Enhancing NLI Models With an Adversarial LLM Approach

Anonymous ACL submission

Abstract

In this paper, we demonstrate that the performance of natural language inference (NLI) models can be enhanced using a novel adversarial approach, in which large language models (LLMs) are used to systematically address NLI models' weaknesses. We first employ the LLMs to adversarially generate challenging NLI examples, looking for instances that are misclassified by the NLI model, effectively creating a dataset. These examples are validated by an ensemble of LLMs to ensure their correctness and are subsequently used to retrain the NLI model, iteratively refining its 014 performance. In our evaluation, the proposed approach demonstrated substantial accuracy improvements on multiple datasets, including 1.43% on the SNLI dataset, 2.75% on the ANLI 017 dataset, and 4.29% on the MultiNLI dataset. 019 Our evaluation highlights the utility of LLMs in adversarial model improvement, providing a pathway toward robust and human-independent enhancements for NLI systems. Additionally, our LLM-based approach can also be used to automate the creation of NLI datasets.

1 Introduction

027

028

032

A fundamental task in natural language processing (NLP), natural language inference (NLI) is performed to determine the relationship between two sentences, ascertaining whether one sentence entails, contradicts, or is neutral to the other. While NLI models have achieved impressive performance, their robustness remains a challenge (Glockner et al., 2018; Carmona et al., 2018). Addressing these weaknesses is crucial for improving the reliability of NLI systems.

Inspired by the methodology used to create the adversarial NLI (ANLI) dataset (Nie et al., 2019), we propose a novel approach for automatically identifying and addressing the weaknesses of NLI models. Our approach leverages large language models (LLMs) to adversarially generate challenging NLI examples that aim to gather instances that are misclassified by the NLI model. These examples are validated by an ensemble of LLMs to ensure their correctness before being used to retrain the NLI model. This iterative process focuses on strengthening the model's ability to handle difficult cases, ultimately improving its performance. 041

042

043

044

045

047

049

051

057

059

060

061

062

063

064

065

066

067

068

069

071

072

073

To evaluate our approach, we trained a leading NLI model using our approach and another data augmentation method, on the same amount of data, using 10 different sets of hyper-parameters. We then evaluated this model on three popular NLI test-sets and observed consistent improvements.

The contributions of our work are as follows: (1) our proposed approach systematically addresses NLI model weaknesses, improving their robustness and accuracy, as demonstrated by performance improvements on the SNLI (Oana-Maria, 2018), ANLI, and MultiNLI (Williams et al., 2018a) datasets; (2) we introduce a fully automated dataset creation process that eliminates the traditional reliance on human annotators; and (3) our approach can scale to generate complete NLI datasets, enabling large-scale training of NLI models.

By combining automation, adversarial examples, and LLMs, our approach represents a significant step forward in enhancing NLI model performance and reliability. Moreover, by applying our method extensively to generate NLI examples, we can assemble a dataset that can be used to train NLI models.

2 Background and Related Work

Improving the robustness and performance of NLI074models remains a significant challenge in natural075language understanding (Glockner et al., 2018; Car-076mona et al., 2018). While traditional approaches077heavily relied on manually created datasets, such078as the Stanford NLI (SNLI) corpus (Oana-Maria,079

111 112

113 114

115

116

117

118

119

120 121

122

123

125

127

129

2018), this labor-intensive process highlighted the need for more efficient alternatives.

Recent advances in LLMs have enabled their use in the creation of NLI datasets, offering a more automated and scalable alternative to current practice. Our methodology leverages state-of-the-art LLMs such as Llama-3.1-70B (Touvron, 2023), Mistral-Large 2 (Jiang et al., 2023), and Mixtral-8x7B (Jiang et al., 2024) to generate and validate NLI examples. These models give our approach the ability to generate high-quality NLI examples and fine-tune NLI models like RoBERTa-Base (Liu et al., 2019), enhancing their robustness and performance.

Several recent studies have explored the use of LLMs for data generation. For example, counterfactual generation (Li et al., 2023) has been used to improve the robustness of the model in various downstream tasks, while paraphrasing bridging and NLI (Klemen and Robnik-Šikonja, 2021) has facilitated the expansion of existing datasets. TextAttack (Morris et al., 2020) is a framework for adversarial attacks and data augmentation, which has proven to be effective in enhancing models.

In the domain of NLI datasets, ANLI (Nie et al., 2019) used a human-and-model-in-the-loop approach to iteratively identify and address model weaknesses by manually creating challenging examples. Similarly, SNLI, with its 570K manually labeled sentence pairs, has become a standard benchmark for evaluating NLI models. Building on SNLI, the MultiGenre NLI (MultiNLI) dataset (Williams et al., 2018b) consists of 433K sentence pairs from various text genres, enhancing the training and evaluation of the models' generalization capabilities and robustness in varied contexts.

Methodology 3

In this section, we describe the four stages in our suggested approach for improving NLI models. The complete flow is presented in Figure 1.

Automated **Hypothesis** Generation Our methodology leverages LLMs to automate hypothesis generation, thus eliminating the need to rely on human annotators. To create diversity in the hypotheses, we begin by inputting premises and their corresponding labels into multiple LLMs. These models are tasked with generating a hypothesis that aligns with the given premise, such that the given label reflects the relation between

them. The pseudocode of the full algorithm is provided in Appendix A.2.

130

131

132

133

134

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

Adversarial Data Filtering Once the hypothesis is generated, it is sent, along with the premise, for classification by a pretrained NLI model, which we try to improve. If the model assigns the correct label for the input pair, both the hypothesis and the premise are discarded. If the model misclassifies the input pair, the pair and its correct label continue on to the validation stage. This is done because we want to gather examples that leading NLI models struggle with, in order to address their weaknesses.

Automated Validation The validity of a hypothesis misclassified by the NLI model is evaluated by an ensemble of three LLMs. These models act as independent judges, using majority voting to ensure robust, unbiased validation.

Iterative Refinement and Retraining If, in the previous stage, the LLMs agree on the validity of the misclassified example, the hypothesis and premise are then used for retraining. This iterative loop is aimed at refining the accuracy of the base NLI model. This process also enhances the training dataset by continually challenging the model and increasing its exposure to complex cases, thereby improving its overall robustness.



Figure 1: Illustration of our approach for improving an NLI model.

3.1 **Dataset Comparison and Semantic** Analysis

To gain insights into the relation between the data generated in out experiment and existing datasets, we examined the 10 most common non-stopwords in each dataset. We also assessed the similarity between the datasets using the TF-IDF and BERTScore F1 metrics (Zhang et al., 2019). The

TF-IDF metric, employing cosine similarity, mea-164 sures lexical overlap to reveal how much vocabu-165 lary and how many syntactic patterns are shared 166 between datasets. The BERTScore metric evaluates 167 semantic similarity using contextual embeddings from transformer language models. 169

Key Findings From the Dataset Analysis 3.1.1 170 In the SNLI train dataset, some of the most frequent words are 'man,' 'woman,' and 'people,' indi-172 cating themes of gender and social interactions. In 173 contrast, the ANLI test dataset focuses on media 174 and chronology with words like 'film' and 'first,' 175 while the MultiNLI test dataset uses more ab-176 stract language. The Generated dataset, containing misclassified examples, consist mainly of spec-178 ulative and gender-focused language. 179

177

180

181

182

186

We also analyzed the hypotheses' length and word counts in the datasets. The hypotheses in the Generated dataset were the longest, whereas SNLI train and SNLI test had similar lengths, suggesting a consistent style. The ANLI test and MultiNLI test datasets had longer hypotheses, highlighting their complexity. A comparison of the text length and word counts in the hypotheses of the examined datasets is provided in Figure 2.



Figure 2: Average text length and word count in the hypothesis column for the examined datasets.



Figure 3: TF-IDF cosine similarity among NLI datasets, including our generated dataset.



Figure 4: BERTScore F1 similarity among NLI datasets, including our generated dataset.

189

190

191

192

193

194

195

196

197

199

200

202

203

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

222

As for the similarity between datasets, Figure 3 presents the TF-IDF cosine similarity between every pair of the datasets' test sets. As can be seen, there is limited lexical overlap, with the greatest expected similarity between the SNLI train and SNLI test datasets and the least similarity between the ANLI test and MultiNLI test datasets. Figure 4 presents the BERTScore similarity; as can be seen, there are notable semantic alignments, particularly between the SNLI train and SNLI test datasets. These insights provide further validation of our approach, confirming that the data generated falls within the range of expected lexical and semantic similarities of existing NLI datasets.

Avoiding Forgetness 3.2

One of the challenges of fine-tuning existing pretrained models is 'forgetness.' Providing a pretrained model with many new training examples from a different distribution may cause the model to overfit the new distribution and degrade its performance on the original distribution on which it was pretrained. To prevent this adverse effect, we added several examples from the original SNLI training set to the new training set we created with the newly generated examples. We experimented with different ratios of generated to original training samples and selected the ratio that maximized accuracy. The different ratios and their corresponding accuracy value are presented in Figure 5. The incorporation of both original and generated train samples also enhances their generalizability. This diversity helps models recognize a broader spectrum of patterns and scenarios, reducing the risk of overfitting and enabling more reliable performance.



Figure 5: Model performance comparison across datasets.

4 Evaluation and Results

For this study, we used the RoBERTa-base-SNLI model from Hugging Face (HuggingFace, 2022) (125M parameters), a popular, open-source NLI model trained on a single dataset, to evaluate our approach. To evaluate our approach, we generated, filtered, and validated thousands of data samples, ending up with 2.5K high-quality samples of NLI data according to our approach. We used Llama-3.1-70B and Mistral-Large 2 (123B) for generation and Mistral-Large 2, Mixtral-8x7B, and Qwen-2.5-72B-Instruct (Qwen) for validation. Then, we finetuned the RoBERTa-base-SNLI on it, along with another 10K samples from the MNLI train set, to maintain our suggested ratio of 1:4. To fine-tune the NLI model, we used a single T4 GPU. We conducted experiments using 10 different sets of hyperparameters to confirm the robustness of our approach. This evaluation demonstrates notable improvements across three different and diverse test sets. In the first experiment, conducted on the SNLI test set, the model trained on our data achieved accuracy of 90.87%, surpassing the RoBERTa-base-SNLI's accuracy of 88.48%. This demonstrates that our approach effectively boosts performance on the dataset that the base model was originally trained on. In the second experiment, using the ANLI test set, our model again outperformed RoBERTa-base-SNLI, achieving an accuracy of 78.38% compared to 75.04%. This result shows that our approach improved the model's ability to handle challenging adversarial examples. Finally, on the MultiNLI dataset, the model trained on our data achieved an accuracy of 59.28%, which is significantly higher than RoBERTa-base-SNLI's accuracy of 54.67%. This emphasizes the enhanced generalization capabilities of our approach across diverse data distributions. For comparison, we fine-tuned the same model on the same amount of data taken from the MNLI train set. We also performed paraphrasing to transform the same amount of samples from MNLI. This approach achieved moderate improvements, with accuracies of 84.73% on SNLI, 72.39% on ANLI, and 50.01% on MultiNLI, but remained below the performance of our proposed method. These results are summarized in Table 1.

260

261

262

263

264

265

266

267

270

271

272

273

274

275

276

277

278

279

281

282

283

285

287

288

290

291

292

293

294

295

296

Dataset	RoBERTa base- SNLI	Additional Data	Para- phrasing	Our Approach
SNLI	88.48%	89.42%	84.73%	90.87% ± 0.58
Adversarial NLI	75.04%	77.07%	72.39%	78.38% ± 0.37
MultiNLI	54.67%	57.61%	50.01%	59.28% ± 0.32

Table 1: Comparison of accuracy on the examined datasets, for RoBERTa-base-SNLI, RoBERTa-base-SNLI fine-tuned with additional data from MNLI, RoBERTa base-SNLI fine-tuned with additional data generated using paraphrasing based on the SNLI train set, and RoBERTa-base-SNLI fine-tuned with additional data generated using our approach.

5 Discussion and Future Research

This study demonstrated the effectiveness of employing LLMs to automatically identify and address NLI models' weaknesses by generating and validating challenging datasets. By targeting model misclassifications, our approach systematically enhances NLI model robustness and accuracy, achieving significant performance improvements on diverse datasets - SNLI, ANLI, and MultiNLI. Our approach represents a major step forward in automating model refinement, reducing reliance on human annotators while preserving data quality and consistency.

Using an ensemble of LLMs for hypothesis validation reduces human biases and errors while enabling a scalable, iterative process for creating complete NLI datasets. This scalability supports both retraining existing models and building comprehensive datasets for future NLI models.

Future research should explore ways to further diversify the data generated by LLMs, incorporating varied linguistic structures and content domains. To explore our approach's potential to further address model weaknesses, its performance when employed on a larger scale and with multiple iterations should be explored. Additionally, applying these techniques to other NLP tasks could examine our approach's utility in other domains.

259

6 Limitations

297

315

320

321

323

324

325

326

327

330

331

336

339

341

342 343

345

346

Our approach's dependence on the initial quality of LLMs and the substantial computational resources required for training and deploying multiple models simultaneously could be prohibitive for some 301 applications. This research was conducted with 302 low-resource computation, which imposed certain constraints, limiting the scale and speed of processing. Additionally, the use of outsourced APIs for model generation introduced a bottleneck, as API response times delayed the generation of nec-307 essary data. These limitations prevented us from generating data at scale and testing our approach by generating hundreds of thousands of examples. We also have not yet examined our approach cyclically, 311 using the model trained with our data as a base 312 313 model for another iteration of data generation. We plan to address these limitations in future research. 314

References

- Vicente Iván Sánchez Carmona, Jeff Mitchell, and Sebastian Riedel. 2018. Behavior analysis of nli models: Uncovering the influence of three factors on robustness. *arXiv preprint arXiv:1805.04212*.
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. Breaking nli systems with sentences that require simple lexical inferences. *arXiv preprint arXiv:1805.02266*.
- HuggingFace. 2022. pepa/roberta-base-snli. Accessed: October 12, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. arXiv preprint arXiv:2401.04088.
- Matej Klemen and Marko Robnik-Šikonja. 2021. Extracting and filtering paraphrases by bridging natural language inference and paraphrasing. *arXiv preprint arXiv:2111.07119*.
- Yongqi Li, Mayi Xu, Xin Miao, Shen Zhou, and Tieyun Qian. 2023. Prompting large language models for counterfactual generation: An empirical study. arXiv preprint arXiv:2305.14791.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.

Roberta: A robustly optimized bert pretraining approach. *arXiv*.

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

400

- John X. Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. *arXiv preprint arXiv:2005.05909*.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2019. Adversarial NLI: A new benchmark for natural language understanding. *arXiv preprint arXiv:1910.14599*.
- Camburu Oana-Maria. 2018. e-snli: Natural language inference with natural language explanations. Advances in Neural Information Processing Systems 31.
- Qwen. Qwen2.5-72b-instruct. https://qwen2.org/ qwen2-5.
- Hugo Touvron. 2023. Llama: Open and efficient foundation language models. *arXiv*.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018a. A broad-coverage challenge corpus for sentence understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1112–1122. Association for Computational Linguistics.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018b. A broad-coverage challenge corpus for sentence understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1112–1122. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

A Appendix

A.1 Model Prompting Procedure for Validation

In Table 2, we present the final prompt used for LLM validation of the NLI dataset. The prompt asks the model if the provided label matches the premise-hypothesis relationship, with the system responding 'Accepted' or 'Not Accepted.' This process is repeated with multiple LLMs to filter challenging and problematic examples. The prompt was designed with detailed instructions, illustrative examples, and a structured response format to ensure consistency and accuracy in the validation process, contributing to the overall quality and robustness of the dataset.

Component	Content		
System	You are a language expert. Your job is to		
Prompt	filter rows of an NLI dataset, which contain		
	some data that may not be good enough.		
	Given a premise and a hypothesis, you		
	should determine whether the label reflects		
	the relationship between them or not.		
User	This is the premise: {premise}.		
Prompt	This is the hypothesis: {hypothesis}.		
	The relationship between them is {label}.		
	Do you accept this relationship? Respond only with 'Accepted' or 'Not Accepted.'		

Table 2: Prompting procedure used to validate the NLI dataset examples.

A.2 Model Prompting Procedure for Generation

The few-shot learning process used in the hypothesis generation stage of our approach is described in the following Algorithm 1. This process uses curated examples to guide the model in generating hypotheses that align with the desired premisehypothesis relationship. Leveraging these examples, the model produces contextually appropriate and accurate hypotheses, ensuring efficiency and consistency.

Algorithm 1 Hypothesis Generation Using Few-Shot Learning

- 1: Shuffle the SNLI train dataset D
- 2: Randomly select *n* observations from *D*: { $(p_1, h_1, l_1), (p_2, h_2, l_2), \dots, (p_n, h_n, l_n)$ }
- 3: for each (p_i, h_i, l_i) , where $i \in \{1, ..., n\}$ do 4: Format the example as:

This is a premise: p_i , this is the hypothesis: h_i , and the label between them is l_i .

- 5: end for
- 6: Provide these *n* formatted examples as fewshot inputs to the model
- 7: After providing the examples, prompt the model with the following instruction:

You are a language expert that helps create an NLI dataset. Given a premise sentence p and a desired label l, generate a one-sentence hypothesis hsuch that the label is relevant to the relation between the premise and the generated hypothesis. Keep the hypothesis short.

- 8: The model generates a one-sentence hypothesis *h* for the given premise *p* and label *l*
- 9: return Generated hypothesis h

We use a few-shot learning approach for hypoth-412 esis generation, providing the model with exam-413 ples from the SNLI train set. This approach lever-414 ages a high-quality examples to guide the model 415 in producing hypotheses that are both contextu-416 ally relevant and accurate. In Table 3, we present 417 the final prompt used, which includes detailed in-418 structions, carefully selected examples, and a struc-419 tured response format. This design ensures that 420 the generated hypotheses align with the desired 421 premise-hypothesis relationship while maintaining 422 consistency and reducing ambiguity in the output. 423

Component	Content			
Few-Shot	This is a premise: {premise}			
Example	This is the hypothesis: {hypothesis}.			
	The label between them is {label}.			
	(eight examples are shown to the model in this			
	format, randomly selected from the SNLI train			
	set.)			
System	You are a language expert that helps create an			
Prompt	NLI dataset. Given a premise and a desired label,			
	your job is to provide a one-sentence hypothe-			
	sis such that the label is relevant to the relation			
	between the given premise and your generated			
	hypothesis. Make sure to keep the hypothesis			
	short and no longer than a sentence.			

Table 3: Prompting procedure used to generate hypotheses for the NLI dataset.

This generation process helps the model create a hypothesis that is aligned with a given premise and label by first showing it several few-shot examples from the SNLI train dataset. After being shown these examples, the model is tasked with generating hypotheses following the same pattern, ensuring relevance and consistency.

A.3 Optimized Hyperparameters for RoBERTa-base-SNLI Model

In this section, we provide the optimized hyperparameters for the RoBERTa-base-SNLI model. After conducting 10 experiments, the best-performing parameters were identified as: a learning rate of 5.31×10^{-6} , a per-device training batch size of 16, a per-device evaluation batch size of 8, one training epoch, and a weight decay of 0.0093. These values were carefully selected to balance training efficiency and model generalization. The learning rate was adjusted to avoid overfitting, while batch sizes optimized the use of available computational resources. The weight decay was included to regularize the model and improve its performance on unseen data.

435

436

437

438

439

440

441

442

443

444

445

446

424

425

426

A.4 Examples of Generated Hypotheses

In Table 4, we provide some examples of the hypotheses generated. Each row contains the original premise, the generated hypothesis, and the original label, highlighting the model's generalization ability.

Premise	Hypothesis (Generated)	Lahel
Dooplo are alaan	A group of individuals are	Dabei
People are clean-	A group of individuals are	0
ing up a street.	picking up trash and debris	
	from the street.	
Swimmers leap	Athletes jump off the start-	0
off the starting	ing blocks into their desig-	
blocks into their	nated lanes at the beginning	
race lanes at an	of a swimming competition.	
indoor pool.		
Two women are	The women are engaged in	1
sitting at a table	a quiet activity.	
working with clay.		
Young man play-	A young man is throwing	0
ing darts in a cur-	darts in a private space.	
tained room.		
A man is bent over	A man is working outside	0
working outside	under construction flags.	
under red, green,		
and yellow flags.		

Table 4: Examples of generated hypotheses with their corresponding original labels.