

# Plausible Extractive Rationalization through Semi-Supervised Entailment Signal

Yeo Wei Jie<sup>1</sup>, Ranjan Satapathy<sup>2</sup>, Erik Cambria<sup>1</sup>

<sup>1</sup>Nanyang Technological University

<sup>2</sup>Institute of High Performance Computing (IHPC),  
Agency for Science, Technology and Research (A\* STAR),  
1 Fusionopolis Way, #16-16 Connexis, 138632, Singapore  
yeow0082@e.ntu.edu.sg

## Abstract

The increasing use of complex and opaque black box models requires the adoption of interpretable measures, one such option is extractive rationalizing models, which serve as a more interpretable alternative. These models, also known as Explain-Then-Predict models, employ an explainer model to extract rationales and subsequently condition the predictor with the extracted information. Their primary objective is to provide precise and faithful explanations, represented by the extracted rationales. In this paper, we take a semi-supervised approach to optimize for the plausibility of extracted rationales. We adopt a pre-trained natural language inference (NLI) model and further fine-tune it on a small set of supervised rationales (10%). The NLI predictor is leveraged as a source of supervisory signals to the explainer via entailment alignment. We show that, by enforcing the alignment agreement between the explanation and answer in a question-answering task, the performance can be improved without access to ground truth labels. We evaluate our approach on the ERASER dataset and show that our approach achieves comparable results with supervised extractive models and outperforms unsupervised approaches by  $> 100\%$ . Code available at [https://github.com/wj210/NLI\\_ETP](https://github.com/wj210/NLI_ETP).

## 1 Introduction

Large language models such as Google’s BERT (Devlin et al., 2018) and OpenAI’s GPT series (Brown et al., 2020) are gaining widespread adoption in natural language processing (NLP) tasks. These models achieved impressive performance in multiple NLP tasks ranging from solving text generation to information extraction (Liu et al., 2023). However, little is known regarding how answers are generated or which portion of the input text the model focuses on. These flaws highlight concerns surrounding trust and fear of

undesirable biases in the model’s reasoning chain. Explainable AI (XAI) is currently an active field of research aimed at addressing these issues (Adadi and Berrada, 2018; Cambria et al., 2023; Yeo et al., 2023). In this work, we focus on extractive rationalizing models (Lei et al., 2016), which are also known as Explain-Then-Predict (ETP) models, and are designed towards producing highlights serving as **faithful** explanations. Faithfulness is defined as serving an explanation that represents the model’s reasoning process for a given decision, while plausibility refers to the level of agreement with humans (Jacovi and Goldberg, 2020). An advantageous characteristic of ETP models is that they concurrently produce the explanation and the task label, eliminating the necessity for an added layer of interpretation.

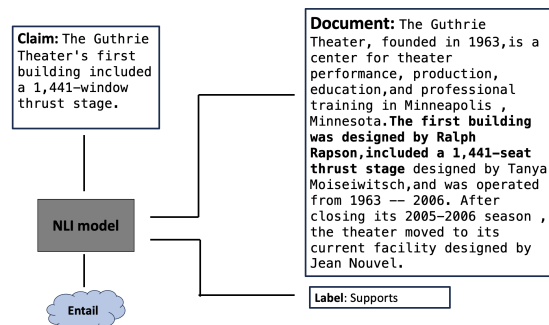


Figure 1: An example from the FEVER dataset, where the bold statement is the annotated rationale. Given the document and claim, the label denotes that the document contains evidence supporting the claim. The NLI predictor interprets this as a form of entailment between the claim and rationale.

This differs from post-hoc techniques such as LIME (Ribeiro et al., 2016) or SHAP (Lundberg and Lee, 2017), specifically tailored to interpret black-box models. Although these techniques are model-agnostic by design, they are computationally expensive and do not guarantee faithfulness nor optimized for plausibility. Chain-of-thought

(CoT) (Wei et al., 2022) is another popular approach, aimed at prompting Large Language Models (LLM) such as OpenAI’s GPT4 to elucidate its own prediction, in the form of reasoning steps which is said to be a form of explanation. However, we note that though the reasoning steps are seemingly plausible and convincing, there is no guarantee of the reasoning being faithful towards the supported output, since there is no constraint on the conditional variables. ETP models instead constrain the predictions on a compressed subset of the input, referred to as rationales, thereby guaranteeing the output to be solely conditioned on the subset, analogous to a binary form of feature relevance.

In our work, we focus on improving the plausibility of rationales, measured via matching human annotations. Several work has established benchmark datasets that consist of both the task label as well as human-annotated rationales (Bao et al., 2018; DeYoung et al., 2019). Current works in extractive rationalization mostly implement a pipeline procedure of training an explainer and a predictor (DeYoung et al., 2019), trained either jointly or separately. The training approach for these models can be bifurcated into two primary methods: supervised or unsupervised rationale extraction. In our methodology, we strike a balance by leveraging a minimal subset of annotated rationales ( $\leq 10\%$ ) to refine an ETP model. This refinement is applied to a separate NLI predictor, functioning as an auxiliary instructor for the explainer in the event of limited annotated rationales. More importantly, the explainer has no access to the annotations, which are exclusively presented to the NLI predictor.

Our approach is inspired by recent work in ensuring factual consistency in abstraction summarization (Roit et al., 2023), which has been found useful in cases of hallucination. Firstly, we create an augmented dataset based on a label transformation algorithm, based on the annotated rationales and NLI classes. This is used to provide further training on the NLI predictor to generate sentence-level rationale annotations.

NLI models are designed to determine whether a hypothesis contradicts, entails, or is neutral to a given premise. As such, they provide useful signals to align a given explanation to the answer produced by the predictor, as shown in subsequent experiments, this can have some desirable effects on the robustness of rationales (Chen et al., 2022). An example is illustrated in Figure 1 on a fact verifica-

tion task, where the purpose of the rationale is to act as evidence to either support or refute the given claim. In summary, the three key contributions of this work are the following:

- A simple yet effective approach that improves the plausibility and robustness of extracted rationales, while simultaneously improving task performance. The approach achieves competitive results against supervised models while outperforming unsupervised models by a large margin ( $>100\%$ ).
- To the best of our knowledge, this is the first work to utilize an auxiliary NLI predictor to generate augmented labels for extractive rationalization.
- Our approach has low resource requirements, using models of  $<300\text{M}$  parameters, and a small set of human-annotated rationales.

## 2 Methodology

### 2.1 Problem setting

Given an input document consisting of  $N$  sentences,  $x_i = \{x_{i,1}, x_{i,2}, \dots, x_{i,N}\}$ . The task objective can be decomposed into two steps, namely rationale extraction, and task prediction. An explainer,  $g_\phi$  takes in the input document and generates a binary mask over the sentences indicating the rationales,  $g_\phi(\hat{z}_i|x_i) \in \{0, 1\}_N$ .

The predictor,  $f_\theta$  can only consider the masked inputs during inference, since the initial reason for extractive rationalization is to present the rationales as a faithful explanation towards the task prediction,  $\hat{y}_i = f_\theta(\hat{z}_i \odot x_i)$ ,  $\odot$  is the element-wise multiplication. As rationales are designed to be a concise representation of the original text, there naturally exists a trade-off between generating a sparse  $z$  and retaining sufficient information to infer the task label accurately. In various studies, optimization strategies are generally consistent, differing mainly in the use of human-annotated labels for training rationale extractors. Our approach, instead employs a semi-supervised method using an auxiliary predictor optimized for NLI, denoted as  $f_{NLI}$ .

### 2.2 Semi-supervised NLI signal

Humans tend to prefer explanations that are aligned with the supported answer, similar to how NLI tasks involve generating the alignment between

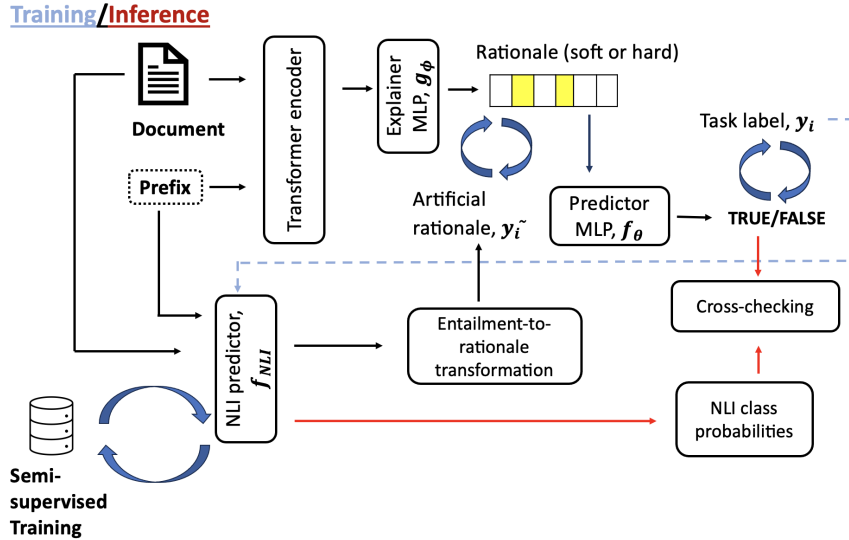


Figure 2: An overview of the proposed approach during training (bold in blue) and inference (bold in red). The NLI predictor only has access to the task label during training. The NLI predictor is initially fine-tuned using a limited set of annotated rationales, before generating artificial targets for the explainer. Cross-checking alignment is conducted during inference against the predictor (explained in Section 2.4)

two sentences. As such, NLI predictors naturally serve as helpful supervision in the absence of annotated rationales. This is especially applicable in a fact-verification scenario where the task is to infer if a given claim is supported by the provided document. For example, given a document containing the following annotated rationale: "Kung Fu Panda opened in 4,114 theaters, grossing \$20.3 million on its opening day" along with a claim: "Kung Fu Panda made more than \$1 million on opening day.". The rationale acts as supporting evidence if the corresponding label,  $y_i = SUPPORT$ , indicates that the claim should be supported given the document and vice versa. The NLI predictor is fine-tuned based on this simple heuristic, to match each sentence in the document against the query. It is trained on the augmented dataset created via a label transformation technique shown in Algorithm 1<sup>1</sup>. The transformation operates under the assumption that there are no contradictory rationales against the label, ie in  $z_{i,j}$  contradicts the claim when the label is *support*. This is a valid assumption since a sample with rationales contradicting the label should be considered an erroneous sample. The human-annotated rationales are used solely in the transformation, thus the explainer does not explicitly see the annotations.

During training, the NLI predictor acts as the

<sup>1</sup>Shown on FEVER dataset, where the supporting facts either support or refute the claim. Similarly applied for true/false settings.

---

### Algorithm 1 Rationale to NLI label transformation

---

**Input:** Annotated rationale,  $z_i$ , task label,  $y_i$

**Output:** NLI label,  $\tilde{z}_i$

- 1: **for**  $z_{i,j}$  in  $z_i$  **do**
  - 2:   **if**  $z_{i,j} = 1 \wedge y_i = SUP$  **then**
  - 3:      $\tilde{z}_{i,j} = \text{entailment}$
  - 4:   **else if**  $z_{i,j} = 1 \wedge y_i = REF$  **then**
  - 5:      $\tilde{z}_{i,j} = \text{contradiction}$
  - 6:   **else**
  - 7:      $\tilde{z}_{i,j} = \text{neutral}$
  - 8:   **end if**
  - 9: **end for**
  - 10: **return**  $\tilde{z}_i$
- 

source of supervision in place of the human-annotated rationales. As the explainer is trained to predict a binary mask, Algorithm 1 can be implemented in reverse to transform the NLI outputs back to rationale labels,  $\tilde{z}$  for the explainer's training, (see Appendix for more details). We note that the above approach is likewise applicable to any binary true/false tasks where the predictor has to indicate if the answer is true or false concerning the question.

### 2.3 Sentence-level training

We utilizes a pipeline approach consisting of a shared encoder, along with separate decoders for the explainer and predictor. The input is first encoded into contextualized hidden states,  $h_{i,1:L} =$

$enc(x_{i,1:L})$ , where  $L$  is at the token level. We follow (Paranjape et al., 2020) and transform the token-level hidden states into sentence-level by concatenating hidden states at the starting and ending positions and feeding it into an explainer to produce rationales,  $\tilde{z}_i = g_\phi(h_i)$ , where  $h_i = MLP(h_{i,s} \oplus h_{i,e})$ ,  $\oplus$  is the concatenation process.

The predictor is conditioned on the rationales and trained using standard cross entropy.

$$L_{f_\theta} = -\mathbb{E}_{z \sim g_\phi(z|x)}[\log(\hat{y}_i | \hat{z} \odot x)] \quad (1)$$

The explainer loss,  $L_{g_\phi}$  is similarly computed with (1), but against the augmented targets,  $\tilde{z}_i = f_{NLI}(\tilde{z}_i|x_i, y_i) \in \{0, 1\}^N$ , instead of the annotated targets. The full training and inference approach is depicted in Figure 2, where the NLI predictor is first fine-tuned before training the ETP model. The choice of a shared encoder allows for a form of dependency between  $e_i$  and  $\hat{y}_i$ , as the encoder has to jointly optimize the representation to infer both the task label and rationales accurately. The final loss is thus a combination of both the predictor and explainer cross-entropy loss,  $L_{total} = L_{f_\theta} + \lambda L_{g_\phi}$ , where  $\lambda$  balances the trade-off between classification and plausibility performance.

The label transformation is only used during training as it requires access to  $y_i$  which is not available at test time. However, we will show how  $f_{NLI}$  can remain useful during inference.

## 2.4 Inference

During inference, the rationales are extracted solely by the trained explainer,  $f_\theta$ . However,  $f_{NLI}$  can act as a cross-checker against the predictor  $g_\phi$  in the event of a distributional shift in  $g_\phi$ . Given  $\hat{z}_i$  and a prefix (claim in fact verification or question-answer pair in Q&A task),  $f_{NLI}$  denotes if  $\hat{z}_i$  contradicts or entails the prefix. We ignore the neutral class and re-weight the NLI class probabilities,  $p(\tilde{y}_i^C)$  before averaging across the  $n$  selected sentences in each instance,

$$p(\tilde{y}_i^C) = \frac{1}{n} \sum_{j=1}^n p(\tilde{y}_{i,j}^C) \quad (2)$$

where  $C$  denotes the NLI class instance. The task probabilities,  $p(\hat{y}_i^C)$  are then scaled with the NLI probabilities.

$$p(\hat{y}_i^C) = p(\tilde{y}_i^C) \cdot p(\tilde{y}_i^C) \quad (3)$$

This is helpful in the case where  $f_\theta$  is less confident around the decision boundary and  $f_{NLI}$  can provide additional support, given the additional training which aligns  $f_{NLI}$ 's decision between supporting rationales and task label.

## 3 Experiments

### 3.1 Datasets

We evaluate our approach against unsupervised and supervised baselines across three benchmark tasks from ERASER. ERASER contains a suite of NLP tasks, extended with human-annotated rationales, to assess plausibility.

- **FEVER**: A fact-verification dataset (Thorne et al., 2018), each instance consists of a claim and a document, where the goal is to determine if the claim is supported or refuted using information from the document.
- **BoolQ**: Question-answering task (Clark et al., 2019), containing a context document from Wikipedia and a question, the answer is either true or false. Due to the long sequence, we select the most relevant portion of the context using TF-IDF scoring similar to (Paranjape et al., 2020).
- **MultiRC**: A multi-hop dataset (Khashabi et al., 2018), requiring reasoning over multiple sentences to infer to correct answer. Multiple answer choices can be associated with a single question and the task is to predict if the answer is true or false.

### 3.2 Experimental Setup

We use RoBERTa-base (Liu et al., 2019) as the shared encoder between the explainer and predictor. The NLI predictor,  $f_{NLI}$  is a DeBERTa-large transformer (He et al., 2021) fine-tuned on multiple NLI datasets, we use the v3 variant. Our approach is agnostic to the choice of the pre-trained transformer for both the backbone encoder and NLI predictor. We selected RoBERTa-base, with its 125M parameters, due to its computational efficiency compared to larger models, while still maintaining high performance. We fine-tune the NLI predictor with 10%<sup>2</sup> of the annotated rationales. We list the full hyperparameter details in A.2. A notable benefit of

<sup>2</sup>The restriction in training samples applies only to rationales, where we train the predictor on the full set of task labels.

our approach is that it does not require an expensive search over objective-related hyperparameters.

### 3.3 Baselines

We evaluate our approach against both supervised and unsupervised settings, along with predictors subjected to full context. We refer to **Full-C** as the predictor-only set up to assess the gap in task performance between using the full context as compared to a subset. **Supervised** trains the explainer against human-annotated labels,  $z_i$ , instead of  $\tilde{z}_i$  in our approach, serving as the upper bound for plausibility.

**IB** is an unsupervised approach from (Paranjape et al., 2020) which optimizes an information-bottleneck objective and selects top  $N\%$  according to pre-defined sparse prior. The author additionally introduces a semi-supervised approach of using 25% of the annotated rationales which we refer to as **IB-25%**. Note that this baseline is subjected to higher supervision compared to ours (10%). We included the reported results for the sake of fairness (**R**). We choose 10% based on empirical results, serving as a good trade-off between minimal resource requirement and performance, albeit a comparable level of supervision (25%) can be referred from Table 4. All evaluated approach implements an ETP-type setup, consisting of an explainer and predictor except for Full-C.

### 3.4 Metrics

We report task performance using classification metrics such as accuracy and F1-score, while the plausibility of extracted rationales is assessed using F1-score (DeYoung et al., 2019) at the sentence level, Sentence-F1. We leave out any faithfulness metrics such as sufficiency as we assume ETP models to be inherently faithful given that the predictor is only subjected to the extracted explanation. We also assess the robustness by exposing the explainer to adversarial inputs (Chen et al., 2022) in Section 4.2.

We employ the following equations (Chen et al., 2022) to compute the normalized discrepancy in task performance,  $\Delta_T$  and plausibility,  $\Delta_P$  between the original and perturbed inputs as an indicator of robustness. Additionally, we utilize the attack rate,  $AR$  to gauge the frequency with which the explainer identifies adversarial sentences.

$$\Delta_T = \frac{1}{N} \sum_{i=1}^N \frac{M_t(\hat{y}_i, y_i) - M_t(\hat{y}_i^A, y_i)}{M_t(\hat{y}_i, y_i)} \quad (4)$$

$$\Delta_P = \frac{1}{N} \sum_{i=1}^N \frac{M_p(\hat{z}_i, z_i) - M_p(\hat{z}_i^A, z_i)}{M_p(\hat{z}_i, z_i)} \quad (5)$$

$$AR = \frac{1}{N} \sum_{i=1}^N \hat{z}_i \cap z^{AS} \quad (6)$$

$M_t$  and  $M_p$  is the scoring function for task and plausibility performance, for which we use the F1 and Sentence-F1 measurement.  $\hat{y}_i^A$ , and  $\hat{z}_i^A$  refer to the generated class label and rationale given the adversarial input.  $z^{AS}$  refers to the position of the adversarial prefix.

## 4 Results

All results are averaged over three runs with different seeds. For Full-C, we do not report plausibility performance since there is no explainer module. In the ERASER benchmark, the number of annotated rationale varies between tasks, where BoolQ features a higher quantity of annotated sentences and also includes larger number of contiguous spans. The main objective of this study is to assess between different unsupervised and supervised approaches in generating plausible and robust rationales, while minimizing negative effects on downstream task performance.

### 4.1 Plausibility and Task Analysis

The task and plausibility performance is shown in Table 1. Judging from the results, our approach achieves highly competitive performance against the gold standard for both task (Full-C) and plausibility (Supervised). In FEVER, it even surpasses the full context approach (94.2 vs 93). It goes to show that ETP-like models can benefit from ignoring spurious noise by conditioning the predictor to only text considered essential for inferring the target class. The additional usage of  $f_{NLI}$  as a cross-checker during inference also provided considerable improvements across all three benchmarks, at little to no cost in computational resources.

In terms of plausibility, our method delivers a plausibility score that is on par with the fully supervised approach across all datasets except BoolQ. We note that a likely reason is that the target rationales are largely inconsistent in length, with instances stretching across as many as six contiguous sentences. Since the NLI predictor is optimized toward matching each sentence with the given query. It may fare worse when individual sentences appear to be unrelated to the query but are nonetheless

Approach	FEVER			MultiRC			BoolQ		
	Task Acc	F1	Plausibility Sent-F1	Task Acc	F1	Plausibility Sent-F1	Task Acc	F1	Plausibility Sent-F1
Full-C	93	91.8	-	<b>76</b>	<b>72</b>	-	65.8	53	-
Supervised	90.1	88.4	<b>83.4</b>	74.3	70.5	<b>64.1</b>	<b>72.4</b>	<b>65.9</b>	<b>76</b>
IB	85.9	85.9	38.9	64.1	63	23.1	64	63.5	10.3
IB w 25%	85.1	85.1	38.4	67.6	67.5	52.7	58.6	52.1	11.4
IB w 25% (R)	-	88.8	63.9	-	66.4	54	-	63.4	19.2
Ours (10%)	<b>93.7</b> <sub>+0.7</sub>	<b>92.6</b> <sub>+0.7</sub>	<b>80.1</b>	72.5 <sub>+0.2</sub>	68.6 <sub>+0.5</sub>	<b>56.4</b>	67.4 <sub>+2.2</sub>	51.4 <sub>+9</sub>	<b>29.6</b>

Table 1: Classification and plausibility performance comparison across the three ERASER tasks. Test results are averaged across 3 seeds. The subscript refers to the case where the NLI predictor is used as a cross checker, in 2.4. Results highlighted in bold refer to the best-performing approach. The supervised approach acts as the upper bound on plausibility performance. **R** is the reported results of the IB approach (Paranjape et al., 2020).

FEVER	MultiRC	BoolQ
100	56.7	20

Table 2: Percentage of extracted over target rationales. BoolQ has the lowest percentage out of all three datasets.

annotated as rationales. Table 2 shows the percentage proportion of sentences annotated as rationales over the target. It’s noteworthy that the explainer marks fewer sentences due to the NLI predictor’s tendency to classify the majority of sentences as neutral, deeming them non-essential for task prediction.

Nonetheless, optimizing the explainer with NLI supervision is proven to be superior compared to the unsupervised information bottleneck objective by (Paranjape et al., 2020). Our approach outperforms the former by large margins in terms of plausibility on FEVER (> 25%) and BoolQ (> 50%), even when provided with a lower amount of supervision (10% vs 25%). The performance gap is even larger when compared to fully unsupervised (IB), with more than twice the scores. The IB method learns a sparse mask over the input document,  $x_i$  by maximizing the mutual information between rationale  $z_i$  and task label  $y_i$  while limiting the extraction budget to a pre-defined prior. However, estimating the prior is difficult and can be detrimental in instances with varying rationale lengths such as in BoolQ. Our approach sidesteps the complicated training yet achieves a better-tuned explainer in extracting plausible rationales.

## 4.2 Robustness

In this section, we evaluate the robustness of ETP models when faced with inputs prefixed with an

adversarial query. The query is unrelated to the document and carries a contrastive meaning with respect to the original. For example, given a claim in FEVER, "Earl Scruggs was a musician who played banjo.", the noun, "Earl Scruggs" and "banjo" is replaced to form the adversarial sentence "manchester archer was a songwriter who played mandolin.". The attack is minimally changed from the query to distract the explainer. A model with limited robustness might interpret the attack as pertinent due to its analogous semantics, thereby influencing the predictor and undermining task performance. The robustness results are reported in Table 3. We found similar findings as compared to (Chen et al., 2022) who note that ETP models exhibit greater robustness compared to predictors subjected to the full context.

In the FEVER dataset, our approach suffers the lowest drop in task and plausibility performance, while having the lowest AR in both datasets. IB has the highest AR, even extracting every adversarial sentence in FEVER. A contributing factor to our approach’s low AR rate is that the NLI signal is derived by verifying if a sentence aligns with the query based on the provided task label. This strengthens the explainer’s proficiency in dismissing instances that don’t satisfy this criterion. On the other hand, the explainer trained with IB is emphasized to maximize the task objective, which can lead to situations where a minimally perturbed sentence is mistakenly perceived as useful. This further proves that training with NLI feedback produces more robust and plausible models.

Approach	FEVER			MultiRC		
	$\Delta_T$	$\Delta_P$	AR	$\Delta_T$	$\Delta_P$	AR
Full-C	11.2	-	-	29.6	-	-
Supervised	10.8	37	54.6	14.3	26.7	68.3
IB (25%)	12	35.8	100	<b>4.9</b>	<b>19.1</b>	93
Ours (10%)	<b>7.8</b>	<b>8.4</b>	<b>32.3</b>	10.2	27.1	<b>67.2</b>

Table 3: In both FEVER and MultiRC, we measure robustness with a preference for lower values. Models considering the full context are evaluated solely based on the difference in task performance as they don’t engage in rationale extraction. All values are normalized percentages drop computed via (4), (5) and (6)

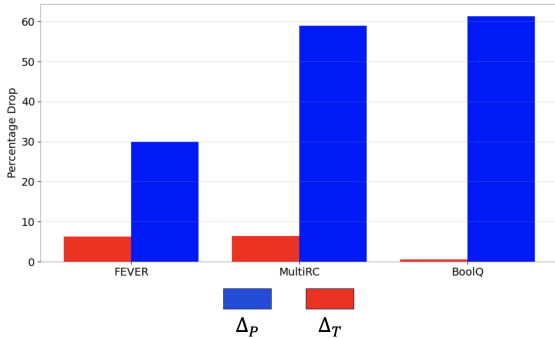


Figure 3: Task and Plausibility performance drop when there is no further fine-tuning on the NLI predictor (10% data). The metrics are computed similarly to robustness using (4) and (5) and are presented in normalized percentages.

## 5 Ablation

### 5.1 Importance of NLI training

In this study, we seek to question the usefulness of introducing further fine-tuning using the limited set of annotations. While the NLI predictor is previously fine-tuned on various NLI tasks, the sentence lengths in its training distribution differ from those in our experimental datasets. Additionally, domain-specific semantics differences can introduce variations in the NLI predictor’s inference process. Consequently, the NLI predictor might not always accurately discern the NLI class, leading to the generation of misleading signals for the explainer. To quantify the effectiveness of further fine-tuning, we compute the drop-in performance on both task and plausibility between an NLI predictor that is fine-tuned, referred to as **FI** and one that is not, **NFI**. The gap in task and plausibility performance is reported in Figure 3.

These results substantiate our initial hypothesis. Without fine-tuning, NFI struggles to provide meaningful feedback to the explainer, primarily

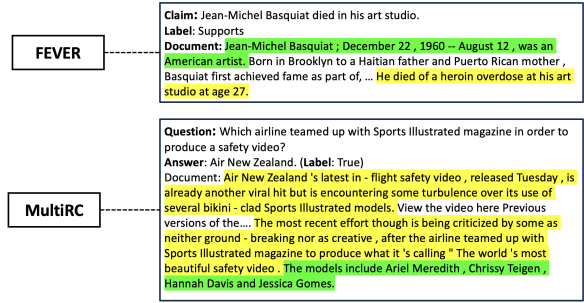


Figure 4: Example of query and input document where the sentences highlighted in green refer to the NLI predictor without fine-tuning. Yellow refers to the annotated rationale as well as extracted by the fine-tuned predictor.

because of its limited capability to accurately determine whether a specific sentence should support or contradict the query based on the given task label. Taking a closer look at sentence classifications in Figure 4 reveals that the NFI tends to mistakenly identify neutral sentences as entailments. In the FEVER example, although the initial sentence shares a noun with the claim, it does not address the death of the noun’s subject yet the NFI incorrectly recognizes it as entailment. Similarly, in the MultiRC instance, the concluding sentence lacks any significant connection to the given question or answer. This may be the reason why despite a significant drop in accurately extracting the correct rationale, 58.9% in MultiRC and 61.3% in BoolQ, the task performance surprisingly does not incur a huge loss ( $< 10\%$ ). A neutral sentence would not drastically change the class probabilities of the predictor as compared to a contradicting sentence. Nevertheless, incorporating additional fine-tuning on the NLI predictor is still essential in filtering out false positives such as sentences with neutral relationships in inferring the output.

### 5.2 Model sizes and NLI supervision

Ablation type	Acc	F1	Sent-F1
Original (base w 10%)	72.5	68.6	56.4
Large	74.3	71.2	58.1
Base w 25%	73.3	71.1	57.9
Base w 50%	73.8	70.9	58
Without NLI Pre-FT	69.2	64.3	52.6

Table 4: Ablation on model size, % NLI supervision and effects of not doing pre-finetuning of  $f_{NLI}$  on general NLI tasks (SNLI, MultiNLI). Implemented on MultiRC.

We carry out further analysis on the effect of both model size and amount of NLI supervision

given to  $f_{NLI}$ . We compare RoBERTA-large (330M) with the original 10% supervision and the base model with increased level of supervision  $\in [25, 50]$ . We additionally compare a DeBERTa encoder without prior fine-tuning on NLI datasets, while similarly fine-tuning on 10% of annotated rationales. The benefits of using a larger encoder and increased NLI supervision for  $f_{NLI}$  can be observed from Table 4. Notably, there is little difference in both accuracy and plausibility scores with higher supervision. Furthermore, our approach remains effective using off-the-shelf encoders without prior fine-tuning. This highlights the strength of our approach which remains effective even in low-resource conditions.

## 6 Related Works

**Linguistic Interpretability:** In recent years, extractive rationalization and attention-based interpretability have emerged as significant approaches within AI research, with numerous studies contributing to the field (Gurrapu et al., 2023; Mohankumar et al., 2020; Serrano and Smith, 2019). However, the efficacy of using attention as an interpretability mechanism has been debated, highlighting a division among researchers. Critics argue that attention scores do not significantly impact model predictions and present challenges in generating counterfactuals (Jain and Wallace, 2019), and may be biased due to their reliance on neighboring token information (Bai et al., 2021; Tutek and Šnajder, 2022). In contrast, proponents suggest that the relevance of attention weights varies with the definition of faithfulness and that multiple weight combinations can yield the same output (Wiegrefe and Pinter, 2019). Efforts to enhance the reliability of attention mechanisms include task-specific attention constraints (Chrysostomou and Aletras, 2021b) and penalties on scores for key words (Chrysostomou and Aletras, 2021a).

**Extractive Rationalization:** The foundational work by (Lei et al., 2016) introduced extractive rationalization using REINFORCE with sparsity regularization for end-to-end training of explainers via predictors’ objectives. The ERASER benchmark, established by (DeYoung et al., 2019), evaluates explainer-predictor models across seven NLP tasks, utilizing a BERT-to-BERT framework and sequential training. Subsequent research, such

as (Atanasova et al., 2022), focused on explainer consistency and confidence, while (Lakhotia et al., 2020) employed Fusion-In-Decoder for rationale extraction in lengthy documents.

The challenge of acquiring supervised rationales has spurred interest in unsupervised methods for generating reliable rationales. Efforts include (Paranjape et al., 2020), targeting rationale conciseness through information bottleneck optimization, and (Ghoshal et al., 2022), which mitigates spurious correlations in QA by adding a question generation objective. (Jain et al., 2020) modularizes the objective, and (Glockner et al., 2020) uses sentence-specific encoding with aggregated loss for rationale selection. Unlike these methods, our approach seeks to optimize all rational sentences without relying on the assumption of pre-defined rationale availability.

**NLI signals:** There have also been works utilizing NLI signals to enhance downstream tasks. (Roit et al., 2023) directly uses entailment scores as a reward signal to optimize factual summarization using RL, (Laban et al., 2022; Kryściński et al., 2019) for mitigating inconsistency in abstractive summarization. (Chen et al., 2023) performs self-rationalization training using a small set of annotated rationales and then annotating the rest of the unlabelled dataset. However, the work is based on abstractive setting, using free-text self-generated rationales, thereby violating the faithfulness property of the explanation. (Golovneva et al., 2022) utilizes NLI as a metric for ensuring semantic correctness in explanations.

## 7 Conclusion

In this paper, we have introduced a simple yet unique way of generating artificial learning signals from an alternative source, to cope with scenarios where human-annotated rationales are scarce. The method harnesses a transformer pre-trained on the NLI task. Through additional fine-tuning, the NLI predictor can produce less biased labels, enhancing the learning process for the explainer.

Through the extensive experiments conducted, we have shown that our work can alleviate the plausibility and robustness of ETP models in a low-resource environment. Notably, with just 10% of the annotated rationale, our method delivers performance on par with fully supervised models and



significantly outperforms both semi-supervised approaches that utilize more annotated data and unsupervised settings. In future directions, we plan to extend this work toward models that generate abstractive explanations, where the NLI signal can act as verification feedback to ensure the mitigation of biased explanations. Another interesting direction is to study how can we extend the NLI predictor’s coverage beyond a single sentence, to capture the correspondence between longer documents.

## 8 Limitations

We only evaluate a singular trait of interpretability: plausibility. We note that multiple other traits of interpretability are equally important and we leave that to further work. The sizes of the encoder models implemented in this work are relatively small, with the biggest consisting of 300M parameters. Though model scaling is the primary objective, we note the importance of extending our work towards larger models given the popularity of NLP research surrounding LLMs.

## References

- Amina Adadi and Mohammed Berrada. 2018. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160.
- Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2022. Diagnostics-guided explanation generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10445–10453.
- Bing Bai, Jian Liang, Guanhua Zhang, Hao Li, Kun Bai, and Fei Wang. 2021. Why attentions may not be interpretable? In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 25–34.
- Yujia Bao, Shiyu Chang, Mo Yu, and Regina Barzilay. 2018. Deriving machine attention from human rationales. *arXiv preprint arXiv:1808.09367*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Erik Cambria, Lorenzo Malandri, Fabio Mercurio, Mario Mezzanzanica, and Navid Nobani. 2023. A survey on XAI and natural language explanations. *Information Processing and Management*, 60(103111).
- Howard Chen, Jacqueline He, Karthik Narasimhan, and Danqi Chen. 2022. Can rationalization improve robustness? *arXiv preprint arXiv:2204.11790*.
- Wei-Lin Chen, An-Zi Yen, Hen-Hsen Huang, Cheng-Kuang Wu, and Hsin-Hsi Chen. 2023. Zara: Improving few-shot self-rationalization for small language models. *arXiv preprint arXiv:2305.07355*.
- George Chrysostomou and Nikolaos Aletras. 2021a. Enjoy the salience: Towards better transformer-based faithful explanations with word salience. *arXiv preprint arXiv:2108.13759*.
- George Chrysostomou and Nikolaos Aletras. 2021b. Improving the faithfulness of attention-based explanations with task-specific information for text classification. *arXiv preprint arXiv:2105.02657*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *NAACL*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C Wallace. 2019. Eraser: A benchmark to evaluate rationalized nlp models. *arXiv preprint arXiv:1911.03429*.
- Asish Ghoshal, Srinivasan Iyer, Bhargavi Paranjape, Kushal Lakhota, Scott Wen-tau Yih, and Yashar Mehdad. 2022. Quaser: Question answering with scalable extractive rationalization. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1208–1218.
- Max Glockner, Ivan Habernal, and Iryna Gurevych. 2020. Why do you think that? exploring faithful sentence-level rationales without supervision. *arXiv preprint arXiv:2010.03384*.
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2022. Roscoe: A suite of metrics for scoring step-by-step reasoning. *arXiv preprint arXiv:2212.07919*.
- Sai Gurrupu, Ajay Kulkarni, Lifu Huang, Ismini Lourentzou, Laura Freeman, and Feras A Batarseh. 2023. Rationalization for explainable nlp: A survey. *arXiv preprint arXiv:2301.08912*.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? *arXiv preprint arXiv:2004.03685*.
- Sarthak Jain and Byron C Wallace. 2019. Attention is not explanation. *arXiv preprint arXiv:1902.10186*.

- Sarthak Jain, Sarah Wiegrefe, Yuval Pinter, and Byron C Wallace. 2020. Learning to faithfully rationalize by construction. *arXiv preprint arXiv:2005.00115*.
- Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: a challenge set for reading comprehension over multiple sentences. In *Proceedings of North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Evaluating the factual consistency of abstractive text summarization. *arXiv preprint arXiv:1910.12840*.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. Summac: Re-visiting nli-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Kushal Lakhota, Bhargavi Paranjape, Asish Ghoshal, Wen-tau Yih, Yashar Mehdad, and Srinivasan Iyer. 2020. Fid-ex: Improving sequence-to-sequence models for extractive rationale generation. *arXiv preprint arXiv:2012.15482*.
- Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2016. Rationalizing neural predictions. *arXiv preprint arXiv:1606.04155*.
- Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. 2023. Summary of chatgpt/gpt-4 research and perspective towards the future of large language models. *arXiv preprint arXiv:2304.01852*.
- 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 preprint arXiv:1907.11692*.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.
- Akash Kumar Mohankumar, Preksha Nema, Sharan Narasimhan, Mitesh M Khapra, Balaji Vasani, and Balaraman Ravindran. 2020. Towards transparent and explainable attention models. *arXiv preprint arXiv:2004.14243*.
- Bhargavi Paranjape, Mandar Joshi, John Thickstun, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. An information bottleneck approach for controlling conciseness in rationale extraction. *arXiv preprint arXiv:2005.00652*.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Paul Roit, Johan Ferret, Lior Shani, Roei Aharoni, Geoffrey Cideron, Robert Dadashi, Matthieu Geist, Serkan Girgin, Léonard Hussenot, Orgad Keller, et al. 2023. Factually consistent summarization via reinforcement learning with textual entailment feedback. *arXiv preprint arXiv:2306.00186*.
- Sofia Serrano and Noah A Smith. 2019. Is attention interpretable? *arXiv preprint arXiv:1906.03731*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. 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 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Martin Tutek and Jan Šnajder. 2022. Toward practical usage of the attention mechanism as a tool for interpretability. *IEEE Access*, 10:47011–47030.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. *arXiv preprint arXiv:1908.04626*.
- Wei Jie Yeo, Wihan van der Heever, Rui Mao, Erik Cambria, Ranjan Satapathy, and Gianmarco Mengaldo. 2023. A comprehensive review on financial explainable ai. *arXiv preprint arXiv:2309.11960*.

## A Appendix

In the main paper, we showed how the label transformation technique is used to transform an annotated rationale into an NLI-associated label, for the purpose of fine-tuning the NLI predictor. We will now show how the reverse is applied to facilitate the training of the explainer.

### A.1 Reverse label transformation

Given a query,  $q$  and each sentence,  $x_i$ , we concatenate the query and sentence as input to the NLI predictor, where the NLI class label is generated as  $\tilde{y}_i = f_{NLI}(q \oplus x_i)$ . This applies to both queries with a single sentence such as the claim in FEVER and BoolQ or double sentences in MultiRC, comprising of both the question and answer.

The NLI class,  $\tilde{y}_i$  is used together with the task label,  $y_i$  to generate  $\tilde{z}_i$ , used in place of  $z_i$  for the semi-supervised explainer. The transformation is detailed in Algorithm 2, applied in reverse to Algorithm 1. C refers to Contradiction, and E to Entailment (example shown on FEVER task). Note that if the  $f_{NLI}$  indicates that the sentence is neutral to the query, the sentence is automatically labeled as a non-rationale. This is similar in the case where if a document is annotated as false with respect to the query, all rationales should be a contradiction and vice versa.

---

**Algorithm 2** Reverse label transformation

---

**Input:** query,  $q_i$ , input document,  $x_i$ , task label,  $y_i$  and NLI predictor,  $f_{NLI}$

**Output:** NLI label,  $\tilde{z}_i$

```

1: for each  $x_{i,j} \in x_i$  do
2:    $\tilde{y}_i \leftarrow f_{NLI}(q_i \oplus x_{i,j})$ 
3:   if ( $\tilde{y}_i = E \wedge y_i = \text{SUP}$ ) or ( $\tilde{y}_i = C \wedge y_i = \text{REF}$ ) then
4:      $z_{i,j} \leftarrow 1$ 
5:   else
6:      $z_{i,j} \leftarrow 0$ 
7:   end if
8: end for
9: return  $\tilde{z}_i$ 

```

---

## A.2 Hyperparameters

We use the AdamW optimizer from (Loshchilov and Hutter, 2017) with  $\epsilon$  set at 1e-8 and fix the batch size at 8. We use a learning rate warm-up scheduler with the final rate capped at 2e-5 and clip all gradient norms at a value of 1.0 while applying a dropout of 0.2 for the explainer. The explainer is a two-layer MLP with ReLU activation. Early stopping is implemented where the training is stopped if the validation loss does not improve after 3 epochs. We run all our experiments for a maximum of 10 epochs, on NVIDIA A6000s, implemented with PyTorch. We do not find much difference between various values of  $\lambda$  and set it to 1.