Measuring Association Between Labels and Free-Text Rationales

Sarah Wiegreffe[♣] Ana Marasović^{†♦}

Noah A. Smith^{†♦}

*School of Interactive Computing, Georgia Institute of Technology

†Allen Institute for Artificial Intelligence

Paul G. Allen School of Computer Science and Engineering, University of Washington saw@gatech.edu, {anam, noah}@allenai.org

Abstract

In interpretable NLP, we require faithful rationales that reflect the model's decision-making process for an explained instance. prior work focuses on extractive rationales (a subset of the input words), we investigate their less-studied counterpart: free-text natural language rationales. We demonstrate that pipelines, models for faithful rationalization on information-extraction style tasks, do not work as well on "reasoning" tasks requiring free-text rationales. We turn to models that jointly predict and rationalize, a class of widely used high-performance models for freetext rationalization. We investigate the extent to which the labels and rationales predicted by these models are associated, a necessary property of faithful explanation. Via two tests, robustness equivalence and feature importance agreement, we find that state-ofthe-art T5-based joint models exhibit desirable properties for explaining commonsense question-answering and natural language inference, indicating their potential for producing faithful free-text rationales.1

1 Introduction

Interpretable NLP aims to better understand predictive models' internals for purposes such as debugging, validating safety before deployment, or revealing unintended biases and behavior (Molnar, 2019). These objectives require faithful rationales—explanations of the model's behavior that are accurate representations of its decision process (Melis and Jaakkola, 2018).

One way towards faithfulness is to introduce architectural modifications or constraints that produce rationales with desirable properties (Andreas et al., 2016; Schwartz et al., 2018; Jiang et al., 2019, *inter alia*). For example, pipeline models (Figure 2) were designed for information extraction (IE) tasks

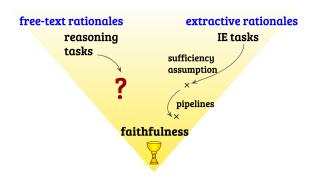


Figure 1: A categorization of interpretable NLP on an illustrative faithfulness spectrum. Two predominant forms of explanation exist that align with two predominant classes of NLP tasks. Unlike models for IE tasks, the desirable properties of interpretable models for reasoning tasks have not been explored. We investigate architectures and tests for explaining reasoning tasks.

for which a rationale can be extracted as a subset of the input and is *sufficient* to make a prediction on its own, without the rest of the input (Lei et al., 2016). Such models approach faithfulness by construction (Jain et al., 2020).

There is a growing interest in tasks that require world and commonsense "knowledge" and "reasoning", such as commonsense question-answering (CommonsenseQA; Talmor et al., 2019) and natural language inference (SNLI; Bowman et al., 2015). Here, extractive rationales necessarily fall short—rationales must instead take the form of free-text natural language to fill in the reasoning or knowledge gap (Camburu et al., 2018; Rajani et al., 2019). In Table 1, for example, the highlighted extractive rationale of the first problem instance lacks at least one reasoning step to adequately justify the answer; the natural language rationale (which is not extractive) fills in the gap.

¹Our code is available at https://github.com/allenai/label_rationale_association.

²We use "free-text" and "natural language" rationales interchangeably. We additionally use the term "rationale" to also mean "explanation"; for a more detailed discussion of terminology see Jacovi and Goldberg (2021); Wiegreffe and Marasović (2021).

Commonsense QA (CoS-E)	Question: While eating a hamburger with friends, what are people trying to do? Answer choices: have fun, tasty, or indigestion Natural language rationale: Usually a hamburger with friends indicates a good time.					
Natural	Premise: A child in a yellow plastic safety swing is laughing as a dark-haired woman stands behind her.					
Language	Hypothesis: A young mother is playing with her daughter in a swing.					
Inference (E-SNLI)	Label choices: <u>neutral</u> , entailment, or contradiction Natural language rationale: Child does not imply daughter and woman does not imply mother.					

Table 1: Examples from the CoS-E v1.0 and E-SNLI datasets (§2). Extractive rationales annotated by humans are highlighted, while human-written free-text rationales are presented underneath the answer/label choices. These examples illustrate that the extractive rationales fail to adequately explain the correct (underlined) label.

We study two distinct model classes: selfrationalizing models, which are fully differentiable and jointly predict the task output with the rationale; and pipelines, which rationalize first and then predict task output with a separate model. We first show that, for CommonsenseQA and SNLI, a selfrationalizing model provides rationales that better indicate the correct label than a pipeline (§3.1). Next, we show that sufficiency is not universally applicable: a natural language rationale on its own does not generally provide enough information to arrive at the correct answer (§3.2). These findings suggest that a faithful-by-construction pipeline is not an ideal approach for reasoning tasks, leading us to ask: is there is a way to achieve faithful freetext rationalization with self-rationalizing models?

We note that there is currently no way to assess the relationship between a prediction and a free-text rationale within the same fully differentiable model. Jacovi and Goldberg (2020) argue for the development of evaluations that measure the *extent* and *likelihood* that a rationale is faithful in practice (illustrated in Figure 1). To do so, we propose two measurements to initiate testing the extent to which predicted labels and explanations are associated within the model that produces them.

The first experiment, *robustness equivalence* (§4.1), analyzes whether a predicted label and generated rationale are similarly robust to noise. The second, *feature importance agreement* (§4.2), analyzes whether the gradient-attributions of the input with respect to the predicted label are similar to those with respect to the predicted rationale. We show that a self-rationalizing finetuned variant of T5 (Raffel et al., 2020; Narang et al., 2020) demonstrates good robustness equivalence and feature importance agreement on the datasets investigated. This result motivates future work on more measurements for testing label-rationale association.

2 Tasks, Datasets, and Models

Before we turn to our analyses we introduce datasets and models used for our experiments.

Tasks and Datasets We explore two large-scale datasets for textual reasoning tasks that contain human-written natural language rationales: E-SNLI (Camburu et al., 2018), an extension of SNLI (Bowman et al., 2015); and CoS-E (Rajani et al., 2019), an extension of CommonsenseQA (Talmor et al., 2019) (both in English). For the former, the task is to infer whether a given hypothesis sentence entails, contradicts, or is neutral towards a premise sentence. For the latter, the task is to select the correct answer from 3 (v1.0) or 5 (v1.11) answer choices for a question. We use both versions of CoS-E in our experiments (see Appendix A.2). Table 1 contains examples and Table 7 (Appendix A.2) data statistics.³

T5 Models All of the models in this work are based on T5, though our methods can in principle be applied to any architecture. The base version of T5 is a 220M-parameter transformer encoder-decoder (Vaswani et al., 2017). To carry out supervised finetuning, T5 is trained by maximizing the conditional likelihood of the correct text output (from annotated data), given the text input.

We finetune five T5-Base models for each dataset, supervising with ground-truth labels and rationales (further details in Appendix A.3-A.4):

- I→R, which maps task inputs to rationales, without ever being exposed to task outputs.
- R→O, which maps rationales to task outputs. The only input elements this model is exposed to are answer choices (for CoS-E).
- I

 OR, which maps inputs to outputs and rationales.

³CoS-E does not contain test set rationale annotations, so we report performance values on the validation set.

Source	CQA	SNLI	SST	AgNews	Evidence Inference	Movie Reviews	MultiRC	LGD	20 News	Amazon Reviews	Beer Reviews	BoolQ	FEVER
True Pipelines (no gradient flow)													
Camburu et al. (2018)		E + NL											
Kumar and Talukdar (2020)		NL											
Rajani et al. (2019)	E + NL												
Jain et al. (2020)			E	E	E	E	E						
Jacovi and Goldberg (2021)			E	E	E	E	E	E	E	E	E		
DeYoung et al. (2020)	E	E			E	E	E					E	E
Lehman et al. (2019)					E								
Discrete Optimization Variants													
Lei et al. (2016)											E		
Bastings et al. (2019)		E	E								E		
Latcinnik and Berant (2020)	NL												
Paranjape et al. (2020)					Е	Е	Е				Е	Е	Е

Table 2: An overview of text-only datasets and rationale types (E for extractive, NL for natural language rationales) used in prior work on pipeline architectures. We focus on the two tasks we believe require a more complex notion of "reasoning" to solve: CommonsenseQA (CQA) and NLI. Unlike the other tasks in the table, prior work for rationalizing these two tasks lacks consensus on (1) the type of rationales best-suited, and (2) the form of the model for these tasks. We argue for natural language rationales, and demonstrate that pipeline models are poorly-suited for CQA and SNLI given this choice. Dataset citations: Appendix A.1.

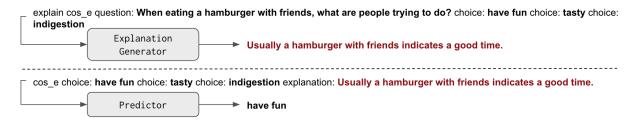


Figure 2: An illustration of a pipeline model (composed of $I \rightarrow R$ and $R \rightarrow O$; §2) for CoS-E v1.0 with a human-written rationale. The dotted line indicates two separate models with no gradient flow.

- IR→O, which maps pairs of inputs and rationales to outputs.
- $I \rightarrow O$, which maps inputs to outputs.

We provide input-output formatting in Table 8 (Appendix A.3). Using these building blocks, we can instantiate two important approaches.

Pipeline Model (I\rightarrowR;R\rightarrowO) This architecture composes I \rightarrow R with R \rightarrow O, each of which is trained entirely separately, for a total of 440M parameters. It is illustrated in Figure 2 and is faithful-by-construction (with caveats; see Jacovi and Goldberg, 2021). The vast majority of prior work using pipelines has focused only on extractive rationales (see Table 2).

Self-Rationalizing Model (I→OR) A joint, self-rationalizing model (Melis and Jaakkola, 2018), illustrated in Figure 3, predicts both a label and rationale. This is the most common approach to free-text rationalization (Hendricks et al., 2016; Kim et al., 2018; Hancock et al., 2018; Camburu et al., 2018; Ehsan et al., 2018; Liu et al., 2019a; Wu and Mooney, 2019; Narang et al., 2020; Do et al., 2020; Tang et al., 2020), but little is understood about model internals. I→OR models are

desirable for their ease-of-use, task-effectiveness, parameter efficiency, and their ability to generate fluent and plausible rationales. We expect models of this kind to play an important role in continuing research on explainable AI for these reasons.

We use the $I\rightarrow OR$ variant of T5 (Narang et al., 2020). Because only one instance of T5 is used to instantiate it, the total number of parameters is half that of the pipeline. We replicate two prior findings (Tables 9–10 in Appendix B): the T5 pipeline does not perform as well as the self-rationalizing model (despite having double the parameters), and T5-Base outperforms pretrained models used in prior work.

Evaluation We do not report BLEU scores (Papineni et al., 2002), because BLEU and related metrics do not measure plausibility (Camburu et al., 2018; Kayser et al., 2021; Clinciu et al., 2021) or faithfulness (Jacovi and Goldberg, 2020). In addition to low correlation with human scores, there can be many valid rationales for a given instance (Miller, 2019); metrics that compare generated rationales to a single ground-truth do not address this and are thus a poor measure of quality.

Human simulatability (Doshi-Velez and Kim, 2017) has a rich history in machine learning interpretability as a reliable measure of rationale quality from the lens of utility to an end-user (Kim et al., 2016; Chandrasekaran et al., 2018; Hase and Bansal, 2020; Yeung et al., 2020; Poursabzi-Sangdeh et al., 2021; Rajagopal et al., 2021, *i.a.*). Rather than computing word-level overlap with a ground-truth explanation, simulatability measures the additional predictive ability towards the predicted label a rationale provides over the input, computed as the difference between task performance when a rationale is given as input vs. when it is not (IR \rightarrow Ô minus I \rightarrow Ô).⁴ Historically, humans have served as the predictors, but recent work has shown that the computation of simulatability can be automated using trained models. Hase et al. (2020) demonstrate that automated metrics for simulatability have moderate to high correlation with human scores in both an expert and a crowdsourced setting. In our experiments, model predictions are often unable to be simulated because they degenerate under high values of noise (§4, A.5). We thus use a variant of this metric that relies on predicting the *gold* labels as our measure of rationale quality: IR \rightarrow O minus I \rightarrow O.⁵ We discuss the effects of this difference in Appendix A.6.

3 Shortcomings of Free-Text Pipelines

We first analyze "faithful-by-construction" pipeline models ($I \rightarrow R; R \rightarrow O$) for free-text rationalization with respect to two properties: quality of generated rationales (§3.1) and appropriateness of the sufficiency assumption (§3.2).

3.1 Joint Model Rationales are More Indicative of Labels

Rationales should be a function of the input *and* the predicted label. To demonstrate why this is the case, consider training an $I \rightarrow R$ model on a dataset with multiple annotation layers, e.g., OntoNotes, that contains word sense, predicate structure, and coreference (Pradhan et al., 2007). Without additional task-specific input, this model would produce the

Source of Rationales	R*	$I \rightarrow R$	I→OR
E-SNLI	97.67 84.84 68.63	89.11	90.52
CoS-E v1.0	84.84	53.47	62.00
CoS-E v1.11	68.63	45.45	53.15

Table 3: Accuracy of the trained $R \rightarrow O$ model evaluated on ground-truth natural language rationales (R^*) and rationales generated from two model architectures: $I \rightarrow OR$ and $I \rightarrow R$ (see §2 for model descriptions).

Source of Rationales	R*	I→R	I→OR
E-SNLI	7.77	-1.63	-0.86
	21.26	-12.11	-6.21
CoS-E v1.11	19.09	-12.77	-6.06

Table 4: Rationale quality scores (§2; higher is better) of ground-truth rationales (R*) and rationales generated from two model architectures: $I\rightarrow OR$ and $I\rightarrow R$. These results demonstrate that rationales generated as a function of the input and the predicted label ($I\rightarrow OR$) are higher quality than those generated as a function of the input alone ($I\rightarrow R$) across datasets (§3.1).

same rationale, regardless of the task being rationalized. Prior work has also critiqued $I \rightarrow R; R \rightarrow O$ models because it is counter-intuitive to generate a rationale before deciding the label to explain (Kumar and Talukdar, 2020; Jacovi and Goldberg, 2021). Therefore, the $I \rightarrow R$ model will first need to implicitly predict a label. But can $I \rightarrow R$ infer the label well, when it is trained without label signal?

To address this question, we study whether $I \rightarrow OR$ rationales are better at predicting the gold labels than $I \rightarrow R$ rationales. We train a $R \rightarrow O$ model on ground-truth rationales (R^*), and evaluate on the following inputs:

- test set ground-truth R* rationales,
- test set rationales generated by I→OR, and
- test set rationales generated by $I \rightarrow R$.

In Table 2, we show that $I\rightarrow OR$ rationales recover 8–9% more ground-truth (R*) performance than R $\rightarrow O$ rationales on both versions of CoS-E, and 1% on E-SNLI. A smaller improvement for E-SNLI could be explained by the fact that E-SNLI has substantially more training examples for each label than CoS-E, which helps a pipeline model learn features predictive of each label. We additionally demonstrate $I\rightarrow OR$ rationales are higher quality than $R\rightarrow O$'s, as measured by our ratio-

⁴The predicted label is from the same system that produced the predicted rationale.

⁵Given the large scale of our analysis (>250K instances evaluated), an automated metric provides coverage, reproducibility and consistency not achievable with human annotation. An author of this paper annotated 60 instances from both E-SNLI and CoS-E v1.11, and found 82.5% agreement between their rationale quality score and the automated metric.

 $^{^6}$ E-SNLI has 549,357 training examples and only 3 labels. In contrast, the number of answer options across all instances in CoS-E v1.0 (v1.11) is 6,387 (12,992), but the training set size is 7,610 (9,741), i.e., \sim 56 (72) times smaller than E-SNLI.

Figure 3: An example of a joint architecture ($I\rightarrow OR$; §2) for CoS-E v1.0 with a human-written rationale. Trained on both task signal and human rationales, these models are effective at generating fluent rationales while retaining good task performance (Table 9-10 in Appendix B).

Model	$R\rightarrow O$ with R^*	IR→O with R*	Δ
E-SNLI	97.67	98.72	+1.05
CoS-E v1.0	84.84	90.42	+5.58
CoS-E v1.11	68.63	80.84	+12.21

Table 5: A comparison of the IR \rightarrow O and R \rightarrow O models (§2) evaluated with ground-truth natural language rationales (R*). In some cases accuracy improves substantially with the addition of the input, indicating that rationales are not always sufficient and pipelines are not always effective.

nale quality metric (Table 4). The fact that the pipeline's strong performance does not generalize to a complex prediction task such as CoS-E empirically demonstrates that training on label signal O is important to generate good-quality rationales and avoid cascading errors.

3.2 Sufficiency is not Universally Valid

"Faithful-by-construction" pipelines rely on the *sufficiency* assumption: the selected rationale must be sufficient to make the prediction without the remaining input. This assumption is suitable for IE tasks for which a subset of the input tokens is predictive of the label. Indeed, humans can serve as $R \rightarrow O$ models on certain IE tasks and make accurate predictions, validating that rationales are sufficient for these tasks (Jain et al., 2020).

To illustrate why sufficiency might not be justified for reasoning tasks, consider the example in Figure 2. The task of the $R \rightarrow O$ model is to select between the answer choices "have fun", "tasty", and "indigestion" given the rationale "Usually a hamburger with friends indicates a good time". The rationale is designed to complement the input question, but the $R \rightarrow O$ model does not see the question, changing the fundamental nature of the task it is solving. We thus wonder: does task obfuscation hurt pipelines' ability to perform the task?

We report the accuracy difference between a $R{\rightarrow}O$ model and a model that receives both the input and rationale (IR \rightarrow O), both trained on R*. We evaluate on test set R*.⁷ In Table 5, the IR \rightarrow O

models on CoS-E have a 5–12% increase in accuracy over R \rightarrow O, indicating that the rationales are not sufficient. The difference is much smaller for E-SNLI (1%), likely due to the fact that E-SNLI was collected by instructing annotators to provide self-contained rationales. However, using dataset-collection to explicitly collect sufficient rationales does not address the unnaturalness of such a task formulation (Wiegreffe and Marasović, 2021). Table 5 indicates that (especially) in the case of CoS-E, sufficiency is not a valid assumption, and the use of $I\rightarrow R;R\rightarrow O$ models is sub-optimal in these cases.

So far, we have highlighted shortcomings of pipelines for reasoning tasks:

- cascading errors caused by low-quality rationales that are not indicative of labels (§3.1),
- missing information due to rationales not being sufficient (§3.2),
- double the number of parameters and more manual labor needed to reach comparable performance to an end-to-end (I→O) model; still often performing worse (§2).

We next turn our focus to self-rationalizing $(I\rightarrow OR)$ models currently in widespread use, which in contrast to pipelines are high-performing, easy to implement via a multi-task loss, and more parameter-efficient (§2).

4 Analyzing Necessary Properties of Joint Models

Despite their popularity and widespread use, the extent to which self-rationalizing models exhibit faithful rationalization has not been studied. To illustrate this point, we reference Narang et al. (2020):

... Much like humans, our approach does not guarantee that the produced explanation actually explains the specific reasons why a model generated its prediction. In other words, the model could potentially just make up a reasonable-sounding

⁷Evaluating on R* instead of generated rationales serves

as an upper-bound on pipeline performance, removing the confounding factor that $I \rightarrow R$ rationales can be poor (§3.1).

⁸Camburu et al. (2018) give an example: the rationale "A woman is not a person" could predict either a contradiction or entailment label depending on the input.

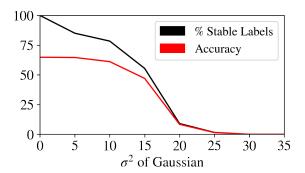


Figure 4: Results of the label portion of the robustness equivalence test for CoS-E v1.0.

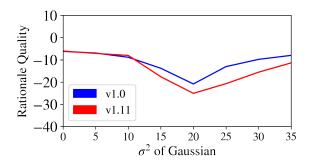


Figure 5: Results of the rationale portion of the robustness equivalence test for both CoS-E datasets.

explanation instead of providing a truly accurate description of its causal decision-making process.

It is not infeasible that a large, overparameterized model trained on both gold-rationale emulation and a labelling task can learn to do both equally well, without having to rely on shared information in its parameters. Therefore, rationales from $I\rightarrow OR$ models cannot be treated as faithful explanations without further investigation.

At minimum, rationales must be implicitly or explicitly tied to the model's prediction. We present two metrics to analyze the association between the mechanisms that produce labels and rationales in a multi-task, $I\rightarrow OR$ model: robustness equivalence (§4.1) and feature importance agreement (§4.2). These experiments serve as a necessary sanity check for the reliability of $I\rightarrow OR$ models' explanations.

4.1 Robustness Equivalence

We aim to analyze whether predicted labels and rationales are similarly or dissimilarly robust to noise applied to the input. The former indicates predicted labels and rationales are strongly associated, while the latter indicates the opposite. Given some amount of noise, there are four possi-

ble cases for a model's output: $\{l_{\text{stable}}, l_{\text{unstable}}\} \times \{r_{\text{stable}}, r_{\text{unstable}}\}$, where l is a label and r is a rationale.

The case where l and r are both stable or both unstable indicates that both tasks are similarly affected by noise. The case where l is unstable but r is stable (or vice versa) is a failure case—if only one output is stable, we conclude the two generation mechanisms cannot be strongly associated within the model.

Method Following related work (Wang et al., 2019; Lakshmi Narayan et al., 2019; Liu et al., 2019b), we add zero-mean Gaussian noise $\mathcal{N}(0,\sigma^2)$ to each input embedding in the I \rightarrow OR encoder at inference time. We measure changes in label prediction as the number of predicted test set labels that flip, i.e., change from their original prediction to something else, alongside changes in accuracy of the I \rightarrow OR model. We measure changes in rationale quality using our rationale quality metric (§2), with details of the metric calculation illustrated in Figure 6a. We report metrics on rationales generated by I \rightarrow OR under different levels of noise, controlled by σ^2 .

An example of noisy outputs for the running CoS-E v1.0 example is presented in Table 6.

Results We present results on the effect of noise on labels in Figure 4 (E-SNLI and CoS-E v1.11 in Figure 8 of Appendix B). As expected, the accuracy of the $I\rightarrow OR$ model (red line) and the percent of labels in the $I\rightarrow OR$ model which have not flipped (black line) are almost identical for all three datasets. We present results on the effect of noise on rationales in Figure 5 for CoS-E (E-SNLI in Figure 9 of Appendix B).

By examining the regions of largest slope, we gain insights into model behavior. On the rationale quality measure, both versions of CoS-E's rationales reach a minimum contribution to task accuracy at $\sigma^2 = 20$ (Figure 5). We similarly observe the largest drop in task accuracy (Figure 4) for CoS-E v1.0 between $\sigma^2 = 15$ to $\sigma^2 = 20.9$ Thus, at lower noise levels (0–15), the model exhibits both stable labels and rationales, and at higher levels (20+), both unstable, indicating robustness equivalence. Similar conclusions can be reached for E-SNLI and CoS-E v1.11; we conclude that the I \rightarrow OR model demonstrates high label-rationale as-

 $^{^9}$ See Appendix A.5 for a note on why the I→OR models achieve worse-than-random accuracy at high values of σ^2 .

σ^2	Predicted Output
0	have fun explanation: having fun is the only thing that people are trying to do.
5	have fun explanation : having fun is the only thing that people are trying to do.
10	have fun explanation : eating a hamburger with friends is fun.
15	have fun explanation : having fun is the only thing people are trying to do while eating a hamburger with friends.
20	<pre><extra_id_0> a hamburger with friends<extra_id_1> are people trying to do? explanation: a hamburger is a hamburger</extra_id_1></extra_id_0></pre>
25	" is the only thing that is """""""""""""""""""""""""""""""""""
30	a indigestion:
35	a ``````````

Table 6: Output of the $I\rightarrow OR$ model for the running CoS-E v1.0 example under differing noise levels. While the rationale changes from variance 0-15, it is still valid for the given (correct) predicted label. At a variance of 20 and beyond, the model fails to generate both the correct label and a valid rationale. The model's predictions for this instance exhibit robustness equivalence.

$$I + \epsilon_{\sigma} \rightarrow OR$$
 $acc(IR \rightarrow O) - acc(I \rightarrow O)$
(a) In §4.1. ϵ_{σ} is the noise.

$$I \setminus I_a \to OR$$

$$acc(IR \to O) - acc(I \to O)$$

(b) In $\S4.2$. I_a are input tokens selected by an attribution method.

Figure 6: An illustration of how rationale quality is calculated in §4.

sociation for all 3 datasets as measured by robustness equivalence.

We expect that rationale quality in Figure 5 does not monotonically decrease because as rationales continue to worsen in quality (see the example in Table 6), the IR \rightarrow O model may ignore them completely and more closely emulate the I \rightarrow O model. For example, in Table 6, at $\sigma^2=30$, the rationale provides signal for an incorrect answer choice ("indigestion") that does not exist at $\sigma^2=35$. Therefore, we consider rationales produced with σ larger than the value at which the metric reaches the minimum as unstable (see further examples corroborating this in Appendix B, Tables 11–13).

4.2 Feature Importance Agreement

If label prediction and rationale generation are associated, input tokens important for label prediction should be important for rationale generation and vice versa. We refer to this property as *feature importance agreement*. To measure to what ex-

tent I→OR models exhibit this property, we use gradient-based attribution (Baehrens et al., 2010; Simonyan et al., 2014) to identify tokens important for label prediction, and the **Remove and Retrain** (**ROAR**) occlusion method (Hooker et al., 2019) to analyze their impact on rationale generation (or vice versa).

Gradient Attribution For a predicted class p, gradient attribution is a function of the gradient of the predicted class' logit l_p with respect to an input token embedding $\mathbf{x}^{(i)} \in \mathbb{R}^d$:

$$a(\mathbf{x}^{(i)}; l_p) = f(\nabla_{\mathbf{x}^{(i)}} l_p) \in \mathbb{R},$$
 (1)

where the function f reduces the gradient to a scalar. Choices for f include L_1 or L_2 norm (Atanasova et al., 2020), or an element-wise sum (Wallace et al., 2019). Intuitively, the gradient measures how much an infinitesimally small change in the input changes the predicted class' logit, using a first-order Taylor series approximation of the logit function. Such methods have been extended to sequence-output models such as neural machine translation (He et al., 2019; Ding et al., 2019; Li et al., 2020) by computing the sum of m decoded logits $\{l_p^{(k)}\}_{k=1}^m$ with respect to the input:

$$a(\mathbf{x}^{(i)}; \{l_p^{(k)}\}_k) = \sum_{k=1}^m a(\mathbf{x}^{(i)}; l_p^{(k)}) \in \mathbb{R}.$$
 (2)

The attribution of a sequence of n input token embeddings, $\mathbf{X} \in \mathbb{R}^{n \times d}$, is a vector $a(\mathbf{X}) = [a(\mathbf{x}^{(1)}), \dots, a(\mathbf{x}^{(n)})] \in \mathbb{R}^n$, where $a(\mathbf{x}^{(i)})$ is shorthand for the value defined in Equation 2.

By decomposing the term in Equation 2 into two parts, we obtain two attribution vectors over the input tokens; one for the predicted label logits \mathcal{L} , and one for the predicted rationale logits \mathcal{R} in the decoded output:

$$a(\mathbf{x}^{(i)}; \{l_p^{(k)}\}_k) = \sum_{k \in \mathcal{L}} a(\mathbf{x}^{(i)}; l_p^{(k)}) + \sum_{k \in \mathcal{R}} a(\mathbf{x}^{(i)}; l_p^{(k)}),$$

$$a(\mathbf{X}) = a(\mathbf{X})_{\mathcal{L}} + a(\mathbf{X})_{\mathcal{R}}.$$
 (3)

Reliability of Gradient Attribution Before we measure feature importance agreement, it is critical to evaluate whether the gradient-attribution scores truly capture token importance, since these methods can be unreliable for certain datasets or architectures (Kindermans et al., 2019). To validate that our attributions are reliable, we perform the ROAR test (Hooker et al., 2019). Using attribution scores, we obtain the top-k% attributed tokens for every instance and occlude them following T5's pretraining procedure and mask tokens. We retrain a model on the occluded training set and evaluate on the occluded test set. We repeat this procedure for $k \in \{10\%, 20\%, 30\%\}$, and compare the drop in performance as k increases to a baseline in which k% random tokens are dropped. To the extent that the occluded model fails to match the random model's performance, we can attribute such degradation to the removal of tokens that the original model finds informative. A large drop in performance indicates that gradient attributions successfully identify important tokens in the input.

We first use this method to select an optimal gradient-attribution method and f function (Figure 11 in Appendix B). We find the L_1 norm of the embedding vector as f to outperform the element-wise sum (which may suffer from dampened magnitudes). Unlike prior work in computer vision (Hooker et al., 2019), we find raw gradients to perform comparably to the input*gradient variant (Shrikumar et al., 2017). Thus, we compute attributions in subsequent experiments following Equation 1 with f equal to the L_1 norm.

We validate that attributions from the label logits, $a(\mathbf{X})_{\mathcal{L}}$, degrade label accuracy when compared to random occlusion (orange vs. blue line in Figure 7, left). The two rationale quality lines (Figure 7, right) for CoS-E v1.0 have an inflection point. We illustrate how the metric is calculated in Figure 6b. Similar to §4.1, we expect this is due to rationales so noisy that IR \rightarrow O ignores them and behaves like I \rightarrow O. If an input attribution degrades rationale quality more than a random attribution, then the line corresponding to that attribution (for values of k for which neither that attribution nor the random attribution have reached the inflection

point) has to be below the "random" line. For values of k for which both attributions have passed the inflection point, the "random" line should be below the attribution line, assuming that after this point, a noisier rationale is more similar is IR \rightarrow O to I \rightarrow O and hence the score is closer to 0. Both criteria hold for attributions from the rationale logits, $a(\mathbf{X})_{\mathcal{R}}$ (green vs. blue line in Figure 7, right) for CoS-E v1.0 and other datasets (see Figure 12 in Appendix B). This reliability check confirms that gradient-attribution works well in our setting.

Agreement Method and Results To measure feature importance agreement—whether tokens important for label prediction are important for rationale generation (and vice versa)—we repeat the same experiment, but measure performance with respect to the *other* output's metric. For attributions computed from label logits, $a(\mathbf{X})_{\mathcal{L}}$, we measure the effect of their occlusion on rationale quality using the rationale quality score. For attributions with respect to rationale logits, $a(\mathbf{X})_{\mathcal{R}}$, we measure the effect of their occlusion on label accuracy. If at least one of these values is notably different from random, we can conclude that the $\mathbf{I} \rightarrow \mathbf{OR}$ model displays feature-importance similarity in a given direction.

Results for CoS-E v1.0 are once again in Figure 7 (and for other datasets in Figure 12 in Appendix B). In Figure 7 (left), we find that removing top-k%tokens by $a(\mathbf{X})_{\mathcal{R}}$ magnitude degrades label performance compared to the baseline (green vs. blue line). Intuitively, this drop is less than token attributions by $a(\mathbf{X})_{\mathcal{L}}$ magnitude (orange line). In Figure 7 (right), we observe that removing top-k% tokens by $a(\mathbf{X})_{\mathcal{L}}$ consistently degrades rationale performance more than random according to the two criteria for comparing the rationale quality lines (orange vs. blue line). This also holds for E-SNLI and CoS-E v1.11 (Figure 12). We conclude that the I→OR model demonstrates label-rationale association as measured by feature importance agreement for the datasets studied.

5 Related Work

Analysis of NLP Models Structural tests for analyzing models' internals include probing (Tenney et al., 2019) and attention analysis (Jain and Wallace, 2019; Serrano and Smith, 2019; Wiegreffe and Pinter, 2019; Tutek and Snajder, 2020). These, along with behavioral tests such as challenge sets (McCoy et al., 2019) and checklists (Ribeiro et al.,

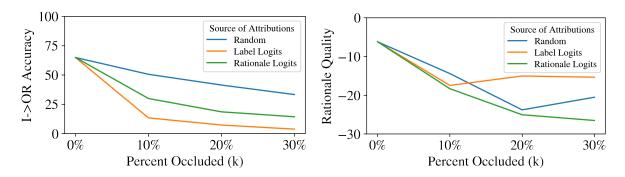


Figure 7: Performance of $I \rightarrow OR$ models trained with the ROAR method on CoS-E v1.0. **Left:** Impact of occlusion by source of attribution on label accuracy. **Right:** Impact of occlusion by source of attribution on rationale quality.

2020), are conceptually similar to our experiments, but study different model properties.

Although gradient-attribution has been extensively studied in NLP, its interplay with free-text rationalization has not. Wu and Mooney (2019) use feature importance agreement to train the explanation module of a VQA model. To the best of our knowledge, we are the first to evaluate gradient-attribution reliability for NLP tasks with the ROAR test (Hooker et al., 2019).

Robustness Analysis Robustness of post-hoc extractive interpretability methods has been studied (Kindermans et al., 2019; Ghorbani et al., 2019; Heo et al., 2019; Zheng et al., 2019; Slack et al., 2020). Zhang et al. (2020) show that saliency maps and model predictions can be independently adversarially attacked in vision and clinical tasks, and conclude this is due to a misalignment between the saliency map generator and model predictor. Such methods have not been tested for models producing natural language (NL) rationales. Future work could include expanding robustness equivalence (§4.1) to model discrete edits of input words.

Analyzing Faithfulness The aim of our work is to initiate placing models that provide NL rationales on the faithfulness spectrum conceptualized by Jacovi and Goldberg (2020). Prior work proposing models (Jain et al., 2020; Schuff et al., 2020; Jacovi and Goldberg, 2021) and evaluations (DeYoung et al., 2020) of faithful explanation focus on extractive rationales and generally rely on the sufficiency assumption. Schuff et al. (2020) propose a regularization term to couple answers and extractive explanations on HotPotQA.

Turning to exceptions that focus on natural language rationales, Latcinnik and Berant (2020) train a differentiable $I \rightarrow R; IR \rightarrow O$ pipeline for Common-

senseQA, controlling the complexity of the IR→O model to increase the likelihood that the model is faithful to the rationale. Kumar and Talukdar (2020) propose an IO→R;IR→O pipeline that generates an explanation for every possible NLI label using label-specific explanation generators—an alternative solution to the problem raised in §3.1 for datasets with a small number of shared labels.

6 Conclusion

After demonstrating the weaknesses that pipeline models exhibit for free-text rationalization tasks, we propose two measurements of label-rationale association in self-rationalizing models. We find that on three free-text rationalization datasets for CommonsenseQA and SNLI, models based on T5 exhibit high robustness equivalence and feature importance agreement, demonstrating that they pass a necessary sanity check for generating faithful free-text rationales.

Future work can expand analysis to more properties. We believe this research direction to be important moving forward due to the advantages of large multi-task explanation models, and as a complement to development of interpretable architectures that can be fickle and task-specific. Although our measurements address only necessary and not sufficient properties, by viewing faithful interpretability as a spectrum, we make a step to quantitatively situate common models on it.

Acknowledgements

We are grateful to Jonathan Berant, Peter Hase, Alon Jacovi, Yuval Pinter, Mark Riedl, Vered Shwartz, Ian Stewart, Swabha Swayamdipta, and Byron Wallace for feedback on the draft. We thank members of the AllenNLP team at the Allen Institute for Artificial Intelligence (AI2), members of the Entertainment Intelligence lab at Georgia Tech, and reviewers for valuable feedback and discussions. We thank Aaron Chan for pointing out an issue (now corrected) in our use of the term "simulatability". This work was done while SW was an intern at AI2.

References

- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. Neural module networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. A diagnostic study of explainability techniques for text classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3256–3274, Online. Association for Computational Linguistics.
- David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. 2010. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11:1803–1831.
- Jasmijn Bastings, Wilker Aziz, and Ivan Titov. 2019. Interpretable neural predictions with differentiable binary variables. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2963–2977, Florence, Italy. Association for Computational Linguistics.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-SNLI: Natural language inference with natural language explanations. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Arjun Chandrasekaran, Viraj Prabhu, Deshraj Yadav, Prithvijit Chattopadhyay, and Devi Parikh. 2018. Do explanations make VQA models more predictable to a human? In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1036–1042, Brussels, Belgium. Association for Computational Linguistics.
- 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 *Proceedings of the 2019 Conference of the North American*

- Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.
- Miruna-Adriana Clinciu, Arash Eshghi, and Helen Hastie. 2021. A study of automatic metrics for the evaluation of natural language explanations. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2376–2387, Online. Association for Computational Linguistics.
- Gianna M Del Corso, Antonio Gulli, and Francesco Romani. 2005. Ranking a stream of news. In Proceedings of the 14th international conference on World Wide Web.
- Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C. Wallace. 2020. ERASER: A benchmark to evaluate rationalized NLP models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4443–4458, Online. Association for Computational Linguistics.
- Shuoyang Ding, Hainan Xu, and Philipp Koehn. 2019. Saliency-driven word alignment interpretation for neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 1–12, Florence, Italy. Association for Computational Linguistics.
- Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. 2020. e-SNLI-VE-2.0: Corrected visual-textual entailment with natural language explanations. In *IEEE CVPR Workshop on Fair, Data Efficient and Trusted Computer Vision*.
- Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning.
- Upol Ehsan, Brent Harrison, Larry Chan, and Mark O Riedl. 2018. Rationalization: A neural machine translation approach to generating natural language explanations. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society.*
- Amirata Ghorbani, Abubakar Abid, and James Zou. 2019. Interpretation of neural networks is fragile. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Braden Hancock, Paroma Varma, Stephanie Wang, Martin Bringmann, Percy Liang, and Christopher Ré. 2018. Training classifiers with natural language explanations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1884–1895, Melbourne, Australia. Association for Computational Linguistics.
- Peter Hase and Mohit Bansal. 2020. Evaluating explainable AI: Which algorithmic explanations help users predict model behavior? In *Proceedings of the*

- 58th Annual Meeting of the Association for Computational Linguistics, pages 5540–5552, Online. Association for Computational Linguistics.
- Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal. 2020. Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language? In *Findings of the Association for Computational Linguistics: EMNLP* 2020, pages 4351–4367, Online. Association for Computational Linguistics.
- Shilin He, Zhaopeng Tu, Xing Wang, Longyue Wang, Michael Lyu, and Shuming Shi. 2019. Towards understanding neural machine translation with word importance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 953–962, Hong Kong, China. Association for Computational Linguistics.
- Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. 2016. Generating visual explanations. In *European Conference on Computer Vision (ECCV)*.
- Juyeon Heo, Sunghwan Joo, and Taesup Moon. 2019. Fooling neural network interpretations via adversarial model manipulation. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. 2019. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4198–4205, Online. Association for Computational Linguistics.
- Alon Jacovi and Yoav Goldberg. 2021. Aligning faithful interpretations with their social attribution. *Transactions of the Association for Computational Linguistics*.
- Sarthak Jain and Byron C. Wallace. 2019. Attention is not Explanation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3543–3556, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sarthak Jain, Sarah Wiegreffe, Yuval Pinter, and Byron C. Wallace. 2020. Learning to faithfully rationalize by construction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4459–4473, Online. Association for Computational Linguistics.

- Yichen Jiang, Nitish Joshi, Yen-Chun Chen, and Mohit Bansal. 2019. Explore, propose, and assemble: An interpretable model for multi-hop reading comprehension. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2714–2725, Florence, Italy. Association for Computational Linguistics.
- Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-vil: A dataset and benchmark for natural language explanations in vision-language tasks.
- 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 the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262, New Orleans, Louisiana. Association for Computational Linguistics.
- Been Kim, Rajiv Khanna, and Oluwasanmi Koyejo. 2016. Examples are not enough, learn to criticize! criticism for interpretability. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 2288–2296.
- Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. 2018. Textual explanations for self-driving vehicles. In *Proceedings of the European conference on computer vision (ECCV)*.
- Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T. Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. 2019. *The (Un)reliability of Saliency Methods*, pages 267–280. Springer International Publishing.
- Sawan Kumar and Partha Talukdar. 2020. NILE: Natural language inference with faithful natural language explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8730–8742, Online. Association for Computational Linguistics.
- Pooja Lakshmi Narayan, Ajay Nagesh, and Mihai Surdeanu. 2019. Exploration of noise strategies in semi-supervised named entity classification. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics* (*SEM 2019), pages 186–191, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *Machine Learning Proceedings* 1995.
- Veronica Latcinnik and Jonathan Berant. 2020. Explaining question answering models through text generation. arXiv:2004.05569.

- Eric Lehman, Jay DeYoung, Regina Barzilay, and Byron C. Wallace. 2019. Inferring which medical treatments work from reports of clinical trials. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3705–3717, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2016. Rationalizing neural predictions. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 107–117, Austin, Texas. Association for Computational Linguistics.
- Jierui Li, Lemao Liu, Huayang Li, Guanlin Li, Guoping Huang, and Shuming Shi. 2020. Evaluating explanation methods for neural machine translation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 365–375, Online. Association for Computational Linguistics.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of lstms to learn syntax-sensitive dependencies. *Transactions of the Association for Computational Linguistics (TACL)*, 4:521–535.
- Hui Liu, Qingyu Yin, and William Yang Wang. 2019a. Towards explainable NLP: A generative explanation framework for text classification. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5570–5581, Florence, Italy. Association for Computational Linguistics.
- Zihan Liu, Jamin Shin, Yan Xu, Genta Indra Winata, Peng Xu, Andrea Madotto, and Pascale Fung. 2019b. Zero-shot cross-lingual dialogue systems with transferable latent variables. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1297–1303, Hong Kong, China. Association for Computational Linguistics.
- Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: Understanding rating dimensions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*.
- Julian McAuley, Jure Leskovec, and Dan Jurafsky. 2012. Learning attitudes and attributes from multiaspect reviews. In 2012 IEEE 12th International Conference on Data Mining.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.

- David Alvarez Melis and Tommi Jaakkola. 2018. Towards robust interpretability with self-explaining neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38.
- Christoph Molnar. 2019. *Interpretable Machine Learning*. https://christophm.github.io/interpretable-ml-book/.
- Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan. 2020. WT5?! training text-to-text models to explain their predictions. arXiv:2004.14546.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of* the 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Bhargavi Paranjape, Mandar Joshi, John Thickstun, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. An information bottleneck approach for controlling conciseness in rationale extraction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1938–1952, Online. Association for Computational Linguistics.
- Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach. 2021. Manipulating and measuring model interpretability. In *Conference on Human Factors in Computing Systems (CHI)*.
- Sameer S Pradhan, Eduard Hovy, Mitch Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2007. Ontonotes: A unified relational semantic representation. In *International Conference on Semantic Computing (ICSC)*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Dheeraj Rajagopal, Vidhisha Balachandran, Eduard Hovy, and Yulia Tsvetkov. 2021. Selfexplain: A self-explaining architecture for neural text classifiers
- Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4932–4942, Florence, Italy. Association for Computational Linguistics.

- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Hendrik Schuff, Heike Adel, and Ngoc Thang Vu. 2020. F1 is Not Enough! Models and Evaluation Towards User-Centered Explainable Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7076–7095, Online. Association for Computational Linguistics.
- Roy Schwartz, Sam Thomson, and Noah A. Smith. 2018. Bridging CNNs, RNNs, and weighted finite-state machines. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 295–305, Melbourne, Australia. Association for Computational Linguistics.
- Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2931–2951, Florence, Italy. Association for Computational Linguistics.
- Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, pages 3145–3153.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Deep inside convolutional networks: Visualising image classification models and saliency maps. In Workshop at International Conference on Learning Representations (ICLR Workshop).
- Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. 2020. Fooling LIME and SHAP: Adversarial attacks on post hoc explanation methods. In *Proceedings of the 2020 AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.

- Zheng Tang, Gus Hahn-Powell, and Mihai Surdeanu. 2020. Exploring interpretability in event extraction: Multitask learning of a neural event classifier and an explanation decoder. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 169–175, Online. Association for Computational Linguistics.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- 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 Snajder. 2020. Staying true to your word: (how) can attention become explanation? In *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 131–142, Online. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Eric Wallace, Jens Tuyls, Junlin Wang, Sanjay Subramanian, Matt Gardner, and Sameer Singh. 2019. AllenNLP interpret: A framework for explaining predictions of NLP models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations*, pages 7–12, Hong Kong, China. Association for Computational Linguistics.
- Dilin Wang, Chengyue Gong, and Qiang Liu. 2019. Improving neural language modeling via adversarial training. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*,.
- Sarah Wiegreffe and Ana Marasović. 2021. Teach me to explain: A review of datasets for explainable natural language processing. In *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks*.
- Sarah Wiegreffe and Yuval Pinter. 2019. Attention is not not explanation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*

Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 11–20, Hong Kong, China. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.

Jialin Wu and Raymond Mooney. 2019. Faithful multimodal explanation for visual question answering. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 103–112, Florence, Italy. Association for Computational Linguistics.

Arnold Yeung, Shalmali Joshi, Joseph Jay Williams, and Frank Rudzicz. 2020. Sequential explanations with mental model-based policies. In *ICML Workshop on Human Interpretability in Machine Learning*.

Omar Zaidan and Jason Eisner. 2008. Modeling annotators: A generative approach to learning from annotator rationales. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 31–40, Honolulu, Hawaii. Association for Computational Linguistics.

Xinyang Zhang, Ningfei Wang, Hua Shen, Shouling Ji, Xiapu Luo, and Ting Wang. 2020. Interpretable deep learning under fire. In 29th USENIX Security Symposium (USENIX Security).

Haizhong Zheng, Earlence Fernandes, and Atul Prakash. 2019. Analyzing the interpretability robustness of self-explaining models. In *ICML* 2019 Security and Privacy of Machine Learning Workshop.

A Additional Information

A.1 Overview of Prior Work on Pipelines

In Table 2, we overview the datasets and types of rationales used in prior work on pipelines. The sources of datasets are: CommonsenseQA (Talmor et al., 2019), SNLI (Bowman et al., 2015), SST (Socher et al., 2013), AgNews (Del Corso et al., 2005), Evidence Inference (Lehman et al., 2019), Movie Reviews (Zaidan and Eisner, 2008), MultiRC (Khashabi et al., 2018), LGD (Linzen et al., 2016), 20 News (Lang, 1995), Amazon Reviews (McAuley and Leskovec, 2013), Beer Reviews (McAuley et al., 2012), BoolQ (Clark et al., 2019), FEVER (Thorne et al., 2018).

A.2 Details of Datasets

We summarize dataset statistics in Table 7. The two versions (v1.0, v1.11) of CoS-E correspond to the first and second versions of the CommonsenseQA dataset. CoS-E v1.11 has some noise in its annotations (Narang et al., 2020). This is our primary motivation for reporting on v1.0 as well, which we observe does not have these issues.

A.3 Details of T5

The T5 model (Raffel et al., 2020) is pretrained on a multi-task mixture of unsupervised and supervised tasks, including machine translation, question answering, abstractive summarization, and text classification. Its inputs and outputs to every task are text sequences; we provide the input-output formatting for training and decoding of our T5 models in Table 8. T5 can provide any word in the vocabulary as an answer.

A.4 Implementation Details

We use Huggingface Datasets¹¹ to access all datasets, and Huggingface Transformers (Wolf et al., 2020) to access pretrained T5 weights and tokenizer. To optimize, we use Adam with $\epsilon = 1e-8$, $\beta_1 = 0.9$, and $\beta_2 = 0.99$. We use gradient clipping to a maximum norm of 1.0 and a dropout rate of 0.1. We train each model on a NVIDIA RTX 8000 GPU (48GB memory) for maximum 200 epochs with a batch size of 64 and a learning rate linearly decaying from 5e-5. Training ends if the validation set loss has not decreased for 10 epochs. Early stopping occurs within 15 epochs for most models. Most CoS-E models train in less than 1 hour and most E-SNLI models in around 30. At inference time, we greedy-decode until an EOS token is generated (or for 200 tokens). Approximating the 64-batch model with a batch-size of 16 and 4 gradient accumulation steps on 8GB memory cloud GPUs, we sweep starting learning rates of 1e-2, 1e-3, 1e-4, 5e-5, and 1e-5. The two largest learning rates never result in good performance. Among the smallest three rates, performance across all model variants ($I \rightarrow R$, $I \rightarrow OR$, $R \rightarrow O$, $I \rightarrow O$, and $IR \rightarrow O$) on E-SNLI and CoS-E v1.0 never varies by more than 1.58% accuracy or 0.34 BLEU.

¹⁰https://github.com/salesforce/cos-e/
issues/2

[&]quot;"https://huggingface.co/docs/datasets/
master/#

A.5 Note on Robustness Equivalence Convergence

Worst-case model performance under large noise values in the Robustness Equivalence experiments (Figures 4 and 8) reaches 0 rather than random accuracy due to structure of the models' output. The $I\rightarrow OR$ model is trained to produce a delimiter to distinguish the label from the rationale in a long string of output tokens. When it fails to produce the delimiter under high noise, we cannot delineate the label from the rationale in the output where multiple answer choices are often mentioned, so we mark the label as incorrect.

A.6 Further Discussion of Rationale Quality Metric

Traditional simulatability is often considered to be lower-bounded at 0, assuming model-predicted explanations are consistent with model-predicted labels, because a model-predicted explanation should not provide negative utility when given as input to a simulator predicting that label, unless the explanation is inconsistent (i.e., it explains a different label). Our rationale quality metric does not have this property, since model-predicted explanations that explain an incorrect label may provide negative utility toward predicting the gold label. As reported in the paper, it is commonly negative in our experiments.

The limitation of our rationale quality metric is that low scores may be due to poor-quality rationales (what we aim to measure) or poor label prediction performance of the model generating the rationales, assuming consistent predicted labels and rationales. Future work may focus on a robust version of simulatability that can both separate these confounders *and* be computed when model predictions are noisy or ill-defined.

B Additional Results

We provide additional results that supplement the main body of the paper:

- Table 9 presents results comparing the selfrationalizing T5 model to baselines.
- Table 10 presents results comparing the selfrationalizing T5 model to its pipeline variant (from §2).
- Figure 8 presents robustness equivalence label results for E-SNLI and CoS-E v1.11 (§4.1).
- Figure 9 presents robustness equivalence rationale results for E-SNLI (§4.1). CoS-E v1.0

- and v1.11 rationale results are included in Figure 5 in main paper.
- Figure 10 presents an example of L_1 -normalized gradient attributions for a single instance (§4.2).
- Figure 11 presents a comparison of attribution methods on the ROAR reliability check (§4.2) for CoS-E v1.0.
- Figure 12 presents results of the ROAR feature importance agreement measure (§4.2) for CoS-E v1.11 and E-SNLI.
- Tables 11-13 contain additional validation set examples (non-cherry-picked) of noised inputs for CoS-E v1.0 (§4.1).

Dataset	Num. Instances Input Length		Extractiv	e Rationale	Natural Language Rationale		
	Train-Val-Test	# Tokens	# Tokens	% of doc.	# Tokens	% of doc.	
E-SNLI	549,367-9,842-9,824	20.27 +/- 6.95	4.01 +/- 3.01	21.30 +/- 15.82	12.39 +/- 6.43	65.67 +/- 35.46	
CoS-E v1.0	7,610-950-none	13.69 +/- 5.97	4.57 +/- 4.16	35.36 +/- 27.22	12.74 +/- 6.99	108.18 +/- 77.26	
CoS-E v1.11	9,741-1,221-none	13.40 +/- 5.77	6.80 +/- 5.79	53.24 +/- 36.05	6.97 +/- 4.14	58.01 +/- 39.60	

Table 7: Statistics on datasets with ground-truth rationales. Results are presented as mean (one standard deviation) on the training set. CoS-E does not contain test set rationale annotations.

The $I \rightarrow OR$ and $I \rightarrow R; R \rightarrow O$ rationale generator's inputs:

explain cos_e question: [question] choice: [choice_0] choice: [choice_1] choice: [choice_2] explain nli hypothesis: [hypothesis] premise: [premise].

The $I \rightarrow R; R \rightarrow O$ pipeline label predictor's input:

cos_e choice: [choice_0] choice: [choice_1] choice: [choice_2] explanation: [free-text rationale]

nli explanation: [free-text rationale]

The $I\rightarrow OR$ models' outputs are trained and decoded as:

[label] explanation: [free-text rationale]

Table 8: T5 input-output formatting.

Dataset	Random Accuracy	Majority-Vote Accuracy	I→O (T5)	I→OR	Δ
E-SNLI	33.33	33.39	90.95 (84.01 [‡])	$90.81 (83.96^{\ddagger}, 90.9^{\dagger})$	-0.14 (-0.05 [‡])
CoS-E v1.0	33.33	33.75	69.16 (63.8§)	64.84	-4.32
CoS-E v1.11	20.0	20.31	61.75	55.61 (59.4 [†])	-6.14

Table 9: Label accuracy on baseline $I\rightarrow O$ T5 models versus their rationalizing $I\rightarrow OR$ variants fine-tuned for each dataset. We observe that adding rationalization results in some loss in accuracy. We also validate that T5-Base models outperform other architectures. Source of prior results in parentheses: \dagger Narang et al. (2020) using T5, \ddagger Camburu et al. (2018) using a bi-directional LSTM, and \S Rajani et al. (2019) using BERT.

Dataset	I→OR	$I \rightarrow R; R \rightarrow O$	$ \Delta $
E-SNLI	90.81 (83.96 [‡])	89.11 (81.71 [‡])	-1.70 (-2.25 [‡])
CoS-E v1.0	64.84	53.47	-11.37
CoS-E v1.11	55.61	45.45	-10.16

Table 10: Label accuracy on the joint self-rationalizing model $I \rightarrow OR$ compared to a pipeline using natural language rationales. We observe that $I \rightarrow OR$ models have stronger task performance. Source of prior results in parentheses: \ddagger Camburu et al. (2018) using bi-directional LSTMs.

σ^2	Predicted Output
0	house explanation: a house is the only place that would have air conditioning.
5	house explanation: a house is the only place that would have air conditioning.
10	house explanation: a house is the only place that would have air conditioning.
15	<pre><extra_id_0> house explanation: a house is the only place that will have air conditioning.</extra_id_0></pre>
20	<extra_id_0> movie theatre explanation: movie theatre is the only option that is not a movie. 911 911 911</extra_id_0>
25	<pre><extra_id_0> explain<extra_id_1> explain<extra_id_2> explain<extra_id_3> explain<extra_id_4> explain<extra_id_5> movie theatre<extra_id_6></extra_id_6></extra_id_5></extra_id_4></extra_id_3></extra_id_2></extra_id_1></extra_id_0></pre>
30	house of house
35	house of house of office

Table 11: Noised output of the $I\rightarrow OR$ model for the CoS-E v1.0 example "A man wants air conditioning while we watches the game on Saturday, where will it likely be installed?" The correct answer is "house".

σ^2	Predicted Output
0	stress explanation: a computer is used to communicate with a granddaughter.
5	stress explanation: a computer is used to communicate with a granddaughter.
10	stress explanation: a computer is used to talk to people.
15	stress explanation: a computer is used to talk to people.
20	<pre><extra_id_0> is using a computer to<extra_id_1> to talk to<extra_id_2> is using a computer to talk to a person is using a computer to talk to a person</extra_id_2></extra_id_1></extra_id_0></pre>
25	<pre><extra_id_0> answer: answ</extra_id_0></pre>
30	<pre><extra_id_0> answer: answ</extra_id_0></pre>
35	office of the of

Table 12: Noised output of the $I\rightarrow OR$ model for the CoS-E v1.0 example "If a person is using a computer to talk to their granddaughter, what might the computer cause for them?" The correct answer is "happiness".

σ^2	Predicted Output
0	transfer of information explanation: transfer of information is the only option that would be appropriate when communicating with a boss.
5	transfer of information explanation: transfer of information is the only option that would be appropriate when communicating with a boss.
10	transfer of information explanation: transfer of information is the only option that would be appropriate when communicating with a boss.
15	transfer of information explanation: transfer of information is the only option that would be appropriate when communicating with my boss.
20	transfer of information: transfer of information is the only thing that is transfer of information. transfer of information is the only thing that is
25	transfer of information: transfer of information is information. transfer of information. transfer of information is information
30	i believe that is the answer of the question. argument is the answer of the question. argument is the answer of the question
35	i can't handle the argument argument argument is the answer of the argument argument is the answer of the argument

Table 13: Noised output of the $I\rightarrow OR$ model for the CoS-E v1.0 example "When communicating with my boss, what should I do?". The correct answer is "transfer of information".

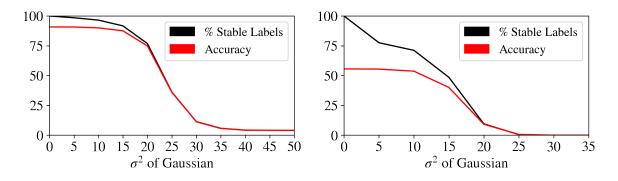


Figure 8: Results of the label portion of the robustness equivalence test for E-SNLI (left) and CoS-E v1.11 (right). Accuracy of the I \rightarrow OR model (red) and % stable labels in the I \rightarrow OR model (black) show that most changes take place in the 10-20 σ^2 range for CoS-E and 15-30 σ^2 for E-SNLI. See §4.1.

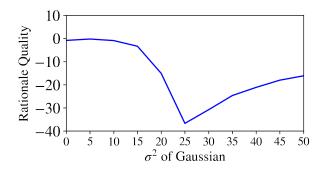


Figure 9: Results of the rationale portion of the robustness equivalence test for E-SNLI.

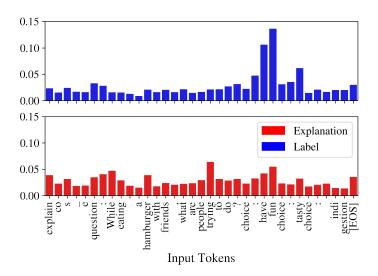


Figure 10: Normalized attributions for the running CoS-E v1.0 example in Figures 2–3. The decoded label is "have fun" and generated rationale "having fun is the only thing that people are trying to do". Important input terms vary across the two loss terms. For example, the predicted label term assigns high importance to the predicted answer choice, "have fun", while the explanation attends more uniformly over the input with peaks on relevant entities and verbs such as "trying". In this example, the explanation- and label-attribution vectors are each L_1 -normalized in order to compare the relative importance of tokens (irrespective of gradient magnitudes).

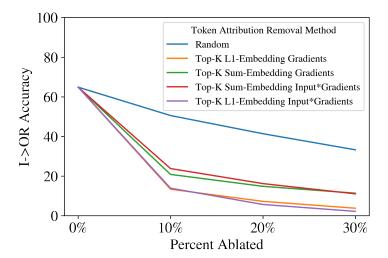


Figure 11: Effect of various gradient attribution methods on the ROAR test at k=10-30% occlusion for the CoS-E v1.0 validation set. We compute attributions with respect to the label logit and measure label accuracy of the resulting model after masking and re-training (see §4.2 for details). The largest drop in performance comes from the L_1 norm embedding-combination method, and raw gradients are not significantly different from input*gradient. On average, input*gradient and raw gradients share 17% of tokens in the top 10%, 24% of tokens in the top 20%, and 31% of tokens in the top 30%.

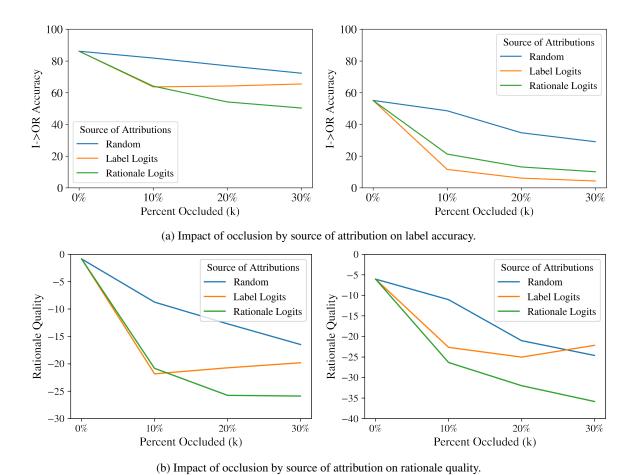


Figure 12: ROAR Feature Importance Agreement results on E-SNLI (left) and CoS-E v1.11 (right). Figure 12a shows label accuracy of the $I\rightarrow OR$ model. Figure 12b shows quality of generated rationales from the $I\rightarrow OR$ model.