

A Novel Computational Modeling Foundation for Automatic Coherence Assessment

Anonymous ACL submission

Abstract

Coherence is an essential property of well-written texts, that refers to the way textual units relate to one another. In the era of generative AI, coherence assessment is essential for many NLP tasks; summarization, long-form question-answering, etc. Current NLP approaches for modeling coherence often rely on a proxy task, specifically sentence reordering. However, such an approach may not capture the full range of factors contributing to coherence. To bridge this gap, in this work we employ the formal linguistic definition of Reinhart (1980) of what makes a discourse coherent, consisting of three conditions — *cohesion*, *consistency* and *relevance* – and formalize these conditions as respective computational tasks. We hypothesize that (i) a model trained on all of these tasks will learn the features required for coherence detection, and that (ii) a joint model for all tasks will exceed the performance of models trained on each task individually. We evaluate this modeling approach on two human-rated coherence benchmarks: one of automatically-generated stories and one of real-world texts. Our experiments confirm that jointly training on the proposed tasks leads to better performance on each task compared with task-specific models, and to better performance on assessing coherence overall, compared with strong baselines. Our formal coherence framework paves the way for advanced, broad-coverage automatic assessment.

1 Introduction

The term *coherence* refers to the quality of texts where sentences and paragraphs flow smoothly and are logically connected, creating a clear and understandable progression of ideas. Coherence detection is crucial for NLP tasks involving text quality measurements such as essay scoring or text quality assessment (Somasundaran et al., 2014; Feng et al., 2014; Lai and Tetreault, 2018). Its importance is

further amplified nowadays in the era of large language models (LLMs). Ensuring coherent LLM outputs is essential for producing meaningful and understandable text for generative tasks as summarization and question answering (Guan et al., 2021; Xu et al., 2018; Yi et al., 2019), to name a few.

The elusive nature of coherence makes it challenging for NLP systems to assess it effectively. While various linguistic theories of coherence exist (e.g., Halliday and Hasan (1976); Joshi and Weinstein (1981); Givon (1995); Hobbs (1979); Dijk (1979); Mann and Thompson (1988)), current approaches often rely on proxy tasks like *sentence reordering* (Lapata, 2003; Miltsakaki et al., 2000). However, this approach oversimplifies coherence, potentially leading to models that struggle with real-world texts (Laban et al., 2021). Furthermore, coherence is multifaceted, varying across genres, contexts, and styles. Proxy tasks often fail to capture this complexity, leading to models unable to generalize across different domains and real-world settings. On top of that, existing NLP models for coherence, such as Barzilay and Lapata (2008)’s sentence reordering, often lack direct evaluation and are instead assessed on downstream tasks like readability (Guinaudeau and Strube, 2013; Mesgar and Strube, 2016, 2018) or essay scoring (Mesgar and Strube, 2018; Somasundaran et al., 2014; Tay et al., 2017). This approach can be expensive, biased towards the downstream task, and may overlook core aspects of coherence.

This work aims to provide a computationally workable definition of coherence, leveraging Reinhart (1980)’s theory which defines three conditions: *cohesion*, *consistency*, and *relevance*. Concretely, we propose to instantiate these conditions as computational tasks, and train them jointly to create a model that captures these properties. We hypothesize that (i) such a model will effectively assess coherence, and (ii) shared information between tasks will improve task-specific performance.

To test our hypotheses, we implement a model trained jointly on tasks capturing Reinhart’s conditions of coherence. The unified model incorporates five tasks: sentence reordering (Lapata, 2003), discourse relation detection (Miltsakaki et al., 2004), natural language inference (Dagan et al., 2005), NP enrichment (Elazar et al., 2022), and irrelevant-sentence detection. Despite its relatively simple architecture, our model achieves SOTA results on most of these tasks. We then evaluate the model’s coherence assessment capability on two human-annotated benchmarks: the *Grammarly Corpus of Discourse Coherence (GCDC)* (Lai and Tetreault, 2018) for real-world texts across four domains, and *CohesSentia* (Maimon and Tsarfaty, 2023) for artificially-generated stories. These benchmarks cover both natural and artificially-generated texts, human-rated for their levels of coherence.

Our empirical findings confirm our hypotheses, showing significant accuracy improvements and producing new SOTA results on both benchmarks. The joint model outperforms standalone models for individual tasks. Our significant performance improvements demonstrate the model’s efficacy in automatic coherence assessment. This framework paves the way for future models to not only identify incoherence, but also analyze its causes. Furthermore, integrating this methodology into text generation may lead to higher-quality outputs.

2 The Proposal: Coherence à la Reinhart

This work aims to provide a computationally workable definition of coherence by adopting Reinhart (1980)’s formalization, which identifies coherence through three criteria: *Cohesion*, *Consistency*, and *Relevance*. According to Reinhart, a text is coherent only if it meets all three criteria. Recently, Maimon and Tsarfaty (2023) used this framework to create a benchmark for coherence scoring of GPT-generated text, with human scores for these criteria. Here we take a different approach, using these conditions as the basis for computational modeling, aiming to predict the ingredients of these three properties via jointly trained tasks.

Cohesion The cohesion condition focuses on the formal elements that link sentences together.¹ Rein-

¹The terms “coherence” and “cohesion” may be confusing to non-linguists. Cohesion relates to the surface forms used (e.g., connectors, pronouns), while coherence pertains to the overall semantics and flow of ideas.

hart states that a text is cohesive if, for every sentence pair, at least one of two conditions is met:

(1) *Referentially linked*: A pair of sentences $\langle S_1, S_2 \rangle$ is referentially linked when S_2 references an entity mentioned in S_1 . A simple example is using a pronominal anaphor:

“Dan is nice. Even Su likes him.”

Here, the underlined entities co-refer. Other types of referential links are prepositional links (Elazar et al., 2022) or bridging anaphora (Hou, 2021).

(2) *Linked by a semantic sentence connector*. A pair of sentences $\langle S_1, S_2 \rangle$ is connected if a discourse relation links them. These connectors indicate semantic relations like cause and effect, comparison, contrast, and more (Prasad et al., 2008). An example of linking by a semantic connector is:

“It was raining. So, we stayed inside.”

The sentences are cohesive due to the existence of the “So” connector. These connectors may be explicit or implicit (Pitler et al., 2009).

Consistency The consistency condition pertains to the formal semantic aspects of a text, ensuring *logical* coherence, which is crucial for interpreting and deriving meaning. Formally, this condition requires that for a set of sentences $\{S_i\}_{i=0}^{n-1}$, the meaning of each sentence S_i must be consistent with all previous sentences $\{S_j\}_{j=0}^{i-1}$. This means all sentences can be true within a single world, not violating this world’s assumptions and restrictions. An example of a violation is shown below:

“My father is dead now. That’s why he has decided to smoke a pipe” (Freeman and Gathercole, 1966)

Despite being cohesive (anaphora & connectors), the passage lacks coherence due to world knowledge violations (a deceased cannot decide).²

Relevance The relevance condition involves *pragmatic* aspects, imposing constraints on the relationships of all sentences $\{S_i\}_{i=0}^{N-1}$ to the discourse topic and other contextual elements. An example of a violation of this condition is as follows:

“I poured some chemical into a beaker. The chemical fell on my hand. The professor immediately took me to the emergency bath. He is a great musician.”

²The consistency condition was further explored by Honovich et al. (2021) to enhance the reliability of automatically generated texts.

The last sentence is cohesive and consistent with the previous sentences but is irrelevant to the overall context and topic of the story.

All in all, Reinhart’s theory outlines conditions encapsulating the fundamental aspects of coherence to determine text coherence. We propose designing NLP tasks to detect these properties.

3 Research Hypotheses and Tasks

At the core of our approach is the implementation of Reinhart’s coherence conditions as computational tasks, using a minimal set of NLP tasks designed to capture the features of cohesion, consistency, and relevance. We hypothesize that a model trained jointly on *all* tasks will detect coherence effectively, and will outperform models trained on each task individually.

To verify this, we define five tasks reflecting these coherence conditions.

The Sentence Reordering (SRO) Task This self-supervised task, proposed by Lapata (2003), involves reordering shuffled sentences to restore their original coherent form. For example, given the following input: “(1) Finally, the parser is evaluated. (2) We develop a useful parser. (3) Then we present our parser. (4) We first describe the older one.” the correct order is (2) → (4) → (3) → (1).

Extending prior work on natural sentence order and coherence (Lin et al., 2011), a model excelling at paragraph reconstruction should capture syntactic and semantic relationships between sentences, reflecting both *cohesion* and *consistency*.

The Discourse-Relation Recognition (DRR) Task Given a pair of sentences (discourse units - DUs), we aim to predict their discourse relation, reflecting notions such as cause and effect, comparison, and contrast. For example, with the following input: “John worked all night. He slept all day today.” the model is expected to detect a relation marker reflecting *contingency* (e.g., ‘so’, ‘hence’).

The discourse relation identification task enhances the model’s ability to connect sentences, addressing the second sub-condition of *cohesion*.

The NP Enrichment (NPE) Task Introduced by Elazar et al. (2022), the NPE task identifies prepositional links between noun phrase (NP) entities. Given NP pairs, it determines the existence of a prepositional relation and identifies the best preposition describing it $p(NP_1 NP_2)$. For a paragraph

with k NP entities, the model outputs the prepositional links for all NP pairs where such a relation exists (or NONE otherwise).

For example, in the paragraph: “Crown Princess Mary of Denmark gives birth to a male child.” there are 4 NPs and thus 12 NP pairs. Sample outputs for these NP pairs are: (1) in(birth, Denmark) and (2) of(birth, male child).

A model trained on this task captures referential links between different parts of the discourse, serving as a proxy for the referential linking sub-condition of *cohesion*.

The Natural Language Inference (NLI) Task The NLI task (Bowman et al., 2015) aims to determine the truth value of a hypothesis based on a given premise. For example, given the premise: “John inspects the uniform of a figure in some East Asian country.” and the hypothesis: “John is sleeping.” the output will be a *contradiction*.

NLI evaluates NLP models’ ability to capture logical relationships between sentences, serving as a proxy for the *consistency* condition.

The Irrelevant Sentence Recognition (ISR) Task We propose a self-supervised task where the model identifies irrelevant sentences in a coherent paragraph. Given a paragraph with N sentences, including one irrelevant sentence s , the model detects and outputs the irrelevant sentence.

For example, given the following input: “(1) Rick is a helpful kid. (2) He does the dishes. (3) He avoids doing his homework. (4) He helps older people.” The irrelevant sentence is (3).

The model is trained to assess sentence relevance to the overall topic and context, acting as a proxy for the *relevance* condition.

Putting It All Together: We propose a Multi-Task Learning (MTL) approach, where a model is jointly trained on these tasks to capture all coherence conditions outlined by Reinhart. This method leverages shared information during training, with the goal of enhancing both overall coherence detection and individual task performance. To assess coherence, we define two types of tasks:

- **The Coherence Scoring Task** To confirm our hypothesis that the proposed model captures coherence, we evaluate its performance on the coherence scoring task, where given a paragraph P the model predicts the coherence score C as a human reader would.

271	• The Coherence Reasoning Task To examine	a coherence classification layer. The SOTA model	319
272	conditions contributing to incoherence (cohe-	for GCDC by Liu et al. (2023) uses a multi-step	320
273	sion, consistency, relevance) beyond a final	approach: identifying document graph structures,	321
274	score, we use the coherence reasoning task	converting subgraphs, constructing corpus-level	322
275	proposed by Maimon and Tsarfaty (2023).	graphs based on shared subgraphs, and encoding	323
276	Given a paragraph P and a new sentence s , the	connections with a GCN.	324
277	model predicts whether s is <i>cohesive</i> , <i>consis-</i>	For CoheSentia Maimon and Tsarfaty (2023)	325
278	<i>tent</i> , or <i>relevant</i> to P using distinct classifiers.	created the SOTA model using a prompt-based ap-	326
279		proach with Flan-T5-large to assess coherence by	327
280	We hypothesize that utilizing the MTL-powered	adding a question at the beginning of each text.	328
281	architecture will improve the results on both coher-		
282	ence scoring and coherence reasoning.		
283		4.2 The Coherence Reasoning Task	329
284		Models: In Classification-Based models, the rea-	330
285	4 Coherence Assessment Setup	soning for each coherence attribute given the para-	331
286		graph P and the new sentence s is predicted using	332
287	Here, we detail the models and experimental setup	a classification head.	333
288	for the coherence assessment tasks we define.	In Generation-Based models, the input includes	334
289		the text with prompts and an output. Prompts and	335
290		outputs for both datasets are in Appendix F.	336
291	4.1 The Coherence Scoring Task		
292	Models: We use two architectures for coherence	Datasets and Evaluation: We evaluate our	337
293	scoring: Classification-Based (BERT (Devlin et al.,	model on the CoheSentia corpus (Maimon and	338
294	2019)) and Generation-Based (T5 (Raffel et al.,	Tsarfaty, 2023), which contains automatically gen-	339
295	2020)). The model predicts for a given text 3-way	erated stories with human annotations for cohesion,	340
296	or 5-way scores, depending on the dataset.	consistency, and relevance. We use precision, re-	341
297	In the Classification-Based models, the coher-	call, and F1 scores for each property.	342
298	ence score C is predicted given the text P using a		
299	classification head.	Baselines: We evaluate our model’s effectiveness	343
300	In Generation-Based models, the input includes	on the CoheSentia dataset by comparing it to the	344
301	the text with dataset-specific prompts and an output.	current SOTA model by (Maimon and Tsarfaty,	345
302	Example prompts and outputs are in Appendix F.	2023), which uses a prompt-based approach with	346
303	Further details on experimental settings are in	Flan-T5-large, adding a question at the beginning	347
304	Appendix B.	of each text to assess coherence.	348
305			
306	Datasets and Evaluation: We evaluate our	5 Task-Specific Experimental Setup	349
307	model on two datasets: GCDC (Lai and Tetreault,	In this section, we elaborate on the modeling of the	350
308	2018) and CoheSentia (Maimon and Tsarfaty,	coherence proxy tasks. For each task, we evalu-	351
309	2023). GCDC includes real-world text from vari-	ate two model variants: Classification-Based and	352
310	ous domains (Clinton emails, Enron emails, Yahoo	Generation-Based, as detailed below. Table 2 sum-	353
311	Answers, Yelp reviews) with coherence scores from	marizes the datasets, evaluation metrics, and key	354
312	1 (not coherent) to 3 (highly coherent). CoheSentia	statistics for each task (see further elaboration in	355
313	features GPT-3 generated stories (fiction and non-	Appendix A). For the Generation-Based models,	356
314	fiction) with scores ranging from 1 to 5. We use	prompts and outputs are detailed in Appendix F.	357
315	the “incremental final score” for CoheSentia stories.		
316	Dataset sizes and splits are detailed in Table 1.	The Sentence Reordering Models For the	358
317	To remain compatible with Lai and Tetreault	Classification-Based models, we adopt the topo-	359
318	(2018), we use accuracy as the metric for evaluating	logical sort architecture from Shrimai Prabhumoye	360
	the final coherence score of the text.	(2020). Each paragraph’s sentence pairs are repre-	361
		sent as triplets $\langle S_i, C_k, S_j \rangle$, indicating whether	362
		S_i precedes or follows S_j . The model has two	363
		stages: a binary classification head that predicts	364
		the pairwise relations and a second stage that pro-	365

Dataset	Split			Per Instance			
	Train	Validation	Test	Max #tokens	Avg #tokens	Max #sent.	Avg #sent.
GCDC	3.6k	800	800	333	156	10	32
CoheSentia	350	75	75	226	150	15	6.5

Table 1: Main Statistics on the Datasets for Coherence Scoring

duces the predicted order using the topological sort algorithm (Tarjan, 1976).

The Generation-Based models use prompts and predict the outputs with the final order.

The Discourse-Relation Recognition Models

In the Classification-Based models, the input consists of a pair of DUs: $\langle DU_1, DU_2 \rangle$. A classification head predicts the discourse relation between them.

In the Generation-Based models, the input is an argument pair, and the model employs a chain-of-thought (CoT) method (Wei et al., 2023) to predict the discourse relation.³ The CoT structure is $\langle \text{connector} \rangle \rightarrow \langle l_1 \text{ relation} \rangle \rightarrow \langle l_2 \text{ relation} \rangle$. That is, the model adopts a three-stage approach to predicting the L_2 relation type. The model first infers the implicit connective, then generalizes it to a broader relation category.

The NP Enrichment Models In the Classification-Based models, we extend the Bi-Affine architecture from Dozat and Manning (2017) to predict preposition relations between NP pairs instead of syntactic dependency labels. NP embeddings are created by pooling tokens representing each NP, and the model head predicts the preposition using the NP’s anchor and complement representations (Figure 4 in Appendix).

The Generation-Based models predict prepositional relations for each NP pair independently. The input consists of the document text and a prompt specifying the NP pair.

The NL Inference Models Classification-Based models predict relations between the premise and hypothesis.

Generation-Based models use prompts and predict outputs for premise-hypothesis pairs.

Irrelevant Sentence Recognition Models The Classification-Based models have two stages. Sentence pairs form triplets $\langle S_i, C_k, S_j \rangle$, where $C_k = 0|1$ indicates relevance. A binary classification head determines the relation. In the second stage,

³CoT detection of discourse relations outperformed simpler prompts in our preliminary experiments.

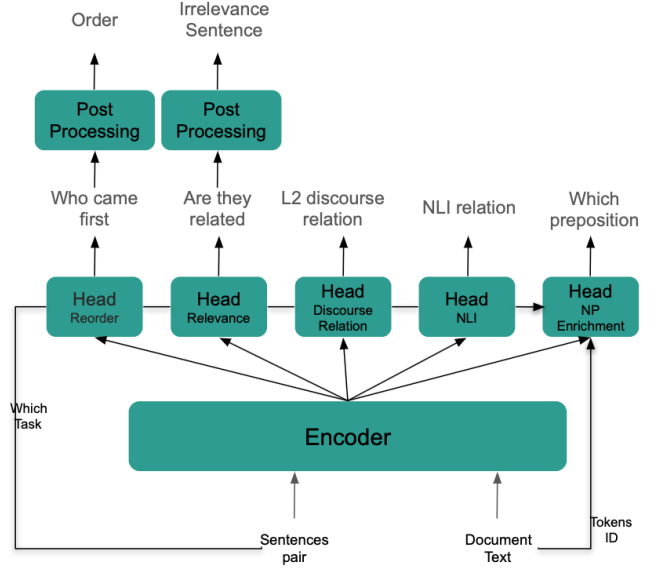


Figure 1: Illustration of the encoder-only model where the input is a pair of sentences (most tasks) or for NPE task the input is a document and the token IDs of different NPs

the sentence with the lowest combined relations score is deemed irrelevant.

Generation-Based models use prompts and predict the irrelevant sentence as the output.

The Overall Joint Architecture To test our hypotheses, we implemented both Classification-Based and Generation-Based models, trained to solve all tasks jointly.

For the Classification-Based variant, we use MTL (Caruana, 1997) with a BERT encoder shared across tasks, each with a unique classification head. Each head predicts task-specific outputs (see Fig. 1). To address forgetting in MTL (Goodfellow et al., 2014), we implement an interleaved training strategy, alternating tasks in batches, ensuring each batch contains samples from a distinct task, effectively mitigating forgetting during MTL training.

For the Generation-Based model, we use the T5 encoder-decoder model, which allows concurrent fine-tuning of multiple tasks using distinct prompts. The prompt structure remains the same as individual task fine-tuning, but batches contain samples from specific tasks, distinct from the previous batch.

Task	Dataset	Metrics	Split			Per Instance			
			Train	Dev	Test	Max #toks	Avg #toks	Max #sent.	Avg #sent.
SRO	RocStories (Mostafazadeh et al., 2016)	PMR (Chen et al., 2016) Acc (Logeswaran et al., 2017)	68k	14k	14k	135	57	5	5
ISR	RocStories (Mostafazadeh et al., 2016)	Accuracy	68k	14k	14k	152	77	6	6
DRR	PDTB3 (Prasad et al., 2019)	Accuracy	17.5k	1.7k	1.5k	556	30	2	2
NPE	TNE (Elazar et al., 2022)	F1, Precision & Recall	3.5k	500	500	284	163	15	6.9
NLI	MNLI (Williams et al., 2018)	Accuracy	393k	7.5k	2.5k	(194,70)	(20,10)	(8,8)	(2,2)

Table 2: The datasets and metrics used for each task and the train/dev/test split size with the max and average number of tokens and sentences. For the NLI task (x,y) refer to the numbers of (premise, hypothesis) respectively

Model	GCDC	CoheSentia
Lai and Tetreault (2018)	57.5	—
SOTA	61.2	35.3
Ours-None (bert-large)	50.2	34.3
Ours-ALL (bert-large)	72.5	55.7
Ours-None (t5-large)	56.3	34.8
Ours-ALL (t5-large)	76.4	62.3
Controlled-nonCoherence (t5-large)	52.8	36.8

Table 3: Accuracy on Coherence Scoring The SOTA for GCDC is by Liu et al. (2023) and for CoheSentia is Maimon and Tsarfaty (2023)

6 Results

We first aim to test the hypothesis that a model jointly trained on tasks reflecting the different coherence conditions will effectively assess coherence. Table 3 shows the coherence scores of our jointly fine-tuned model (Ours-ALL) on the GCDC and CoheSentia datasets, compared to current SOTA models on either dataset. Compared to these models, our jointly fine-tuned model shows significant improvements in coherence scoring, especially on CoheSentia. We observe a 15% and 27% accuracy gain for GCDC and CoheSentia respectively, demonstrating that our proposed approach and selected tasks effectively contribute to capturing fundamental aspects of coherence.

We further analyze the contribution of the proxy tasks (Ours-All) by comparing it to a model without such fine-tuning (Ours-None) to isolate performance gains. As evident in Table 3 these tasks dramatically enhance performances. These results are supported by further qualitative analysis in Appendix D.

Next we analyze the MTL model’s success on assessing the coherence conditions (cohesion, consistency, relevance), by fine-tuning on the coherence reasoning task. We compare the results to SOTA from Maimon and Tsarfaty (2023), who fine-tuned the Flan-T5 model with a simple prompt, and to our model without initial fine-tuning on the coherence proxy tasks (Ours-None). Table 4 summarizes the

coherence reasoning task results for all attributes and metrics. Our model achieves SOTA performance across all coherence conditions, demonstrating the efficacy of our approach.

Finally, We evaluate the task-specific performance of models trained with either individual or joint fine-tuning on the proxy tasks, using both the Classification-Based and Generation-Based variations. Results are in Table 5, alongside comparisons to current SOTA on these benchmarks. Our findings consistently show that Generation-Based models outperform Classification-Based ones. More importantly, joint fine-tuning across all tasks consistently surpasses individual fine-tuning, particularly in the SRO, ISR, and DRR tasks, where it leads to significant performance improvements and even surpasses SOTA benchmarks. For the NPE task, joint fine-tuning achieves substantial recall gains, though precision falls short of SOTA results, offering a more balanced performance. The exception is the NLI task, where our model’s performance is lower than SOTA.

In summary, our MTL model outperforms single-task models on all tasks, achieving SOTA results except for NLI, in line with our hypothesis on the benefits of the joint architecture. Moreover, our MTL model, jointly trained on coherence proxy tasks, significantly improves performance, enhancing coherence scoring for both datasets and excelling in coherence reasoning, in line with the second part of our hypothesis.

7 Analysis

7.1 The Effect of Different Tasks on Coherence Scoring

This section examines how fine-tuning on diverse subsets of coherence proxy tasks affects coherence scoring. We fine-tune models on various combinations of these tasks, then perform final fine-tuning and evaluation on the coherence scoring task.

Model	Cohesion			Consistency			Relevance		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
SOTA	72.4	72.1	72.2	59.6	67.5	63.3	56.4	74.6	59.5
Ours-None (bert-large)	66.4	59.4	62.7	60.4	56.5	59.6	49.2	49.9	49.5
Ours-ALL (bert-large)	74.7	70.5	72.5	70.6	68.2	69.3	59.8	61.1	60.4
Ours-None (t5-large)	81.1	80.3	80.7	60.4	62.6	61.5	48.1	49.6	48.8
Ours-ALL (t5-large)	83.1	83.2	83.1	78.5	80.3	79.4	70.8	76.9	73.7

Table 4: Results for Coherence Reasoning Task. The SOTA is by [Maimon and Tsarfaty \(2023\)](#)

Model	SRO		ISR	DRR	NPE			NLI
	PMR	ACC	Accuracy	Accuracy	F1	P	R	Accuracy
SOTA	81.9	90.8	-	64.7	64.0	80.5	53.1	92.0
Ours-Individual (bert-large)	51.8	69.5	60.4	60.0	53.1	67.1	44.0	87.4
Ours-ALL (bert-large)	67.1	83.2	78.6	65.7	64.4	79.8	54.2	90.2
Ours-Individual (t5-large)	75.7	87.8	80.4	64.8	59.8	68.5	53.1	89.9
Ours-ALL (t5-large)	83.8	92.1	82.2	67.3	76.7	76.7	76.7	91.5

Table 5: Results for all proxy tasks compared to SOTA performances. The SOTA model for SRO is ReBART ([Basu Roy Chowdhury et al., 2021](#)), for DRR is Contrastive Learning ([Long and Webber, 2023](#)), for NPE is TNE ([Elazar et al., 2022](#)) and for NLI T5-11B ([Raffel et al., 2020](#))

Figure 2 shows the impact of fine-tuning proxy coherence tasks on coherence scoring performance. Models fine-tuned on any coherence proxy task outperform those without fine-tuning (Ours-None), highlighting their effectiveness. Performance generally improves with more tasks, especially beyond three, indicating cumulative benefits.

Interestingly, NLI fine-tuning significantly enhances performance, likely due to its role in improving the model’s ability to capture consistency, crucial for coherence assessment. Additionally, ISR fine-tuning is more impactful when combined with other tasks. These findings underscore the importance of task selection and task interaction during fine-tuning for optimal coherence scoring.

7.2 Impact of Non-Coherence Tasks Fine-Tuning

In this Section we aim to empirically refute a possible hypothesis that the joint ALL model outperforms the NONE model simply due to its complexity, regardless of the nature of the tasks used (i.e., tasks reflecting coherence conditions).

To this end, we compare the performance of our fine-tuned on coherence-tasks model (Ours-ALL) with a model fine-tuned on three tasks orthogonal to coherence, followed by fine-tuning on the coherence scoring task: (i) Machine Translation (MT): We sample 15k instances from the WMT14 dataset ([Bojar et al., 2014](#)). (ii) Named-Entity Recognition (NER): We use the Conll2003 dataset ([Tjong Kim Sang and De Meulder, 2003](#)) containing annotations for 14k instances. (iii) Part-of-Speech (POS), using the same Conll2003 dataset containing the POS tags as well.

Model In these experiments, we used the T5-large model as the basis, employing specific prompt and output designs for each task. For the NER and POS tasks, we adapted the “Sentinel + Tag” architecture by [Raman et al. \(2022\)](#). Detailed prompts and sample outputs are provided in Appendix F.

Results Following the same procedure as in our main experiments, fine-tuned models were assessed for coherence scoring using the GCDC and CoheSentia benchmarks (detailed in Table 3).

Fine-tuning on tasks orthogonal to coherence yielded minimal to no improvements over the baseline (Ours-None) and significantly underperformed compared to our final MTL model (Ours-ALL). This highlights the importance of coherence-specific proxy tasks for effective coherence detection, as unrelated tasks can hinder performance.

7.3 Cross-Domain Generalization

Since GCDC and CoheSentia present different domains and writing styles, we evaluate model generalizability by fine-tuning on one dataset and assessing coherence on the other. Table 6 presents the results for our MTL model (Ours-ALL) and the non-coherence fine-tuned model (Ours-None) under three settings: fine-tuning on CoheSentia only, GCDC only, and both combined.

Results demonstrate performance gains across domains, highlighting the generalizability of our method. Combining data improves performance, with Ours-ALL showing a 12% and 14% error reduction on CoheSentia and GCDC, respectively, compared to in-domain scenarios, underscoring the utility and transferability of the learned features.

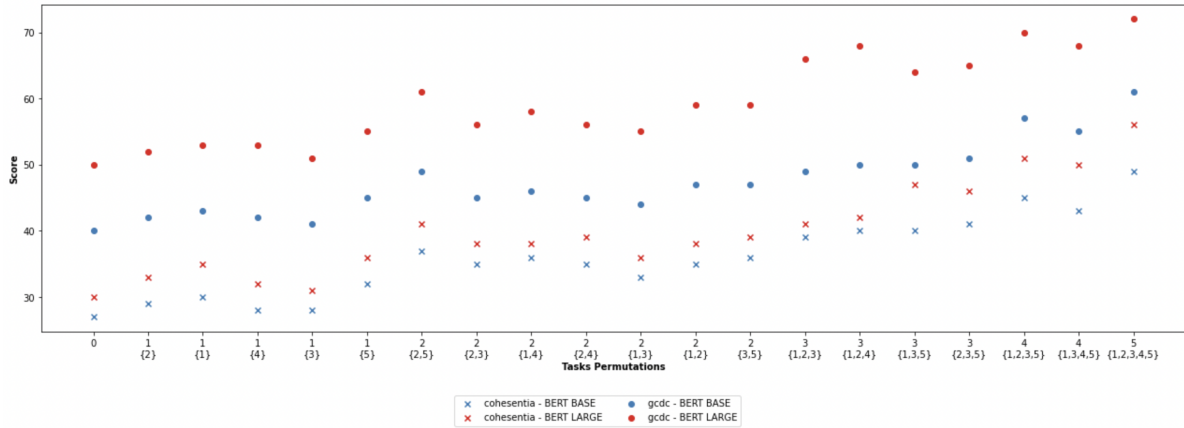


Figure 2: Accuracy for Coherence Scoring Task for both GCDC and CoheSentic with different proxy coherence task-subsets. The labels are tasks IDs (1-SRO, 2-ISR, 3-DRR, 4-NPE, 5-NLI)

Model	GCDC	CoheSentic
Ours-None-CoheSentic	52.8	34.8
Ours-None-GCDC	56.3	28.5
Ours-None-Both	57.5	35.4
Ours-ALL-CoheSentic	71.8	62.3
Ours-ALL-GCDC	76.4	59.5
Ours-ALL-Both	79.8	66.7

Table 6: Accuracy on coherence scoring on both datasets when fine-tuned based on T5-model on only one dataset

7.4 Effects of Different Tasks on One Another

We investigate the impact of fine-tuning on various coherence task subsets on individual task performance. The model was trained on different task combinations with increasing numbers of tasks and evaluated on each task separately.

Figure 3 shows consistent performance gains in the Sentence Reordering (SRO) task for BERT models as more tasks are jointly fine-tuned (see Appendix E for other tasks). This supports our hypothesis that shared information among coherence tasks enhances individual task performance.

The impact of specific tasks varies; for instance, DRR minimally affects SRO, likely due to limited training data. Notably, NLI significantly influences the performance of various tasks. The ISR task notably improves performance on other tasks, suggesting its effectiveness in capturing relevance errors, crucial for coherence assessment. We thus emphasize the introduction of this self-supervised ISR task and advocate for its exploration in future research to enhance coherence assessment.

The overall performance trends are similar for both BERT-base and BERT-large models, indicating that the impact of specific tasks is consistent regardless of model size.

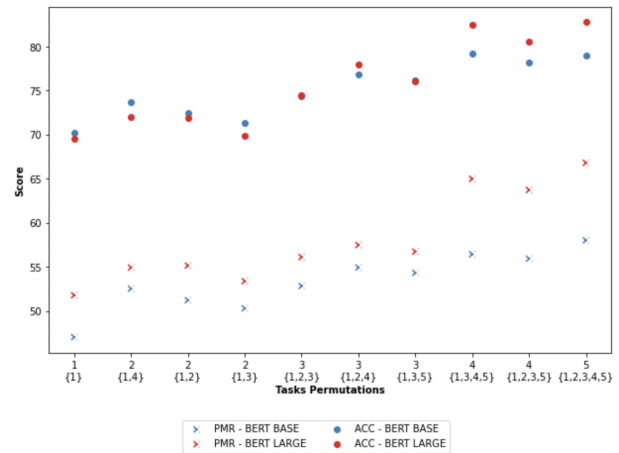


Figure 3: Results for SRO task, for different subsets of coherence tasks fine-tuned upon. The labels are the number of tasks and in curly brackets which tasks (1 - SRO, 2 - ISR, 3 - DRR, 4 - NPE, 5 - NLI)

8 Conclusion and Future Work

In this paper we propose a new coherence modeling method, based on Reinhart (1980)’s theory which defines the conditions needed for coherence: cohesion, consistency, and relevance. We use five key NLP tasks as proxies for these conditions, and train an MTL model on them jointly. Our unified coherence model achieves SOTA results on these individual tasks, and moreover it excels in coherence scoring for both real-world and generated texts. We conjecture that this framework will enhance NLP systems’ ability to quantify and evaluate text quality automatically. Future follow-up research will focus on using these conditions for improved coherent-text generation, and for detecting particular causes of incoherence automatically. Our code and models are publicly available, to encourage further research on coherence scoring and coherence reasoning.

612 Limitations

613 While this work advances the modeling and auto-
614 matic evaluation of coherence, limitations exist that
615 suggest promising avenues for future research.

616 **Dataset Limitations** Existing coherence evalua-
617 tion datasets like GCDC and CoheSentia, along
618 with datasets for our proxy tasks, primarily fo-
619 cus on relatively short texts. To address this, we
620 analyzed the performance of our MTL models
621 (Ours-ALL) and the non-coherence version (Ours-
622 None) on GCDC and CoheSentia across various
623 text lengths after fine-tuning for coherence scoring
624 (see Figure 7 in the Appendix). As expected, for
625 both models and datasets, accuracy decreased with
626 longer texts, highlighting the increased difficulty of
627 assigning coherence scores for complex passages.
628 This observation aligns with recent work suggest-
629 ing that while LLMs can handle longer texts, their
630 reasoning abilities might decline with increasing
631 text length (Levy et al., 2024).

632 **Focus on Short Texts** Our current study focused
633 on short texts (≤ 512 tokens). The effectiveness
634 of our approach on longer documents remains an
635 open question for future exploration. We hypothe-
636 size that incorporating coherence proxy tasks could
637 benefit the model’s performance on longer texts,
638 but further investigation is necessary.

639 Acknowledgements

640 References

641 Hongxiao Bai and Hai Zhao. 2018. [Deep enhanced rep-](#)
642 [resentation for implicit discourse relation recognition.](#)
643 *Preprint*, arXiv:1807.05154.

644 Regina Barzilay and Mirella Lapata. 2008. Modeling
645 local coherence: An entity-based approach. *Compu-*
646 *tational Linguistics*, 34(1):1–34.

647 Somnath Basu Roy Chowdhury, Faeze Brahman, and
648 Snigdha Chaturvedi. 2021. [Is everything in order? a](#)
649 [simple way to order sentences.](#) In *Proceedings of the*
650 *2021 Conference on Empirical Methods in Natural*
651 *Language Processing*, pages 10769–10779, Online
652 and Punta Cana, Dominican Republic. Association
653 for Computational Linguistics.

654 Ondrej Bojar, Christian Buck, Christian Federmann,
655 Barry Haddow, Philipp Koehn, Johannes Leveling,
656 Christof Monz, Pavel Pecina, Matt Post, Herve Saint-
657 Amand, Radu Soricut, Lucia Specia, and Ale s Tam-
658 chyna. 2014. [Findings of the 2014 workshop on](#)
659 [statistical machine translation.](#) In *Proceedings of the*
660 *Ninth Workshop on Statistical Machine Translation*,
661 pages 12–58, Baltimore, Maryland, USA. Associa-
662 tion for Computational Linguistics.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, 663
and Christopher D. Manning. 2015. [A large anno-](#) 664
[tated corpus for learning natural language inference.](#) 665
Preprint, arXiv:1508.05326. 666

Rich Caruana. 1997. Multitask learning. *Machine* 667
learning, 28:41–75. 668

Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 669
2016. [Neural sentence ordering.](#) *Preprint*, 670
arXiv:1607.06952. 671

Ido Dagan, Oren Glickman, and Bernardo Magnini. 672
2005. [The pascal recognising textual entailment chal-](#) 673
[lenge.](#) pages 177–190. 674

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 675
Kristina Toutanova. 2019. [BERT: Pre-training of](#) 676
[deep bidirectional transformers for language under-](#) 677
[standing.](#) In *Proceedings of the 2019 Conference of* 678
the North American Chapter of the Association for 679
Computational Linguistics: Human Language Tech- 680
nologies, Volume 1 (Long and Short Papers), pages 681
4171–4186, Minneapolis, Minnesota. Association for 682
Computational Linguistics. 683

Teun A Van Dijk. 1979. Pragmatic connectives. *Dis-* 684
course processes, 3:447–456. 685

Timothy Dozat and Christopher D. Manning. 2017. 686
[Deep biaffine attention for neural dependency pars-](#) 687
[ing.](#) *Preprint*, arXiv:1611.01734. 688

Yanai Elazar, Victoria Basmov, Yoav Goldberg, and 689
Reut Tsarfaty. 2022. [Text-based np enrichment.](#) 690
Preprint, arXiv:2109.12085. 691

Vanessa Wei Feng, Ziheng Lin, , and Graeme Hirst. 692
2014. [The impact of deep hierarchical discourse](#) 693
[structures in the evaluation of text coherence.](#) *Pro-* 694
ceedings of COLING 2014, the 25th International 695
Conference on Computational Linguistics: Technical 696
Papers, 1:940–949. 697

Thomas Freeman and CE Gathercole. 1966. Persevera- 698
tion—the clinical symptoms—in chronic schizophre- 699
nia and organic dementia. *The British Journal of* 700
Psychiatry, 112(482):27–32. 701

T Givon. 1995. Coherence in text vs. coherence in mind. 702
Coherence in spontaneous text, pages 31–59. 703

Ian J. Goodfellow, Mehdi Mirza, Xia Da, Aaron C. 704
Courville, and Yoshua Bengio. 2014. [An empirical](#) 705
[investigation of catastrophic forgetting in gradient-](#) 706
[based neural networks.](#) In *2nd International Confer-* 707
ence on Learning Representations, ICLR 2014, Banff, 708
AB, Canada, April 14-16, 2014, Conference Track 709
Proceedings. 710

Jian Guan, Xiaoxi Mao, Changjie Fan, Zitao Liu, Wen- 711
biao Ding, and Minlie Huang. 2021. [Long text gener-](#) 712
[ation by modeling sentence-level and discourse-level](#) 713
[coherence.](#) *Preprint*, arXiv:2105.08963. 714

715	Camille Guinaudeau and Michael Strube. 2013. Graph-based local coherence modeling . In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 93–103, Sofia, Bulgaria. Association for Computational Linguistics.	770
716		771
717		772
718		773
719		
720		
721	M.A.K. Halliday and Ruqaiya Hasan. 1976. <i>Cohesion in English</i> . Longman, Dallas, Texas.	
722		
723	Jerry R Hobbs. 1979. Coherence and coreference. <i>Cognitive science</i> , 3(1):67–90.	
724		
725	Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. 2021. q^2: Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering . <i>Preprint</i> , arXiv:2104.08202.	
726		
727		
728		
729		
730	Yufang Hou. 2021. End-to-end neural information status classification . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1377–1388, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
731		
732		
733		
734		
735	Yangfeng Ji and Jacob Eisenstein. 2015. One vector is not enough: Entity-augmented distributed semantics for discourse relations . <i>Transactions of the Association for Computational Linguistics</i> , 3:329–344.	
736		
737		
738		
739	Aravind K Joshi and Scott Weinstein. 1981. Control of inference: Role of some aspects of discourse structure-centering . <i>IJCAI</i> , pages 385–387.	
740		
741		
742	Philippe Laban, Luke Dai, Lucas Bandarkar, and Marti A. Hearst. 2021. Can transformer models measure coherence in text? re-thinking the shuffle test . <i>Preprint</i> , arXiv:2107.03448.	
743		
744		
745		
746	Alice Lai and Joel Tetreault. 2018. Discourse coherence in the wild: A dataset, evaluation and methods . In <i>Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue</i> . Association for Computational Linguistics, 1:214–223.	
747		
748		
749		
750		
751	Mirella Lapata. 2003. Probabilistic text structuring: Experiments with sentence ordering. <i>EMNLP/ICNLP</i> , pages 2273–2283.	
752		
753		
754	Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. Same task, more tokens: the impact of input length on the reasoning performance of large language models . <i>Preprint</i> , arXiv:2402.14848.	
755		
756		
757		
758	Li Liang, Zheng Zhao, and Bonnie Webber. 2020. Extending implicit discourse relation recognition to the PDTB-3 . In <i>Proceedings of the First Workshop on Computational Approaches to Discourse</i> , pages 135–147, Online. Association for Computational Linguistics.	
759		
760		
761		
762		
763		
764	Ziheng Lin, Hwee Tou Ng, , and Min-Yen Kan. 2011. Automatically evaluating text coherence using discourse relations . <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> . Association for Computational Linguistics, 1:997–1006.	
765		
766		
767		
768		
769		
	Wei Liu, Xiyan Fu, and Michael Strube. 2023. Modeling structural similarities between documents for coherence assessment with graph convolutional networks . <i>Preprint</i> , arXiv:2306.06472.	774
		775
		776
		777
	Xin Liu, Jiefu Ou, Yangqiu Song, and Xin Jiang. 2020. On the importance of word and sentence representation learning in implicit discourse relation classification . <i>Preprint</i> , arXiv:2004.12617.	
	Lajanugen Logeswaran, Honglak Lee, and Dragomir Radev. 2017. Sentence ordering and coherence modeling using recurrent neural networks . <i>Preprint</i> , arXiv:1611.02654.	778
		779
		780
		781
	Wanqiu Long and Bonnie Webber. 2023. Facilitating contrastive learning of discourse relational senses by exploiting the hierarchy of sense relations . <i>Preprint</i> , arXiv:2301.02724.	782
		783
		784
		785
	Aviya Maimon and Reut Tsarfaty. 2023. Cohesentia: A novel benchmark of incremental versus holistic assessment of coherence in generated texts . <i>CoRR</i> , abs/2310.16329.	786
		787
		788
		789
	William C Mann and Sandra A Thompson. 1988. Rhetorical structure theory: Toward a functional theory of text organization. <i>Text Interdisciplinary Journal for the Study of Discourse</i> , 1:243–281.	790
		791
		792
		793
	Mohsen Mesgar and Michael Strube. 2016. Lexical coherence graph modeling using word embeddings . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1414–1423, San Diego, California. Association for Computational Linguistics.	794
		795
		796
		797
		798
		799
		800
	Mohsen Mesgar and Michael Strube. 2018. A neural local coherence model for text quality assessment . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4328–4339, Brussels, Belgium. Association for Computational Linguistics.	801
		802
		803
		804
		805
		806
	Miltsakaki, Eleni, Kukich, and Karen. 2000. Automated evaluation of coherence in student essays.	807
		808
	Eleni Miltsakaki, Rashmi Prasad, Aravind Joshi, and Bonnie Webber. 2004. The Penn Discourse Treebank . In <i>Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)</i> , Lisbon, Portugal. European Language Resources Association (ELRA).	809
		810
		811
		812
		813
		814
	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and evaluation framework for deeper understanding of commonsense stories . <i>Preprint</i> , arXiv:1604.01696.	815
		816
		817
		818
		819
		820
	Emily Pitler, Annie Louis, and Ani Nenkova. 2009. Automatic sense prediction for implicit discourse relations in text . In <i>Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th</i>	821
		822
		823
		824

825				
826				
827				
828				
829	Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Milt-	Adina Williams, Nikita Nangia, and Samuel Bowman.		880
830	sakaki, Livio Robaldo, Aravind Joshi, and Bonnie	2018. A broad-coverage challenge corpus for sen-		881
831	Webber. 2008. The Penn Discourse TreeBank 2.0.	tence understanding through inference. In <i>Proceed-</i>		882
832	In <i>Proceedings of the Sixth International Conference</i>	<i>ings of the 2018 Conference of the North American</i>		883
833	<i>on Language Resources and Evaluation (LREC'08),</i>	<i>Chapter of the Association for Computational Lin-</i>		884
834	Marrakech, Morocco. European Language Resources	<i>guistics: Human Language Technologies, Volume 1</i>		885
835	Association (ELRA).	<i>(Long Papers)</i> , pages 1112–1122. Association for		886
836	Rashmi Prasad, Bonnie Webber, Alan Lee, and Aravind	Computational Linguistics.		887
837	Joshi. 2019. Penn Discourse Treebank Version 3.0.			
838	Lianhui Qin, Zhisong Zhang, and Hai Zhao. 2016. Shal-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien		888
839	low discourse parsing using convolutional neural net-	Chaumond, Clement Delangue, Anthony Moi, Pier-		889
840	work. In <i>Proceedings of the CoNLL-16 shared task,</i>	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtow-		890
841	pages 70–77, Berlin, Germany. Association for Com-	icz, Joe Davison, Sam Shleifer, Patrick von Platen,		891
842	putational Linguistics.	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,		892
843	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Teven Le Scao, Sylvain Gugger, Mariama Drame,		893
844	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	Quentin Lhoest, and Alexander M. Rush. 2020. Hug-		894
845	Wei Li, and Peter J. Liu. 2020. Exploring the limits	gingface's transformers: State-of-the-art natural lan-		895
846	of transfer learning with a unified text-to-text trans-	guage processing. <i>Preprint</i> , arXiv:1910.03771.		896
847	former. <i>Preprint</i> , arXiv:1910.10683.			
848	Karthik Raman, Iftekhar Naim, Jiecao Chen, Kazuma	Wei Xiang, Bang Wang, Lu Dai, and Yijun Mo. 2022.		897
849	Hashimoto, Kiran Yalasangi, and Krishna Srinivasan.	Encoding and fusing semantic connection and lin-		898
850	2022. Transforming sequence tagging into a seq2seq	guistic evidence for implicit discourse relation rec-		899
851	task. <i>Preprint</i> , arXiv:2203.08378.	ognition. In <i>Findings of the Association for Compu-</i>		900
852	Tanya Reinhart. 1980. <i>Conditions for text coherence.</i>	<i>tational Linguistics: ACL 2022</i> , pages 3247–3257,		901
853	<i>Poetics Today 1(4): 16t-180</i> , volume 1.	Dublin, Ireland. Association for Computational Lin-		902
854	Alan W Black Shrimai Prabhumoye, Ruslan Salakhut-	guistics.		903
855	dinov. 2020. Topological sort for sentence ordering.			
856	Swapna Somasundaran, Jill Burstein, and Martin	Jingjing Xu, Xuancheng Ren, Yi Zhang, Qi Zeng,		904
857	Chodorow. 2014. Lexical chaining for measuring	Xiaoyan Cai, and Xu Sun. 2018. A skeleton-		905
858	discourse coherence quality in test-taker essays. <i>Pro-</i>	based model for promoting coherence among sen-		906
859	<i>ceedings of COLING 2014, the 25th International</i>	tences in narrative story generation. <i>Preprint</i> ,		907
860	<i>Conference on Computational Linguistics: Technical</i>	arXiv:1808.06945.		908
861	<i>Papers</i> , 1:950–961.	Sanghyun Yi, Rahul Goel, Chandra Khatri, Alessan-		909
862	Robert Endre Tarjan. 1976. Edge-disjoint spanning	dra Cervone, Tagyoung Chung, Behnam Hedayatnia,		910
863	trees and depth-first search. <i>Acta Informatica,</i>	Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-		911
864	6(2):171–185.	Tur. 2019. Towards coherent and engaging spoken		912
865	Yi Tay, Minh C. Phan, Luu Anh Tuan, and Siu Cheung	dialog response generation using automatic conversa-		913
866	Hui. 2017. Skipflow: Incorporating neural coher-	tion evaluators. <i>Preprint</i> , arXiv:1904.13015.		914
867	ence features for end-to-end automatic text scoring.			
868	<i>Preprint</i> , arXiv:1711.04981.	A Tasks Specific Experimental Settings		915
869	Erik F. Tjong Kim Sang and Fien De Meulder.			
870	2003. Introduction to the CoNLL-2003 shared task:	In this section, we further elaborate on the datasets		916
871	Language-independent named entity recognition. In	and evaluation metrics used for each one of the		917
872	<i>Proceedings of the Seventh Conference on Natural</i>	coherence proxy tasks.		918
873	<i>Language Learning at HLT-NAACL 2003</i> , pages 142–			
874	147.	A.1 The Sentence Reordering Task Setup		919
875	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten			
876	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	Topological Sort: A topological sort (Tarjan,		920
877	Denny Zhou. 2023. Chain-of-thought prompting elic-	1976) linearly orders vertices in a DAG. The al-		921
878	its reasoning in large language models. <i>Preprint</i> ,	gorithm is presented in Algo A.1.		922
879	arXiv:2201.11903.	[t] Input: A digraph G with n vertices		923
		Output: A topological ordering v_1, v_2, \dots, v_n of G.		924
		L \leftarrow Empty list that will contain the sorted nodes		925
		S \leftarrow Set of all nodes with no incoming edge S is		926
		not empty remove a node n from S add n to L each		927
		node m with an edge e from n to m remove edge		928
		e from the graph m has no other incoming edges		929
		insert m into S graph has edges error (graph has at		930
		least one cycle) L (a topologically sorted order)		931

Dataset: We use the ROCStories (Mostafazadeh et al., 2016) dataset (Licence ID is CC-BY 4.0.) which contains 5-sentence stories. We use the standard 85:15 train/test split and randomly select a subset of the train for validation.

Evaluation: We use two common evaluation metrics for the reordering task:⁴

- Perfect Match Ratio (PMR): Chen et al. (2016) calculate the percentage of samples for which the entire sequence was correctly predicted.

$$PMR = \frac{1}{N} \sum_{i=1}^N 1\{\hat{O}^i = O^i\}$$

- Sentence Accuracy (Acc): Logeswaran et al. (2017) calculate the percentage of sentences for which their absolute position was correctly predicted.

$$Acc = \frac{1}{N} \sum_{i=1}^N \frac{1}{v_i} \sum_{j=1}^{v_i} 1\{\hat{O}_j^i = O_j^i\}$$

A.2 The Discourse-Relation Recognition Task Setup

Dataset: We use the Penn Discourse TreeBank3 (PDTB3) Level 2 dataset (Miltsakaki et al., 2004; Prasad et al., 2008; Liang et al., 2020). We only used labels with more than 100 instances, which leaves us with 14 senses from L_2 . The variability of data splits used in the literature is substantial, therefore, we follow earlier work by Ji and Eisenstein (2015); Bai and Zhao (2018); Liu et al. (2020); Xiang et al. (2022) using Sections 2-20, 0-1 and 21-22 for training, validation and testing respectively. When multiple annotated labels are present, we adopt the approach described by Qin et al. (2016) and consider them as distinct instances during the training phase. During testing, if a prediction matches any of the reference labels, it is considered correct.

Evaluation: We use the accuracy metric on the number of sentence pairs the model correctly predicted the L_2 discourse relation:

$$Accuracy = \frac{1}{N} \sum_{i=1}^N 1\{\hat{R}^i = R^i\}$$

A.3 The NP Enrichment Task Setup

Token Classification Head: Figure 4 is an illustration of the token classification head for the NPE task.

⁴There are 5 metrics, we used the most common 2.

Dataset: We use the TNE dataset (Elazar et al., 2022) (Licence Free) which contains documents and relations between every noun pair in it (with a total number of nouns of 190k and a total number of NP relations of 1M). There are 28 possible relations (including ‘no relation’). This dataset’s advantage is that it contains real-world long paragraphs. As in the original publication split the data at the document level.

The distribution of the possible preposition between pair of nouns in TNE dataset is in Figure 5

Evaluation: We report precision, recall & F1 on NP pairs with prepositional links between them.

A.4 The NL Inference (NLI) Task Setup

Dataset and Evaluation: We use the MNLI dataset (Williams et al., 2018) (Licence ID CC-BY-3.0). with the accuracy metric on the amount of hypothesis-premise pairs that the model correctly predicts their relation R :

$$Accuracy = \frac{1}{N} \sum_{i=1}^N 1\{\hat{R}^i = R^i\}$$

A.5 The Irrelevant Sentence Recognition Task Setup

Dataset: We again use ROCStories as in sentence reordering. Each story within the ROCStories dataset was augmented with a single, randomly inserted sentence. The irrelevant sentence for each story was randomly selected from the entire ROCStories dataset, with the sole constraint that it contained entities present in the target story. Both this and the sentence reordering task leverage the same benchmark, retaining the same train/dev/test splits.

Evaluation: We use the accuracy metric on the percentage of paragraphs where the model correctly detected the irrelevant sentence S :

$$Accuracy = \frac{1}{N} \sum_{i=1}^N 1\{\hat{S}^i = S^i\}$$

A.6 Overall Experimental Settings

We trained each model three times, reporting the mean performance. Training utilized multiple Tesla V100 GPUs (up to 4) with 32GB memory each. For each architecture, the settings are:

1. Classification-Based: BERT (base and large) served as the encoder with fine-tuning across

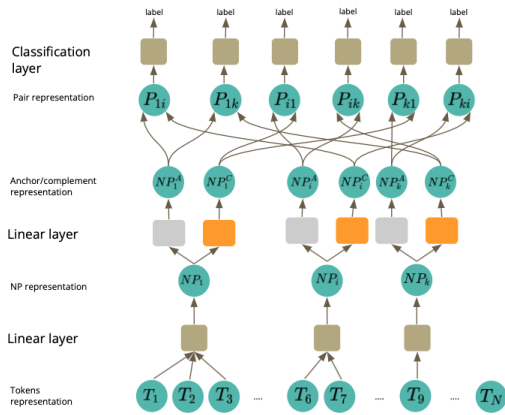


Figure 4: Illustration of the token head which contains several stages: starting with (1) embedding for each token in the text, (2) creating an embedding for each NP when it acts as the complement and the anchor separately, (3) a representation for each NP pair and finally (4) a classification layer

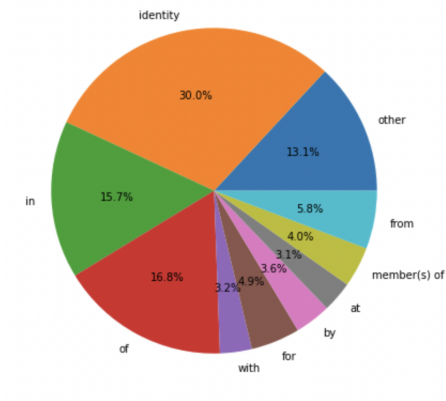


Figure 5: Distribution of the main prepositions in the NP Enrichment test set

all layers. We used Adam optimizer with a learning rate of $5e-5$ and a dropout of 0.5. For tasks requiring classification (SRO, ISR, DRR, NLI), we employed a linear classification head with 512 hidden dimensions and 0.3 dropouts. The NPE utilized a different head structure (details omitted for brevity). Cross-Entropy loss was used for all datasets.

2. Generation-Based: T5 (base and large) models were used as the backbone. Training employed Adam optimizer with a learning rate of $5e-5$. Models were trained with task-specific prompts and corresponding ground truth labels for supervised learning.

Both architectures shared the following hyper-parameters: fine-tuning for 3 epochs with early stopping, batch size of 4, and gradient accumulation steps of 2. The hyper-parameters were chosen using parameters-grid. Our code is based on the Huggingface library (Wolf et al., 2020).

B Coherence Assessment Experimental Settings

For each architecture, the settings are:

1. Classification-Based (BERT base and large): Encoder with fine-tuning across all layers, Adam optimizer (learning rate $5e-4$), dropout (0.3). Each dataset used a linear classification head (512 hidden dimensions, 0.1 dropout). Cross-Entropy loss was used.
2. Generation-Based (T5 base and large): Encoder-decoder architecture, Adam optimizer (learning rate $2e-5$). Inputs included prompts specific to each dataset (GCDC or CoheSenta) and the paragraph text.

The models share hyperparameters: 50 epochs with early stopping (accuracy), batch size of 4, and gradient accumulation steps of 2. We employed 10-fold cross-validation on both datasets (following Lai and Tetreault (2018)) using a single Tesla V100 GPU with 32GB memory. The hyper-parameters were chosen using parameters-grid. Our code is based on the Huggingface library (Wolf et al., 2020).

C Text Length vs. Coherence Score

The accuracy of the models on both coherence datasets based on different lengths is in Figure 7.

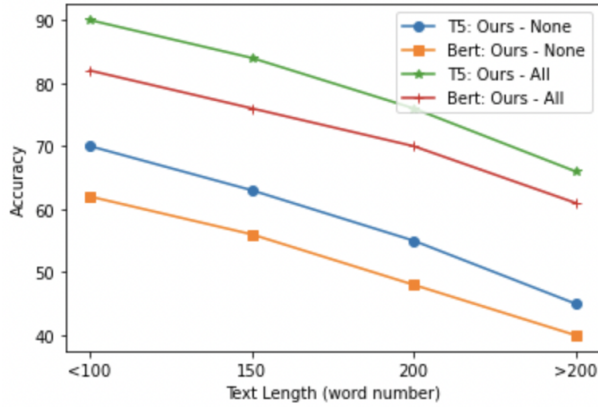


Figure 6: Accuracy For GCDC based on number of words

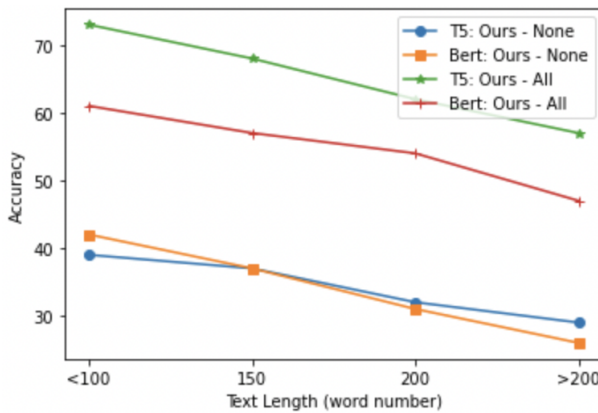


Figure 7: Accuracy For CoheSentia based on number of words

D Qualitative Analysis

D.1 Qualitative Analysis

To gain qualitative insights, we sampled 50 misclassified examples by SOTA models, from CoheSentia and GCDC. We then assessed these examples on various models, including our MTL model (Ours-ALL) and the non-coherence fine-tuning version (Ours-None).

For CoheSentia, the previous SOTA models favor extreme scores, likely due to training data imbalance. Our model exhibits greater robustness, predicting a more balanced distribution of scores. Figure 8a and Table 7 present an example of a text from the CoheSentia dataset and the predictions of the models. In this example, the base model (Ours-None) failed on coherence prediction, while our final model (Ours-ALL) succeeded. Figure 8b presents an example of text from GCDC dataset and Table 8 the predictions of different models on the coherence scoring task. This example highlights

Model	Prediction
Ground Truth	Medium
SOTA	High
Ours-None (BERT-large)	High
Ours-None (T5-large)	High
Ours-ALL (BERT-large)	Medium
Ours-ALL (T5-large)	Medium

Table 7: Predicted Coherence scores for the text in Figure 8a

Model	Prediction
Ground Truth	Low
SOTA	Medium
Ours-None (BERT-large)	Medium
Ours-None (T5-large)	Medium
Ours-ALL (BERT-large)	Low
Ours-ALL (T5-large)	Low

Table 8: Predicted Coherence scores for the text in Figure 8b

a complex case with cohesion and relevance violations. Both the baseline and ISR-trained models missed this issue, while our MTL model achieved accurate prediction.

E Results for Subsets of Tasks

Figures 9a, 9b, 9c, 10a and 10b visualize the performance of coherence proxy tasks across fine-tuning settings for BERT-base and BERT-large models. It highlights how subsets of tasks impacts target task performance.

F T5 Prompts and Outputs for Different Tasks

In Table 9 we detail the various prompts used for fine-tuning T5 models on all explored tasks in this work.

In Table 10 we detail the various outputs used for fine-tuning T5 models on all explored tasks in this work.

‘Shed been a widow for over two years and was starting to lose hope of ever seeing her husband again. One day, she received an email from him asking if he could come out for fun at her funeral. She skeptically agreed but soon found herself enjoying his company more than she could have imagined. As they went around the Neapolitan town where she belonged, it quickly became clear that their bond was even stronger then before. they laughed and danced together like teenagers on celebrated days like this one. It seemed equitable that he should be there too. as long as he didnt mind being the man in attendance at her burial pyre.’

(a) CoheSentia

‘Guy from Mexico is in NY and is cooperating. Discussions with him continue this am. Since he is cooperating, no move to court or to presentment scheduled yet. \n\nMexican support has been excellent throughout. Alice has call sheet for Espinosa — call can take place whenever its convenient for you later this morning (Espinosa is apparently out on West Coast, but Ops could confirm time difference).

Holding off for now on other calls that rest of us would make (Saudis, et al), pending further developments in NY.
Will let you know as soon as we have more.’

(b) GCDC

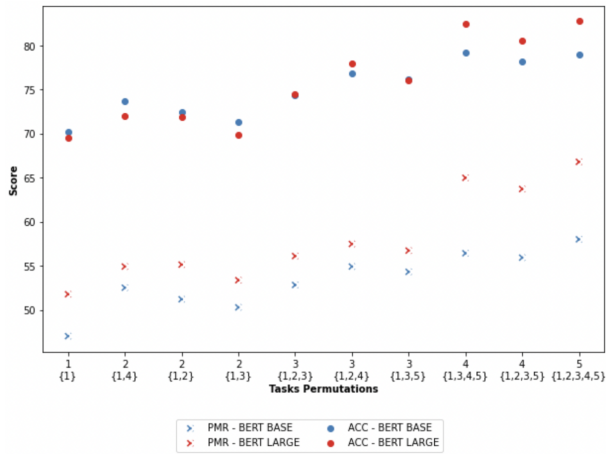
Figure 8: Sample Texts for coherence scoring tasks: GCDC & CoheSentia benchmarks

Task Name	Dataset Name	Prompt
SRO	ROCStories	“reorder: what is the order of the sentences so that the paragraph is coherent? sentence 1: $\langle S_1 \rangle$ sentence 2: $\langle S_2 \rangle$... $\langle S_N \rangle$ ”
ISR	ROCStories	“relevance: what is the irrelevant sentence in the text? sentence1: $\langle S_1 \rangle$ sentence2: $\langle S_2 \rangle$ sentence3: ... $\langle S_N \rangle$ ”
DRR	PDTB3	“discourse relation: what is the discourse relation between $\langle DU_1 \rangle \langle DU_2 \rangle$ ”
NPE	TNE	“coreference text: what are the preposition relations between $\langle NP_i \rangle$ and $\langle NP_j \rangle$? text: $\langle P \rangle$ ”
NLI	MNLI	“mli: does this hypothesis contradict, entail, or neutral with the premise? hypothesis: $\langle H \rangle$ premise: $\langle P \rangle$ ”
Coherence Scoring	GCDC	“GCDC coherence: what is the coherence score of the text (3 - high, 1 - low)? text: $\langle P \rangle$ ”
Coherence Scoring	CoheSentia	“CoheSentia coherence: what is the coherence score of the text (5 - high, 1 - low)? title: $\langle T \rangle$ text: $\langle P \rangle$ ”
MT	WMT14	“Machine Translation: what is the translation of the next text from language $\langle source_language \rangle$ to $\langle target_language \rangle$? text in source language”
NER	Conll2003	“NER task: what is the entity recognition tagging of each token in the next text? $\langle extra_id_0 \rangle$ token1 $\langle extra_id_1 \rangle$ token2 ...”
POS	Conll2003	“POS task: What is the part of speech tagging of each token in the next text? $\langle extra_id_0 \rangle$ token1 $\langle extra_id_1 \rangle$ token2 ...”
Cohesion Reasoning	CoheSentia	“Cohesion reasoning: previous data: $\langle d_i \rangle$ new sentence: $\langle si \rangle$. Task: is the new sentence cohesive in regard to the previous data? give a yes or no answer to each item ”
Consistency Reasoning	CoheSentia	“Consistency reasoning: previous data: $\langle d_i \rangle$ new sentence: $\langle si \rangle$. Task: is the new sentence consistent in regard to the previous data? give a yes or no answer to each item ”
Relevance Reasoning	CoheSentia	“Relevance reasoning: previous data: $\langle d_i \rangle$ new sentence: $\langle si \rangle$. Task: is the new sentence relevant in regard to the previous data? give a yes or no answer to each item ”

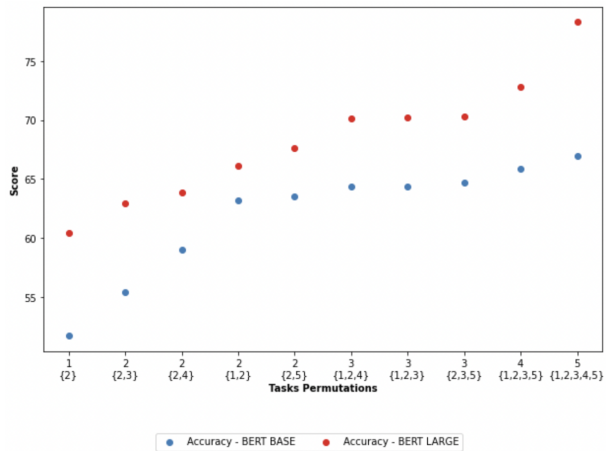
Table 9: Prompts for all tasks in this paper when using T5 model as the backbone model

Task	Dataset	Outputs
SRO	ROCStories	list of position markers $[Y_1, Y_2, \dots, Y_N]$ (Y_i -position of the i_{th} sentence of the corresponding ordered sequence S_i in the shuffled input)
ISR	ROCStories	the index of the irrelevant sentence in the paragraph
DRR	PDTB3	" $\langle \text{connector} \rangle \rightarrow \langle l_1 \text{ relation} \rangle \rightarrow \langle l_2 \rangle$ "
NPE	TNE	the preposition
NLI	MNLI	Contradict / Entails / Neutral
Coherence scoring	GCDC	the score
Coherence scoring	CoheSentia	the score
MT	WMT14	the translated text
NER	Conll2003	" $\langle \text{extra_id_0} \rangle \text{ner_tag_token1} \langle \text{extra_id_2} \rangle \text{ner_tag_token2} \dots$ "
POS	Conll2003	" $\langle \text{extra_id_0} \rangle \text{pos_tag_token1} \langle \text{extra_id_2} \rangle \text{pos_tag_token2} \dots$ "
Cohesion reasoning	CoheSentia	Yes / No
Consistency reasoning	CoheSentia	Yes / No
Relevance reasoning	CoheSentia	Yes / No

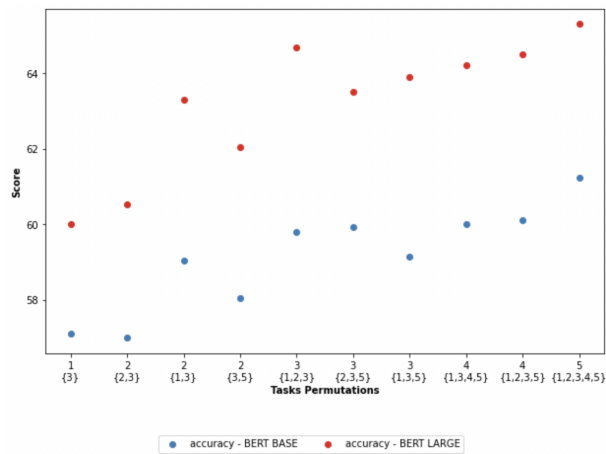
Table 10: Outputs for all tasks in this paper when using T5 model as the backbone model



(a) SRO

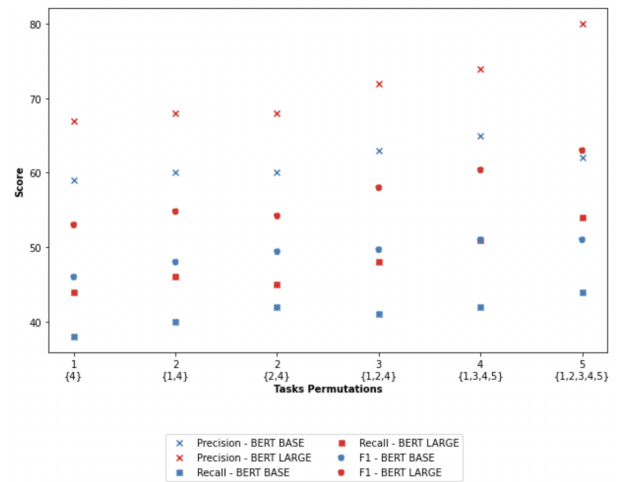


(b) ISR

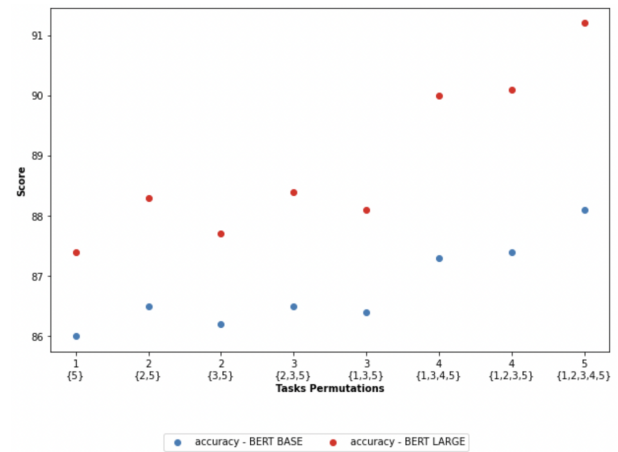


(c) DRR

Figure 9: Results for all tasks, for different permutations of tasks fine-tuned upon. The labels are the number of tasks and in curly brackets which tasks (1 - SRO, 2 - ISR, 3 - DRR, 4 - NPE, 5 - NLI)



(a) NPE



(b) NLI

Figure 10: Results for all tasks, for different permutations of tasks fine-tuned upon. The labels are the number of tasks and in curly brackets which tasks (1 - SRO, 2 - ISR, 3 - DRR, 4 - NPE, 5 - NLI)