

Probing Discourse Structure in Dialogue: Evaluating and Fine-Tuning RoBERTa and BART on Sentence Ordering and Next Sentence Prediction

Anonymous ACL submission

Abstract

This study investigates how large language models represent discourse structure in dialogue through two probing tasks: sentence ordering and next sentence prediction (NSP). Using the STAC corpus of multi-party conversational data, RoBERTa and BART are evaluated in both pre-trained and fine-tuned settings. Fine-tuning yields substantial gains across both tasks, with BART’s generative architecture proving more effective for sentence ordering ($\rho = 0.473$) while RoBERTa excels at NSP classification (accuracy = 0.697). Focusing on four core discourse relations, the analysis finds that models handle frequent, surface-cued relations (Question-Answer Pairs, Comment) effectively but struggle with relations requiring deeper semantic dependencies (Continuation, Elaboration). These findings highlight both the capabilities and limitations of current models in capturing discourse-level coherence in dialogue.

1 Introduction

Large language models (LLMs) have revolutionized natural language processing but face challenges in capturing discourse-level coherence in dialogue (Holtzman et al., 2020). While excelling at surface-level language tasks, their ability to maintain logical flow across multi-turn conversations remains limited. This work extends the probing framework of (Koto et al., 2021) from monologic texts to dialogue, evaluating discourse coherence through sentence ordering and next sentence prediction (NSP) tasks.

The contributions are threefold: (1) demonstrating the necessity of fine-tuning for dialogue discourse understanding, (2) revealing architectural strengths—BART for generation, RoBERTa for classification, and (3) providing focused relation-specific analysis on four core discourse relations showing that models excel at surface-cued relations but struggle with semantic dependencies.

Using the STAC corpus (Asher et al., 2016) with Segmented Discourse Representation Theory (SDRT) annotations (Asher and Lascarides, 2003), this research offers insights into discourse modeling capabilities of transformer architectures (Vaswani et al., 2017).

2 Related Work

2.1 Discourse Structure and Dialogue

Discourse coherence frameworks like SDRT (Asher and Lascarides, 2003) model conversations through relations between Elementary Discourse Units (EDUs). While resources like RST Discourse Treebank (Carlson et al., 2001) and Penn Discourse Treebank (Prasad et al., 2008) focus on written texts, the STAC corpus (Asher et al., 2016) provides fine-grained discourse annotations for multi-party game dialogues, capturing strategic conversations with non-linear discourse links essential for probing dialogue coherence.

2.2 Probing Language Models

(Koto et al., 2021) established probing tasks including sentence ordering and NSP to evaluate discourse awareness, finding RoBERTa and BART most sensitive to discourse phenomena. However, their study focused on monologic texts, leaving dialogue coherence unexplored. This work extends this approach to conversational data with focused relation-specific evaluation on four core discourse relations.

2.3 Sentence Ordering and NSP

Sentence ordering serves as a coherence proxy (Barzilay and Lapata, 2008), with neural approaches demonstrating pattern recognition (Logeswaran et al., 2018; Gong et al., 2016). BART’s sentence permutation pre-training (Lewis et al., 2020) aligns naturally with this task, while RoBERTa’s contextual representations (Liu et al.,

2019) suit ranking approaches. NSP, though questioned for pre-training utility (Liu et al., 2019), remains valuable for probing local coherence and has been used in dialogue response selection (Zhang et al., 2018; Whang et al., 2020).

2.4 Differences from Previous Work

This study differs from (Koto et al., 2021) in three key aspects: (1) focusing on dialogue rather than monologic texts, (2) employing relation-specific analysis across discourse types, and (3) comparing architectural preferences across generative and discriminative paradigms. Unlike dialogue response selection studies (Zhang et al., 2018; Whang et al., 2020), this work probes inherent discourse understanding rather than response generation performance.

3 Methodology

3.1 Dataset and Tasks

The STAC corpus (Asher et al., 2016) contains multi-party online game dialogues annotated with SDRT relations. After preprocessing — including cleaning malformed text, handling token limits, and segmenting utterances into Elementary Discourse Units (EDUs) — dialogues with fewer than three EDUs were excluded to ensure meaningful evaluation of ordering. The study uses two probing tasks:

Sentence Ordering: Models reconstruct the original EDU sequence from shuffled input. The task is evaluated by Spearman’s ρ , calculated after aligning predicted EDU sequences to gold sequences using the Ratcliff–Obershelp algorithm (Ratcliff and Obershelp, 1988), considering only overlapping units.

Next Sentence Prediction: A binary classification task on pairs of EDUs. Positive pairs are adjacent EDUs connected by a discourse relation; negative pairs are non-adjacent EDUs. Balanced accuracy is used to account for class distribution across relation types.

The analysis focuses on four discourse relations with substantial data support: Question-Answer Pair, Comment, Continuation, and Elaboration.

3.2 Models and Evaluation

BART (Lewis et al., 2020): Encoder-decoder fine-tuned for sequence generation (sentence ordering) and binary classification (NSP).

RoBERTa (Liu et al., 2019): Encoder-only model used for reranking (sentence ordering) and classification (NSP).

Both models were fine-tuned using standard hyperparameters with relation-aware special tokens. Specifically, relation-type tokens (e.g., [REL:QAP], [REL:CONT]) were prepended to each EDU to make discourse relations explicit in the input. Evaluation uses Spearman’s ρ for sentence ordering and balanced accuracy for NSP. Statistical significance of improvements was tested using paired Wilcoxon signed-rank tests (Wilcoxon, 1945), with confidence intervals reported via bootstrap resampling (Cumming, 2014).

3.3 Experimental Setup

Fine-tuning used consistent hyperparameters across both models (Table 1): learning rate of $2e-5$, batch size 16, and 3 epochs. Data splits included cross-validation (sentence ordering) and 80/10/10 splits (NSP). All splits were performed at the dialogue level to ensure that EDUs from the same dialogue do not appear in both training and test sets, preventing data leakage. Relation-aware special tokens were added during tokenization. BART ordering used beam search (size 5) with Ratcliff–Obershelp alignment; RoBERTa reranking evaluated 6 permutations per dialogue (inputs truncated to 512 tokens). All runs used early stopping (patience 2) and a fixed random seed.

Table 1: Hyperparameter settings for fine-tuning experiments.

Parameter	BART	RoBERTa
Learning rate	$2e-5$	$2e-5$
Batch size	16	16
Training epochs	3	3
Warmup steps	500	500
Weight decay	0.01	0.01
Max sequence length	512	512

4 Results

4.1 Overall Performance

Fine-tuning yields substantial improvements across both tasks (Table 2). For sentence ordering, BART achieves $\rho = 0.473$ (vs. 0.130 baseline) while RoBERTa reaches $\rho = 0.467$ (vs. 0.149 baseline). For NSP, RoBERTa achieves higher accuracy (0.697) than BART (0.637), reflecting architectural advantages for classification tasks.

All improvements are statistically significant ($p < 0.001$).

Table 2: Overall performance across tasks and models. SO: Sentence Ordering (ρ), NSP: Next Sentence Prediction (accuracy). Scores are means over test instances. All improvements from fine-tuning are statistically significant ($p < .001$, paired Wilcoxon signed-rank tests). Confidence intervals for the improvements are provided in Section 4.1.1.

Task	Model	n	Base.	Fine.
SO	BART	712	0.130	0.473
	RoBERTa	188	0.149	0.467
NSP	BART	1,776	0.499	0.637
	RoBERTa	293	0.510	0.697

4.1.1 Statistical Analysis of Improvements

The effects of fine-tuning were quantified using mean gains (Δ) with 95% confidence intervals (CIs) derived from paired bootstrap resampling (10,000 iterations). For sentence ordering, BART shows $\Delta\rho = 0.343$ (95% CI: [0.309, 0.377]), while RoBERTa shows $\Delta\rho = 0.318$ (95% CI: [0.201, 0.434]). For NSP, the gains are $\Delta\text{acc} = 0.137$ (95% CI: [0.105, 0.171]) for BART and $\Delta\text{acc} = 0.187$ (95% CI: [0.155, 0.219]) for RoBERTa. All fine-tuned models perform significantly above chance level (one-sided Wilcoxon test: $p < .001$ for $\rho > 0$ or $\text{acc} > 0.5$). Complete results for all 15 discourse relations analyzed (including Acknowledgement, Explanation, and Contrast) are consistent with this pattern and are provided in the supplementary material. All experiments used a fixed random seed, with results based on single runs and statistics aggregated over the test set.

4.2 Relation-Specific Analysis

Performance is analyzed across four core discourse relations to understand model capabilities at different levels of discourse complexity (Table 3).

Question-Answer Pairs show the strongest performance across both tasks, with both models achieving high scores in ordering and NSP. **Comment** relations demonstrate solid performance in sentence ordering but more modest gains in NSP. **Continuation** proves challenging for BART in ordering but both models handle it well in NSP. **Elaboration** shows the largest architecture-dependent variation, with RoBERTa substantially outperforming BART in NSP.

Table 3: Performance by discourse relation type (fine-tuned models). Values are mean scores over test instances. Confidence intervals and full statistical details are available upon request.

Relation	Sentence Ordering (ρ)		NSP (Accuracy)	
	BART	RoBERTa	BART	RoBERTa
Question-Answer	0.498	0.487	0.675	0.701
Comment	0.450	0.431	0.544	0.630
Continuation	0.384	0.552	0.678	0.699
Elaboration	0.395	0.342	0.548	0.691

5 Discussion

5.1 Discourse Awareness Through Fine-Tuning

The significant performance gap between baseline and fine-tuned models demonstrates that pre-training alone is insufficient for reliable discourse coherence modeling. While pretrained models show basic syntactic and semantic knowledge, domain-specific fine-tuning is essential for capturing dialogue-specific discourse patterns, particularly for the four focal relations.

5.2 Architectural Alignment with Task Demands

The results reveal clear architectural preferences: RoBERTa’s encoder-only design excels at classification tasks (NSP), while BART’s encoder-decoder framework benefits generation tasks (sentence ordering). This pattern is most pronounced for Elaboration relations, where RoBERTa achieves 0.691 NSP accuracy compared to BART’s 0.548, suggesting RoBERTa’s discriminative approach better captures semantic dependencies.

5.3 Relation-Specific Capabilities and Limitations

The focused analysis reveals a clear hierarchy of model capabilities across the four relations:

Question-Answer Pairs represent the most accessible relation type, benefiting from strong adjacency patterns and lexical complementarity within the game negotiation context. Models reliably learn to preserve question-response sequences after fine-tuning.

Comment relations show intermediate performance, with models successfully grouping related evaluative utterances but sometimes failing to capture pragmatic ordering conventions.

Continuation exposes limitations in tracking

241	procedural logic across turns. While models identify related terms in negotiations, they often break the sequence of proposals and counteroffers, indicating reliance on lexical coherence over logical flow.	289
242		290
243		291
244		292
245		
246	Elaboration demonstrates the potential for learning semantic dependencies, particularly with RoBERTa’s strong NSP performance. However, the relation remains challenging for sentence ordering, suggesting difficulties with global coherence modeling.	
247		
248		
249		
250		
251		
252	5.4 Error Analysis and Case Studies	
253	Qualitative analysis reveals systematic error patterns. The following examples illustrate how models, particularly BART, can confuse <i>Continuation</i> (a sequential, procedural relation) with <i>Elaboration</i> (a semantic dependency relation) due to overlapping lexical cues.	
254		
255		
256		
257		
258		
259	Example 1 (Continuation misidentified as Elaboration):	
260		
261	<i>Gold order:</i> “I need wood.” → “I can give you clay for it.” → “No, I need wood for a settlement.”	
262		
263		
264	<i>Predicted order (BART):</i> “I need wood.” → “No, I need wood for a settlement.” → “I can give you clay for it.”	
265		
266		
267	Here, the model groups the two “need wood” utterances together (treating them as <i>Elaboration</i>) and detaches the offer (“clay for it”), breaking the negotiation <i>Continuation</i> . The model prioritizes lexical repetition over procedural logic.	
268		
269		
270		
271		
272	Example 2 (Elaboration with spurious causal chain):	
273		
274	<i>Gold:</i> “He stole my sheep.” → “That’s why I’m attacking next turn.”	
275		
276	<i>Predicted (BART):</i> “I’m attacking next turn.” → “He stole my sheep.”	
277		
278	The model reverses the cause–effect order, placing the action before its explanation, thereby losing the <i>Elaboration</i> link. RoBERTa’s NSP classifier correctly identifies the true order in this case, reflecting its stronger grasp of semantic dependency.	
279		
280		
281		
282		
283	For <i>Continuation</i> relations, models frequently misorder negotiation sequences, placing counteroffers before initial proposals when they share lexical overlap. In <i>Elaboration</i> examples, BART often generates plausible but incorrect causal chains, while RoBERTa more reliably identifies	
284		
285		
286		
287		
288		
	true explanatory relationships. A notable failure case involves nested relations, where models correctly handle local coherence but fail to maintain global discourse structure across multiple turns.	289
		290
		291
		292
	5.5 Impact of Input Length Constraints	293
	The 512-token input constraint, particularly for RoBERTa, limited the evaluation to shorter dialogues (n=188 for ordering vs. BART’s n=712). This truncation may remove contextual cues necessary for global coherence modeling in longer conversations, potentially underestimating the models’ capabilities on relations like <i>Continuation</i> and <i>Elaboration</i> that span multiple turns. Future work should explore long-context architectures or hierarchical encoding to better capture extended dialogue structure.	294
		295
		296
		297
		298
		299
		300
		301
		302
		303
		304
	6 Conclusion	305
	This work extends discourse probing to dialogue through focused analysis of four core relations. The results demonstrate that fine-tuning is essential for discourse coherence modeling and reveal clear architectural preferences: RoBERTa for classification, BART for generation. The relation-specific hierarchy—from easily learned Question-Answer Pairs to challenging <i>Continuation</i> and <i>Elaboration</i> relations—provides a nuanced understanding of current capabilities and limitations. Future research could explore specific hybrid architectures, such as a pipeline where BART generates candidate sequences and RoBERTa reranks them, or incorporate graph neural networks to explicitly model SDRT structures. Additionally, expanding to more diverse dialogue domains beyond game negotiations would test the generalizability of these findings.	306
		307
		308
		309
		310
		311
		312
		313
		314
		315
		316
		317
		318
		319
		320
		321
		322
		323
	Limitations	324
	This study focuses on four discourse relations from the STAC game dialogues, which may limit generalizability to other relation types and domains. The 512-token input constraint particularly affects RoBERTa’s sentence ordering evaluation, limiting analysis of longer conversational sequences. Future work should expand to include other dialogue domains and address input length constraints.	325
		326
		327
		328
		329
		330
		331
		332
		333

334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387

References

Nicholas Asher, Julie Hunter, Mathieu Morey, Farah Benamara, and Yannick Le Braud. 2016. [Discourse structure and dialogue acts in multiparty dialogue: The STAC corpus](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 2721–2727.

Nicholas Asher and Alex Lascarides. 2003. *Logics of Conversation*. Cambridge University Press.

Regina Barzilay and Mirella Lapata. 2008. [Modeling local coherence: An entity-based approach](#). *Computational Linguistics*, 34(1):1–34.

Lynn Carlson, Daniel Marcu, and Mary Ellen Okurowski. 2001. [Building a discourse-tagged corpus in the framework of rhetorical structure theory](#). In *Proceedings of the Second SIGDIAL Workshop on Discourse and Dialogue*, pages 1–10.

Geoff Cumming. 2014. [The new statistics: Why and how](#). *Psychological Science*, 25(1):7–29.

Yeyun Gong, Hang Luo, and Jun Zhang. 2016. [Natural language inference over interaction space](#). *arXiv preprint arXiv:1607.06952*.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text de-generation](#). In *Proceedings of the 8th International Conference on Learning Representations*.

Fajri Koto, Cedric Shin Shen Wu, Ling Liu, and Timothy Baldwin. 2021. [Discourse probing of pre-trained language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 5412–5424, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A robustly optimized BERT pretraining approach](#). *arXiv preprint arXiv:1907.11692*.

Lajanugen Logeswaran, Honglak Lee, and Yoshua Bengio. 2018. [Sentence ordering and coherence modeling using recurrent neural networks](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pages 5285–5292, New Orleans, Louisiana. AAAI Press.

Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Milt-sakaki, Livio Robaldo, Aravind Joshi, and Bonnie

Webber. 2008. [The Penn discourse treebank 2.0](#). In *Proceedings of the Sixth International Conference on Language Resources and Evaluation*, pages 2961–2968, Marrakech, Morocco. European Language Resources Association. 388
389
390
391
392

John W. Ratcliff and David E. Osherson. 1988. [Pattern matching by gestalt](#). *Software: Practice and Experience*, 18(11):1045–1058. 393
394
395

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. 396
397
398
399
400

Taesun Whang, Dongyub Lee, Dongsuk Oh, Chan-hee Lim, Sangyun Lee, Kijong Park, Jaechoon Lee, and Heuisook Heo. 2020. [Response selection for multi-party conversations with dynamic topic tracking](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6586–6598, Online. Association for Computational Linguistics. 401
402
403
404
405
406
407
408

Frank Wilcoxon. 1945. [Individual comparisons by ranking methods](#). *Biometrics Bulletin*, 1(6):80–83. 409
410

Zhuosheng Zhang, Jiangtong Li, Pengfei Zhu, and Hai Zhao. 2018. [Modeling multi-turn conversation with deep utterance aggregation](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3740–3752, Santa Fe, New Mexico, USA. Association for Computational Linguistics. 411
412
413
414
415
416
417