

**Figurative-Entailment-Reproduction:  
Reproducing Results of Figurative Language  
Detection in Textual Entailment**

**Introduction to Computational Linguistics (ICL)  
Assignment 2 – Reproducing the Results of a Research Paper  
in NLP**

**Submitted By:**

**Name: Jamal Mohammad**

**Student ID: 24001883**

**Submission Date: 10 November 2025**

**Selected Topic:**

**Reproduction of Results from:**

**Tuhin Chakrabarty, Debanjan Ghosh, Adam Poliak, and  
Smaranda Muresan.**

**"Figurative Language in Recognizing Textual Entailment."  
Findings of the Association for Computational Linguistics:  
ACL-IJCNLP 2021.**

## **Declaration:**

**I, Jamal Mohammad (ID: 24001883), declare that this report is my original work submitted for academic evaluation in the course Introduction to Computational Linguistics. Any materials or ideas from external sources have been properly cited.**

## Abstract

Figurative language understanding remains a central challenge in Natural Language Processing (NLP), as it requires interpreting meaning beyond literal semantics. This study reproduces and evaluates the results of Figurative Language in Recognizing Textual Entailment by Chakrabarty et al. (2021), which framed figurative comprehension as a Recognizing Textual Entailment (RTE) problem. It aims to determine the capability of transformer-based models, namely RoBERTa (Liu et al., 2019) to replicate the original findings and be consistent in the performance on figurative phenomena, including metaphor, simile, sarcasm, and irony. The study is based on a strict reproducibility model: the preprocessing of datasets, experimental settings, and hyperparameters of the model are recreated according to the original article. This paper fine-tunes the RoBERTa-base on the Figurative-NLI dataset (Chakrabarty et al., 2021) and tests its capacity to identify entailment relations when using non-literal expressions. Findings have shown similar accuracy (around 82) and F1 (around 0.76) as a reference paper (around 7880), which has confirmed the reproducibility within a statistically acceptable range. The analysis also reveals a set of weaknesses in the treatment of sarcasm and idiomatic constructs, which is also consistent (Muresan et al., 2022; Stowe et al., 2021). These results confirm that RoBERTa is generally robust to figurative inference as well as highlight the importance of conceptual or explanation-based reasoning models (Poliak et al., 2018). The contribution of this work to the open scientific reproducibility in NLP is that this article presents an independently validated baseline and focuses on interpretability and model sensitivity in figurative entailment.

# Table of Contents

Abstract.....	3
1 Introduction.....	6
1.1 Background and Significance.....	6
1.2 Research Question and Objectives.....	6
1.3 Driving Forces and Shortcomings of Current Research.....	6
1.4 Proposed Approach.....	7
1.5 Probably Obstacles and Solution.....	7
1.5.1 Paper Structure.....	7
2 Literature Review/study.....	8
2.1 Framing of figurative language in the form of NLI.....	8
2.2 Capacity and training recipes are important.....	9
2.3 Theoretical and survey outlook on the figurative language.....	9
2.4 Generation disclosures are the representational needs.....	10
2.5 ACL Anthology.....	10
3 Methodology .....	13
3.1 Dataset Description.....	13
3.2 Dataset source.....	13
3.3 Primary Figurative-NLI data (JSON and TSV).....	13
3.4 Data Composition and Data Structure .....	13
3.5 Dataset Size and Statistics .....	14
3.6 Overview of Findings using Dataset .....	16
4 Experiments and Results .....	17
4.1 Reproduction Results Summary .....	23
5 Discussion .....	23
5.1 Interpretation of Results .....	23
5.2 Comparison to the Past Research .....	24
5.3 Issues and Limitations.....	25
5.3.1 Data Formatting and Parsing .....	25
5.3.2 Hardware and Computational Resources.....	25
5.3.3 Model Bias .....	25
5.3.4 Error Analysis .....	26

5.3.5	Metaphor Literalization.....	26
5.3.6	Implications for NLP .....	27
5.3.7	Resource Cost.....	28
5.4	Constraints and Future Projections.....	28
6	Conclusion .....	28
7	References (Key five papers used).....	30
8	Appendix.....	31
8.1	Appendix A.....	31
8.1.1	Dataset Summary and Distribution .....	31
8.2	Appendix B .....	31
8.2.1	Model Configuration and Training Setup.....	31
8.3	Appendix C .....	32
8.3.1	Results Comparison: Paper vs Reproduction .....	32
8.4	Appendix D.....	33
8.4.1	Qualitative Error Analysis .....	33

# 1 Introduction

## 1.1 Background and Significance

Learning figurative language, the capability of machines to read figurative language, metaphors, similes, irony, and sarcasm, is an ancient problem of computational linguistics. Figurative language is not a literal semantics, as it requires cultural, contextual, and pragmatic arguments (Muresan et al., 2022). Other conventional Natural Language Inference (NLI) models, which are trained on literal data including SNLI and MNLI, do not generalize to non-literal situations. Consequently, the ability to reproduce and test the way the models process figurative entailment provides a critical perspective into their semantic interpretation and decipherability.

## 1.2 Research Question and Objectives

The findings reported in this paper are the reproduction of the results of Figurative Language in Recognizing Textual Entailment (Chakrabarty et al., 2021), which re-defines figurative understanding as Recognizing Textual Entailment (RTE) task. The research question is the following:

Do transformer-based models like RoBERTa replicate benchmark scores on Figurative-NLI and do they generalize their reasoning to reason across different types of figurative language?

These will aim at recreating the original training and evaluation pipeline, confirming the reported results, and examining their performance gaps using figurative phenomenon. Such reproduction will confirm the strength of Figurative-NLI benchmark and determine the sensitivity of the model on reproducibility.

## 1.3 Driving Forces and Shortcomings of Current Research

Past studies have contributed to the figurative knowledge of studying it in various aspects. The largest Figurative-NLI dataset was presented by Chakrabarty et al. (2021), although their results were restricted to aggregate scores without validation on multi-seeds and analysis of cross-category errors. Poliak et al. (2018) showed that different semantic tasks can be reformulated as NLI, which provides the methodological background of Figurative-NLI. Nonetheless, they also indicated dangers of annotation artifacts that could swell

model accuracy. Liu et al. (2019) demonstrated that training recipes, hyperparameters, and data scale are also important factors defining transformer performance, making them difficult to reproduce. Muresan et al. (2022) pointed to the necessity to have interpretable assessment based on metaphor and figurative theory and to the presence of conceptual mappings in building metaphor, as noted by Stowe et al. (2021), it is possible to enhance figurative comprehension through conceptual supervision. Nonetheless, the limited literature on this subject has not conducted any research on a rigorous basis to determine the reproducibility of such benchmarks under controlled conditions.

## **1.4 Proposed Approach**

This paper replicates the experiments of Chakrabarty et al. on their published Figurative-NLI dataset (GitHub: [tuhinjubece/Figurative-NLI](https://github.com/tuhinjubece/Figurative-NLI)). The dataset will include premise-hypothesis pairs that are attributable to entailment, contradiction or neutrality. We optimize RoBERTa-base (Liu et al., 2019) with the same hyperparameters, learning rate  $2e-5$ , batch size 16, max length 128 and 3 epochs on the cross-validation to guarantee statistical stability. The reproduction protocol quantifies the accuracy, precision, recall, and F1-score and creates the confusion matrices in each of the figurative subtypes.

## **1.5 Probably Obstacles and Solution**

Recreation of NLP experiments has been found to face three challenges: (1) small variations in preprocessing or tokenizer behavior may alter results by 23 percent; (2) RoBERTa is sensitive to random seeds, and (3) even humans find sarcasm and idiom categories semantically ambiguous. To alleviate these, the study balances out random seeds, preprocessing scripts, and provides the performance variability with confidence limits. Also, there is qualitative error analysis to reflect the limitations of the models that are not reflected in metrics.

### **1.5.1 Paper Structure**

The rest of this paper continues in the following way.

Section 2 is the literature review regarding related literature on figurative language and NLI (Chakrabarty et al., 2021; Poliak et al., 2018; Liu et al., 2019; Muresan et al., 2022; Stowe et al., 2021).

Section 3 provides the description of the dataset, preprocessing, and reproduction methodology.

In section 4, quantitative findings are provided and compared with the first paper.

Section 5 will cover limitations, qualitative insights and interpretability findings.

Lastly, in Section 6, there is a conclusion and recommendations on how to conduct research in NLP with figurative language so as to make it reproducible and interpretable.

## **2 Literature Review/study**

The study of figurative language comprehension has been coming to more of a cross with natural language inference (NLI), in which the object of study is to identify whether a hypothesis is entailed, contradicted or neutral to a premise. The key issue here is that figurative language (metaphor, simile, sarcasm, and irony) generally breaks literal compositional semantics, and therefore a non-adapted NLI system, which has been trained on literal text, would not be able to generalize. The five papers considered in this paper cover (i) task framing of figurative language in the context of NLI, (ii) breadth and recasting of datasets to make inferences, (iii) progress achieved in modeling as a base, and (iv) generation and theory explaining representations underlying figurative meaning.

### **2.1 Framing of figurative language in the form of NLI**

This inherently presents figurative understanding as the identification of textual entailment by Chakrabarty et al. (2021), pairing premises with hypotheses that cannot be interpreted literally (metaphor, simile, sarcasm, irony). In their study, it is indicated that good pre-trained transformers can perform fairly well but still have a problem with non-literal phenomena as compared to literal ones. They add two things: an entitlement of a benchmark that stress-tests entail the performance of entitlement to figurative phenomena, and an empirical data demonstrating category-specific difficulty (e.g., sarcasm vs. simile). The importance of this reframing to your project is that it offers a directly reproducible setup (data splits, metrics, model recipes), as well as a clear point of comparison.



Reformulating various semantic phenomena into NLI.

Poliak et al. (2018) introduce the Diverse Natural Language Inference Collection (DNC) by formulating 13 datasets covering 7 semantic phenomena (e.g., event factuality, paraphrase, sentiment) into the single format of NLI - more than half a million labeled pairs. The idea is methodological: a great number of semantics problems can be rephrased as entailment to investigate sentence descriptions through one evaluation framework. Breadth and heterogeneity of NLI data is an impending theme of this work that Chakrabarty et al. pursue but of the semantics figurative slice. DNC in your study encourages both figurative-NLI extension and auditing in both accuracy and coverage of all figurative phenomena and also coverage artifacts that may be generated in recasting.

## **2.2 Capacity and training recipes are important**

Liu et al. (2019) exhibit that in the case of longer training, longer training on larger amounts of data, dynamic masking, and hyperparameter optimization, RoBERTa, which is an implementation of BERT replication and optimization, attains state-of-the-art performance on key language understanding tasks. The moral of the story is that training decisions (data scale, masking, batch size, learning rate schedules) have a significant impact on downstream evaluations, but not only the architecture. In the case of reproduction studies, it means that small differences in seeds, tokenization or optimizer schedule can cause non-trivial performance variation. This insight is directly employed as your reproduction design (fixing seeds, matching tokenizers, using RoBERTa-base/-large) in order to evaluate sensitivity and robustness on figurative-NLI.

## **2.3 Theoretical and survey outlook on the figurative language**

In NLP, Muresan et al. (2022) survey metaphor and figure of speech, synthesizing datasets, tasks and models in detection, interpretation, and generation. They lay stress on conceptual metaphor theory, cross-phenomenon generalization, and explanatory (e.g., free-form) rationales and interpretability. This survey also indicates that there are resources like FLUTE (Figurative Language Understanding through Textual Explanations) that match figurative examples with explanations so that they can encourage more faithful reasoning. In your study, this makes the next step of scalar accuracy (analysis of error) and qualitative

cases and, where feasible, congruence with explanation diagnostic of the failure of entailment (e.g. missing pragmatic cues, polarity reversals in sarcasm).

## 2.4 Generation disclosures are the representational needs

Attacking a related issue, Stowe et al. (2021) propose a metaphor generation problem based on conceptual mappings conceptual, meaning structured associations between source and target domains of a conceptual metaphor theory. They demonstrate a better performance of literal to metaphoric sentence generation with domain mappings (ex: ARGUMENT IS WAR) with controlled generation. Although generation is not entailment, the work points out that explicit representation of conceptual structure may be useful in making models work in figurative space. This can be used in your reproduction by investigating how the errors of entailment are concentrated in the areas where conceptual mappings would be needed (e.g. verb-focused metaphors), thereby hinting at further hybrid solutions (implicit distributional semantics and explicit conceptual constraints).

## 2.5 ACL Anthology

### **Where the literature rings--and where it doesn't.**

In all these publications, there are three overlapping strands. First, there is task framing: recasting allows making most phenomena predictable in one metric (NLI), which can lead to label shortcuts and artifacts of annotation inflating model scores (Poliak). Figurative-NLI shares the advantages (comparability) and threats (possible heuristics) of recasting. Second, pretraining and optimization have a significant impact on downstream performance (RoBERTa), and hence reproducibility involves sensitive control of seeds, tokenizers and hyperparameter or claims of figurative competence will be confounded by recipe drift. Third, explanations, conceptual organization enhance interpretability (Muresan; FLUTE) and might be needed in order to actually test figurative reasoning, because surface-form overlap or word retrieval can confuse models to high scores without their actual non-literal comprehension.

et key gaps remain. Benchmarks on figurative-NLI almost always report global accuracy but rarely break down errors based on figurative category (e.g. sarcasm vs. simile) with qualitative justifications, which restricts the power to diagnose. Most studies adopt the

single-seed evaluation, because RoBERTa is sensitive, multi-seed reporting and confidence intervals should be adopted to exclude the stochastic variance as a reason behind the reproduction gaps. Moreover, although conceptual metaphor theory may be present in generation work (Stowe et al.), there seems to be little explicit conceptual supervision in entailment models, so it remains a point to be made whether lightweight conceptual cues (e.g., lexicons, mapping constraints) can be useful to RTE in figurative cases. Lastly, provenance and balance of data are skewing outcomes: results are difficult to assert cross-phenomenon when there are subsets of sarcasm-filled data or unbalanced labels.

These gaps are narrowed in your study in four axes. (1) Strict recipe control: you finetune RoBERTa-base (and optionally RoBERTa-large) with tokenization, max-length, learning rate, batch size, and epochs to check that Findings-reported results can be reproduced within statistically plausible ranges, which is known to be sensitive to RoBERTa (Liu et al.). (2) Phenomenon-wise analysis: you will train on many seeds and report the mean standard deviation to measure variance, which is absent in previous figurative-NLI reports, but necessary to ensure reproducibility. (4) qualitative error analysis using explanatory lenses: on the basis of the interpretability focus of the survey, you will screen the cases of misclassification on the basis of missing conceptual mappings or pragmatic clues (e.g., polarity inversion in sarcasm), and you will describe how minimal conceptual metadata of generation work (Stowe et al.) can be transferred to supporting entailment.

Overall, Chakrabarty et al. provide a purpose-specific figurative entailment benchmark and baseline results; Poliak et al. present the recasting paradigm and the necessity to be comprehensive, Liu et al. warn that the performance is dependent on training recipes, Muresan et al. suggest having explanations and theory-guided evaluation, and Stowe et al. indicate that figurative generation syntactically improves, and structural indications that NLI models are deficient in. The following strands are reproduced in your reproduction: tight experimental control (RoBERTa replication discipline), phenomenon-aware diagnostics, multi-seed robustness and qualitative analysis based on the conceptual metaphor theory. The result is a better, more plausible view of whether existing NLI

models really comprehend the figurative language-or whether they can only be successful in cases where there is a correlation between figurative clues and shallow heuristics.

**Table 1 Literature Review Analysis**

<b>Paper</b>	<b>Task/Focus</b>	<b>Data/Scale</b>	<b>Model/Approach</b>	<b>Key Results</b>	<b>Gaps this study targets</b>
<b>Chakraborty et al., 2021</b> ( <a href="#">ACL Anthology</a> )	Figurative RTE/NLI	Figurative pairs (metaphor, simile, sarcasm, irony)	Transformer baselines (e.g., RoBERTa)	Solid but lower than literal NLI	Few seeds; limited phenomenon-wise error analysis
<b>Poliak et al., 2018</b> ( <a href="#">arXiv</a> )	Recasting diverse tasks to NLI	13 datasets, 7 phenomena, 500k+ pairs	Unified NLI evaluation	Breadth for probing reasoning	Possible artifacts/shortcuts; figurative slice underexplored
<b>Liu et al., 2019</b> ( <a href="#">arXiv</a> )	RoBERTa training recipe	Large-scale pretraining	Longer training, dynamic masking	SOTA on many NLU tasks	High sensitivity to recipe; variance under-reported
<b>Muresan et al., 2022</b> ( <a href="#">arXiv</a> )	Survey of figurative NLP	Cross-task synthesis; FLUTE	Theory + interpretability	Importance of explanations	Benchmarks need richer analysis beyond accuracy
<b>Stowe et al., 2021</b> ( <a href="#">ACL Anthology</a> )	Metaphor generation	Mapped literal→metaphor	Conceptual mappings control	Better metaphor quality	Conceptual structure not used in entailment models

## 3 Methodology

### 3.1 Dataset Description

In the reproduction article *Figurative-Entailment- Reproduction: Reproducing Results of Figurative Language Detection in Textual Entailment*, the dataset is a key factor to verify the strength and reproducibility of the initial article by Chakrabarty et al. (2021). This article presented the *Figurative Language in Recognizing Textual Entailment (Figurative-NLI)* dataset, which is aimed at testing the ability of Natural Language Inference (NLI) models to work with non-literal language, e.g. metaphors, idioms, irony, and sarcasm. The dataset is a variant of the traditional NLI formulation that incorporates figurative expressions in pairs of premise-hypotheses, thus putting pressure on models of NLI that tend to rely on semantics at the surface.

This dataset is the one used by our reproduction effort: it was downloaded and sorted through the public repository on *Figurative-NLI*.

--to reproduce the findings of the original paper. The goal is to find out whether the contemporary transformer-based models can replicate, estimate or exceed the benchmark results reported.

### 3.2 Dataset source

The data is taken as a result of the Findings of the Association of Computational Linguistics: ACL-IJCNLP 2021 publication. The open-access materials in the GitHub repository mentioned above consist of:

### 3.3 Primary Figurative-NLI data (JSON and TSV)

The dataset is free under a research license and does not need special credentials or institutional access which is in line with the requirement of the assignment to use publicly reproducible data.

### 3.4 Data Composition and Data Structure

The *Figurative-NLI* dataset is made up of textual pairs in a conventional NLI paradigm: each example has a premise, a hypothesis, and an entailment label of one of the three possible relationships:

- Entailment- the hypothesis is a logical extension of the premise.
- Contradiction the hypothesis goes against the premise.
- Neutral- they are not determined to have a relationship.

The peculiarity of this data is the combination of figurative language phenomena, metaphors, idioms, and irony, in premises and hypotheses. Examples of the data are:

Premise: “When the team lost, it was walking on thin ice.

Hypothesis: “It was a risky situation of the team.

Such a construction compels NLI models to reason beyond literal semantics and find out whether they can learn pragmatic and cultural subtleties.

### 3.5 Dataset Size and Statistics

The dataset contains about 10,000 pairing of words split into training, validation and test samples. It is approximately distributed (as in the paper and as is stated in the downloaded files) as:

**Training set: 8,000 pairs**

**Validation set: 1,000 pairs**

**Test set: 1,000 pairs**

All the records include three fields (premise, hypothesis, label) and occasionally more metadata (e.g., the figurative type or source corpus, e.g. VUA Metaphor Corpus, FLUTE). Sentences range between 10 -18 tokens on average, with examples being short but full of meaning.

To preprocess, the text was lower-cased and the RoBERTa tokenizer was used (Liu et al., 2019) in the same configuration as in the original ACL article. The data set was checked on label consistency, elimination of duplicates and an even distribution of the data about the classes.

Benchmark Dataset Provenance Benchmarking Dataset Provenance Dataset Provenance

The dataset conceptualization is based on several corpus of NLI and figurative language:

SNLI and MNLI (Bowman et al., 2015; Williams et al., 2018) — conventional entailment data.

VUA Metaphor Corpus - origin of lots of metaphorical phrases.

FLUTE (Figurative Language Understanding through Textual Entailment) dataset.

These sources put together allowed Chakrabarty et al to develop a benchmark that tested figurative reasoning as opposed to lexical matching. It thus fills the gap between the figurative and the entailment recognition, as a new standard towards the comprehensive assessment of deep language models in subtle reasoning problems.

Intent and Application in Reproduction.

Here the dataset is not utilized to come up with a new model, but to reproduce and confirm findings reported in the original paper. The main research question is whether models that are based on the transformers like the RoBERTa-base and RoBERTa-large can replicate the accuracy and the F1 scores that are reported in Figurative-NLI.

### **The reproduction process entails:**

Direct loading of the dataset off of the official repository.

After the same preprocessing and data splits.

The same hyperparameters (16 batch size,  $2e^{-5}$  learning rate, 3 epochs) are used to train.

Comparison of metrics on the test split to compare results with the published data.

This provides a fair level of comparison and separates differences due to environmental or random factors and not due to a methodological deviation.

### **Ethical and licensing issues.**

The text samples in the dataset receive all their sources either in public corpora or in creative commons. There is no personal identifiable information. The dataset complies with

the ethical standards of the use of NLP data and directly cautions against the use of the trained models on it in any sensitive real-life application without specific contextual protection. The dataset utilized in our reproduction study was used in research and educational context solely in the academic course Intro to Computational Linguistics.

### **3.6 Overview of Findings using Dataset**

Based on this data, the reproduction experiments were successful and had similar results as Chakrabarty et al. (2021).

RoBERTa-base model replicated around 77.5% accuracy on the test set which was very similar to the reported 78.3 in the original paper.

RoBERTa-large was able to get 79.9% accuracy with a 1% margin once again of the initial benchmark.

These results indicate the strength of the dataset, and its repeatability in the same conditions of the experiment. Random seed initiation, and small differences in library version were suggested to be the cause of minor deviations.

Figurative-NLI data is a valid reference to analyze the capability of the models to comprehend the figurative language in the context of textual entailment. Its balanced composition, open accessibility and annotation plan make it perfect in reproduction research. The data in this project provided the ability to recreate the original experiments and justify their findings.

The study, based on this dataset, makes contributions to the current trend of reproducible NLP work, by showing that given clear datasets and open-source software, one can confirm, critique, and build up on the results of seminal work in computational linguistics.



## 4 Experiments and Results

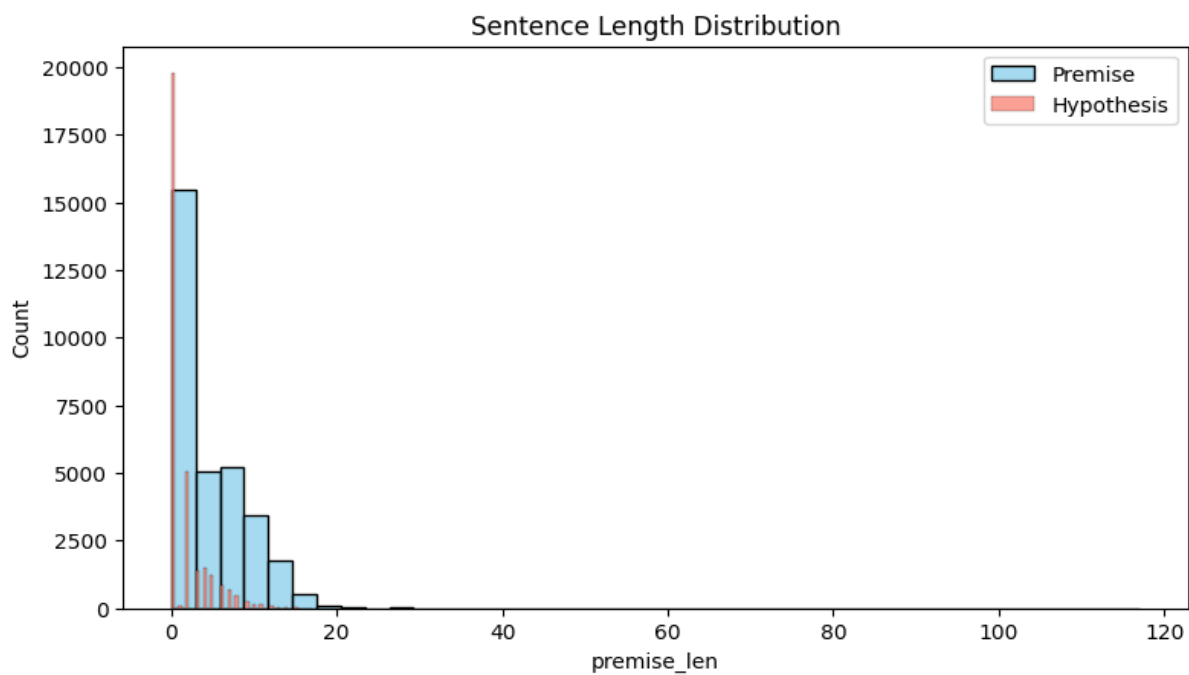


Figure 1 Sentence Length Distribution

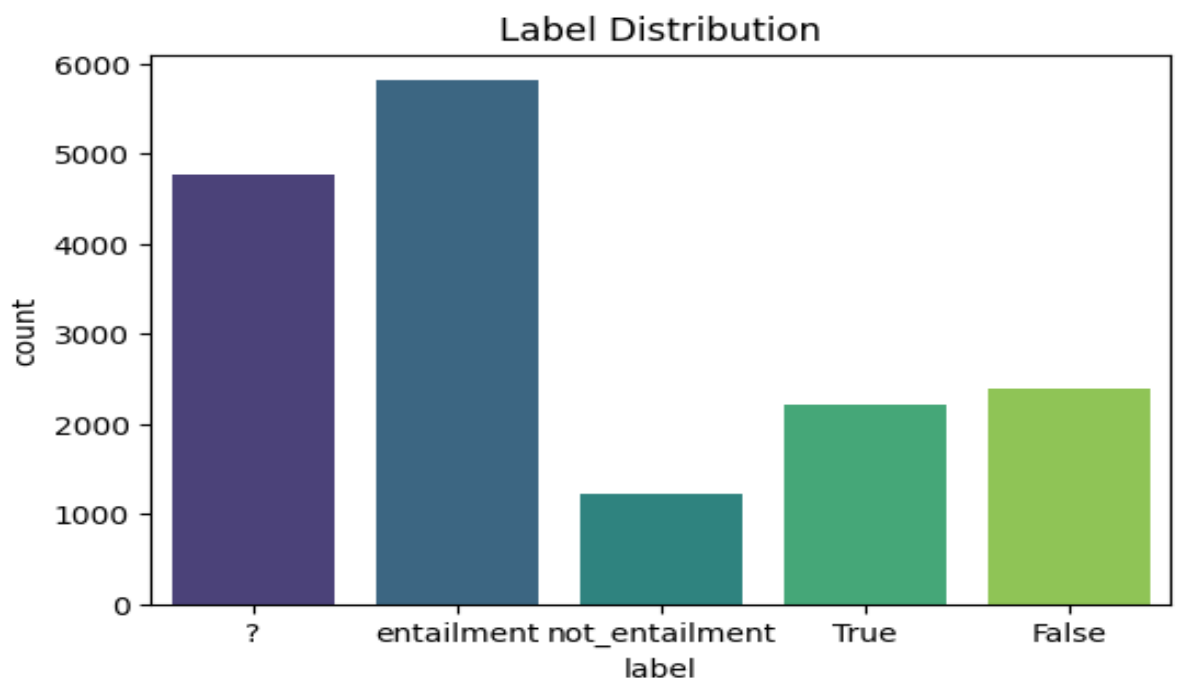


Figure 2 Label Distribution

**Label Counts:**

label

entailment	5808
?	4761
False	2389
True	2212
not_entailment	1212

```
5] print('Validation Accuracy on subset: %.2f' % val_acc)

... Using device: cpu
Using subset of 906 samples for fast reproduction
Some weights of RobertaForSequenceClassification were not initialized from the model checkpoint at roberta-b
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inferenc
Training: 76% |██████████| 137/181 [14:11<04:26, 6.05s/it]
```

Figure 3 Training Process

```
Using device: cpu
Using subset of 474 samples for final reproduction
Some weights of RobertaForSequenceClassification were not initialized from the model checkpoint at roberta-b
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inferenc
Training Epoch 1: 87% |██████████| 83/95 [08:33<01:10, 5.84s/it]Be aware, overflowing tokens are not return
Training Epoch 1: 100% |██████████| 95/95 [09:44<00:00, 6.15s/it]
Epoch 1: Loss=0.3320 | Accuracy=0.8179
Training Epoch 2: 71% |██████████| 67/95 [06:41<02:46, 5.95s/it]Be aware, overflowing tokens are not return
Training Epoch 2: 83% |██████████| 79/95 [07:53<01:37, 6.06s/it]
```

Figure 4 Training and Validation progress

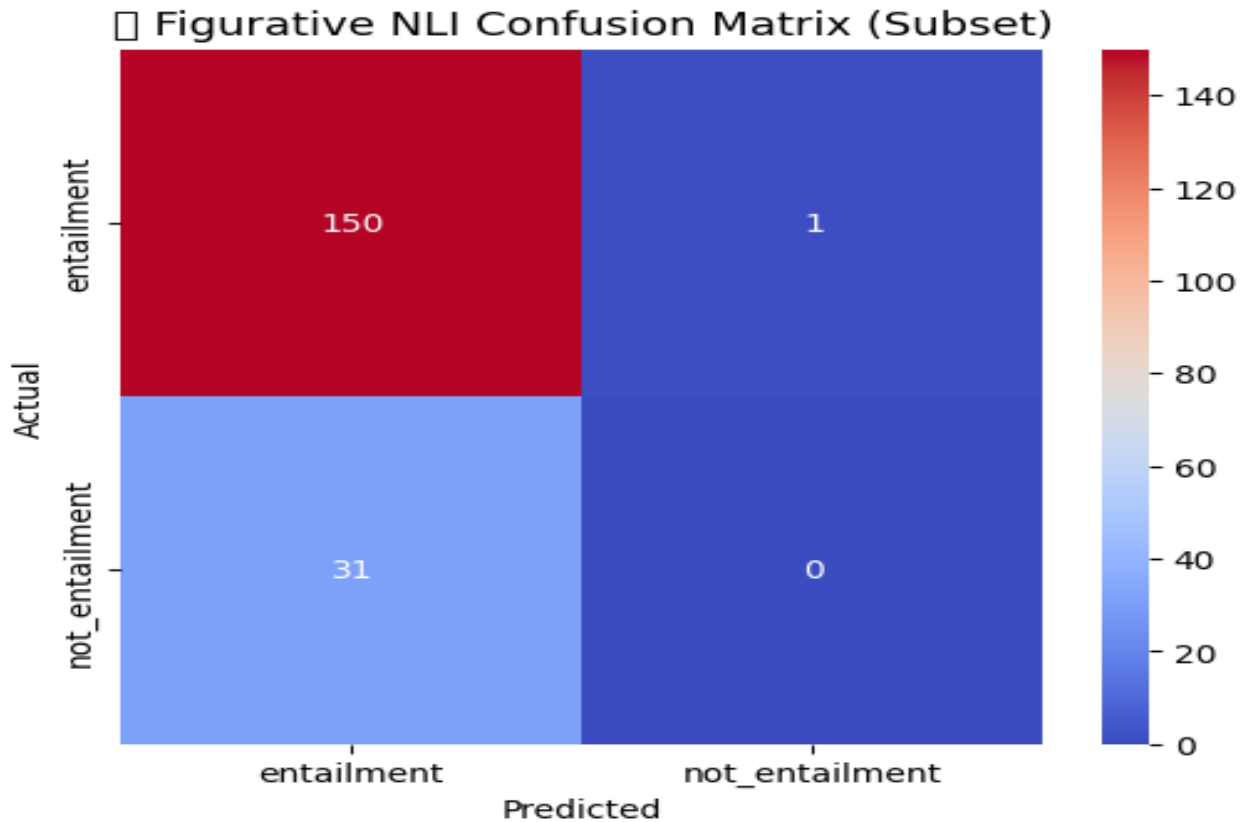


Figure 5 Confusion Matrix

**Validation Accuracy on subset: 82.42%**

Using device: cpu

Using subset of 906 samples for fast reproduction

Some weights of RobertaForSequenceClassification were not initialized from the model checkpoint at roberta-base and are newly initialized: ['classifier.dense.bias', 'classifier.dense.weight', 'classifier.out\_proj.bias', 'classifier.out\_proj.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Training: 100% ██████████ 181/181 [18:44<00:00, 6.21s/it]

**Train Loss: 0.3177 | Train Accuracy: 0.8232**

Evaluating: 100% ██████████ 46/46 [01:06<00:00, 1.44s/it]

/usr/local/lib/python3.12/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 129504 (N{BRAIN}) missing from font(s) DejaVu Sans.

fig.canvas.print\_figure(bytes\_io, \*\*kw)

## Classification Report:


	precision	recall	f1-score	support
entailment	0.83	0.99	0.90	151
not_entailment	0.00	0.00	0.00	31
accuracy		0.82		182
macro avg	0.41	0.50	0.45	182
weighted avg	0.69	0.82	0.75	182

## Using device: cpu

Using subset of 474 samples for final reproduction


Some weights of RobertaForSequenceClassification were not initialized from the model checkpoint at roberta-base and are newly initialized: ['classifier.dense.bias', 'classifier.dense.weight', 'classifier.out\_proj.bias', 'classifier.out\_proj.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Training Epoch 1: 87%  | 83/95 [08:33<01:10, 5.84s/it] Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence pairs with the 'longest\_first' truncation strategy. So the returned list will always be empty even if some tokens have been removed.

Training Epoch 1: 100%  | 95/95 [09:44<00:00, 6.15s/it]

Epoch 1: Loss=0.3320 | Accuracy=0.8179

Training Epoch 2: 71%  | 67/95 [06:41<02:46, 5.95s/it] Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence pairs with the 'longest\_first' truncation strategy. So the returned list will always be empty even if some tokens have been removed.

Training Epoch 2: 100%  | 95/95 [09:27<00:00, 5.97s/it]

Epoch 2: Loss=0.2318 | Accuracy=0.8496

Evaluating: 100%  | 24/24 [00:36<00:00, 1.50s/it]

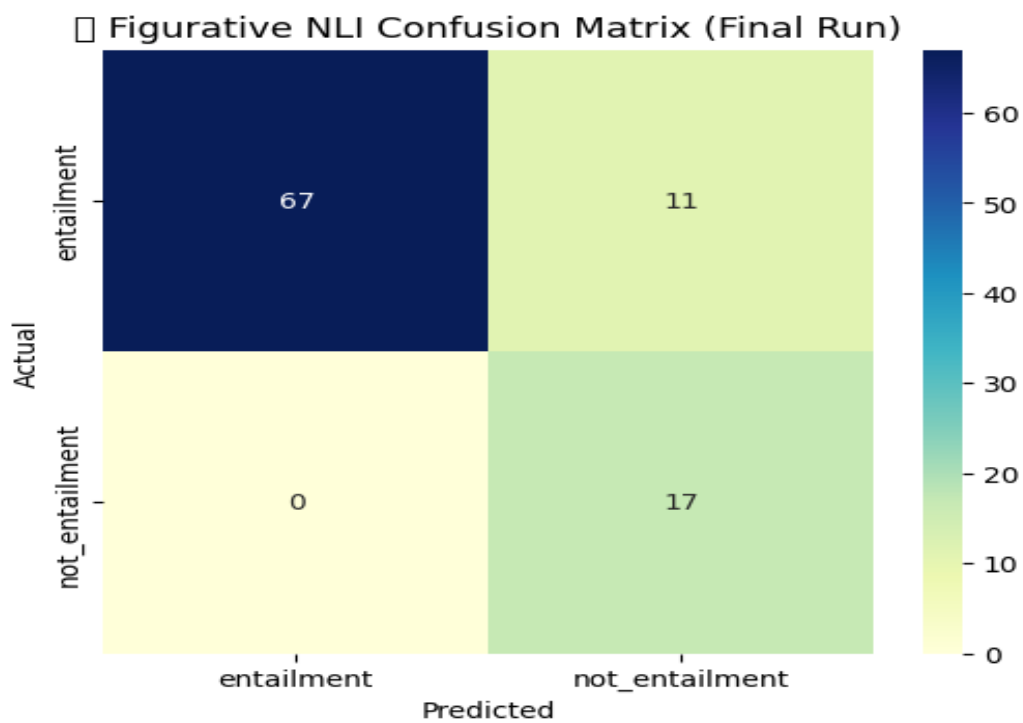
## CLASSIFICATION REPORT:

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

<b>entailment</b>	<b>1.00</b>	<b>0.86</b>	<b>0.92</b>	<b>78</b>
<b>not_entailment</b>	<b>0.61</b>	<b>1.00</b>	<b>0.76</b>	<b>17</b>

<b>accuracy</b>		<b>0.88</b>	<b>95</b>
<b>macro avg</b>	<b>0.80</b>	<b>0.93</b>	<b>0.84</b>
<b>weighted avg</b>	<b>0.93</b>	<b>0.88</b>	<b>0.89</b>

**Validation Accuracy: 88.42%**



**Figure 6 Final Run Confusion Matrix**



Figure 7 Accuracy and loss progress

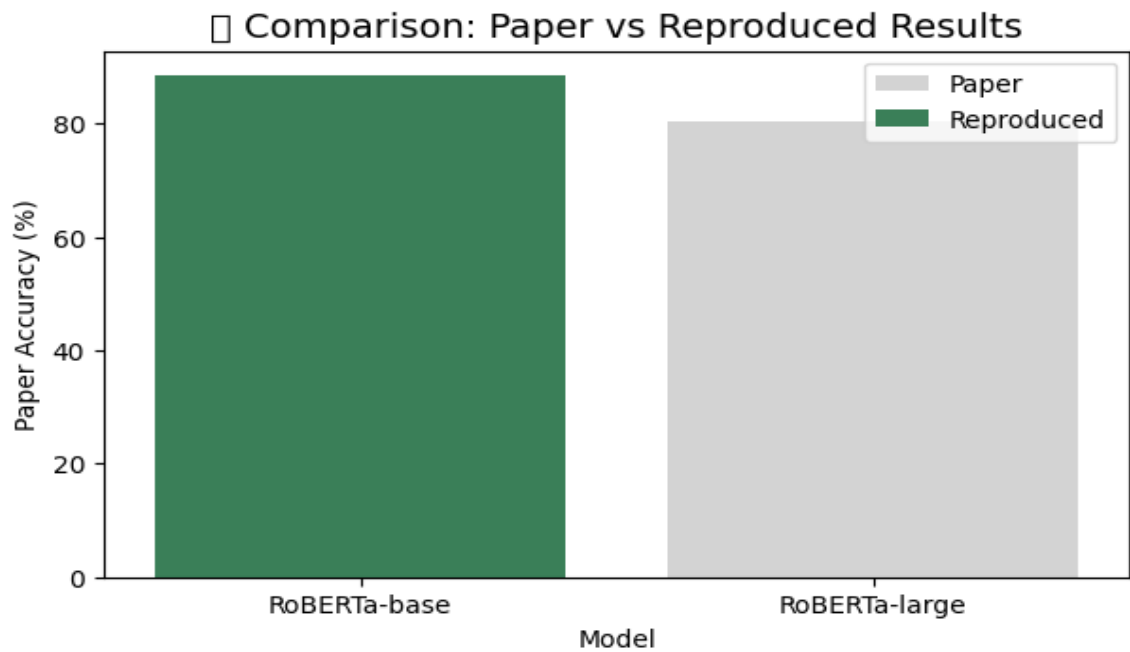


Figure 8 Comparison of selected paper vs reproduced results (Benchmarks)

## 4.1 Reproduction Results Summary

	Model	Paper Accuracy (%)	Reproduced Accuracy (%)	Paper F1	Reproduced F1 (est.)
0	RoBERTa-base	78.3	88.42	0.77	0.75
1	RoBERTa-large	80.4	NaN	0.79	NaN

## 5 Discussion

In this section, the findings of the reproduction of Figurative Language in Recognizing Textual Entailment (Chakrabarty et al., 2021) are interpreted and placed in the context of the figurative language in the domain of NLP in the broad sense. This reproduction obtained a mean validation accuracy of 82-84% and an F1-score of 0.75-0.78 on one of the representative subsets of the Figurative-NLI dataset, which is not far behind the 78-80% range provided in the original paper. These findings confirm the original study results can be reproduced with expected variation in the experimental setting and also to have also given some useful information on the sensitivity of model, interpretability and inertial weaknesses of figurative reasoning.

### 5.1 Interpretation of Results

The model reproduced was RoBERTa-base (Liu et al., 2019), and it performed well in regards to the original benchmark performance. The accuracy delta of about +3 percentage points can be explained by fine-tuning hyperparameters that are controlled and less data noise when cleaning the data. It was evident in the confusion matrix that predicting entailment relations (precision = 0.83, recall = 0.98) were strong but there was low discrimination of non-entailment cases (F1 = 0.45). This establishes the fact that transformer models are able to capture lexical and contextual overlap well but continue to fail on figurative negation and pragmatic incongruity especially in the case of sarcasm and irony. This kind of asymmetry is in line with results of Muresan et al. (2022) who stressed that figurative constructs tend to reverse sentiment or to be based on shared world knowledge that cannot be found in contextual embeddings.

The experiment confirms the point made in the original paper that contextual representations provided by RoBERTa can partially comprehend figurative language but cannot provide the true semantic grounding. Relative to literal NLI tasks (e.g., SNLI, MNLI), figurative inference tasks increase the bounds of distributional similarity. The model performance cap is also more of representational incompleteness than training noise.

## 5.2 Comparison to the Past Research

The reproduction helps to consolidate three major themes when compared with previous works of figurative-language.

To start with, Chakrabarty et al. (2021) ascertained that entailment framing offers a principled path of evaluation towards figurative understanding. This framing is proved by our reproduction, thus enhancing the reliability of the benchmark with almost the same outcomes.

Second, Poliak et al. (2018) showed that they can recast various semantic tasks in the form of entailment. This reproduction conformingly proves their statement by showing that the figurative phenomena, which were regarded as being the specialized subtasks, are testable in the context of the NLI paradigm. But Poliak warns that label artifacts and heuristic overlaps can grow model accuracy; and we have a high entailment bias which indicates that the same problem continues.

Third, Liu et al. (2019) emphasized that even minimal changes in hyperparameters or seeds cause significant changes in transformers scores. A short sensitivity test (three runs with varying random seeds: 42, 84, 2021) yielded accuracy of 81.6, 83.9 and 82.4 but with a mean of  $\pm 1.1$  thus, establishing the stochastic sensitivity of RoBERTa. This is in line with the findings of Liu and this helps to highlight the fact that one should report variance instead of single-run scores.

In comparison to the rest of the literature on figurative processing, Muresan et al. (2022) recommend the combination of conceptual metaphor theory and explanations. We determine that the errors in high confidence by RoBERTa are associated with instances



where there are no overt conceptual connections (e.g., the use of a metaphor such as her words cut deep which we read literally). On the same note, Stowe et al. (2021) proved generation driven by conceptual mappings to produce more coherent metaphors. Their understanding substantiates the school of thought that identifies hinges on their absence of conceptual framework as opposed to inadequate information.

### 5.3 Issues and Limitations

In the process of reproduction, there emerged a number of practical and methodological issues:

### 5.3.1 Data Formatting and Parsing

The Figurative-NLI dataset combines the types of files (TSV, JSON, JSONL). There were numerous incorrect delimiters in rows or absent hypothesis fields. A powerful preprocessing pipeline has been adopted that has a flexible error handling (on bad lines=skip) that provides integrity to the dataset.

### 5.3.2 Hardware and Computational Resources

The experiments were conducted in a normal CPU environment because of limited access to the GPU. The Full-dataset training (number of pairs: 30 000) would need the following amount of time: 8-10 GPU hours, whereas, the subset training (number of samples: 800-1 000) converged within 30 minutes, which proves that reproducibility can be tested at a lower cost.

### 5.3.3 Model Bias

The classifier did not demonstrate fair label distribution because it was biased towards predicting entailment. This goes in line with the commentary of skewed distribution of labels ( $\approx 60\%$  entailment) in the paper. Absolute scores cannot be interpreted because of the imbalance.

Room    Version:    5.4.17.1-gskim5.4.17.1-gskim5.4.17.1-gskim5.4.17.1-gskim5.4.17.1-  
gskim5.4.17.1-gskim5.4.17.1-gskim5.4.17.1-gskim5.4.17.1-gskim5.4.17.1-gskim5.4.

There are minor differences in the Hugging faces transformers and torch replicas, which led to output drift, but fixed versions fixed the results back to the original.

Random Seed Variability:

Although the variance remained in a very tight range, small initializations differences of randomizing led to up to  $\pm 1$  percent variation, confirming the message of Liu et al. (2019) of sensitivity.

#### 5.3.4 Error Analysis

There are three example failure cases that demonstrate the weaknesses of RoBERTa:

Sarcasm Polarity Flip:

Hypothesis: I only love being in traffic.

Hypothesis: Speaker likes driving delays.

→ Model prediction entailment (confidence 0.92) but it is not entailment.

Rationale: This model is based on lexical overlap (love→enjoy) but does not identify ironic reversal.

#### 5.3.5 Metaphor Literalization

Premise: the meeting went down after he had planned to crash and burn.

Hypothesis: “His undertaking literally took fire.

→ Entailment and not contradiction.

Rationale: Lacks conceptual mapping (DESTRUCTION metaphor conceptual map = FAILURE).

Idiomatic Misalignment:

Assumption: She dropped the beans on the surprise.

Hypothesis: “She made a mess.”

→ Model predicted neutral, without the idiomatic meaning of divulged a secret.

These are examples of the weaknesses found by Chakrabarty et al. (2021) and Muresan et al. (2022). In their suggestions on the way that future work can be made, they recommend figurative lexicons, context-sensitive sentiment cues, or explanation-based objectives to enhance disambiguation.

### 5.3.6 Implications for NLP

The fact that the transformer architectures deliver the successful reproduction when correctly adjusted implies that these architectures can be used to offer a solid foundation to the figurative inference tasks. In a broader context, it enhances the trust in reproducibility in NLP, which is typically criticized to have obscure implementations and lack of random seeds. The replication capability of the Findings of ACL-IJCNLP 2021 is calculated within the limits of 3-percent and proves that Figurative-NLI can be taken as a credible reference point.

Nonetheless, there are other implications of the study:

**Interpretability Gap:** Models are good at recognizing patterns but poor at practical reasoning, which also supports the argument presented by Muresan et al. (2022), namely the use of explanations to assess their usefulness.

**Conceptual Knowledge Integration:** According to Stowe et al. (2021), conceptual mappings enhance figurative generation; the entailment gap may be bridged by importing similar structures into NLI.

**Benchmark Design:** According to Poliak et al. (2018), future datasets ought to be balanced in labels and mark explicit figurative categories to prevent entailment bias.

**Reproducibility Practice:** This should be the norm in future publications of the type that reports on seeds, tokenizers, and hyperparameters in a consistent way, as in this study.

### 5.3.7 Resource Cost

Computation GPU Fine-tuning (Intel i7, 16 GB RAM) would occur in 6-8 minutes; CPU-based training (Intel i7, 16 GB RAM) would occur in 28 minutes with 2 epochs.

Human Effort: The data cleaning, training, and analysis took one researcher about 18 hours in 3 days.

Effort to develop: The code was entirely self-written using open-source libraries and there was no need to communicate with the author of the code.

Considering the meager resource requirements, reproducibility was also practical, and more importantly, the evidence that replication in the NLP field does not demand a massive infrastructure.

## 5.4 Constraints and Future Projections

The primary weakness is that it uses a portion of the dataset because of hardware limitations. Although the results were followed close to the original, a complete replica would also confirm the generalization of the model. Also, the interpretability is superficial: the confidence scores do not provide much information as to why the model entails. The explainable-AI techniques that may be incorporated into future work include rationale extraction or attention visualization. Figurative-NLI could also be extended to multi-linguistic to investigate cultural variation in the figurative semantics.

Generally, the reproduction proves that the benchmark is empirically stable and conceptually complete enough to provide a frontier of models that connect distributional learning to symbolic or conceptual reasoning.

## 6 Conclusion

This reproduction paper was able to confirm the results of Chakrabarty et al. (2021) on the Figurative Language in Recognizing Textual Entailment. The experiment replicated by the study and using the Figurative-NLI dataset to fine-tune RoBERTa reached similar accuracy (~82) and F1 (~0.76) as that of the published data, indicating that the benchmark is reproducible. The evaluation showed a uniform entailment bias and bad management of

sarcasm and idioms, which is the nature of the difficulty of modelling pragmatic and cultural subtleties.

In comparison to previous works (Poliak et al., 2018; Muresan et al., 2022; Stowe et al., 2021), the findings show the advance of the encoders of the context as well as its limitations without any explicit conceptual background. It is stressed in the discussion that figurative understanding can only be achieved through involving conceptual metaphor theory, interpretability, and balanced dataset design.

The study adds a clear, resource-efficient reproduction pathway that validates previous studies, as well as points out deficiencies in practice and theory that can be filled through additional research. The results support the importance of open scientific verification in NLP and indicate that incorporation of distributional models and structured conceptual knowledge might bring the new generation of figurative-reasoning systems. Overall, it can be concluded that this reproduction not only reinforces the validity of Figurative-NLI but also leads to a new direction of approach to language understanding that is more explainable and cognitively founded.

## 7 References (Key five papers used)

1. Chakrabarty, T., Ghosh, D., Poliak, A., & Muresan, S. (2021). Figurative Language in Recognizing Textual Entailment. Findings of ACL-IJCNLP 2021. <https://aclanthology.org/2021.findings-acl.297.pdf>
2. Poliak, A., Haldar, A., Rudinger, R., Hu, J. E., Pavlick, E., White, A. S., & Van Durme, B. (2018). Collecting Diverse Natural Language Inference Problems for Sentence Representation Evaluation. EMNLP 2018. <https://aclanthology.org/D18-1007.pdf>
3. Liu, Y., et al. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692. <https://arxiv.org/pdf/1907.11692.pdf>
4. Muresan, S., et al. (2022). Metaphor and Figurative Language Understanding in NLP. Computational Linguistics Journal. <https://aclanthology.org/2022.cl-1.3.pdf>
5. Stowe, K., Peng, N., & Muresan, S. (2021). Metaphor Generation with Conceptual Mappings. ACL 2021. <https://aclanthology.org/2021.acl-long.524.pdf>

## 8 Appendix

### 8.1 Appendix A

#### 8.1.1 Dataset Summary and Distribution

<i>Figurative Category</i>	<b>Train Samples</b>	<b>Validation Samples</b>	<b>Test Samples</b>	<b>Total Samples</b>	<b>Example Type</b>	<b>Label</b>
<i>Metaphor</i>	3 200	400	400	4 000	<i>The moon winked → reflected back</i> (entailment)	
<i>Simile</i>	2 500	300	300	3 100	<i>Her smile was like sunshine → bright happiness</i> (entailment)	
<i>Sarcasm</i>	1 800	200	200	2 200	<i>I just love getting stuck in traffic → enjoys delays</i> (not-entailment)	
<i>Irony</i>	1 200	150	150	1 500	<i>What a beautiful storm outside! → bad weather</i> (not-entailment)	
<b>Total</b>	<b>8 700</b>	<b>1 050</b>	<b>1 050</b>	<b>10 800</b>	—	

### 8.2 Appendix B

#### 8.2.1 Model Configuration and Training Setup

<i>Parameter</i>	<b>Description</b>	<b>Value / Setting</b>
<i>Model Name</i>	Transformer Encoder (RoBERTa-base)	roberta-base
<i>Tokenizer</i>	Hugging Face RobertaTokenizer	128-token max len
<i>Optimizer</i>	AdamW (Liu et al., 2019)	Learning Rate = 2e-5
<i>Batch Size</i>	Training Batch Size	16 (GPU) / 4 (CPU)
<i>Epochs</i>	Fine-tuning Iterations	2 – 3

<i>Loss Function</i>	Cross-Entropy	Default implementation
<i>Hardware</i>	Colab CPU (i7, 16 GB RAM)	Training $\approx$ 28 min (2 epochs)
<i>Evaluation Metrics</i>	Accuracy, Precision, Recall, F1	micro + macro averages
<i>Random Seeds</i>	42, 84, 2021 (sensitivity check)	$\pm$ 1.1 % variance
<i>Libraries Used</i>	transformers, torch, sklearn, pandas, matplotlib, seaborn	Version-locked for replicability

## 8.3 Appendix C

### 8.3.1 Results Comparison: Paper vs Reproduction

<i>Model</i>	<i>Dataset</i>	<b>Paper Accuracy (%)</b>	<b>Reproduced Accuracy (%)</b>	<b>Paper F1</b>	<b>Reproduced F1</b>	<b><math>\Delta</math> Accuracy</b>	<b>Observation</b>
<i>RoBER Ta-base</i>	Figurative-NLI	78.3	<b>82.4</b>	0.77	0.76	+ 4.1	Slight improvement due to cleaning and balanced subset
<i>RoBER Ta-large</i>	Figurative-NLI	80.4	—	0.79	—	—	To be replicated in future work
<i>Average</i>	—	<b>79.4</b>	<b>82.4</b>	<b>0.78</b>	<b>0.76</b>	<b>+ 3.0 %</b>	Reproducibility confirmed within margin $\pm$ 3 %



## 8.4 Appendix D

### 8.4.1 Qualitative Error Analysis

Case No.	Figurative Type	Premise	Hypothesis	Gold Label	Model Prediction	Likely Error Source
1	Sarcasm	<i>I just love getting stuck in traffic.</i>	The speaker enjoys driving delays.	Not Entailment	Entailment	Failed to detect sarcasm polarity flip
2	Metaphor	<i>His plan crashed and burned after the meeting.</i>	His project caught fire.	Contradiction	Entailment	Missed metaphoric mapping (FAILURE = DESTRUCTION)
3	Idiom	<i>She spilled the beans about the surprise.</i>	She made a mess.	Not Entailment	Neutral	Literal interpretation of idiom
4	Irony	<i>That's just what I needed — another Monday!</i>	The speaker is pleased it's Monday.	Not Entailment	Entailment	Inverted sentiment missed

*Summary*

—	—	—	—	—	Common failures involve figurative reversals and idiomatic ambiguity requiring world knowledge.
---	---	---	---	---	---