

| Model             | Validation |       |       |       | Test |      |      |       |
|-------------------|------------|-------|-------|-------|------|------|------|-------|
|                   | R@1        | R@10  | R@50  | MRR   | R@1  | R@10 | R@50 | MRR   |
| <b>subtask-1</b>  |            |       |       |       |      |      |      |       |
| ESIM              | 43.76      | 71.70 | 95.84 | 53.24 | 50.1 | 78.3 | 95.4 | 59.34 |
| T-ESIM            | 52.62      | 76.46 | 96.08 | 60.57 | 61.9 | 82.2 | 96.6 | 69.09 |
| T-ESIM-CR         | 54.46      | 79.26 | 97.92 | 62.73 | 63.4 | 84.2 | 98.5 | 70.69 |
| T-ESIM-Sampled    | 53.16      | 78.5  | 96.46 | 61.54 | 62.8 | 83.4 | 96.6 | 69.7  |
| T-ESIM-Sampled-CR | 55.46      | 81.98 | 98.2  | 64.12 | 64.3 | 84.7 | 97.3 | 71.25 |
| <b>subtask-2</b>  |            |       |       |       |      |      |      |       |
| ESIM              | 11.12      | 31.5  | 58.6  | 18.59 | 12.8 | 28.5 | 36.5 | 18.43 |
| T-ESIM            | 18.72      | 35.9  | 61.26 | 25.13 | 21.6 | 36.0 | 44.1 | 26.68 |
| <b>subtask-4</b>  |            |       |       |       |      |      |      |       |
| ESIM              | 40.16      | 76.18 | 96.36 | 53.43 | 43.5 | 82.1 | 96.2 | 57.96 |
| T-ESIM            | 47.76      | 77.22 | 96.4  | 58.43 | 52.5 | 82.3 | 97.1 | 63.6  |
| <b>subtask-5</b>  |            |       |       |       |      |      |      |       |
| K-ESIM            | 44.82      | 72.74 | 96.4  | 54.52 | 50.1 | 78.3 | 96.3 | 60.2  |
| TK-ESIM           | 53.10      | 75.88 | 96.26 | 60.88 | 60.9 | 80.2 | 96.6 | 67.93 |
| TK-ESIM-CR        | 54.84      | 79.26 | 97.96 | 62.98 | 62.3 | 83.4 | 97.8 | 69.56 |

Table 1: Performance of models on the Ubuntu validation and test datasets. R@k refers to Recall at position k in 100 candidates, denoted as R@1, R@10 and R@50. MRR refers to the Mean Reciprocal Rank.

| Model              | Validation |      |      |       | Test |      |      |       |
|--------------------|------------|------|------|-------|------|------|------|-------|
|                    | R@1        | R@10 | R@50 | MRR   | R@1  | R@10 | R@50 | MRR   |
| ESIM (subtask-1)   | 17.2       | 47.6 | 88.8 | 27.5  | 14.8 | 46.2 | 86.6 | 25.43 |
| ESIM (subtask-3)   | 10.2       | 47.6 | 87.6 | 22.08 | 18.6 | 60.2 | 92.6 | 31.62 |
| ESIM (subtask-4)   | 22.2       | 57.2 | 91.8 | 33.89 | 17.0 | 72.8 | 91.2 | 30.14 |
| K-ESIM (subtask-5) | 16.4       | 50.4 | 85.6 | 27.45 | 11.6 | 49.2 | 88.2 | 23.02 |

Table 2: Performance of models on the Advising validation and test datasets. R@k refers to Recall at position k in 100 candidates, denoted as R@1, R@10 and R@50. MRR refers to the Mean Reciprocal Rank.

has moderate workload, large class size, 4 credits, has a discussion, the classes are on Thursday, Tuesday afternoon”. This representation is used as external knowledge input to our proposed model K-ESIM for the given dialog.

### Proposed model: T-ESIM

In a dialog corpus, similar conversations can appear many times. For example, in a customer care scenario, a common problem might arise for multiple users and many similar dialogs would be present in the corresponding dialog corpus. (?) proposed Exemplar Encoder-Decoder (EED) architecture that makes use of similar conversations for the generation-based dialog system and achieved better results than models such as HRED (?) and VHRED (?). We adopt a similar approach and propose a new training strategy to incorporate the information from similar dialogs present in the available training data for the next utterance selection task. We refer to the new training strategy as T-ESIM (Targeted ESIM), where additional information in terms of probable target responses is added to the contextual information.

For our T-ESIM implementation, we use text-based similarity technique to identify relevant dialogs, similar to K-ESIM. Each dialog in the training data is split at multiple

points to create a larger pool of dialogs, which are called sub-dialogs. The sub-dialogs are then converted to TF-IDF vector representations. The current dialog is matched against these sub-dialogs (excluding its children) to identify similar sub-dialogs to select top-k similar sub-dialogs<sup>6</sup>. The corresponding response(s) for the top-k similar sub-dialogs are concatenated to the current dialog context as a new turn in the partial conversation<sup>7</sup>. The core motivation here is that the model can learn to use responses for similar dialogs present in the training data, to get improved performance on the next utterance selection task. We also explore additional training strategies: T-ESIM-Sampled and T-ESIM-CR, which are described below. We evaluate these strategies on the Ubuntu dialog corpus, as shown in Table 1.

**T-ESIM-Sampled:** When the candidate set for each dialog contains 100 utterances, 1 candidate is the correct utterance and the remaining 99 candidates are incorrect responses. To speed up the training, we randomly sample 9 utterances from the 99 incorrect utterances. We refer to this training strategy as *T-ESIM-Sampled*.

**T-ESIM-CR:** The Ubuntu and Advising corpus are constructed from real-world human-to-human conversations.

<sup>6</sup>We use k=3 during training and k=1 during evaluation

<sup>7</sup>The training data is increased by a factor of k

This makes them unique, as different people can answer the same question in different ways. The different answers could theoretically have the same information but would differ in terms of natural language. Therefore, during evaluation, we employ the Candidate Reduction (CR) trick to use the presence of unique responses in the dataset. For evaluation on the validation set, we reduce the total number of candidates in the candidate set by removing the candidates which are present as correct responses in the training data and similarly, for test data, we remove the candidates which are present in the training and validation data.

## Experiments and Results

Our results for the baseline model ESIM and our proposed models: K-ESIM and T-ESIM for the Ubuntu dataset are given in Table 1 and for the Advising dataset are given in Table 2. The models are evaluated on two metrics - Recall@k, which refers to recall at position k in the set of the 100 candidates and MRR (mean reciprocal rank). We observe that the baseline ESIM model achieves 50.1 R@1 on the Ubuntu test set and 14.8 R@1 on the Advising test set for subtask 1. For Advising dataset subtask 5, we observe that the K-ESIM model performance is slightly below the baseline ESIM model. We believe that our external knowledge representation for the Advising dataset is not suited for the task. For the Ubuntu dataset, we also observe that our proposed models: K-ESIM and T-ESIM perform better than the baseline ESIM model. K-ESIM achieves 44.82 R@1 and 0.5452 MRR on Ubuntu subtask 5 validation set, compared to 43.76 R@1 and 0.5324 MRR for ESIM. K-ESIM also performs slightly better than ESIM on MRR on the test set. T-ESIM performs significantly better than the baseline ESIM model on all Ubuntu subtasks and achieves 61.9 R@1 on subtask 1. Our proposed techniques T-ESIM-Sampled and T-ESIM-CR perform well and achieve 64.3 R@1 score on the Ubuntu subtask 1. These results show that our proposed models and training strategies perform well.

For Ubuntu Subtask 2, the size of global pool of candidates is 120000. For training purposes, we reduce the candidate set by randomly sampling 99 incorrect responses from the global pool. These 99 responses, in addition to the correct response, construct our candidate set of 100 responses per dialog, similar to Subtask 1. During evaluation on the validation and test sets, we first employ the CR technique mentioned above. Then, we shortlist the number of candidates to 100, by selecting the top-100 candidates from the reduced candidate global pool using IR-based methods similar to knowledge extraction for K-ESIM and T-ESIM.

## Conclusion and Future Work

In this paper, we introduced two knowledge incorporating end-to-end dialog systems for retrieval-based goal-oriented dialog, by extending the ESIM model. Evaluation based on the Ubuntu dataset show that our methods are effective to improve performance by incorporating additional external knowledge sources and leveraging information from similar dialogs. Although our proposed model K-ESIM shows improvement on the Ubuntu subtask 5, we observe a slight decrease in performance on the Advising subtask 5 as explained

in the previous section. In our future work, we plan to explore the following areas to improve our proposed K-ESIM and T-ESIM models: a) improve the knowledge representation for course information, b) investigate attention mechanisms over a KB (?) and c) explore neural approaches, instead of TF-IDF, for extracting relevant external information (man pages) and identifying similar dialogs for T-ESIM.

## Appendix: Model Training and Hyperparameter Details

In Word Representation Layer, we used 300-dimensional Glove pre-trained vectors<sup>8</sup> ((?)), 100-dimensional word2vec vectors (?) and 80-dimensional character-composed embedding vectors for generating the representation of a word. For training word2vec vectors, we use the `gensim.models.Word2Vec` API with the following hyper-parameters: `size=100`, `window=10`, `min_count=1` and `epochs=20`. The final prediction layer is a 2-layer fully-connected feed-forward neural network with ReLU activation. We use sigmoid function and minimize binary cross-entropy loss for training and updating the model. The baseline model was implemented in Tensorflow (?) and we used the source code released by ? (?)<sup>9</sup> for the baseline model. We generated word2vec word embeddings from scratch on the DSTC7 datasets as mentioned in Algorithm-1 from ? (?). We used Adam (?) with a learning rate of 0.001 and exponential decay with a decay rate of 0.96 decayed every 5000 steps. Batch size used was 128. The number of hidden units for BiLSTM in both the context representation layer and the matching aggregation layer was 200. For the prediction layers, we used 256 hidden units with ReLU activation.

Hic provident veritatis quae impedit iure expedita facere libero aut modi, magnam minima autem id porro, illum distinctio nobis fuga perspiciatis sapiente doloreque, adipisci magni pariat ut inventore doloribus sit. Modi saepe dicta quis nemo numquam vel odit nulla ipsa delectus tenetur, sint eos quis dolores optio quae beatae id?Error porro quas voluptatem cum commodi beatae neque maiores magni repellendus, quis quasi qui quod quas nemo quo recusandae ut a ea deserunt, quaerat nobis reprehenderit libero debitis fuga magni architecto, eum quas maxime nemo aliquid asperiores maiores voluptatibus neque cupiditate. Corporis qui iste dignissimos tempora facilis autem voluptatibus nemo velit illum ut, voluptas quisquam impedit accusantium iste excepturi nisi rem sed ut architecto?Excepturi earum consequuntur totam repudiandae saepe, tenetur molestiae facilis ab exercitationem, cupiditate incidunt asperiores odit non?Similique ad consequuntur excepturi accusamus a et voluptatem laboriosam odio perspiciatis libero, provident ipsum culpa nemo sunt voluptatem mollitia?Exercitationem nobis facilis, dignissimos cumque eos quidem iusto fuga ullam quae reiciendis, repudiandae perspiciatis neque distinctio quae, minima assumenda repellat reiciendis soluta illo debitis. Assumenda rem possimus totam nam tempora

<sup>8</sup>glove.42B.300d.zip : <https://nlp.stanford.edu/projects/glove/>

<sup>9</sup>source code released by ? (?): [https://github.com/jdongca2003/next\\_utterance\\_selection](https://github.com/jdongca2003/next_utterance_selection)