Annotating Training Data for Conditional Semantic Textual Similarity Measurement using Large Language Models

Anonymous ACL submission

Abstract

002

006

016

017

022

024

035

040

042

043

Semantic similarity between two sentences depends on the aspects considered between those sentences. To study this phenomenon, Deshpande et al. (2023) proposed the Conditional Semantic Textual Similarity (C-STS) task and annotated a human-rated similarity dataset containing pairs of sentences compared under two different conditions. However, Tu et al. (2024) found various annotation issues in this dataset and showed that manually re-annotating a small portion of it leads to more accurate C-STS models. Despite these pioneering efforts, the lack of large and accurately annotated C-STS datasets remains a blocker for making progress on this task as evidenced by the subpar performance of the C-STS models. To address this training data need, we resort to Large Language Models (LLMs) to correct the condition statements and similarity ratings in the original dataset proposed by Deshpande et al. (2023). Our proposed method is able to re-annotate a large training dataset for the C-STS task with minimal manual effort. Importantly, by training a supervised C-STS model on our re-annotated training data we achieve a novel state-of-theart (SoTA) for C-STS, thereby validating the accuracy of our dataset. The re-annotated dataset is submitted anonymously to ARR and will be publicly released upon paper acceptance to expedite the progress of C-STS research.

1 Introduction

Semantic Textual Similarity (STS) is a fundamental Natural Language Processing (NLP) task to evaluate the semantic similarity between two given sentences (Agirre et al., 2012). However, the focus on the sentences can vary and affects the judgment of the similarity. To address this, Deshpande et al. (2023) introduced a novel C-STS task, which measures the similarity between two sentences under a specified condition. In the C-STS dataset, each sentence pair has two conditions – a condition c_{low}



Figure 1: An example C-STS instance. The two sentences are compared under two different conditions, focusing on different aspects, resulting in a high (score of 5), and a lower (score of 1) semantic similarities.

producing a low semantic similarity, and a condition c_{high} a high semantic similarity, as shown in Figure 1. The similarity under each condition is rated on an ordinal scale from 1 (low similarity) to 5 (high similarity).

While the C-STS task brings greater specificity to the aspects of sentences being compared, Tu et al. (2024) observed that both the conditions and human similarity ratings suffer from issues such as ambiguity and inaccuracy, introducing label noise into the task. Although recent methods (Li et al., 2024; Liu et al., 2025; Yoo et al., 2024) have advanced the modeling of C-STS, their performance is still limited by the dataset quality, with Spearman correlations generally remaining below 0.5. To reduce those identified annotation errors, Tu et al. (2024) re-annotated the validation portion of the dataset with the help of human annotators. However, as later discussed in §2.1, in addition to annotation errors in similarity ratings, we find that the conditions themselves can be problematic, such as expressing varying granularities and a high-level of subjectivity, further impacting the reliability of the dataset. Moreover, the validation data re-annotated by (Tu et al., 2024) consists of only a small proportion (25%) of the C-STS dataset. Although it would be ideal to manually re-annotate the full C-STS dataset providing better training data for the C-STS prediction models, it is a significantly

labour intensive and costly annotation task.

074

075

084

091

100

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

To address this data cleansing task, we use LLMs to (1) modify the conditions, and (2) re-annotate the similarity ratings between two sentences under the modified conditions, requiring minimum manual effort. LLMs have been successfully used to generate synthetic training data and to provide judgements for several related NLP tasks (Peng et al., 2023; Patel et al., 2024; Wei et al., 2024). It is noteworthy that prior work (Deshpande et al., 2023) using LLMs such as GPT-4 (OpenAI et al., 2023) and Flan-T5 (Chung et al., 2024) to predict C-STS have reported suboptimal performance where they observed numerous issues including semantically similar sentence pairs being incorrectly assigned with low similarity scores. While we also use LLMs to correct the conditions and similarity ratings, we aim to improve the effectiveness of the C-STS training data by increasing both the number of instances and the accuracy of the annotations such that better C-STS models can be trained.

Our contributions in this paper are three-fold: (a) We first correct the errors in the condition statements in the C-STS dataset (§2.1). (b) Next, we use an ensemble consisting of two LLMs (i.e. GPT-40 and Claude-3.7-Sonnet) to independently obtain C-STS ratings, which we then combine with the original human ratings via a voting scheme (??). We found that both of those LLMs demonstrate a high level of agreement with the human C-STS ratings, resulting in Spearman correlations of 62% and 66%, respectively. (c) To evaluate the usefulness of our LLM-cleansed dataset, we train a supervised C-STS model on it following the method proposed by Zhang et al. (2025). The trained model obtains a Spearman correlation of 74% against the humanrated test data, thereby establishing a novel SoTA for C-STS.

2 C-STS Training Data Cleansing

Our data cleansing method for C-STS consists of two steps. In the first step (§2.1), we identify common issues with the conditions and use GPT-40 to refine those. In the second step (??), we re-annotate the labels using both GPT-40 and Claude-3.7-Sonnet, due to their high performance on natural language understanding demonstrated by Chatbot Arena leaderboard (Zheng et al., 2023).¹ Finally, we use a voting method to aggregate the human ratings in the original dataset with

Issue	Condition
Imbalanced Condition	number of # type of # color of #
Subjective Condition	The age of person. The color of animal. The number of people.
Inconsistent Phrasing Style	The all are food. Where the dog is visible from. The amount of stoves/ ovens. Type of room. The person's age.
Varying Granularity	The absence of tomato. The place of the object. The species of the one who's in the room.
Redundant Expression	The fact that they're both girls. String instrument being played. The players move to the position.
Grammatical Issue	The thing that fly.

Table 1: Common stand-alone condition issues.

the two sets of LLM ratings.

2.1 Modifying the Conditions

We identify multiple issues in the conditions that impact the accuracy of the human annotations. These issues fall into two categories: (1) conditions that are inherently ambiguous or misleading in their own, and (2) conditions that are misleading when interpretting the sentence semantics. 123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

2.1.1 Stand-alone Condition Issues

Imbalanced Conditions: Certain condition types occur far more frequently than the others, resulting in a highly imbalanced distribution (see Appendix A), biasing model training and evaluation. **Subjective Conditions:** Some conditions introduce discrepancies and conflicts in human similarity ratings because annotators can interpret the conditions differently. As a result, different annotators can assign contradicting similarity scores to the same sentence-pair. Appendix B presents examples and explanations highlighting such subjectivity and inconsistency in human similarity judgments. The frequently occurring subjective conditions in the original C-STS dataset introduce noise and reduce the reliability of model evaluation.

Inconsistent Phrasing Styles: The phrasing of some conditions is inconsistent, ranging from

¹https://lmarena.ai/

Issue	Sentence Pair	Condition
Ambiguous Condition	A climber with a yellow backpack walks along the ridge of a snowy mountainside. A person in a red hat with a huge backpack going hiking.	The climber.
Invalid Condition	A man wearing yellow and blue is riding a large, bucking bull. A bull rider, in full padding and wearing a helmet, rides a large brown and white bull.	Color of bull.
Unrelated Condition	Three hotdogs on buns with whole slices of relish sit on a white plate. A hot dog on a bun with a drop of ketchup on the table.	The number of <u>dogs</u> .

Table 2: Common condition issues that cause the judgment divergence related to sentences.

148full sentences to fragmented sentences or phrases.149Moreover, they lack uniformity in both stopword150usage and their grammatical structure.

Varying Granularity: Conditions range from
very general to overly specific. This divergence
affects how the models interpret those conditions.
Redundant Expressions: Conditions can some-

times include redundant words or phrases.

155

156

157

158

160

161

164

165

166

168

169

170

171

172

Grammatical Issues: Obvious English grammatical errors exist in some of the conditions. Table 1 shows examples of the above issues.

159 2.1.2 Sentence-dependent Condition Issues

Ambiguous Conditions: Tu et al. (2024) found that conditions presented as singletons without associated entity features to be ambiguous, lacking a clear specification of the aspects being compared.
Invalid Conditions: Tu et al. (2024) showed that some of the conditions to be invalid, as they required information that cannot be inferred from the sentences based on those conditions.

Unrelated Conditions: Some conditions contain typos or imprecise expressions. Although comprehensible by humans, it could mislead LLM judges. Table 2 shows examples of abovementioned issues.

To standardise the condition expressions and im-173 prove their specificity and accuracy to reduce am-174 biguity, we use GPT-40 to refine the conditions. 175 Specifically, we instruct GPT-40 using a prompt 176 that provides explicit guidelines and constraints. 177 The prompt requires that conditions to be clear, 178 specific, and semantically grounded, discouraging 180 vague references (e.g., "animal") in favour of more precise formulations (e.g., "species of animal"). 181 We also remove unnecessary stopwords (e.g., "the") and maintain a uniform phrasing style across all conditions. Additionally, the prompt requests a jus-184

tification for any substantive modifications. The complete prompt, along with examples before/after the modified conditions, is provided in Appendix C.

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

207

209

210

211

212

213

214

215

216

217

218

219

220

2.2 Re-annotating the Similarity Ratings

After refining the conditions, we use LLMs to re-annotate the similarity ratings in the training set. Specifically, we use GPT-40 and Claude-3.7-Sonnet with a few-shot prompt, providing five examples covering similarity ratings (1-5), each accompanied by a human-written justification. We also require LLMs to give corresponding justifications for their similarity ratings. This design serves two purposes: (1) it helps the LLM to understand the scoring rubric in a conditional STS context; and (2) it encourages the generation of not only a similarity rating but also a justification, which serves as a self-check mechanism to reduce hallucinations and improve the annotation quality. We use the same five-point rating scale proposed by Deshpande et al. (2023) and instruct the LLMs to only return a JSON-formatted object instead of a natural language commentary. The complete prompt, along with examples before/after re-annotating the similarity ratings under the modified conditions is provided in Appendix D.

To increase the reliability of the annotations, we use a voting approach to combine original human ratings with the LLM-predicted ratings. Specifically, for each instance, we compute the arithmetic mean of the original human-annotated similarity rating (y^{human}) , the predicted rating by GPT-40 $(y^{\text{GPT-40}})$, and the predicted rating by Claude-3.7-Sonnet (y^{Claude}) , and round to the nearest integer. As shown in Appendix E, combining ratings from both LLMs results in the best performance, justifying the proposed voting scheme.

Train	Test	Spearman
ReVal	ReTest	61.28
ReVal-Mod w/o	ReTest-Mod w/o	64.25
ReVal-Mod w/	ReTest-Mod w/	66.89

Table 3: Comparison of condition modification, evaluated using the supervised non-linear projection. *w/* and *w/o* denote condition modification with and without stopword removal, respectively.

3 Experiments

253

254

259

262

221

We define the following dataset naming conventions. **Train-Orig** is the original training set from Deshpande et al. (2023). **Train-Mod** applies condition modifications to Train-Orig, and **Train-Mod-Reanno** further includes LLM-generated ratings. **Val-Orig** denotes the original validation set, and **Val-Reanno** is the *human* re-annotated version introduced by Tu et al. (2024). Val-Reanno is the most accurate human-verified C-STS data that exists. We split **Val-Reanno** into **ReVal** (randomly selected 70%) as our validation set and **ReTest** (remaining 30%) as our test set. We construct **ReVal-Mod** and **ReTest-Mod** by applying condition modifications to ReVal and ReTest, respectively.

To evaluate the effectiveness of a particular training dataset, we first use it to train a supervised Multi-Head Non-Linear Projection (**MH**) model following the SoTA method proposed by Zhang et al. (2025), and then measure the improvement of C-STS task performance on the same humanlabelled test data (ReTest). Details of this supervised model architecture are provided in Appendix F. Spearman's correlation coefficient with human similarity ratings is the standard evaluation metric for C-STS, where a high correlation indicates an accurate C-STS model. We use an NVIDIA RTX A6000 GPU with PyTorch 2.0.1 and CUDA 11.7 for our experiments.

To evaluate the effectiveness of condition modification, we train **MH** models on ReVal and evaluate on ReTest as shown in Table 3. Further effect of stopword removal from the modified conditions is also considered. We see that the best performance is reported by the LLM-based condition modification with stopword removal (i.e. ReText-Mod w/). Stopwords often contribute little or no semantic distinctions to the conditions, and removing them helps the model to attend to content words.

Following these findings, we apply condition modification with stopword removal and follow **??** to re-annotate the similarity ratings in the condition-

Train	Test	Spearman
Train-Orig	ReTest	68.54
Train-Orig	ReTest-Mod	69.68
Train-Mod	ReTest-Mod	69.39
Train-Mod-Reanno	ReTest-Mod	74.56

Table 4: Spearman correlation coefficients obtained by training a **MH** model on different training datasets

263

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

283

284

285

287

289

290

291

292

293

295

296

297

299

300

301

302

303

304

modified C-STS training set. To measure the consistency of LLM-generated annotations, we randomly select 100 instances from Train-Mod and repeat the annotation process five times using Claude-3.7-Sonnet with our few-shot prompt. We measure the agreement of the five sets of annotations using the Krippendorff's Alpha (Hayes and Krippendorff, 2007). We find a high level of annotation consistency, indicated by a resulting Krippendorff's Alpha of 0.865. Moreover, we manually reviewed 100 randomly selected instances with the modified conditions and the re-annotated similarity ratings to assess the data quality, confirming an overall improvement in the training data, as further elaborated in Appendix G.

To evaluate the ability of our LLM modified conditions and re-annotated similarity ratings for improving C-STS measurement, we train MH models using different training datasets in Table 4. Compared to training C-STS models on Train-Orig, we see that doing so on Train-Mod-Reanno results in the best performance. This is a 6% statistically significant improvement over the best biencoder C-STS performance reported by Zhang et al. (2025). This shows that, keeping the model architecture and all other training settings fixed, our re-annotated C-STS training data alone can improve the SoTA peformance of C-STS. We believe that our re-annotated C-STS training data will fill the gap for large-scale accurate training data for C-STS, and facilitate the future progress of C-STS research.

4 Conclusion

We identify key issues in the condition definitions and human-annotated similarity ratings in the original C-STS dataset. To address these, we propose an efficient LLM-based data cleansing approach that improves dataset quality through condition modification and re-annotation of similarity scores. By integrating this with human-annotated data, our cleansed dataset significantly advanced the performance of a previously proposed C-STS method.

5 Limitations

305

307

311

313

314

316

317

318

320

321

322

323

324

325

330

331

332

335

336

337

339

340

341

343

346

347

350

353

354

There is a large number of LLMs developed and made publicly available. However, it is practically infeasible to use multiple LLMs for the C-STS data re-annotation due to the costs involved. Therefore, we selected two highly popular and accurate models (GPT-40 and Claude-3.7-Sonnet) to balance performance and cost-effectiveness. Although we modified the conditions, certain stand-alone condition issues such as imbalanced conditions still exist, as the overall distribution of condition types has not changed.

This study was conducted using C-STS datasets for English, which is a morphologically limited language. However, this choice is based on the availability of C-STS datasets. To the best of our knowledge, C-STS datasets are not publicly available for languages other than English. We consider it to be an important task for future work to develop multilingual C-STS datasets to study the languagespecific issues pertaining to this task.

6 Ethical Concerns

LLMs have been shown to exhibit social biases, such as those related to age and gender (Gallegos et al., 2024). Such social topics exist in the conditions for the C-STS task. Using LLMs for annotation may further propagate such biases into the dataset. The influence of whether the LLM-based annotation process impacts the data quality with respect to social bias is not evaluated. Additionally, LLM-based condition-aware sentence embeddings could encode unfair social biases. Therefore, it is important to evaluate social bias amplifications (if any) due to training C-STS models on our proposed training dataset before deploying those models in downstream NLP applications.

References

- Eneko Agirre, Daniel Matthew Cer, Mona T Diab, and Aitor Gonzalez-Agirre. 2012. SemEval-2012 task
 6: A pilot on semantic Textual Similarity. *SemEval*, pages 385–393.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, and 16 others. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.

Ameet Deshpande, Carlos E Jimenez, Howard Chen, Vishvak Murahari, Victoria Graf, Tanmay Rajpurohit, Ashwin Kalyan, Danqi Chen, and Karthik Narasimhan. 2023. CSTS: Conditional semantic textual similarity. *Empir Method Nat Lang Process*, pages 5669–5690. 355

356

358

361

362

363

364

365

366

367

370

371

372

373

374

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

- Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097– 1179.
- Andrew F Hayes and Klaus Krippendorff. 2007. Answering the call for a standard reliability measure for coding data. *Communication methods and measures*, 1(1):77–89.
- Baixuan Li, Yunlong Fan, and Zhiqiang Gao. 2024. Seaver: Attention reallocation for mitigating distractions in language models for conditional semantic textual similarity measurement. In *Findings of the Association for Computational Linguistics: EMNLP* 2024, pages 78–95.
- Xinyue Liu, Zeyang Qin, Zeyu Wang, Wenxin Liang, Linlin Zong, and Bo Xu. 2025. Conditional semantic textual similarity via conditional contrastive learning. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 4548–4560.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2023. GPT-4 Technical Report. *arXiv* [cs.CL].
- Ajay Patel, Colin Raffel, and Chris Callison-Burch. 2024. Datadreamer: A tool for synthetic data generation and reproducible llm workflows. *arXiv preprint arXiv:2402.10379*.
- Letian Peng, Yuwei Zhang, and Jingbo Shang. 2023. Generating efficient training data via llm-based attribute manipulation. *arXiv preprint arXiv:2307.07099*.
- Jingxuan Tu, Keer Xu, Liulu Yue, Bingyang Ye, Kyeongmin Rim, and James Pustejovsky. 2024. Linguistically conditioned semantic textual similarity. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1161–1172, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Hui Wei, Shenghua He, Tian Xia, Fei Liu, Andy Wong, Jingyang Lin, and Mei Han. 2024. Systematic evaluation of llm-as-a-judge in llm alignment tasks: Explainable metrics and diverse prompt templates. *arXiv preprint arXiv:2408.13006*.

- 410 411
- 412 413
- 414 415
- 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

- Young Yoo, Jii Cha, Changhyeon Kim, and Taeuk Kim. 2024. Hyper-CL: Conditioning sentence representations with hypernetworks. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 700-711, Bangkok, Thailand. Association for Computational Linguistics.
- Gaifan Zhang, Yi Zhou, and Danushka Bollegala. 2025. Case – condition-aware sentence embeddings for conditional semantic textual similarity measurement. Preprint, arXiv:2503.17279.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. Preprint, arXiv:2306.05685.

Supplementary Materials

A **Imbalanced Condition**

By analysing the distribution of condition types in the C-STS training dataset, we observe a significant imbalance. As shown in Table 5, two broad condition categories, number of # and type of #, dominate the dataset, accounting for 16.7% and 16.6% of all conditions, respectively.

With respect to specific conditions, we present the 15 most frequent ones in Table 6. The most common conditions include The number of people., The type of animal., and The sport. However, these frequently occurring conditions often introduce problems such as ambiguity and subjectivity in the evaluation process.

Condition Type	Count	Percentage
number of #	1892	16.7%
type of #	1886	16.6%
color of #	664	5.9%
action	357	3.1%
position of #	88	0.8%

Table 5: Counts of general condition types (top 5) in the original C-STS training dataset.

B Subjectivity in Human Annotations

Human annotators can give contradictory ratings to some similar instances in the dataset. We show subjectivity in human ratings for the conditions The number of people, Age of person and Gender of person in the original C-STS training dataset as examples. Table 7 lists some examples of instances

Condition	Count
The number of people.	520
The type of animal.	254
The sport.	249
The name of the place.	162
The animal.	154
The color of the shirts.	123
The number of people visible.	103
The action.	94
The type of food.	87
The number of animals.	85
The type of clothing.	85
The number of people in the image.	72
The location.	65
The color of the clothing.	64
The number of objects.	62

Table 6: Counts of specific conditions (top 15) in the original C-STS training dataset.

that show subjectivity. We explain them one by one as follows.

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

Considering the condition The number of people:

In the instance that Sentence 1: A man and woman sitting in a booth together and smiling., Sentence 2: Three people sitting at a table at a restaurant., Rating: 4, there are 2 people in Sentence 1, and 3 people in Sentence 2. While the number of people differs (2 vs. 3), annotators still rated the pair as highly similar. This suggests that some annotators perceive small differences in number (such as 2 versus 3) as relatively minor.

In the instance that Sentence 1: A baseball player swings to hit the ball as another player catches., Sentence 2: A man in a white and black uni- form is attempting to swing a baseball bat., Rating: 4, there are 2 people in Sentence 1 and 1 person in Sentence 2. Human annotators give this small difference in number a score of high similarity.

However, in another instance that Sentence 1: A person is diving into blue water on a rocky coast., Sentence 2: Two males on a rock over water, one in midair about to dive., Rating: 1, there are 1 person in Sentence 1 and 2 people in Sentence 2. The number of people is also different, but similar in number (same case as the previous example). Some annotators interpret it as a stronger signal of dissimilarity.

Additionally, in the instance that Sentence 1: A

Sentence 1	Sentence 2	Condition	Rating
A man and woman sitting in a booth together and smiling.	Three people sitting at a table at a restaurant.	The number of people.	4
A baseball player swings to hit	A man in a white and black uni-	The number of people.	4
the ball as another player catches.	form is attempting to swing a		
	baseball bat.		
A person is diving into blue water	Two males on a rock over water,	The number of people.	1
on a rocky coast.	one in midair about to dive.		
A person is doing a trick in the	Person performing a move on a	The number of people.	1
air on a bike near some buildings.	mountain bike with two people		
	watching.		
A young girl with a sippy cup	A child is making a ridiculous	The number of people.	4
swings on a swing.	face with an open mouth.		
The boy on the bike is wearing	A man dressed in bicycle gear is	Age of person.	1
safety glasses and a red helmet.	riding through a course.		
Two images show a man reach-	A boy in black shorts jumps and	Age of person.	3
ing out to hit a tennis ball with a	holds his tennis racket out in		
racket.	front of him.		
A very happy child sits on a chair	A child is bouncing on a trampo-	Age of person.	3
on top of some rocks.	line that is near a house.		
A man in a red and yellow outfit	A woman is riding a bike with a	Gender of person.	1
is riding a bicycle on one wheel.	basket of flowers.		
A woman with a red scarf around	A man in a black hat looks very	Gender of person.	4
her neck is smiling.	happy.		
A little girl is brushing her teeth	A woman is brushing her teeth in	Gender of person.	1
in a bathroom.	a bathroom mirror.		
A man is skateboarding on the	A girl is rollerblading on a path.	Gender of person.	4
sidewalk.			

Table 7: Examples of sentence pairs under the conditions "The number of people", "Age of person", and "Gender of person" with subjective similarity ratings by human annotators in the original C-STS training set.

person is doing a trick in the air on a bike near some buildings., Sentence 2: Person performing a move on a mountain bike with two people watching., Rating: 1, there are 1 person in Sentence 1 and 3 people in Sentence 2. Human annotators can regard this mismatch in number as dissimilarity.

Moreover, in the instance that Sentence 1: A young girl with a sippy cup swings on a swing., Sentence 2: A child is making a ridiculous face with an open mouth., Rating: 4, both sentences have 1 person. Human annotators give a high similarity score of 4, even though the numbers are exactly the same.

Considering the condition *Age of person*:

In the instance that Sentence 1: *The boy on the bike is wearing safety glasses and a red helmet* and Sentence 2 is: *A man dressed in bicycle gear is riding through a course*, the rating is 1. The perceived age difference between "boy" and "man"

leads to a low similarity rating. Some annotators may weigh age references heavily when evaluating similarity.

In contrast, in the instance that Sentence 1: *Two images show a man reaching out to hit a tennis ball with a racket* and Sentence 2 is: *A boy in black shorts jumps and holds his tennis racket out in front of him*, the rating is 3. While the age difference between "man" and "boy" still exists, annotators give a moderate similarity score.

In another instance that Sentence 1: A very happy child sits on a chair on top of some rocks. and Sentence 2 is: A child is bouncing on a trampoline that is near a house, the rating is 3. Both sentences have description about the "child", which should be a higher similarity score of 4. At least, the label should be different with the previous example which compares the age of "man" and "child".

599

600

601

602

603

604

605

606

607

608

610

611

612

Considering the condition *Gnender of person*:

518

519

520

522

523

524

526

527

529

530

532

533

534

535

536

537

538

540

541

542

544

548

549

550

551

552

553

554

555

556

In the instance that Sentence 1: A man in a red and yellow outfit is riding a bicycle on one wheel and Sentence 2: A woman is riding a bike with a basket of flowers, the rating is 1. Some annotators view gender as a central feature for this condition, leading to a low similarity rating despite shared activity.

However, in the instance that Sentence 1: A woman with a red scarf around her neck is smiling and Sentence 2: A man in a black hat looks very happy, the rating is 4. Even though the genders differ, the facial expressions and emotional tone are similar, suggesting that some annotators focus more on affective similarity than gender cues, which is inaccurate.

In the instance that Sentence 1: A little girl is brushing her teeth in a bathroom. and Sentence 2: A woman is brushing her teeth in a bathroom mirror, the rating is 1. The gender is both sentences is female. Human annotators should not give a dissimilar score based on gender. When gender information matches across two sentences, it should not contribute to a higher dissimilarity rating.

In the instance that Sentence 1: A man is skateboarding on the sidewalk. and Sentence 2: A girl is rollerblading on a path., the rating is 4. The gender is male in Sentence 1, but the gender is female in Sentence 2. Humman annotators should not give a high similarity score of 4 to this mismatching gender information.

C Prompt Used for Modifying the Conditions

Figure 2 shows the full prompt for condition modification. Table 8 provides examples of how our prompt effectively refines various types of problematic conditions.

D Prompt Used for Similarity Annotations

Figure 3 shows the complete prompt for assigning
similarity ratings using LLMs. Table 9 provides
examples of the original and our re-annotated ratings, showing the improvement in the accuracy of
C-STS scores. Selected examples are based on the
conditions of the same semantic focus (conditions
modified only with stopword removal).

E Evaluating the Voting Method

Table 10 reports the average performance across different rating aggregation strategies. We use NV-Embed-v2 (NV) to first generate condition-aware sentence embeddings and then train the supervised multi-head non-linear projection as described in §3. The projection model is fixed with a hidden dimensionality of 768, two heads (each with output dimensionality 256), and a dropout rate of 0.1. Results show that combining human ratings with annotations from both LLMs yields the highest performance.

F Multi-Head Non-Linear Projection

The multi-head non-linear projection consists of concatenated non-linear MLP proposed by Zhang et al. (2025). These supervised models are Siamese bi-encoders tailored for the C-STS task which have proven high performance (Deshpande et al., 2023; Yoo et al., 2024). Each model takes as input two condition-aware embeddings corresponding to sentence 1 and sentence 2 with the condition, respectively.

Zhang et al. (2025) propose that input conditionaware sentence embeddings are generated from LLM-based models, using the prompt "Retrieve semantically similar texts to the [CONDITION], given the Sentence: [SENTENCE]." They show that the LLM-based embeddings work better than the Masked Language Model (MLM)-based embeddings. To improve the condition-specific relevance, a post-processing step of subtracting the corresponding embeddings of the conditions is applied after generating the condition-aware sentence embeddings. Here, the embeddings of the conditions are generated using the prompt "Retrieve semantically similar texts to a given Sentence: [CON-DITION]."

Let us denote the resulting LLM-generated condition-aware sentence embeddings by e_1, e_2 for each instance. The **Multi-Head Non-Linear Projection** (MH) is then defined as follows:

Given two input embeddings $\mathbf{e}_1, \mathbf{e}_2 \in \mathbb{R}^d$, the model computes multiple parallel non-linear transformations and aggregates them via concatenation. Each head independently applies a two-layer feed-forward network with ReLU activation and dropout:

 $\mathbf{t}_i = \text{Dropout}\left(\text{ReLU}\left(\mathbf{W}_1\mathbf{e}\right)\right) \tag{1}$

$$\mathbf{h}_{i} = \text{Dropout}\left(\text{ReLU}\left(\mathbf{W}_{2}\mathbf{t}_{i}\right)\right)$$
(2)

Condition issue	Before	After
Ambiguous Condition	The animal. The sport.	type of animal presence of vehicles
Unrelated Condition	The name of the game.	type of sport
Inconsistent Phrasing Style	What the person is holding.	object being held
Varying Granularity	The setting. Specific areas of the home.	urban environment areas of home
Redundant Expression	If a tv is present.	presence of television
Grammatical Issue	The food with plate. The the size of the room.	food on plate size of room

Table 8: Examples of conditions before and after using our condition modification prompt.

Sentence 1	Sentence 2	Condition	Before	After
A room that has white walls and a window shade up has	A bed appears to have nothing else on it except two pillow in	type of room	2	4
a double unmade bed on the	a bedroom.			
floor.				
A deep dish pizza in a metal	The margarita pizza is on a	type of pizza	5	3
pan topped with several kinds of toppings.	plate, and ready to be cut and served.			
Older men sitting on wooden	There are people looking at a	gender of people	5	3
benches on a sidewalk to-	booth and a woman and man			
gether, with scooters parked in	in a wheelchair on the side-			
the street and stores across the	walk.			
street.				
a man sitting on a couch with	A computer desk topped with	number of people	4	2
a silver laptop in a living room	a monitor and a keyboard next			
	to a mouse.			
A person flying a kite at the	Three people standing on the	action of people	5	3
beach while two others walk	shore of a sandy beach in front			
past him	of waves			
A colorful purple airplane sits	A white and gray passenger	type of vehicle	2	4
on the runway with a darkened	plane has just landed or is			
sky in the background.	about to take off.			
Two elephants are bathing in	A group of people stand on the	Name of animal	2	4
deep water as a person sits on	shore while watching an ele-			
one of their backs.	phant in the water.			

Table 9: Examples of ratings with modified condition before and after using our re-annotation prompt.

where $i \in \{1, \ldots, H\}$, H denotes the number of heads, and $W_{i,1} \in \mathbb{R}^{d \times h}$, $W_{i,2} \in \mathbb{R}^{h \times d'}$ are learned parameters for the *i*-th head.

The final projected embedding is obtained by concatenating the outputs from all heads:

$$\mathbf{z} = [\mathbf{h}_1; \mathbf{h}_2; \dots; \mathbf{h}_H] \in \mathbb{R}^{H \cdot d'}$$
(3)

619 Hyperparameters are tuned on our validation set

ReVal-Mod. We fix the batch size to 512 and the learning rate to 10^{-3} in all of the experiments on non-linear MLP and multi-head non-linear projection. For non-linear projection, we use a dropout rate of 0.15 and a hidden dimensionality of 1024. For multi-head non-linear projection, set the hidden dimensionality to 1024 or 768. Additionally, we can set the number of heads to 2 with output dimen-

Rating Data	Spearman
y ^{GPT-40}	70.88
y^{Claude}	71.95
$V(y^{\text{GPT-4o}} + y^{\text{Claude}})$	72.21
$V(y^{\text{GPT-4o}} + y^{\text{human}})$	71.11
$V(y^{Claude} + y^{human})$	72.74
$V(y^{human} + y^{GPT-4o} + y^{Claude})$	73.10

Table 10: Average Spearman Correlation based on rating data across different aggregation strategies. V() denotes taking the arithmetic mean and rounding to the nearest integer.

Model	Non-linear Multi-layer Perception (MLP)	Linear MLP
NV	69.30	69.95
SFR	62.85	59.22
GTE	64.16	56.10
E5	62.12	47.03
SimCSE_large	56.67	45.96
SimCSE_base	56.60	39.54

Table 11: Spearman correlation of embedding models based on supervised MLPs with reduced dimensionality 512.

sionality of 256, and 4 with output dimensionality of 128. The dropout rate is set to 0.1.

Zhang et al. (2025) found that LLM-based models work better than MLM-based models such as SimCSE for the C-STS task. Although a direct comparison with prior C-STS methods is challenging due to issues in the test sets and lack of implementation details (e.g., Tu et al. (2024) do not release their hyperparameters or test/validation splits), we include a comparison table to highlight the performance improvements achieved using the method proposed by Zhang et al. (2025). Table 11 shows the performance of different embedding models. Three are LLM-based: NV-Embed-v2 (NV), SFR-Embedding-Mistral (SFR), gte-Qwen2-7B-instruct (GTE). Three are MLMbased: Multilingual-E5-large-instruct (E5), supsimcse-roberta-large (SimCSE large), and supsimcse-bert-base-uncased (SimCSE base).² NV achieves the highest Spearman correlation, significantly outperforming all other models. Therefore, we select NV as the base model for evaluating dataset cleansing effectiveness in our study.

G Manual Investigation of Annotations

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

To validate the LLM modified conditions and reannotated similarity ratings, we randomly selected 100 instances from our dataset to conduct a manual verification. From this manual investigation, we found that overall, the LLM modified condition statements to be more precise compared to the original conditions. Importantly, we did not find any conditions that degrade in quality or their meaning altered significantly by the LLM-based modification process.

Investigating the re-annotated similarity ratings, we found that the similarity ratings to accurately reflect the true conditional semantic textual similarity in most cases. For example, cases where similar sentence pairs were previously labelled as dissimilar were correctly assigned higher similarity ratings during this re-annotation process. However, a small proportion of instances (approximately 9%) still remain slightly inaccurate. The positive aspect is that the disagreement in these cases is relatively minor, typically differing by only 1 rating point (recall that the similarity ratings are in [1,5]) from the human rating.

²All models are available at https://huggingface.co/ spaces/mteb/leaderboard and https://huggingface. co/princeton-nlp

This is a Conditional STS task: Evaluate the similarity between the two sentences, with respect to the condition.

Sentence pair has a label (score) between 1 and 5 as follows:

- 1. The two sentences are completely dissimilar with respect to the condition.
- 2. The two sentences are dissimilar, but are on a similar topic with respect to the condition.
- 3. The two sentences are roughly equivalent, but some important information differs or is missing with respect to the condition.
- 4. The two sentences are mostly equivalent, but some unimportant details differ with respect to the condition.
- 5. The two sentences are completely equivalent with respect to the condition.

Check and modify the provided condition if it is inaccurate or ambiguous, following these guidelines strictly:

- Conditions must be clear and specific. (e.g., instead of "animal", specify clearly such as "species of animal".)
- Remove stopword from conditions (e.g., "the").
- Conditions must accurately match human-annotated labels.
- Provide conditions concisely, without context-specific details. Good examples: color of clothing, type of event, intention of travel.
- Do NOT overly specify the condition more narrowly than the original meaning.

Return a JSON object with two fields:

'improved_condition': the improved condition, 'justification': a single sentence explaining why you update the condition. Give empty str this if only stopword 'the' is removed.

Figure 2: Prompt for modifying conditions

Definition: Evaluate the similarity between the two sentences, with respect to the condition. Assign the pair a score between 1 and 5 as follows: 1. The two sentences are completely dissimilar with respect to the condition. 2. The two sentences are dissimilar, but are on a similar topic with respect to the condition. 3. The two sentences are roughly equivalent, but some important information differs or is missing with respect to the condition. 4. The two sentences are mostly equivalent, but some unimportant details differ with respect to the condition. 5. The two sentences are completely equivalent with respect to the condition. Return a JSON object with two fields: "rating": the similarity rating (between 1 to 5 as defined above), "justification": a single sentence explaining why you gave that similarity rating. Do not return anything else other than this JSON object. Do not use code blocks. ## Example 1 Sentence1: A close up of a giraffe laying on a ground near many large rocks. Sentence2: A giraffe reaches up his head on a ledge high up on a rock. Condition: animal's posture {"rating": 1, "justification": "In Sentence1 the giraffe is lying down, while in Sentence2 the giraffe is stretching its head upward."} ## Example 2 Sentence1: This bathroom stall has toilet tissue on the floor while the toilet is raised. Sentence2: A full trashcan is beside the commode in a public restroom toilet that needs to be cleaned. Condition: location of trash {"rating": 2, "justification": "Sentence2 does not clearly state that there is any trash outside the trashcan."} ## Example 3 Sentence1: A large red and blue boat sitting on top of a lake next to other boats. Sentence2: Part of a ship sits in the shallow end of the bay next to a city. Condition: body of water type {"rating": 3, "justification": "The two sentences mention lake and bay and are roughly equivalent, but Sentence2 does not clarify whether it is a bay within a lake."} ## Example 4 Sentence1: A monkey mug in front of a computer with a stuffed penguin beside it. Sentence2: A laptop computer sitting on top of a table next to two computer monitors. Condition: name of the device {"rating": 4, "justification": "Both sentences mention computers, but Sentence1 does not specify the type, while Sentence2 explicitly mentions a laptop."} ## Example 5 Sentence1: This bathroom stall has toilet tissue on the floor while the toilet is raised. Sentence2: A full trashcan is beside the commode in a public restroom toilet that needs to be cleaned. Condition: room function {"rating": 5, "justification": "Both sentences describe a room functioning as a restroom or toilet."}

Figure 3: Few-shot prompt for conditional sentence similarity evaluation