1145/3547138.
- [29] R. O. Ogundokun, R. Maskeliunas, S. Misra, and R. Damaševičius. "Improved CNN based on batch normalization and adam optimizer". In: International Conference on Computational Science and Its Applications. Springer. 2022, pp. 593–604. DOI: 10.1007/978-3-031-10548-7\_43.
- V. Agarwal, S. D. Shivnikar, S. Ghosh, H. 708 [30]709 Arora, and Y. Saini. "Lidsnet: A lightweight on-device intent detection model using deep 710 siamese network". In: 2021 20th IEEE In-711 ternational Conference on Machine Learn-712 ing and Applications (ICMLA). IEEE. 2021, 713 pp. 1112-1117. DOI: 10.1109/ICMLA52953. 714 2021.00182. 715
- [31] Y. Wu, H. Li, L. Zhang, C. Dong, Q. Huang, and S. Wan. "Joint intent detection model for task-oriented human-computer dialogue system using asynchronous training". In: ACM Transactions on Asian and Low-Resource Language Information Processing 22.5 (2023), pp. 1–17. DOI: 10.1145/3558096.
- [32] X. Cheng, Z. Yao, Z. Zhu, Y. Li, H. Li, and Y.
  Zou. "C 2 A-SLU: cross and contrastive attention for improving ASR robustness in spoken
  language understanding". In: *Proc. of INTER- SPEECH*. 2023. DOI: 10.21437/Interspeech.
  2023-93.
- [33] C. Xia, C. Xiong, and P. Yu. "Pseudo siamese network for few-shot intent generation". In: *Proceedings of the 44th International ACM SI-GIR Conference on Research and Development in Information Retrieval.* 2021, pp. 2005–2009.
  DOI: 10.1145/3404835.3462995.