

Findings on Conversation Disentanglement

Rongxin Zhu, Jey Han Lau, Jianzhong Qi

School of Computing and Information System, The University of Melbourne

Introduction

1.1 Problem Definition

Conversation disentanglement aims at identifying separate threads in multi-party conversations, acting as an important preprocessing step for high-level tasks of multi-party conversations such as conversation summarization and response generation. The figure below shows two threads in different colors, with reply-to links between utterances.

[12:05] <ydnar> for what reason would a dvd not play if i have libdvcss2 installed?

[12:05] <gourdin> we will be able to access an edgy repo ?

[12:05] <Ng> ydnar: what are you using to play it?

[12:06] <Anfangs> Edgy Eft is the next codename for Ubuntu dapper+1. See <https://ubuntu.com/0064.html>.

[12:06] <holycow> because it couldn't crack the encoding for the particular portion of the dvde

[12:06] <ydnar> tried vlc. holycow, do you have any

[12:06] <gourdin> I don't think the link works

1.2 Limitation of previous methods

- transformer-based models are not systematically compared with respect to performance, memory consumption and speed
- previous methods don't leverage **dialogue history** when measuring the similarity between UOI and a candidate, or introduce **noise** in context expansion
- greedy decoding algorithm recovers threads by finding the parent utterance for each utterance of interest (UOI) **independently**

Methodology

2.1 Pairwise Model

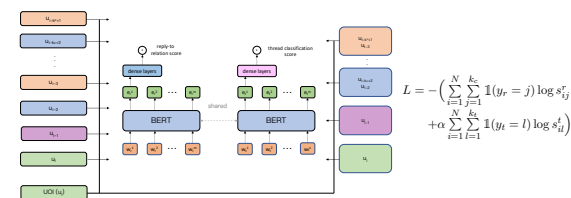
We conduct empirical study to compare rule-based, feature-based and transformer models that measures the similarity between UOI and each candidate (we regard k=50 past utterances, including UOI, as the candidate pool). Experiments show that BERT [1] +MF (manual features) is still a strong baseline.

Model	Link Prediction		F1	Ranking		Clustering	F
	Precision	Recall		R@1	R@5	1-1	VI
Last Mention	37.1	35.7	36.4			21.4	60.5
GLOVE+MF	71.5	68.9	70.1	70.2	95.8	98.6	76.1
MF	71.1	68.5	69.8	70.2	94.0	97.3	75.0
Poly-Batch	71.3	69.3	70.3	70.2	94.0	97.3	75.0
POLY-INLINE	42.2	40.7	41.4	42.8	70.8	81.3	62.0
ALBERT	46.1	44.4	45.3	46.8	77.3	88.4	68.6
BERT	48.2	46.4	47.3	48.3	75.4	84.7	74.3
BERT+TJ	67.9	65.3	66.6	66.9	90.2	95.3	76.0
BERT+MF	73.9	71.3	72.6	73.9	95.8	98.6	77.0

Model	GPU Mem (GB)	Speed (ms)
BERT	18.7	9.4
ALBERT	14.6	9.4
POLY-INLINE	9.9	16.8
POLY-BATCH	5.1	36.4

We compare the GPU memory consumption and speed of four transformer models, finding that Poly-encoder [2] is the fastest and most memory efficient one, with a sacrifice of performance.

2.2 Context Expansion using Multi-task Learning



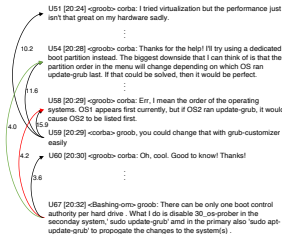
Multi-task Learning: we conduct utterance-to-utterance classification (reply-to relation identification) and utterance-to-thread classification (assign each UOI to an existing thread or create a new thread) at the same time, using a shared BERT model with separate dense layers for two tasks. In training, we use ground truth reply-to relations and thread labels. During inference, we only do reply-to relation identification, which does not introduce noise if the predicted threads are not completely correct (e.g., a predicted thread may contain utterances from other threads). The loss function is a weighted sum of two loss terms for reply-to relation identification and thread classification, respectively.

Model	Link Prediction			Ranking		Clustering		
	Precision	Recall	F1	R@1	R@5	1-1	VI	F
BERT	48.2	46.4	47.3	48.8	75.4	84.7	74.3	89.3
BERT+MF	73.9	71.3	72.6	73.9	95.8	98.6	77.0	92.0
MULTI (alpha = 1)	65.6	63.2	64.4	66.7	91.8	95.6	64.6	87.7
MULTI (alpha = 5)	66.9	64.5	65.7	65.4	91.8	95.6	68.7	88.8
MULTI (alpha = 10)	65.2	62.9	64.0	64.4	91.4	95.6	70.3	89.5
MULTI (alpha = 20)	64.7	62.4	63.5	63.9	91.0	95.0	68.3	88.8
MULTI+MF (alpha = 1)	72.8	70.2	71.5	71.9	94.0	96.4	76.3	91.8
MULTI+MF (alpha = 5)	73.3	70.7	72.0	72.4	94.0	96.5	72.8	90.8
MULTI+MF (alpha = 10)	72.2	69.6	70.8	70.4	93.4	96.4	71.8	90.2
MULTI+MF (alpha = 20)	70.8	68.2	69.5	69.4	93.4	97.3	73.2	90.6

Our multi-task learning framework is helpful when manual features are not available. It doesn't outperform BERT+MF pairwise model.

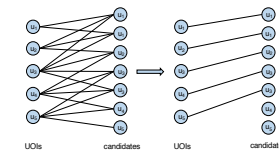
2.3 Bipartite Graph Matching for Conversation Disentanglement

Motivation: pairwise models achieve a high recall@5, which means the correct parents are ranked high but not at the top sometimes. Bipartite matching-based algorithm recovers threads by identifying the parent utterance of a set of UOIs jointly. When conflict occurs (multiple UOIs choose the same parent), some UOI may fall back to choose its second-best candidate as parent.



U₆₇ chooses U₅₄ as parent in global decoding algorithm but chooses U₅₈ in greedy algorithm.

We first build a bipartite graph with two sets of nodes, containing all UOIs and all candidates, respectively. Edges are created between each UOI, and its top-5 ranked candidates from pairwise models. Then we frame conversation disentanglement as a **maximum-weight bipartite matching** [3] problem. The aim is to find a subset of edges, representing the predicted reply-to relations.



	Precision	Recall	F1
Oracle	88.4	85.2	86.8
Rule-Based	73.7	70.9	72.3
FFN	73.8	71.0	72.3
BERT+FFN	72.9	70.3	71.5

The bipartite graph matching-based algorithm has the potential to outperform greedy approaches. Oracle means we use ground truth node frequencies for all candidate nodes. We tried to predict node frequencies, but the results are not ideal. More effective methods to predict node frequencies are needed.

Conclusion

- BERT combined with manual features is still a strong baseline for conversation disentanglement
- The multi-task learning framework that conducts utterance-to-utterance and utterance-to-thread classification at the same time outperforms pairwise models when manual features are not available
- Bipartite graph matching-based conversation disentanglement shows potential to outperform greedy approaches.

References

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2019. Poly-encoders: Architectures and pre-training strategies for fast and accurate multi-sentence scoring. In International Conference on Learning Representations. Y. Zhang, C. Bajaj, B.-S. Sohn. Adaptive and Quality 3D Meshing from Imaging Data, ACM Symposium on Solid Modeling and Applications, pp. 286–291, Seattle, June 2003.
- AMH Gerards. 1995. Matching. Handbooks in operations research and management science, 7:135–224.