# **TROTR:** A Framework for Evaluating the Recontextualization of Text

Anonymous ACL submission

#### Abstract

Current approaches for *detecting* text reuse do not focus on *recontextualization*, i.e., how the new context(s) of a reused text differs from 004 its original context(s). In this paper, we propose a novel framework called TROTR that relies on the notion of topic relatedness for 007 evaluating the diachronic change of context in which text is reused. TROTR includes two NLP tasks: TRiC and TRaC. TRiC is designed to evaluate the topic relatedness between a pair of recontextualizations. TRaC is designed to evaluate the overall topic variation within a set of recontextualizations. We also provide a curated TROTR benchmark of biblical text reuse, human-annotated with topic relatedness. The 015 benchmark exhibits an inter-annotator agree-017 ment of .811. We evaluate multiple, established SBERT models on the TROTR tasks and find that they exhibit greater sensitivity to textual similarity than topic relatedness. Our experiments show that fine-tuning these models can mitigate such a kind of sensitivity.

# 1 Introduction

026

027

As individuals, we often *reuse* someone else's words for diverse reasons and in various ways. This linguistic choice transcends cultural and temporal boundaries, representing an interesting phenomenon to study in Linguistics (Bois, 2014). For instance, linguistic scholars have investigated theories of Reception (Thompson, 1993; Hohendahl and Silberman, 1977) and Resonance (McDonnell et al., 2017; Dimock, 1997) to understand how individuals and communities interpret and reuse historical texts many years after they were written.

With the advent of digitization, recent years have seen a growing interest in computational methods for studying *text reuse*, i.e., "the reuse of existing written sources in the creation of a new text" (Clough et al., 2002). Existing methods focus on the main task of Text Reuse Detection (TRD). In TRD, text reuses are all assumed as "topically related to the source" (Hagen and Stein, 2011; Chiu et al., 2010), the boundaries of reused text are unknown, and the goal is to *detect* text reuse across a diachronic corpus (Seo and Croft, 2008). Whether and how the topic(s) or context(s) of a reused text differs from the source is generally overlooked. Thus, new methods are needed for modeling *recontextualization*, i.e., "the dynamic transfer-andtransformation of a text from one discourse/text-incontext to another" (Connolly, 2014; Linell, 1998). 041

042

043

044

045

049

053

055

059

060

061

062

063

064

065

066

067

069

072

073

074

075

076

077

078

079

083

In this paper, we propose a framework, called Topic Relatedness of Text Reuse (TROTR), to evaluate computational methods for capturing the different recontextualizations of text reuse. In TROTR, the boundaries of reused text are known and the goal is to distinguish reuses of the same text according to their different, latent (i.e., unlabeled) topics. As an example, consider three recontextualizations of the biblical passage *John 15:13* (in bold):

- (1) It's the wonderful pride month!! ♥ ♥ ♥ ♥ ♥
  Honestly pride is everyday! Love is love don't forget I love you ♥. Remember this! John 15:12-13: "My command is this: Love each other as I have loved you. Greater love has no one than this: to lay down one's life for one's friends"
- (2) At a large Crimean event today Putin quoted the Bible to defend the special military operation in Ukraine which has killed thousands and displaced millions. His words "There is no greater love than if someone gives soul for their friends". And people were cheering him. Madness!!!
- (3) "Freeing people from genocide is the reason, motive & goal of the military operation we started in the Donbas & Ukraine", Putin says, then quotes the Bible: "There is no greater love than to lay down one's life for one's friends." It's like Billy Graham meets North Korea

In this example, the biblical passage is incorporated within three texts with different topic recontextualizations. In particular, the text (1) has a different topic with respect to text (2) and (3), while the texts (2) and (3) are topic related. In TROTR, we

178

179

180

182

132

support the recognition of such a kind of recontextualizations by leveraging the notion of topic relatedness. TRoTR represents a new opportunity in Natural Language Processing (NLP) and can be used to distinguish recontextualizations of any kind of text reuse (e.g., proverbs, Ghosh and Srivastava, 2022), to investigate phenomena such as the use of misquotations (Porrino et al., 2008) and dogwhistles (Hertzberg et al., 2022), as well as to provide in-context interpretation to vague utterances, with special focus on enhancing the LLMs' capabilities to this end (DeVault and Stone, 2004).

### Our original contribution.

086

090

100

101

104 105

107

109

110

111

112

113

114

115

116 117

118

- We introduce a novel framework, called TROTR, with two NLP tasks called Text Reuse in-Context (TRiC) and Topic variation Ranking across Corpus (TRaC).
- We provide TROTR with a benchmark containing gold labels derived by human judgements of topic relatedness. The judgements show an inter-annotator agreement of .811, calculated by the average pairwise correlation on assigned assessments.
- We propose a novel annotation process to model topics through topical relatedness in context pairs.
- We evaluate 36 SBERT models by considering 4 settings. Our results reveal that these models reach high performance (correlation of .600 .800), but are more sensitive to textual similarity rather than topic relatedness.

### 2 Related work

Works related to TROTR are about text reuse and recontextualization, semantic textual similarity and relatedness, and topic modeling and annotation.

Text reuse and recontextualization. Although 119 multiple facets of text reuse have been investi-120 gated, such as historical (Büchler et al., 2014), 121 cross-lingual (Muneer and Nawab, 2022), allu-122 sive (Manjavacas et al., 2019), explicit (Franzini 123 et al., 2018), non-literal (Moritz et al., 2016), and 124 local (Seo and Croft, 2008), computational ap-125 proaches primarily focuses on *detecting* instances 126 127 of text reuse. To the best of our knowledge, studies extending beyond mere TRD often leverage text 128 metadata to analyze reuse within temporal and spa-129 tial graphs (Khritankov et al., 2015; Smith et al., 130 2013; Xu et al., 2014). However, these studies do 131

not specifically focus on capturing how the reused text is recontextualized, thereby leaving a gap in the current literature.

Among recent advancements in NLP, some works are related to the recontextualization of text. Wilner et al. (2021) focus on Narrative Analysis by investigating how the recontextualization of events across whole stories impacts word embeddings. Ghosh and Srivastava (2022) introduce a benchmark for evaluating the LLMs' capability of generating proverbs in-context of narratives.

Over the past few years, there has been growing interest in quotations, i.e. "well known phrases or sentences that we use for various purposes such as emphasis, elaboration, and humor" (Lee et al., 2016). This interest extends to various forms of quotations spanning from epigraphs (Bond and Matthews, 2018) to biblical references (Moritz et al., 2016). In particular, there has been a surge of attention in recommendation systems that offers off-the-shelf quotations based on provided context (Wang et al., 2023, 2022, 2021).

Semantic textual similarity and relatedness. In NLP, a possible option for assessing text recontextualization is to use semantic (textual) similarity. However, semantic similarity is traditionally used as a metric to assess paraphrases or entailment equivalence between two texts (Hercig and Kral, 2021; Konopík et al., 2017; Cer et al., 2017; Agirre et al., 2016, 2015, 2014, 2013, 2012); thus, it is not suitable for TROTR. Semantic (textual) relatedness has been long recognized as a core aspect in understanding the meaning of texts (Halliday and Hasan, 2014), and encompasses a multitude of intricate relationships, such as sharing a common topic, expressing similar viewpoints, or originating from the same temporal period (Abdalla et al., 2023). However, there is no universally accepted linguistic theory or set of guidelines for evaluating relatedness. Its assessment is inherently more complex than semantic similarity, as two texts may lack semantic similarity but still be semantically related through some textual relationship (see Table 5).

**Topic modeling and annotation.** Topic models can be useful tools to discover latent topics in collections of documents (Abdelrazek et al., 2023), either as probability distributions like LDA (Blei et al., 2003) or clustering of embeddings like BERTopic (Grootendorst, 2022). When applied, the derived topics need to be carefully evaluated against benchmarks containing manually derived

279

232

ground truth. As topics represent vague concepts, different guidelines for deriving ground truth use different topic definitions tailored to the specific interests of analysis (Orita et al., 2014). Generally, these guidelines result in manual annotations of topic labels that typically differ across annotators and thus require post-processing techniques to be uniform and standardized (Poursabzi-Sangdeh and Boyd-Graber, 2015). For example, annotators can use different wording to express the same concept.

183

184

185

188

189

190

191

192

194

195

196

197

198

199

201

202

203

208

210

211

212

213 214

215

216

217

218

219

As a result, there is no well-established guideline for annotating topics. However, common to different guidelines is a definition of topic that relies on the notion *what the text is about* (Bauwelinck and Lefever, 2020; Hovy and Lin, 1998).

### **3** The TROTR framework

The TROTR framework consists of two tasks, called Text Reuse in-Context (TRiC) and Topic variation Ranking across Corpus (TRaC). TRiC and TRaC are grounded on human judgments of a specific facet of semantic relatedness (see Section 2) that considers the extent to which two texts share a common *topic*. We call this facet **topic relatedness** (see Table 5 for an example). In our study, the definition of **topic** follows the popular notion of *what the text is about*.

When dealing with complex problems, such as recontextualization, a general approach involves starting with a smaller sub-problem to establish a focused foundation before further expanding. Thus, we first present TRiC as a context-pair level task. Then, we present TRaC as a more complex corpuslevel task that must be addressed to identify potential varying targets for real, in-depth analysis.

#### 3.1 Tasks

In the TROTR tasks, instances of text reuse are presented within different contexts, each representing a new recontextualization of the original text.

221**Text Reuse in-Context**frames a text reuse t222within two different contexts  $c_1$  and  $c_2$ . The goal is223to assess the topic relatedness of  $c_1$  and  $c_2$ . TRiC224includes two subtasks, namely binary classifica-225tion and ranking. These subtasks resemble the226structure of the Word-in-Context task (Loureiro227et al., 2022; Martelli et al., 2021; Liu et al., 2021;228Raganato et al., 2020; Pilehvar and Camacho-229Collados, 2019) and the Graded Word Similarity230in Context task (Armendariz et al., 2020), respec-231tively. However, while they focus on distinguishing

the different meanings words can have in different contexts, TRiC focuses on distinguishing different topics in which text are reused.

Each TRiC instance is associated with a binary label  $l \in \{0, 1\}$  and a continuous score  $1 \le s \le 4$ .

- Subtask 1 *binary classification*: the task is to identify, for each instance, whether the contexts  $c_1$  and  $c_2$  share roughly the same topic (i.e., l = 1) or not (i.e., l = 0).
- Subtask 2 *ranking*: the task is to rank the TRiC instances according to the degree of topic relatedness *s* of the contexts *c*<sub>1</sub> and *c*<sub>2</sub>.

Topic variation Ranking across Corpus frames a text reuse t within a corpus C that includes various contexts  $c_i$  where t occurs. TRaC resembles the structure of the Lexical Semantic Change (LSC) detection task defined by (Schlechtweg et al., 2018; Kutuzov and Pivovarova, 2021). However, while this focuses on assessing the semantic change of a word, TRaC focuses on assessing the topic variation of a reused text. Each TRaC instance is associated with a continuous score  $s \in [0, 1]$  of topic variation that indicates the variability in topic usages for a target text reuse t across the corpus C. Specifically, a score of 1 indicates that a target is associated with a high number of topics, while a score of 0 indicates that a target is associated with a single topic.

Given a set of target text reuses  $t \in T$ , the task is to rank the text reuses by the degree of topic variation across the corpus C.

### 3.2 Annotation process

The TROTR annotation process is enforced to collect human judgments topic relatedness. (see Table 5). In our study, we sidestep the need for annotating topics explicitly using a well-established paradigm adopted for modeling word meaning. Our intuition is that annotating topic relatedness, instead of relying on explicit topic labels, closely mirroring recent work exemplified in the Wordin-Context task (Pilehvar and Camacho-Collados, 2019), which relies on annotating meaning relatedness rather than explicit sense labels.

Annotators are asked to evaluate the *topic relat*edness of different text reuse instances  $\langle t, c_1, c_2 \rangle$ , where t is a target text reuse, and  $c_1$  and  $c_2$  are two different contexts in which t occurs.

The topic relatedness is evaluated by utilizing the four-point DURel relatedness scale (see Table 6), with annotators following instructions inspired by the guidelines from Erk et al. (2013), as well as those provided for SemEval-2020 Task 1 (Schlechtweg et al., 2020) and the PLATOS project (Bauwelinck and Lefever, 2020). The annotation guidelines for TRoTR, along with its benchmark, and our code, will be publicly available<sup>1</sup>.

# **4** The **TRoTR** benchmark

The TROTR benchmark is composed of humanannotated instances of text reuse. Specifically, we first manually collected and curated tweets containing biblical text reuse instances. We then incorporated gold labels derived by human annotations.

#### 4.1 Data

283

287

289

290

297

302

303

Tweets were collected through a manual search process. This involved inputting a biblical reference or its corresponding passage into the Twitter search bar and collecting 30 distinct tweets where the passage occurs from the search results.<sup>2</sup> A set of 42 target passages was curated by experts to provide a list of popular biblical quotations commonly used on social networks. On the one hand, the manual search of tweets limited our capacity to retrieve a large number of text reuse instances due to its time-consuming nature. On the other hand, it gave us a rigorous control over the identification of text reuse instances, thereby bypassing a TRD phase and its validation.

Data cleaning. Biblical passages are used with slightly different wording e.g., depending on the 310 version of the quoted Bible. For instance, Table 1 il-311 lustrates various English versions of the same Bible verse. As a results, two tweets that reuse a biblical 313 314 passage may lack a significant common substring, e.g. the longest substring between the first and 315 the fourth passage are the words is kind. To estab-316 lish a controlled setting and ensure that the subsequent evaluation of computational models were not influenced by slight variations in wording (see 319 Section 5), we chose to manually clean the data. 320 Specifically, we removed double spaces, corrected 321 typos, and replaced any outdated or rephrased version of a verse with its contemporary and standard-323 324 ized counterpart, ensuring that the change did not disrupt the flow of the sentence. 325

Text	<b>Bible version</b>	Year
Love is patient, love is kind. Love	Christian Standard Bible	
does not envy, is not boastful, is not	(CSB)	2017
arrogant,	(CSB)	
Love is patient and kind; love does	English Standard Version	2001
not envy or boast; it is not arrogant	(ESV)	2001
Love suffers long and is kind; love	Now King James Varsion	
does not envy; love does not parade		1982
itself, is not puffed up;	$(\mathbf{N}\mathbf{K}\mathbf{J}\mathbf{V})$	
Charity suffereth long, and is kind;	King Jamas Varsian	
charity envieth not; charity vaunteth		1611
not itself, is not puffed up,		

Table 1: Different versions of the passage 1 Cor 13:4

#### 4.2 Human judgments

We collected judgments according to Section 3.2. Specifically, we recruited four native English speakers as annotators; however, during the annotation campaign, two annotators dropped out early. Nevertheless, the majority of instances received three annotations. Annotators were trained and tested on a small set of instances in an online tutorial.

For each target sequence t, we randomly sampled of 150 unique context pairs  $\langle t, c_1, c_2 \rangle$  without replacement from the full set of possible combinations. These were presented to annotators in randomized order to be judged for topic relatedness. The outcome of our annotation pipeline is a dataset of 6,300 annotated context pairs. We measured inter-annotator agreement on judgments using Krippendorff's  $\alpha$  coefficient (Krippendorff, 2018) and the weighted mean of Spearman correlations (Spearman, 1904) between annotator pairs. Table 4 in Appendix provides a summary of our agreement scores. Similar to previous studies that reported Krippendorff's  $\alpha$  of .439 (Loureiro et al., 2022) and weighted mean of Spearman correlation between annotator judgments ranging from .550 to .680 (Erk et al., 2013; Schlechtweg et al., 2018), we obtained a comparable Krippendorff's  $\alpha$  score of .420 and Spearman correlation of .506.

In order to further increase the quality of our data, we employ distinct filtering criteria for annotation instances in both TRiC and TRaC. Specifically, for TRiC, we enforced a two-step filtering process: (i) we filtered out all instances where the difference in judgments between the maximum and minimum judgment exceeds 1, e.g. an instance with three different judgments (e.g., 1, 2, 3) from three annotators. In step (ii), to enforce a more clear-cut separation, we filtered out all the instances  $\langle t, c_1, c_2 \rangle$  where the average judgment score between annotators is greater than 2 and lower than 3.

This filtering results in a more refined dataset of 3,821 annotated context pairs, characterized by a

<sup>&</sup>lt;sup>1</sup>We will insert our GitHub link after acceptance. **Our data** and software are submitted as supplementary material.

<sup>&</sup>lt;sup>2</sup>The Twitter search matching goes beyond exact matches.

Krippendorff's  $\alpha$  agreement of **.709** and a weighted average pairwise Spearman agreement of **.811**.

367

368

373

375

378

384

388

390

394

396

400

401

402

403

404

405

406

407

408

409

410

411

412

For TRaC, we refrained from using the identical filtering, as it could bias the random sample of use pairs. Keeping the sample random and large allows to assume that topic variation inferred on the sample generalizes to the full set of usages for each target quotation corpus. Thus, we adopted a distinct filtering approach at the level of targets. Specifically, we filtered out all the  $\langle t, c_1, c_2 \rangle$  instances for a target t where the weighted average pairwise Spearman agreement was below .150. This criterion led to the exclusion of 2 targets.

**TRiC labels.** We directly use the judgments of each annotated instance to derive binary labels and continuous scores for Subtask 1 and Subtask 2. In particular, our guidelines define a score of 1 as Unrelated topic and a score of 4 as Identical topic. Thus, we assume that pairs with judgments of 3 and 4 are more likely to belong to the same topic, while judgments of 1 and 2 are more likely to belong to different topics. Therefore, for Subtask 1, we aggregate the judgments of all annotators for each instance by averaging, resulting in a score s representing topic relatedness for each instance. We then binarize s as 1 if  $s \ge 2.5$  or as 0 if s < 2.5and associate each instance with the corresponding binary label. A threshold of 2.5 is a reasonable choice, as it represents a midpoint split on the judgment scale. An alternative choice influences the granularity of topics to distinguish, offering flexibility for adjustment in the future. Overall, our benchmark includes a total of 2,621 examples with label 0 and a total of 1,200 examples with label 1.

For Subtask 2, we directly utilize the continuous score s for each instance.

**TRaC labels.** We use a judgment summary measure similar to the DURel EARLIER/LATER measures introduced by Schlechtweg et al. (2018) in the field of LSC (Montanelli and Periti, 2023; Tahmasebi et al., 2021; Kutuzov et al., 2018). Given a target t, this simply involves computing the average of annotator judgments over all instances  $\langle t, c_1, c_2 \rangle$ . As the resulting scores s range between 1 and 4, we normalize them using the following formula:

$$1 - \frac{1}{3}(s - 1)$$
 (4)

This normalization ensures that higher scores correspond to greater topic variation, while lower scores are associated with less topic variation (Appendix C).

### **5** Evaluation setup

We use the TROTR tasks and benchmarks to evaluate the ability of sequence-level models to capture topic relatedness and variation in different text recontextualizations. For this purpose, we chose Sentence-BERT (SBERT) models since they are recognized to be the state-of-the-art architecture for addressing sequence-level tasks (Reimers and Gurevych, 2019). Notably, the SBERT architecture is a modification of the pre-trained BERT network (Devlin et al., 2019) tailored for sequencelevel embeddings and textual similarity. In the following, we introduce the SBERT models and architectures as well as the considered evaluation settings for both TRiC and TRaC. 416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

# 5.1 SBERT models

We consider 36 SBERT models trained on a wide range of tasks including Paraphrasis, Semantic Similarity, and Question Answering. We evaluate each SBERT model in its pre-trained version (baseline) and three different settings, namely:

- +MASK: given an instance (t, c<sub>1</sub>, c<sub>2</sub>), we mask the text-reuse excerpt t in the contexts c<sub>1</sub> and c<sub>2</sub> to prevent that the topic estimate of topic relatedness is influenced by the common t in c<sub>1</sub> and c<sub>2</sub>. To this end, we replace t in c<sub>1</sub> and c<sub>2</sub> with a dash (i.e., "-");
- +*FT*: we fine-tune the pre-trained model on TRiC instances using the *contrastive loss* (Hadsell et al., 2006). This loss minimizes the distance between embeddings of similar sentences and maximizes the distance for dissimilar sentences;
- +*FT*+*MASK*: we combine both the +FT and +MASK settings, meaning that we fine-tune the model and then evaluate it by considering contexts where targets are masked.

### 5.2 SBERT architectures

We evaluate models trained on two SBERT architectures:

Bi-Encoder models are designed to produce a sequence embedding for an input text sequence. Given an instance (t, c<sub>1</sub>, c<sub>2</sub>), we independently feed a Bi-Encoder model with the sequence c<sub>1</sub> and c<sub>2</sub> to obtain the corresponding sequence embeddings u and v. Similar to Abdalla et al. (2023), we use the cosine similarity between u and v as an estimate of the topic relatedness between c<sub>1</sub> and c<sub>2</sub>.

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

558

559

560

561

562

563

564

515

516

Cross-Encoder models are designed to produce an output value that indicates the similarity of two input sequences. Thus, given an instance (t, c<sub>1</sub>, c<sub>2</sub>), we simultaneously pass the sequences c<sub>1</sub> and c<sub>2</sub> to the Cross-Encoder model and use the output value as an estimate of the topic relatedness between c<sub>1</sub> and c<sub>2</sub>.

### 5.3 TRiC evaluation

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

502

503

504

505

508

509

510

511

512

513

514

Similar to the WiC tasks (e.g., Pilehvar and Camacho-Collados, 2019), we consider a supervised scenario where we split the TROTR benchmark into three distinct partitions, namely training set (Train), development set (Dev), and test set (Test), comprising approximately 80%, 10%, and 10% of the instances, respectively. To strengthen the robustness of the evaluation, ten randomized Train-Dev-Test splits were generated (see Appendix B). We consider the average performance across all the splits as a reference for comparison. Additionally, inspired by Raganato et al. (2020), we include the evaluation of target text reuse t that are unseen during fine-tuning. The goal is to evaluate the ability of models to generalize the assessment of topic relatedness. Specifically, we fine-tune each considered model on the Train set and we evaluate it on two different Test sets: i) the standard Test set, containing instances  $\langle t, c_1, c_2 \rangle$  whose target t was either seen or unseen during fine-tuning; and ii) the Out-of-Vocabulary (OOV) Test set, containing only instances  $\langle t, c_1, c_2 \rangle$  whose target t was not seen during fine-tuning. OOV Test set represents half of the Standard Test set.

For TRiC Subtask 1, we need to define a threshold to determine instances  $\langle t, c_1, c_2 \rangle$  where  $c_1$  and  $c_2$  share roughly the same topic or not. Thus, given a model, we tune a threshold-based classifier on the Dev set. Specifically, for each instance  $\langle t, c_1, c_2 \rangle$ in Dev, we use the model to predict the topic relatedness between  $c_1$  and  $c_2$ . Then, we determine the optimal threshold that maximized the Weighted F1 (Harbecke et al., 2022) score over the Dev set. Finally, we apply this threshold to both the Train and Test sets. Due to the unbalanced distribution of gold binary labels, we evaluate models using the F1 metric. Precision (PR) and Recall (RE) for each individual class are also reported for completeness.

For TRiC Subtask 2, given a model, we directly use its predictions as estimates of topic relatedness. Then, we evaluate the model using Spearman correlation (SP) with continuous gold scores.

#### 5.4 TRaC evaluation

Similar to the LSC tasks (e.g., Schlechtweg et al., 2020), we consider an *unsupervised* scenario. In particular, motivated by the limited number of targets (i.e., 42), we do not split the benchmark into Train-Dev-Test partitions with the aim to mitigate the potential evaluation impact of a small Test set. Without training instances, the configurations with +FT and +FT+MASK are not applicable to TRaC.

To quantify the topic variation of a target, we adopted the same approach used for determining the gold scores. Thus, given a model, the topic variation of a target t is calculated as the average prediction of topic relatedness across all the annotated  $\langle t, c_1, c_2 \rangle$  pairs. We then evaluate models using Spearman correlation (SP) with gold scores.

#### **6** Evaluation results

First, we evaluated an extensive set of pre-trained SBERT models on the TRiC task (see Table 8 in Appendix). Then, for simplicity, we opted to consider and fine-tune a smaller set of models, precisely the top-five models by SP over the Train sets. Since we did not perform any training over the models, the Train sets act as a larger sets for testing the models. Specifically, we chose: all-distilroberta-v1 (ADR), distilusebase-multilingual-cased-v1 (DBM), paraphrasemultilingual-MiniLM-L12-v2 (PAM), paraphrasemultilingual-mpnet-base-v2 (PAR), and multi-gampnet-base-cos-v1 (MQA). In particular, ADR and DBM are Bi-Encoders for English. PAM and PAR are multilingual Bi-Encoders fine-tuned on paraphrase pairs. Similarly, MQA is a multilingual Bi-Encoder fine-tuned on question-answer pairs.

As a general remark on our initial evaluation, we note that Bi-Encoder models consistently exhibit superior performance compared to Cross-Encoder models in both TRiC Subtask 1 and Subtask 2. This finding aligns with the recent comparisons by Ishihara and Shirai (2022) and Cassotti et al. (2023) for News Article Similarity and LSC, challenging the idea that the use of cross-attention benefits Cross-Encoder architectures in sequence-level tasks (Lee et al., 2023; Thakur et al., 2021). In the following, we first present the results of our evaluation by comparing the use of pre-trained and fine-tuned models (+FT); then, we discuss the results in the masking settings (+MASK, +FT+MASK). We report in Table 2 and 3 the overall results for TRiC and TRaC, respectively.

	Standard Test set								Out-of-vocabulary (OOV) Test set							
		Label 0			Label 1		All			Label 0		Label 1			All	
Models	PR	RE	F1	PR	RE	F1	F1	SP	PR	RE	F1	PR	RE	F1	F1	SP
ADR	.95±.03	.47±.13	.62±.11	.42±.11	.93±.04	.57±.10	.61±.10	.55±.09	.94±.07	.45±.20	.58±.20	.38±.19	.93±.06	.51±.18	.58±.16	.48±.20
+FT	.95±.03	.61±.15	.73±.11	.50±.14	.93±.03	.64±.10	.71±.10	.66±.07	.91±.12	.49±.24	.61±.22	.40±.21	.91±.06	.52±.18	.61±.18	.51±.22
+MASK	.89±.05	.87±.07	.87±.03	.70±.14	.72±.12	.69±.07	.82±.03	.67±.06	.90±.07	.85±.10	.87±.05	.62±.21	.71±.18	.63±.14	.82±.05	.62±.15
+FT+MASK	.90±.07	.89±.07	.89±.03	.75±.12	.76±.12	.74±.05	.85±.04	.71±.05	.87±.11	.88±.09	.87±.06	.66±.20	.70±.15	.65±.09	.82±.06	.63±.15
DBM	.96±.02	.26±.12	.40±.14	.35±.09	.97±.03	.51±.09	.43±.12	.54±.09	.96±.08	.21±.19	.31±.23	.31±.14	.97±.05	.45±.16	.38±.18	.44±.23
+FT	.97±.02	.46±.17	.60±.15	.43±.10	.96±.03	.58±.09	.61±.13	.64±.07	.93±.15	.34±.23	.46±.26	.34±.14	.95±.05	.49±.15	.50±.19	.48±.29
+MASK	.87±.07	.88±.07	.87±.03	.72±.14	.66±.16	.66±.09	.81±.03	.64±.04	.88±.09	.88±.09	.87±.05	.66±.23	.64±.25	.58±.19	.82±.04	.58±.12
+FT+MASK	.88±.06	.89±.07	.88±.04	.74±.11	.70±.13	.70±.04	.83±.03	.66±.04	.85±.12	.87±.09	.85±.08	.63±.19	.58±.20	.57±.13	.80±.08	.58±.14
PAM	.96±.02	.46±.09	.61±.08	.41±.09	.96±.02	.57±.08	.61±.07	.58±.08	.96±.04	.43±.17	.57±.16	.37±.15	.95±.05	.52±.15	.59±.12	.49±.22
+FT	.95±.03	.59±.12	.72±.11	.48±.09	.92±.04	.63±.08	.70±.09	.66±.06	.90±.18	.45±.21	.57±.23	.37±.13	.92±.06	.51±.13	.59±.17	.51±.22
+MASK	.89±.05	.88±.06	.88±.03	.71±.10	.72±.10	.70±.05	.83±.03	.67±.04	.89±.09	.86±.09	.87±.06	.65±.19	.71±.18	.65±.12	.83±.05	.60±.13
+FT+MASK	.90±.05	.90±.03	.90±.03	.76±.07	.77±.06	.76±.03	.86±.03	.69±.04	.88±.10	.89±.05	.88±.06	.68±.13	.73±.11	.69±.07	.84±.06	.60±.12
PAR	.95±.03	.40±.10	.56±.09	.39±.09	.95±.04	.55±.08	.56±.07	.56±.09	.93±.11	.35±.18	.49±.19	.34±.15	.95±.06	.49±.16	.52±.15	.47±.25
+FT	.95±.05	.60±.10	.73±.08	.49±.10	.93±.05	.63±.08	.71±.07	.66±.06	.91±.17	.46±.21	.58±.21	.38±.16	.91±.08	.51±.15	.59±.18	.53±.24
+MASK	.89±.05	.85±.07	.87±.04	.69±.10	.75±.11	.70±.05	.83±.03	.68±.03	.90±.08	.83±.13	.86±.07	.63±.19	.75±.17	.65±.10	.82±.05	.62±.11
+FT+MASK	.89±.06	.91±.05	.90±.03	.78±.09	.73±.11	.74±.05	.86±.03	.70±.04	.87±.11	.90±.07	.88±.06	.68±.16	.66±.18	.64±.11	.83±.07	.61±.14
MQA	.94±.03	.42±.11	.58±.11	.40±.10	.94±.03	.55±.09	.58±.09	.55±.09	.94±.09	.39±.19	.53±.20	.36±.19	.96±.03	.50±.18	.55±.16	.49±.21
+FT	.96±.03	.61±.13	.74±.10	.50±.10	.94±.04	.65±.08	.72±.09	.68±.06	.92±.15	.47±.22	.60±.24	.39±.16	.94±.05	.53±.15	.61±.19	.54±.21
+MASK	.88±.05	.87±.07	.88±.04	.71±.10	.71±.12	.69±.06	.83±.04	.68±.05	.89±.07	.86±.10	.87±.06	.63±.18	.69±.16	.63±.13	.83±.05	.62±.13
+FT+MASK	.90±.05	.91±.04	.90±.03	.77±.08	.76±.09	.76±.05	.86±.03	.72±.04	.88±.10	.90±.04	.88±.06	.67±.16	.69±.16	.65±.11	.84±.06	.63±.13

Table 2: **TRiC evaluation** on Subtask 1 and Subtask 2 for both Test and OOV Test sets. For Subtask 1, precision (PR), recall (RE), and Weighted -F1 scores (F1) are reported for both label 0 (i.e., different topics) and label 1 (i.e., roughly identical topics). For Subtask 2, Spearman correlation (SP) is reported on the overall set of instances. Standard deviations (±) across the 10 Test splits are presented for comparative analysis. For each metric, the best performance of the comparison between pre-trained/fine-tuned models is highlighted in **bold**. Results for masking settings are reported in *italic*.

Madala	ADR	DBM	PAM	PAR	MQA
widdels	+MASK	+MASK	+MASK	+MASK	+MASK
Spearman	.72	.66	.66	.73	.65
	.84	.80	.81	.76	.80

Table 3: **TRaC evaluation** using the pre-trained models alone and in the +MASK setting (*italic*).

#### 6.1 TRiC: pre-trained vs. fine-tuned

Across the overall *standard* Test sets, when *pre-trained* models are used for Subtask 1, we observe high precision (PR) values, ranging from .93 to .96, and low recall (RE) values ranging from .21 to .47 for label 0 (i.e., different topics). Conversely, for label 1 (i.e., roughly identical topics), we observe an inverse trend of performance, with PR values ranging from .31 to .42 and RE values ranging from .93 to .97. Such results suggest that SBERT models face difficulties in distinguishing different recontex-tualization. For Subtask 1, we observe a moderate F1-score (F1) ranging from .43 to .61; for Subtask 2, we observe only moderate Spearman correlation coefficients (SP) ranging from .54 to .58.

Additional results for the *OOV* Test sets are reported in Table 2. We note that the results for the OOV Test sets are lower in performance, while being associated to higher standard deviations. For pre-trained models, we attributed this drop to (1) the unbalance number of instances and labels available for each target; (2) that the inter-annotator agreements differ between targets. If target words with small number of instances or lower inter-annotator agreement fall in the OOV Test sets, then the performance will be much lower. Finally, (3)

the size of the OOV Test sets is smaller because it splits the standard Test sets in two halves.

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

When the pre-trained models are *fine-tuned* on TRiC instances (i.e., +FT), we observe a significant improvement in performance for both Subtask 1 and Subtask 2 on both the standard Test set and the OOV Test set. This observation indicates that fine-tuning SBERT models on TRiC instances enhances their capability to contextualize a sequence in-context. In particular, the improvement is more pronounced on the standard Test sets than on the OOV Test sets. We attribute this discrepancy to the limited size of our benchmark that includes a small number of target quotations sufficient for testing purposes. A larger number of targets will further improve the models' generalization capability. For Subtask 1, we observe a F1 ranging from .61 to .72 (standard) and from .50 to .61 (OOV); for Subtask 2, we observe SP coefficients ranging from .64 to .68 (standard) and .51 to .54 (OOV).

#### 6.2 TRiC and TRaC: masking settings

When pre-trained and fine-tuned models are used in the masking settings (i.e., +MASK and +FT+MASK), we observe a significant improvement in performance for both TRiC and TRaC. Notably, this improvement for TRiC is substantially larger compared to the one observed in the prior comparison (pre-trained vs. fine-tuned), with +FT+MASK exhibiting slightly superior performance to +MASK. We attribute this improvement to the fact that, in the masking settings, models are

588

589

590

565

566

671

673

622

compelled to pay more attention to the surrounding contexts of reused texts, thereby fostering a more comprehensive understanding of topic relatedness.

When SBERT models are used in the +MASK and +FT+MASK settings for TRiC, we observe the following performance. For Subtasks 1, we observe a F1 ranging from .81 to .83 and from .82 to .86 for +MASK and +FT+MASK, respectively. For Subtask 2, we observe a SP coefficients ranging from .60 to .68 and from .60 to .72 for +MASK and +FT+MASK, respectively.

Similar to the previous results, we observe pronounced discrepancy even when we compare the use of pre-trained models as baselines and their use in the +MASK setting for TRaC (see Table 3). When pre-trained models are used for TRaC, we observe SP coefficients ranging from .65 to .73. Conversely, when pre-trained models are used in the +MASK setting, SP coefficients exhibit a substantial improvement, ranging from .76 to .84.

These results further underscore the difficulty of SBERT models in distinguishing different text recontextualizations. As a matter of fact, pre-trained models exhibit a bias towards their typical pretraining focus, namely *semantic similarity*, while demonstrating only a superficial understanding of topic relatedness. Although the masking settings seem to offer a valuable workaround to sidestep the problem, we claim that their use is generally undesirable in real scenario involving text reuse. First, because masking may disrupt the natural flow of sentences precluding to obtain optimal performance. Second, because the boundaries of text reuse are often nuanced or unbalanced in different recontextualizations, when considering a form of text reuse broader than explicit quotation that implicitly reuses text in-context. In such cases, masking may result in the removal of crucial contextual information.

Consequently, to provide a more accurate modeling of text-reuse *in-context*, we argue that there is a clear imperative to develop or fine-tune novel models specifically tailored on topic relatedness. In this regard, TROTR represents a valuable framework for evaluating language models that extend existing benchmarks on sentence-pair regression tasks, such as Semantic Textual Similarity (Agirre et al., 2012) and Semantic Textual Relatedness (Abdalla et al., 2023). While current benchmarks rely on a notion of *similarity* or *relatedness*, they overlook the potential impact of shared substrings, such as text-reuse excerpts, on computational estimates.

#### 7 Concluding remarks and future work

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

In this work, we relied on the notion of topic relatedness to introduce a novel framework named Topic Relatedness of Text Reuse (TROTR). This framework is designed to assess computational methods in distinguishing diverse text reuse re-Our framework comprises contextualizations. two NLP tasks, namely Text Reuse in-Context (TRiC) and Topic variation Ranking across Corpora (TRaC). The framework also comprises a human annotated benchmark of biblical text reuse extracted from Twitter (now X). In the TROTR framework and benchmark, we consider a synchronic scenario that involves various recontextualizations of a literal reused text (i.e., explicit quotation of a biblical passage), overlooking the original contextualization (e.g., religious sources like the Bible). However, as text reuse is inherently *diachronic*, we argue that the TROTR framework is applicable to address the recontextualization problem across time, space, or domain, by also encompassing both literal and non-literal recontextualizations as well as the original contextualization(s).

We comprehensively evaluate SBERT models on the TRiC and TRaC tasks, and find that the models exhibit a greater sensitivity to semantic *similarity* rather than topic *relatedness*. Fine-tuning on text reuse instances can mitigate such sensitivity.

To the best of our knowledge, this work represents a first pioneering effort in the computational modeling of recontextualization to support Linguistic Recycling and Reception studies. Our ongoing and future work is about advancing this work by extending the current benchmark. This involves: i) carefully selecting a diachronic corpus spanning multiple time periods with numerous text-reuse instances; ii) implementing and rigorously validating a Text Reuse Detection (TRD) pipeline; and, finally, iii) conducting a larger annotation campaign on text-reuse instances. By leveraging a more extensive benchmark, we aim to enhance the TROTR framework by modeling topics in recontextualizations through an extension of the Word Usage Graphs (WUGs, Schlechtweg et al., 2021) paradigm. Specifically, by representing semantic proximity judgments on text reuse pairs as a weighted graph (contexts as nodes, judgments as edges), topics can be inferred using a graph clustering algorithm, sidestepping the need for explicit topic labeling.

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

772

773

774

775

776

777

778

779

# 8 Limitations

724

726

727

728

731

733

734 735

736

737

740

741

742

743

744

745

746

747

749

750

753

754

755

762

767

769

771

The main limitations of this work pertain to the benchmark, including the data collection and processing:

• *Manual tweet search*: we conducted a manual search of tweets by leveraging the Twitter search bar. This allowed us to sidestep a Text Reuse Detection phase and its validation. However, manually checking the suitability of retrieved tweets is extremely time consuming, thus limiting our ability to collect a large amount of tweets. Moreover, due to the Twitter ranking of matching results, the topic distribution of recontextualizations may be biased.

- *Manual processing*: data collection, processing, and cleaning have been executed by two human annotators. For data processing, the annotators adhered to shared guidelines and they replaced any outdated or rephrased version of a biblical verse with one of its contemporary and standardized counterparts, ensuring that the change did not disrupt the flow of the sentence. For data cleaning, the annotators worked on the removal of double spaces, typos, Twitter mentions, and biblical citations from the tweets' content. However, they did not follow formal guidelines, and discrepancies on the treatment of these elements can be observed as a result.
  - Randomization of the annotation instances: in generating the pairs of tweets to compare for human judgement, we randomized the order of ⟨t, c<sub>1</sub>, c<sub>2</sub>⟩ instances. However, we did not randomize the order of the two contexts within a pair. The ordering of c<sub>1</sub> and c<sub>2</sub> in ⟨t, c<sub>1</sub>, c<sub>2</sub>⟩ was fixed and determined by their IDs. If item order influences annotator judgments, this may have created a bias towards certain orderings.
  - *Human judgments*: we discarded a significant number of judgments from human annotators to ensure high-quality of annotation results. This implied a high degree of imbalance in the distribution of TRiC labels for Subtask 1. We addressed and discussed this imbalance in the experimental results (see Section 5.3 and Appendix B).

As a further limitation, the TROTR benchmark contains English tweets only with literal text reuse (i.e., explicit quotations). However, the benchmark can be extended to consider multi-language corpora and implicit text reuse.

As this work is the first of its kind to phrase a new problem, recontextualization of text-reuse, create a human-annotated benchmark, and attempt to solve the problem using computational tools, we do not claim our work to be exhaustive.

# 9 Ethical considerations

The authors have carefully considered the ethics associated with the TROTR benchmark. The benchmark data, extracted from Twitter (now X), and annotations have been used while respecting the privacy and confidentiality of both users and annotators. For users, we made an effort to anonymize publicly available tweets' content by removing tweet mentions and users. For human annotators, we explicitly notified them prior to the annotation that some instances of text reuse might encompass discriminatory language against people or communities. We encourage the research community to approach our benchmark with a critical perspective, recognizing the potential ethical implications of working with data from social media platforms.

The annotation campaign was conducted with Native English speakers who were reached through email broadcasts. Compensation details, set in advance, were based on an hourly rate of  $\leq 12$ . Each annotator spent a total of 53 hours on the annotation process, resulting in an overall compensation of  $\leq 636$ . This fixed compensation was determined according to our time estimation. As per our contract terms, annotators received payment at the conclusion of the annotation campaign.

### References

- Mohamed Abdalla, Krishnapriya Vishnubhotla, and Saif Mohammad. 2023. What Makes Sentences Semantically Related? A Textual Relatedness Dataset and Empirical Study. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 782–796, Dubrovnik, Croatia. Association for Computational Linguistics.
- Aly Abdelrazek, Yomna Eid, Eman Gawish, Walaa Medhat, and Ahmed Hassan. 2023. Topic Modeling Algorithms and Applications: A Survey. *Information Systems*, 112:102131.

- 821 822
- 82 82
- 82
- 82
- 83
- 831
- 833
- 835 836
- 838
- 839
- 8
- 8
- 8
- 847
- 8
- 850 851
- 8
- 854 855 856
- 8! 8!
- 86

- 8
- 8
- 8
- 870
- 872 873

874 875

- 876
- 877 878
- 878 879

- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larraitz Uria, and Janyce Wiebe. 2015. SemEval-2015 Task 2: Semantic Textual Similarity, English, Spanish and Pilot on Interpretability. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 252–263, Denver, Colorado. Association for Computational Linguistics.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2014. SemEval-2014 Task 10: Multilingual Semantic Textual Similarity. In *Proceedings of the* 8th International Workshop on Semantic Evaluation (SemEval 2014), pages 81–91, Dublin, Ireland. Association for Computational Linguistics.
  - Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2016. SemEval-2016 Task 1: Semantic Textual Similarity, Monolingual and Cross-Lingual Evaluation. In *Proceedings of the* 10th International Workshop on Semantic Evaluation (SemEval-2016), pages 497–511, San Diego, California. Association for Computational Linguistics.
- Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. 2012. SemEval-2012 Task 6: A Pilot on Semantic Textual Similarity. In \*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012), pages 385– 393, Montréal, Canada. Association for Computational Linguistics.
- Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. 2013. SEM 2013 shared task: Semantic Textual Similarity. In Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity, pages 32–43, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Carlos Santos Armendariz, Matthew Purver, Matej Ulčar, Senja Pollak, Nikola Ljubešić, and Mark Granroth-Wilding. 2020. CoSimLex: A Resource for Evaluating Graded Word Similarity in Context. In *Proc. of LREC*, pages 5878–5886, Marseille, France. ELRA.
- Nina Bauwelinck and Els Lefever. 2020. Annotating Topics, Stance, Argumentativeness and Claims in Dutch Social Media Comments: A Pilot Study. In *Proceedings of the 7th Workshop on Argument Mining*, pages 8–18, Online. Association for Computational Linguistics.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet Allocation. *J. Mach. Learn. Res.*, 3(null):993–1022.

John W. Du Bois. 2014. Towards a dialogic syntax. *Cognitive Linguistics*, 25(3):359–410.

880

881

882

883

884

885

886

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

- Francis Bond and Graham Matthews. 2018. Toward An Epic Epigraph Graph. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Marco Büchler, Philip R. Burns, Martin Müller, Emily Franzini, and Greta Franzini. 2014. *Towards a Historical Text Re-use Detection*, pages 221–238. Springer International Publishing, Cham.
- Pierluigi Cassotti, Lucia Siciliani, Marco DeGemmis, Giovanni Semeraro, and Pierpaolo Basile. 2023. XL-LEXEME: WiC pretrained model for cross-lingual LEXical sEMantic changE. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 1577–1585, Toronto, Canada. Association for Computational Linguistics.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. SemEval-2017 Task 1: Semantic Textual Similarity Multilingual and Crosslingual Focused Evaluation. In Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017), pages 1–14, Vancouver, Canada. Association for Computational Linguistics.
- Stanford Chiu, Ibrahim Uysal, and W. Bruce Croft. 2010. Evaluating Text Reuse Discovery on the Web. In Proceedings of the Third Symposium on Information Interaction in Context, IIiX '10, page 299–304, New Brunswick, New Jersey, USA. Association for Computing Machinery.
- Paul Clough, Robert Gaizauskas, Scott SL Piao, and Yorick Wilks. 2002. Measuring text reuse. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 152–159.
- John H. Connolly. 2014. Recontextualisation, Resemiotisation and Their Analysis in Terms of an FDGbased Framework. *Pragmatics*, 24(2):377–397.
- David DeVault and Matthew Stone. 2004. Interpreting Vague Utterances in Context. In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*, pages 1247–1253, Geneva, Switzerland. COLING.
- 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.
- Wai Chee Dimock. 1997. A Theory of Resonance. *PMLA*, 112(5):1060–1071.

Katrin Erk, Diana McCarthy, and Nicholas Gaylord. 2013. Measuring Word Meaning in Context. *Computational Linguistics*, 39(3):511–554.

935

936

937

938

940

947

949

951

953

954

955

957

960

961

962

963

964

965

966

967

968

969

972

974

976

977

978

979

980

981

990

- Greta Franzini, Marco Passarotti, Maria Moritz, and Marco Büchler. 2018. Using and Evaluating TRACER for an Index Fontium Computatus of the Summa contra Gentiles of Thomas Aquinas. In Proceedings of the Fifth Italian Conference on Computational Linguistics (CLiC-it 2018), Torino, Italy, December 10-12, 2018, volume 2253 of CEUR Workshop Proceedings. CEUR-WS.org.
- Sayan Ghosh and Shashank Srivastava. 2022. ePiC: Employing Proverbs in Context as a Benchmark for Abstract Language Understanding. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3989–4004, Dublin, Ireland. Association for Computational Linguistics.
  - Maarten Grootendorst. 2022. BERTopic: Neural Topic Modeling with a Class-based TF-IDF Procedure.
  - Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2006), 17-22 June 2006, New York, NY, USA, pages 1735–1742. IEEE Computer Society.
  - Matthias Hagen and Benno Stein. 2011. Candidate Document Retrieval for Web-Scale Text Reuse Detection. In *String Processing and Information Retrieval*, pages 356–367, Berlin, Heidelberg. Springer Berlin Heidelberg.
  - Michael Alexander Kirkwood Halliday and Ruqaiya Hasan. 2014. *Cohesion in English*. 9. Routledge.
  - David Harbecke, Yuxuan Chen, Leonhard Hennig, and Christoph Alt. 2022. Why only Micro-F1? Class Weighting of Measures for Relation Classification. In Proceedings of NLP Power! The First Workshop on Efficient Benchmarking in NLP, pages 32–41, Dublin, Ireland. Association for Computational Linguistics.
  - Tomáš Hercig and Pavel Kral. 2021. Evaluation Datasets for Cross-lingual Semantic Textual Similarity. In Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021), pages 524–529, Held Online. INCOMA Ltd.
  - Niclas Hertzberg, Robin Cooper, Elina Lindgren, Björn Rönnerstrand, Gregor Rettenegger, Ellen Breitholtz, and Asad Sayeed. 2022. Distributional properties of political dogwhistle representations in Swedish BERT. In Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH), pages 170–175, Seattle, Washington (Hybrid). Association for Computational Linguistics (ACL).
  - Peter Uwe Hohendahl and Marc Silberman. 1977. Introduction to Reception Aesthetics. *New German Critique*, 10:29–63.

Eduard Hovy and Chin-Yew Lin. 1998. Automated Text Summarization and the Summarist System. In *TIPSTER TEXT PROGRAM PHASE III: Proceedings* of a Workshop held at Baltimore, Maryland, October 13-15, 1998, pages 197–214, Baltimore, Maryland, USA. Association for Computational Linguistics. 991

992

993

994

995

996

997

998

999

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

- Shotaro Ishihara and Hono Shirai. 2022. Nikkei at SemEval-2022 task 8: Exploring BERT-based biencoder approach for pairwise multilingual news article similarity. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-*2022), pages 1208–1214, Seattle, United States. Association for Computational Linguistics.
- Anton S. Khritankov, Pavel V. Botov, Nikolay S. Surovenko, Sergey V. Tsarkov, Dmitriy V. Viuchnov, and Yuri V. Chekhovich. 2015. Discovering text reuse in large collections of documents: A study of theses in history sciences. In 2015 Artificial Intelligence and Natural Language and Information Extraction, Social Media and Web Search FRUCT Conference (AINL-ISMW FRUCT), pages 26–32.
- Miloslav Konopík, Ondřej Pražák, and David Steinberger. 2017. Czech Dataset for Semantic Similarity and Relatedness. In Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017, pages 401–406, Varna, Bulgaria. INCOMA Ltd.
- Klaus Krippendorff. 2018. Content Analysis: An Introduction to Its Methodology. SAGE Publications.
- Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Velldal. 2018. Diachronic Word Embeddings and Semantic Shifts: a Survey. In *Proceedings of COLING*, pages 1384–1397, Santa Fe, New Mexico, USA. ACL.
- Andrey Kutuzov and Lidia Pivovarova. 2021. Rushifteval: a shared task on semantic shift detection for russian. *Komp'yuternaya Lingvistika i Intellektual'nye Tekhnologii: Dialog conference.*
- Hanbit Lee, Yeonchan Ahn, Haejun Lee, Seungdo Ha, and Sang-goo Lee. 2016. Quote recommendation in dialogue using deep neural network. In Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '16, page 957–960, New York, NY, USA. Association for Computing Machinery.
- Hyun Seung Lee, Seungtaek Choi, Yunsung Lee, Hyeongdon Moon, Shinhyeok Oh, Myeongho Jeong, Hyojun Go, and Christian Wallraven. 2023. Cross encoding as augmentation: Towards effective educational text classification. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2184–2195, Toronto, Canada. Association for Computational Linguistics.
- Per Linell. 1998. Approaching Dialogue: Talk, Interac-<br/>tion and Contexts in Dialogical Perspectives. John1044<br/>1045Benjamins.1046

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

Qianchu Liu, Edoardo Maria Ponti, Diana McCarthy, Ivan Vulić, and Anna Korhonen. 2021. AM2iCo: Evaluating word meaning in context across lowresource languages with adversarial examples. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 7151–7162, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

1047

1048

1049

1051

1055

1056

1057

1058

1059

1060

1061

1063

1067

1068

1069 1070

1071

1072

1073

1074

1075

1076

1077

1081

1082

1083

1084

1085

1087

1089

1091

1092

1093

1094

1095

1096

1097

1098 1099

1100

1101

1102

1103

- Daniel Loureiro, Aminette D'Souza, Areej Nasser Muhajab, Isabella A. White, Gabriel Wong, Luis Espinosa-Anke, Leonardo Neves, Francesco Barbieri, and Jose Camacho-Collados. 2022. TempoWiC: An Evaluation Benchmark for Detecting Meaning Shift in Social Media. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3353–3359, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
  - Enrique Manjavacas, Brian Long, and Mike Kestemont. 2019. On the Feasibility of Automated Detection of Allusive Text Reuse.
  - Federico Martelli, Najla Kalach, Gabriele Tola, and Roberto Navigli. 2021. SemEval-2021 task 2: Multilingual and cross-lingual word-in-context disambiguation (MCL-WiC). In Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021), pages 24–36, Online. Association for Computational Linguistics.
  - Terence E. McDonnell, Christopher A. Bail, and Iddo Tavory. 2017. A Theory of Resonance. *Sociological Theory*, 35(1):1–14.
  - Stefano Montanelli and Francesco Periti. 2023. A Survey on Contextualised Semantic Shift Detection. *arXiv preprint arXiv:2304.01666.*
  - Maria Moritz, Andreas Wiederhold, Barbara Pavlek, Yuri Bizzoni, and Marco Büchler. 2016. Non-Literal Text Reuse in Historical Texts: An Approach to Identify Reuse Transformations and its Application to Bible Reuse. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1849–1859, Austin, Texas. Association for Computational Linguistics.
  - Iqra Muneer and Rao Muhammad Adeel Nawab. 2022. Cross-Lingual Text Reuse Detection at sentence level for English–Urdu language pair. *Computer Speech* & *Language*, 75:101381.
  - Naho Orita, Naomi Feldman, Jordan Boyd-Graber, and Eliana Vornov. 2014. Quantifying the Role of Discourse Topicality in Speakers' Choices of Referring Expressions. In Proceedings of the Fifth Workshop on Cognitive Modeling and Computational Linguistics, pages 63–70, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. 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 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.

- Jack A. Porrino, Virak Tan, and Aaron Daluiski. 2008. Misquotation of a Commonly Referenced Hand Surgery Study. *The Journal of Hand Surgery*, 33(1):2.e1–2.e9.
- Forough Poursabzi-Sangdeh and Jordan Boyd-Graber. 2015. Speeding Document Annotation with Topic Models. In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop, pages 126–132, Denver, Colorado. Association for Computational Linguistics.
- Alessandro Raganato, Tommaso Pasini, Jose Camacho-Collados, and Mohammad Taher Pilehvar. 2020. XL-WiC: A multilingual benchmark for evaluating semantic contextualization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7193–7206, Online. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Dominik Schlechtweg, Barbara McGillivray, Simon Hengchen, Haim Dubossarsky, and Nina Tahmasebi.
  2020. SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection. In *Proceedings of SemEval*, pages 1–23, Barcelona. ICCL.
- Dominik Schlechtweg, Sabine Schulte im Walde, and Stefanie Eckmann. 2018. Diachronic Usage Relatedness (DURel): A Framework for the Annotation of Lexical Semantic Change. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 169–174, New Orleans, Louisiana. Association for Computational Linguistics.
- Dominik Schlechtweg, Nina Tahmasebi, Simon Hengchen, Haim Dubossarsky, and Barbara McGillivray. 2021. DWUG: A large resource of diachronic word usage graphs in four languages. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7079–7091, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Dominik Schlechtweg, Shafqat Mumtaz Virk, Pauline1157Sander, Emma Sköldberg, Lukas Theuer Linke,<br/>Tuo Zhang, Nina Tahmasebi, Jonas Kuhn, and<br/>Sabine Schulte im Walde. 2023. The durel anno-<br/>tation tool: Human and computational measurement1160

- 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1179 1180 1181 1182 1183 1184 1185 1186 1187 1188 1189 1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1202 1203 1204 1205 1206 1207 1208 1210 1211 1212

1163

1164

1213 1214 1215

1216 1217 1218 of semantic proximity, sense clusters and semantic change.

- Jangwon Seo and W. Bruce Croft. 2008. Local Text Reuse Detection. In Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '08, page 571–578, New York, NY, USA. Association for Computing Machinery.
- David A. Smith, Ryan Cordell, and Elizabeth Maddock Dillon. 2013. Infectious texts: Modeling text reuse in nineteenth-century newspapers. In 2013 IEEE International Conference on Big Data, pages 86–94.
- Charles Spearman. 1904. The proof and measurement of association between two things. American Journal of Psychology, 15:88–103.
- Nina Tahmasebi, Lars Borin, and Adam Jatowt. 2021. Survey of Computational Approaches to Lexical Semantic Change Detection. In Computational approaches to semantic change. Language Science Press.
- Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. 2021. Augmented SBERT: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 296-310, Online. Association for Computational Linguistics.
- Martyn P. Thompson. 1993. Reception Theory and the Interpretation of Historical Meaning. History and Theory, 32(3):248-272.
- Lingzhi Wang, Xingshan Zeng, and Kam-Fai Wong. 2021. Quotation Recommendation and Interpretation Based on Transformation from Queries to Quotations. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 754–758, Online. Association for Computational Linguistics.
- Lingzhi Wang, Xingshan Zeng, and Kam-Fai Wong. 2022. Learning when and what to quote: A quotation recommender system with mutual promotion of recommendation and generation. In Findings of the Association for Computational Linguistics: EMNLP 2022, pages 3094–3105, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Lingzhi Wang, Xingshan Zeng, and Kam-Fai Wong. 2023. Quotation Recommendation for Multi-Party Online Conversations Based on Semantic and Topic Fusion. ACM Trans. Inf. Syst., 41(4).
- Sean Wilner, Daniel Woolridge, and Madeleine Glick. 2021. Narrative Embedding: Re-Contextualization Through Attention. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 1393-1405, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Shaobin Xu, David Smith, Abigail Mullen, and Ryan 1219 Cordell. 2014. Detecting and Evaluating Local Text 1220 Reuse in Social Networks. In Proceedings of the 1221 Joint Workshop on Social Dynamics and Personal 1222 Attributes in Social Media, pages 50–57, Baltimore, 1223 Maryland. Association for Computational Linguistics.

# Appendix

1228

1226

# A Inter-annotator agreement

Id	Target sequence	Agreement	Agreement	# Instances	
Bible reference	Reused text	Krippendorf's $\alpha$	Avg. pairwise Spearman's $\rho$	Number of pairs	
Overall	Overall	.420 / .709	.506 / .811	6,300 / 3,821	
(Matthew 18:22)	Seventy times seven	.118 / .764	.619 / .857	150 / 95	
(John 17:21)	That all may be one	.494 / .677	.183 / .782	150 / 79	
(Matthew 5:39)	Turn the other cheek	036 / .193	.510 / .405	120/132	
(Matthew 7:7)	Seek and you will find	.210 / .097	.558 / .475	150 / 125	
(Psalm 23:1)	The Lord is my shepherd	.213 / .138	.476 / .431	150 / 117	
(John 8:32)	The truth will set you free	.250 / .217	.347 / .368	150 / 10	
(1 Corinthians 13:4)	Love is patient, love is kind	.282 / .798	.382 / .834	150 / 61	
(Matthew 7:1)	Judge not, that ye be not judged	.472 / .450	.555 / .469	150/96	
(Ecclesiastes 3:1)	For everything there is a season	.263 / .369	.443 / .557	150 / 104	
(Romans 8:28)	All things work together for good	.110 /030	.543 / .430	150 / 128	
(2 Corinthians 5:7)	For we walk by faith, not by sight	.383 / .823	.437 / .866	150 / 79	
(Psalm 121:7)	The Lord will keep you from all harm	.169 / .178	.213 / .166	150 / 103	
(Mark 12:17)	Give to Caesar what belongs to Caesar	.196 / .117	.460 / .522	150 / 106	
(Proverbs 27:5)	Better is open rebuke than hidden love	.431 / .828	.481 / .911	150/95	
(Exodus 20:3)	You shall have no other gods before me	.259 / .506	.331 / .646	150 / 80	
(Genesis 1:1)	In the beginning God created the heaven	.151 / .177	.219 / .277	150 / 66	
(Romans 12:10)	Love one another with brotherly affection	.410 / .566	.571 / .657	150 / 101	
(Leviticus 20:13)	If a man lies with a male as with a woman	.315 / .492	.339 / .517	150 / 86	
(Joshua 1:9)	Be strong and courageous. Do not be afraid	.223 / .822	.281 / .833	150/69	
(Mark 9:23)	Everything is possible for one who believes	.081 / .509	.118 / .557	150 / 81	
(Philippians 4:13)	I can do all things through Christ who strengthens me	.128 / .624	.349 / .787	150 / 73	
(Ephesians 5:25)	Husbands, love your wives, as Christ loved the church	.487 / .802	.507 / .802	150/95	
(Matthew 5:44)	Love your enemies and pray for those who persecute you	.073 /004	.389 / .532	150 / 104	
(John 15:12)	My command is this: Love each other as I have loved you	.288 / .589	.426 / .790	150/97	
(Isaiah 43:4)	You are precious in my eyes and honored, and I love you	.421 / .439	.625 / .730	150/110	
(Matthew 7:25)	The rain came down, the streams rose, and the winds blew	.166 / .625	.406 / .802	150 / 71	
(1 Timothy 2:12)	But I suffer not a woman to teach, nor to usurp authority	.302 / .232	.315/.217	150/71	
(Proverbs 10:12)	Hatred stirs up conflict, but love covers over all wrongs	.172 / .148	.377 / .517	150 / 100	
(Hosea 8:7)	They have sown the wind, and they shall reap the whirlwind	.261 / .621	.423 / .668	150 / 55	
(Proverbs 12:25)	Anxiety weighs down the heart, but a kind word cheers it up	.093 / .518	.253 / .777	150/81	
(1 John 4:8)	Whoever does not love does not know God, because God is love	.329 / .355	.373 / .393	150/121	
(Solomon 4:7)	You are altogether beautiful, my darling: there is no flaw in you	.385 / .782	.537 / .878	150/92	
(Leviticus 18:22)	You shall not lie with a male as with a woman: it is an abomination	.423 / .648	.458 / .753	150/101	
(Psalm 118:24)	This is the day that the Lord has made: let us rejoice and be glad in it	.294 / .847	.492 / .884	150 / 77	
	A wife of noble character who can find? She is worth far more than				
(Proverbs 31:10)	rubies	.108 / .775	.261 / .797	150/74	
	Greater love has no one than this: to lay down one's life for one's				
(John 15:13)	friends	.347 / .355	.640 / .737	150/125	
	Come to me, all you who labour and are overburdened, and I shall				
(Matthew 11:28)	give you rest	007 /024	.268 / .355	150 / 101	
	The heart is deceitful above all things, and desperately wicked; who				
(Jeremiah 17:9)	can know it?	.164 / .432	.311 / .591	150 / 75	
	Now faith is confidence in what we hope for and assurance about				
(Hebrews 11:1)	what we do not see	.3917.853	.4107.870	150772	
	Take heed to yourselves. If your brother sins against you, rebuke him;	011 / 0/7	1211105	150 / 01	
(Luke 17:3)	and if he repents, forgive him	011 / .267	.124 / .485	150/91	
00:4: 51	Therefore, if anyone is in Christ, he is a new creation. The old has	207 1 127	414 / 5/4	150 / 25	
(2 Corinthians 5:17)	passed away; behold, the new has come	.307 / .655	.414 / .744	150/75	
(1.6	The Lord does not look at the things people look at. People look at the	422 1 665	500 / 722	150 / 00	
(1 Samuel 16:7)	outward appearance, but the Lord looks at the heart	.4327.005	.508 / ./33	150 / 80	

Table 4: Biblical passages included in TROTR and their inter-annotator agreement agreement. We report data using the x / y format, where x denotes the data on the entire set of instance pairs, and y denotes the data post-filtering process.

Text 1	Text 2	Semantic Textual Similarity	Semantic Textual Relatedness	Semantic Textual <i>Topic</i> Relatedness		
It's the wonderful pride month!! ♥ ♥♥♥♥♥ Honestly pride is every- day! Love is love don't forget I love you ♥. Remember this! John 15:12- 13: "My command is this: Love each other as I have loved you. Greater love has no one than this: to lay down one's life for one's friends"	Happy Pride Month! ♥ Remember, pride isn't just for a month—it's a daily celebration! Love knows no boundaries, and I want you to know that I cherish you every single day. ♥ Let's always remember these power- ful words from John 15:12-13: "My command is this: Love each other as I have loved you. Greater love has no one than this: to lay down one's life for one's friends"	√ paraphrase	related in some aspects	related in topic		
"Freeing people from genocide is the reason, motive & goal of the military operation we started in the Donbas & Ukraine", Putin says, then quotes the Bible: "There is no greater love than to lay down one's life for one's friends." It's like Billy Gra- ham meets North Korea	At a large Crimean event today Putin quoted the Bible to defend the special military operation in Ukraine which has killed thousands and displaced millions. His words "There is no greater love than if someone gives soul for their friends". And people were cheering him. Madness!!!	× neither paraphrases nor entailment	√ related in some aspects	related in topic		
It's the wonderful pride month!! ♥ ♥ ♥ ♥ ♥ Honestly pride is every- day! Love is love don't forget I love you ♥. Remember this! John 15:12- 13: "My command is this: Love each other as I have loved you. Greater love has no one than this: to lay down one's life for one's friends"	At a large Crimean event today Putin quoted the Bible to defend the special military operation in Ukraine which has killed thousands and displaced millions. His words "There is no greater love than if someone gives soul for their friends". And people were cheering him. Madness!!!	× neither paraphrases nor entailment	related in some aspects	× unrelated in topic		
You are altogether beautiful, my dar- ling; there is no flaw in you. Charm is deceitful, and beauty is vain, but a woman who fears the Lord is to be praised	At a large Crimean event today Putin quoted the Bible to defend the special military operation in Ukraine which has killed thousands and displaced millions. His words "There is no greater love than if someone gives soul for their friends". And people were cheering him. Madness!!!	× neither paraphrases nor entailment	× unrelated in any aspects	× unrelated in topic		

Table 5: Examples of *semantic textual similarity, semantic textual relatedness*, and *topic relatedness*. The first and last pair of sentences are examples of paraphrases and semantically unrelated content, respectively. Most people will agree that the second pair of sentences is more related in topic than the third pair of sentences. However, some people may still consider the third pair as semantically related due to the presence of the same quotation.

1232

1233

1234

1235

1236

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1258

1259

1260

1261

1263

1265

1266

1267

1268

1271

1272 1273

1274

1275

# **B** Train-Dev-Test partitions

For each randomized split, we use the filtered instances (see Section 4.2) to create the Train-Dev-Test partitions, comprising approximately 80%, 10%, and 10% of the instances, respectively. In the creation of the Train set of a split, we exclude the  $\langle t, c_1, c_2 \rangle$  instances associated to four targets t(i.e., 10% of the benchmark's targets). We include these instances in Dev and Test to enforce the Outof-Vocabulary (OOV) evaluation. Specifically, we include in Dev the instances associated with two targets, and in Test the instances of the remaining excluded targets.

Notably, we ensure that each partition has a distinct set of OOV targets, such that the intersection of the OOV sets for each split is empty.

# C Model evaluation

We evaluate almost all the pre-trained models available at https://www.sbert.net/index. html. Specifically, we considered only pre-trained models trained on tasks based on textual similarity and excluded those trained on other tasks (e.g., models for Image Search). Table 8 reports results for all the evaluated models.

For the sake of transparency and completeness, we have included the computation of Precision (PR) and Recall (RE) for each considered class. Specifically, for label 1, PR and RE are calculated as  $\frac{TP}{(TP+FP)}$  and  $\frac{TP}{(TP+FN)}$  respectively. Similarly, for label 0, PR and RE are computed as  $\frac{TN}{(TN+FP)}$ . In scientific literature, these latter metrics are also known as Negative Predictive Value and Sensitivity. For the sake of clarity, we preferred using PR and RE for *label 0* and *label 1* instead of distinguishing between Precision (PR), Recall (RE), Negative Predictive Value (NPV), and Specificity (SP).

As for TRaC label (see Equation 4), we emphasize that, given our use of the Spearman rank order correlation for evaluation, alternative normalization formulas, such as  $1 - \frac{s}{S_{max}}$ , can be employed without altering our results.

# D Fine-tuning

For each randomized split, we fine-tuned each considered model on the Train set and subsequently validated its performance on the Dev set. To do this, we employed the AdamW optimizer, coupled with a linear learning rate warm-up applied to the first 10% of the Train set. We used grid search to 1277 optimize hyper-parameters, with a particular focus 1278 on fine-tuning the learning rate by testing values 1279 from the set {1e-6, 2e-6, 5e-6, 1e-5, 2e-5}. We do 1280 not use weight decay, since our initial experiments did not yield any additional benefits. During the 1282 training, we leveraged an early stopping strategy. 1283 In particular, we fine-tuned each pre-trained model 1284 on TRiC instances using the contrastive loss (Had-1285 sell et al., 2006). This loss minimizes the distance 1286 between embeddings of similar sentences and maximizes the distance for dissimilar sentences. We 1288 finally ceased training when there was no further 1289 improvement observed on the Dev set. Details on 1290 the setup of hyper-parameters are shown in Table 7. 1291

# **E** Annotation

Annotating topic relatedness, instead of relying on 1293 explicit topic labels, closely resembles recent work exemplified in the Word-in-Context task (Pilehvar 1295 and Camacho-Collados, 2019), which relies on an-1296 notating word meaning relatedness rather than ex-1297 plicit sense labels. The methodology underlying 1298 this approach is thoroughly elucidated in our guidelines, submitted as supplementary material along 1300 with our paper. The topic relatedness is evaluated 1301 by using the four-point DURel relatedness scale in 1302 Table 6. Annotator were trained in a 30-minute on-1303 line session and tested on a small set of 25 instances 1304 (tutorial). In particular, we ensured that each an-1305 notator achieved a minimum agreement (measured 1306 by Spearman correlation) of at least .550 with the 1307 tutorial judgments. We interpreted these results 1308 as reliable, and consequently, we proceeded with 1309 the annotation of our benchmark. Then, we derive 1310 TRiC and TRaC labels after conducting an empiri-1311 cal analysis of the agreement of each level of our 1312 topic relatedness scale (see Section 4.2). 1313

- 4: Identical
- 3: Closely related
- 2: Distantly related
- 1: Unrelated

Table 6: The DURel relatedness scale proposedby Schlechtweg et al. (2023, 2018).

# **F** Hyper-parameters

Models	Learning Rate
all-distilroberta-v1(ADR)	1e-05
distiluse-base-multilingual-cased-v1(DBM)	1e-05
paraphrase-multilingual-MiniLM-L12-v2 (PAM)	2e-05
paraphrase-multilingual-mpnet-base-v2(PAR)	5e-06
multi-qa-mpnet-base-cos-v1(MQA)	1e-05

Table 7: Models learning rates.

1314

		Standard Test Set						Out-of-vocabulary (OOV) Test set								
Models	PR	Label 0	F1	PR	Label I RE	FI	E1	.ll SP	PR	Label 0	FI	PP	Label I RE	FI	E1 A	ll SP
paraphrase-multilingual-MiniLM-L12-v2 (PAM)	.96±.02	.46±.09	.61±.08	.41±.09	.96±.02	.57±.08	.61±.07	.58±.08	.96±.04	.43±.17	.57±.16	.37±.15	.95±.05	.52±.15	.59±.12	.49±.22
+MASK	.89±.05	.88±.06	.88±.03	.71±.10	.72±.10	.70±.05	.83±.03	.67±.04	.89±.09	.86±.09	.87±.06	.65±.19	.71±.18	$.65 \pm .12$	.83±.05	.60±.13
multi-qa-mpnet-base-cos-v1 (MQA) +MASK	.94±.03 88+.05	.42±.11 87+ 07	.58±.11 88+ 04	.40±.10 71+10	.94±.03 71+12	.55±.09 69±.06	.58±.09 83+ 04	.55±.09 68±.05	.94±.09 89±.07	.39±.19 86+ 10	.53±.20 87±.06	.36±.19	.96±.03 69+16	.50±.18 63+13	.55±.16 83±.05	.49±.21
all-distilroberta-v1 (ADR)	.95±.03	.47±.13	.62±.11	.42±.11	.93±.04	.57±.10	.61±.10	.55±.09	.94±.07	.45±.20	.58±.20	.38±.19	.93±.06	.51±.18	.58±.16	.48±.20
+MASK	.89±.05	.87±.07	.87±.03	.70±.14	.72±.12	.69±.07	.82±.03	.67±.06	.90±.07	.85±.10	.87±.05	.62±.21	.71±.18	.63±.14	.82±.05	.62±.15
all-mpnet-base-v2 +MASK	.93±.03	.48±.14 .84+.09	.62±.13	.42±.12	.91±.03	.5/±.10 .67+.04	.61±.11	.53±.10	.93±.09	.44±.22 .82+.11	.56±.21	.38±.20	.94±.05	.51±.18	.5/±.18 .81+.05	.48±.20
paraphrase-multilingual-mpnet-base-v2 (PAR)	.95±.03	.40±.10	.56±.09	.39±.09	.95±.04	.55±.08	.56±.07	.56±.09	.93±.11	.35±.18	.49±.19	.34±.15	.95±.06	.49±.16	.52±.15	.47±.25
+MASK	.89±.05	.85±.07	.87±.04	.69±.10	.75±.11	.70±.05	.83±.03	.68±.03	.90±.08	.83±.13	.86±.07	.63±.19	.75±.17	.65±.10	.82±.05	.62±.11
aii-MiniLM-L12-V2 +MASK	.95±.05 .88±.05	.40±.12 .87±.08	.55±.11 .87±.03	.39±.11 .70±.13	.94±.04 .72±.11	.54±.10 .69±.06	.55±.10 .82±.03	.52±.08 .68±.04	.94±.03 .89±.08	.3/±.18 .85±.10	.50±.18 .86±.05	.55±.18	.93±.06	.48±.18 .62±.14	.52±.16 .82±.04	.4/±.1/ .62±.13
multi-qa-distilbert-cos-v1	.96±.03	.33±.11	.48±.11	.37±.10	.97±.02	.53±.09	.50±.10	.53±.09	.97±.06	.29±.18	.42±.19	.33±.16	.97±.05	.47±.17	.46±.16	.47±.21
+MASK	.88±.06	.86±.06	.87±.03	.68±.11	.73±.10	.69±.05	.82±.03	.68±.05	.89±.10	.85±.08	.86±.05	.61±.19	.71±.14	.63±.11	.82±.05	.62±.14
+MASK	.92±.03	.46±.08	.86±.03	.69±.12	.65±.19	.63±.09	.80±.03	.63±.05	.91±.15 .87±.09	.45±.19	.85±.06	.62±.23	.63±.20	.57±.13	.80±.05	.40±.22
all-MiniLM-L6-v2	.96±.02	.40±.10	.55±.10	.39±.10	.95±.04	.55±.09	.56±.08	.53±.09	.97±.03	.37±.17	.51±.19	.35±.17	.95±.07	.49±.18	.54±.14	.44±.23
+MASK	.88±.05	.88±.06	.88±.03	.72±.12	.70±.12	.69±.06	.83±.03	.67±.05	.89±.07	.88±.09	.88±.05	.67±.22	.66±.19	.62±.14	.83±.04	.61±.16
+MASK	.87±.02	.20±.12 .88±.07	.40±.14 .87±.03	.72±.14	.97±.03	.51±.09	.45±.12 .81±.03	.54±.09	.90±.08 .88±.09	.21±.19 .88±.09	.31±.25 .87±.05	.51±.14	.64±.25	.43±.10 .58±.19	.38±.18 .82±.04	.44±.23 .58±.12
distiluse-base-multilingual-cased-v2	.96±.03	.26±.08	.40±.10	.34±.09	.97±.03	.50±.09	.43±.09	.54±.10	.96±.08	.21±.16	.32±.20	.30±.15	.96±.08	.44±.17	.38±.16	.44±.25
+MASK	.87±.06	.89±.07	.87±.03	.72±.14	.66±.14	.66±.09	.82±.03	.65±.04	.88±.08	.88±.10	.87±.05	.66±.24	.64±.23	.60±.18	.82±.05	.59±.12
muiti-qa-distribert-dot-v1 +MASK	.95±.04 .85±.05	.40±.12 .87±.08	.55±.11 .85±.03	.39±.09	.92±.05	.54±.09 .61±.08	.56±.09	.51±.09 .62±.05	.92±.12 .86±.09	.30±.10 .87±.09	.50±.16 .86±.05	.54±.15	.92±.07	.48±.16 .55±.16	.55±.11 .80±.03	.43±.19 .57±.14
paraphrase-albert-small-v2	.96±.02	.36±.09	.52±.09	.38±.09	.96±.02	.54±.09	.53±.07	.53±.09	.95±.10	.32±.16	.46±.18	.33±.14	.97±.04	.48±.16	.50±.12	.43±.25
+MASK	.88±.06	.84±.07	.86±.03	.65±.11	.70±.14	.66±.07	.80±.02	.65±.05	.88±.08	.82±.12	.84±.07	.56±.19	.67±.20	.58±.14	.80±.05	.57±.14
multi-qa-MiniLM-L6-cos-v1 +MASK	.95±.03 .88±.05	.3/±.09 .88±.04	.52±.09 .88±.02	.38±.10 .70±.09	.95±.04 .69±.09	.53±.10 .68±.06	.53±.08 .83±.02	.52±.10 .66±.04	.91±.14 .87±.09	.34±.18 .87±.07	.48±.19 .87±.05	.34±.17	.94±.08 .64±.19	.4/±.1/ .60±.14	.50±.16 .82±.04	.42±.25 .60±.15
cross-encoder/stsb-roberta-large	.33±.08	.99±.02	.49±.09	.97±.03	.19±.10	.30±.13	.36±.11	.52±.07	.29±.14	.99±.03	.42±.16	.94±.15	.11±.15	.18±.19	.28±.15	.42±.20
+MASK	.70±.13	.68±.15	.66±.07	.87±.06	.87±.08	.87±.03	.81±.03	.66±.04	.62±.27	.64±.28	.57±.21	.87±.10	.86±.11	.86±.06	.80±.06	.62±.08
parapnrase-MiniLM-L3-V2 +MASK	.95±.05	.28±.09	.45±.10	.35±.08	.96±.03	.51±.09	.40±.08	.49±.11	.96±.05	.23±.19	.34±.21	$.51\pm.14$ $.61\pm.21$	.97±.05	.45±.16	.41±.17	.40±.27
msmarco-distilbert-dot-v5	.93±.04	.36±.10	.51±.09	.37±.09	.93±.03	.52±.08	.52±.08	.47±.08	.92±.09	.31±.18	.43±.20	.32±.14	.92±.08	.46±.15	.48±.13	.38±.19
+MASK	.87±.05	.91±.04	.89±.03	.75±.10	.66±.08	.69±.06	.84±.03	.64±.04	.87±.09	.90±.05	.88±.06	.67±.18	.60±.16	.61±.15	.83±.06	.58±.10
msmarco-MiniLM-L12-cos-v5 +MASK	.91±.04 .85±.05	.44±.09 .88±.06	.59±.08 .86±.03	.39±.09 .68±.11	.90±.05 .60±.10	.54±.08 .62±.06	.58±.07 .80±.03	.44±.08 .59±.04	.91±.09 .85±.10	.44±.17 .88±.08	.58±.16 .86±.06	.36±.16	.88±.10 .55±.21	.49±.16 .53±.16	.59±.12 .79±.05	.38±.19 .54±.12
multi-qa-MiniLM-L6-dot-v1	.89±.07	.54±.07	.67±.06	.42±.09	.84±.08	.55±.08	.64±.05	.46±.10	.87±.16	.51±.15	.63±.15	.37±.16	.83±.11	.49±.15	.62±.12	.37±.26
+MASK	.83±.07	.86±.06	.84±.03	.61±.13	.56±.12	.56±.07	.76±.04	.53±.07	.82±.12	.86±.08	.83±.07	.53±.23	.50±.20	.47±.17	.76±.08	.45±.18
+MASK	.95±.05 .85±.06	.41±.10 .87±.07	.36±.10 .86±.04	.59±.09 .67±.10	.92±.00	.62±.07	.30±.08 .79±.03	.44±.10 .59±.04	.95±.07 .85±.11	.36±.18	.32±.18 .85±.07	.54±.10	.58±.24	.48±.17 .55±.16	.79±.05	.57±.22
msmarco-distilbert-base-tas-b	.93±.04	.36±.13	.51±.13	.38±.09	.93±.05	.53±.09	.52±.11	.45±.10	.92±.10	.32±.22	.44±.21	.33±.15	.92±.10	.47±.16	.48±.17	.36±.23
+MASK	.86±.07	.86±.08	.86±.03	.67±.14	.64±.14	.63±.07	.80±.03	.62±.05	.86±.11	.87±.11	.85±.06	.61±.23	.59±.26	.53±.20	.79±.07	.56±.14
+MASK	.55±.08	.96±.04 .61±.15	.49±.08 .61±.07	.94±.06 .85±.07	.25±.10 .86±.09	.35±.12 .85±.04	.40±.10 .78±.04	.43±.08 .59±.04	.29±.14 .58±.22	.96±.06	.43±.15 .51±.17	.89±.21 .85±.11	.1/±.16 .84±.12	.27±.19 .84±.07	.34±.15 .77±.07	.50±.21 .55±.08
msmarco-distilbert-cos-v5	.94±.03	.30±.09	.45±.11	.36±.09	.95±.03	.51±.09	.48±.09	.42±.09	.91±.12	.26±.14	.38±.17	.31±.14	.94±.06	.44±.15	.43±.13	.34±.17
+MASK	.88±.05	.84±.06	.85±.03	.64±.10	.71±.11	.66±.07	.80±.02	.62±.03	.88±.08	.82±.08	.84±.05	.56±.19	.67±.16	.59±.15	.80±.04	.56±.09
+MASK	.32±.09	.98±.03	.48±.10	.96±.03	.10±.13	.20±.10	.33±.14	.41±.07	$.29\pm.14$ .61+.23	.97±.05	.43±.17	.77±.39	.13±.18 .85+.11	.19±.25	.28±.20	.54±.19
cross-encoder/stsb-roberta-base	.31±.08	.98±.02	.47±.09	.95±.05	.13±.07	.22±.10	.30±.08	.42±.07	.28±.14	.97±.05	.41±.16	.91±.15	.10±.10	.16±.15	.26±.13	.33±.20
+MASK	.68±.10	.64±.15	.64±.08	.86±.06	.87±.07	.86±.03	.80±.04	.63±.06	.57±.21	.57±.26	.52±.20	.86±.11	.86±.10	.85±.06	.78±.08	.57±.11
+MASK	.93±.03	.32±.10 .90±.05	.4/±.11 .88±.03	.74±.11	.94±.03	.51±.09	.49±.09	.43±.09 .65±.03	.91±.07 .86±.09	.20±.19	.38±.21 .88±.05	.51±.14	.92±.08	.43±.10 .58±.17	.43±.13 .82±.05	.55±.24 .58±.11
cross-encoder/ms-marco-TinyBERT-L-2-v2	.32±.08	.97±.02	.48±.09	.93±.06	.17±.11	.28±.14	.34±.12	.34±.10	.29±.14	.97±.03	.43±.16	.78±.30	.13±.19	.20±.23	.29±.19	.26±.20
+MASK	.67±.15	.64±.14	.63±.07	.86±.06	.86±.09	.85±.04	.79±.03	.60±.06	.60±.23	.61±.24	.55±.17	.87±.10	.86±.12	.85±.05	.79±.06	.55±.15
cross-encoder/ms-marco-MiniLM-L-2-V2 +MASK	.52±.08 .67±.14	.97±.02 .61±.13	.48±.09 .62±.07	.94±.05 .85±.06	.10±.12 .87±.07	.20±.15 .85±.03	.33±.13 .79±.03	.50±.10 .57±.08	.29±.14 .58±.22	.97±.05 .55±.25	.43±.16 .50±.18	.91±.15 .85±.11	.13±.20 .85±.10	.19±.24 .84±.06	.29±.20 .77±.07	.20±.23 .51±.16
cross-encoder/ms-marco-MiniLM-L-4-v2	.32±.08	.95±.03	.47±.09	.89±.04	.18±.10	.29±.13	.35±.11	.31±.10	.29±.14	.93±.09	.42±.17	.91±.10	.16±.16	.24±.20	.32±.17	.24±.22
+MASK	.63±.13	.64±.13	.62±.07	.86±.06	.83±.08	.84±.04	.78±.04	.56±.07	.56±.21	.62±.25	.53±.16	.87±.11	.81±.14	.83±.07	.77±.06	.52±.15
cross-encoder/quora-roberta-base +MASK	.51±.08 .63±.12	.55±.02	.46±.09 .58±.08	.90±.04 .83±.04	.10±.05 .87±.03	.18±.07 .85±.03	.2/±.0/ .78±.03	.32±.08 .47±.09	.28±.14 .58±.30	.98±.05 .47±.18	.41±.17 .49±.20	.78±.39 .84±.08	.09±.10 .88±.08	.15±.15 .85±.05	.25±.13 .79±.04	.23±.17 .41±.16
cross-encoder/quora-roberta-large	.31±.08	.97±.05	.46±.09	.26±.40	.09±.15	.13±.21	.23±.17	.31±.10	.28±.14	.97±.06	.41±.17	.25±.39	.08±.19	.11±.23	.22±.21	.22±.19
+MASK	.40±.20	.76±.37	.40±.10	.22±.34	.29±.44	.25±.38	.30±.26	.48±.08	.35±.25	.73±.41	.32±.20	.23±.36	.29±.44	.25±.39	.30±.28	.42±.14
cross-encoder/ms-marco-MiniLM-L-6-v2 +MASK	.33±.09 .63±.14	.91±.05 .62±.13	.48±.09 .60±.08	.8/±.06 .85±.06	.25±.11 .84±.07	.5/±.12 .84±.03	.41±.11 .78±.03	.30±.09 .55±.07	.29±.15 .55±.24	.88±.12 .57±.26	.42±.17 .49±.20	.89±.10 .86±.11	.23±.15 .83±.12	.34±.16 .83±.05	.39±.14 .77±.05	.21±.18 .50±.15
cross-encoder/ms-marco-MiniLM-L-12-v2	.33±.09	.79±.07	.46±.09	.81±.07	.35±.09	.49±.09	.49±.07	.24±.09	.30±.17	.78±.17	.41±.18	.82±.16	.34±.15	.47±.16	.48±.12	.20±.16
+MASK	.58±.11	.58±.13	.56±.05	.83±.06	.81±.09	.82±.05	.75±.04	.48±.05	.47±.19	.53±.25	.45±.18	.83±.11	.80±.12	.81±.08	.74±.06	.44±.09
cross-encoder/quora-distilroberta-base +MASK	.31±.08 .39±.20	.9/±.06	.46±.09	.18±.36	.08±.18 .20±.39	.11±.23	.22±.18	.25±.10 .34±.10	.28±.14 .34+.19	.98±.05 .81+.37	.42±.17 .36±.21	.19±.38	.07±.20 .20±.40	.09±.23	.21±.21	.16±.20 .30±.18
cross-encoder/qnli-electra-base	.33±.10	.45±.12	.36±.08	.74±.08	.63±.12	.67±.08	.58±.07	.04±.11	.31±.18	.49±.18	.34±.16	.78±.14	.64±.14	.68±.09	.60±.11	.07±.18
+MASK	.41±.12	.36±.12	.35±.07	.74±.08	.77±.14	.74±.08	.63±.08	.07±.08	.40±.23	.38±.19	.32±.11	.77±.14	.78±.16	.75±.10	.64±.12	.11±.14
cross-encoder/qnli-distilroberta-base +MASK	.31±.09 .46±.24	.30±.18 .31±.15	.35±.10 .30±.11	.73±.07	.33±.18 .77±.16	.00±.11 .74±.08	.53±.08 .61±.06	.05±.06 .13±.09	.30±.18 .32±.27	.48±.19 .26±.15	.52±.14 .24±.13	.75±.13	.55±.19 .78±.19	.39±.14 .74±.11	.54±.12 .62±.10	.02±.10 .15±.13

Table 8: **TRiC evaluation** using various SBERT models on Subtask 1 and Subtask 2. Results are presented for each model using pre-trained models and the +MASK setting (*italic*). For Subtask 1, precision (PR), recall (RE), and Weighted -F1 scores (F1) are reported for both label 0 (i.e., different topics) and label 1 (i.e., roughly identical topics). For Subtask 2, Spearman correlation (SP) is reported on the overall set of instances. The reported metrics include standard deviations (±) across the 10 Test splits for comparative analysis. The superior performance for each metric between pre-trained models is highlighted in **bold**. Results for both Test and OOV Test sets are provided for completeness.