

# Few-Shot Self-Rationalization with Natural Language Prompts

Anonymous ACL submission

## Abstract

Self-rationalization models that predict task labels and generate free-text elaborations for their predictions could enable more intuitive interaction with NLP systems. These models are, however, currently trained with a large amount of human-written free-text explanations for each task which hinders their broader usage. We propose to study a more realistic setting of self-rationalization using few training examples. We present FEB—a standardized collection of four existing English-language datasets and associated metrics. We identify the right prompting approach by extensively exploring natural language prompts on FEB. Then, by using this prompt and scaling the model size, we demonstrate that making progress on few-shot self-rationalization is possible. We show there is still ample room for improvement in this task: the average plausibility of generated explanations assessed by human annotators is at most 51%, while plausibility of human explanations is 76%. We hope that FEB and our proposed approach will spur the community to take on the few-shot self-rationalization challenge.

## 1 Introduction

Models constrained to be more understandable to people are easier to troubleshoot and more useful in practice (Rudin et al., 2021). For instance, constraining a model that answers the question “Which linguist invented the lightbulb?” with “none” to also provide the reason—“Thomas Edison is the inventor of the lightbulb and he was not a linguist”—makes the model easier to control and interact with (Kim et al., 2021). Models that jointly predict task labels and generate *free-text explanations* for their predictions (as in the previous example) are known as *self-rationalization models* (Wiegreffe et al., 2021). Their explanations are arguably more faithful and stable than post-hoc explanations since they are intrinsic to the model (Melis and Jaakkola, 2018). The free-text format is essential for explaining tasks requiring reasoning about unstated

knowledge such as commonsense (Marasović et al., 2020), and it makes explanations more intuitive to people compared to highlights of individual words (Camburu et al., 2018). Despite these benefits, self-rationalization models are not widely used, in part because their training currently requires an abundance of human-authored explanations for each task (Narang et al., 2020). A possible solution is few-shot learning, which has shown promising results in recent years. To help the research community begin tackling self-rationalization with only a few examples, we present (i) FEB—a standardized collection of four existing English-language datasets and associated metrics, and (ii) the first approach for the task established through an extensive evaluation of natural language prompts.<sup>1</sup>

One approach to few-shot learning is *prompt-based finetuning* with *natural language prompts*. Such prompts are produced by formatting finetuning instances using a format similar to that used in pretraining, based on the idea that finetuning examples that look similar to pretraining ones will be more informative in the fewshot setting. A few prompts are then used for finetuning. In this paper, we explore whether prompt-based finetuning can be extended to induce few-shot self-rationalization behavior in addition to few-shot prediction. To measure our progress, we first introduce FEB as benchmark dataset consisting of human authored free-text explanations across four distinct end tasks including natural language inference and commonsense tasks (§2). Since finding appropriate prompts is often challenging (Gao et al., 2021), we then extensively explore natural language prompts for few-shot self-rationalization. In our experiments, we fine-tune the T5 and UNIFIEDQA pretrained encoder-decoder transformers (Raffel et al., 2020; Khashabi et al., 2020), and show that versatile question-answering prompts (defined in §3.1) outperform prompts based on span infilling by 8.73

<sup>1</sup>Few Explanations Benchmark (FEB)

accuracy points, as well as prompts designed by following the most similar T5’s supervised pretraining task by 3.21.

We then study the impact of model size on few-shot self-rationalization to investigate whether the quality of generated explanations scales with the size as good as the accuracy of predicting task labels. To this end, we also evaluate GPT-3’s (Brown et al., 2020) self-rationalization behavior. Our experiments show that explanation plausibility scored by human annotators and end-task accuracy improve with increasing model size, despite models being overparametrized. Specifically, the difference in plausibility scores between the BASE and 3B model ranges from [6.24, 24.85] (on average 14.85). The average plausibility across datasets is 43.36 (UNIFIEDQA-3B) and 50.58 (GPT-3). While encouraging, our results show that there is still a large gap between model and human performance (25.75 for GPT-3), and we hope this work will help enable the research community to take on the few-shot self-rationalization challenge.

Our code for producing data splits, prompt construction, model training/evaluation, and human evaluation templates will be publicly available.

## 2 FEB Benchmark

There has been an explosion of interest in generating free-text explanations and in few-shot learning in the last 1–2 years. However, appropriate datasets and metrics for few-shot self-rationalization have not yet been established. We thus introduce the FEB benchmark—a suite of existing English-language datasets with human-authored free-text explanations and associated metrics for few-shot self-rationalization. We expect that FEB will simplify future model comparison and lower barriers to entry for those interested in working on this task.

**Datasets in FEB** To identify available datasets suitable for few-shot self-rationalization, we start with a recent overview of datasets with free-text explanations (Wiegrefe and Marasović, 2021) and filter them according to the following criteria: (i) the input is textual, (ii) the explanation consists of one sentence or 2–3 simple sentences, (iii) the task has a fixed set of possible labels, (iv) the explanation is human-authored, and (v) the dataset has at least 389 instances. We use the second and third criteria to narrow the scope to easier self-rationalization since we expect that few-shot self-rationalization is very challenging. The last requirement is introduced to

FEB Tasks		# Shots
E-SNLI (Camburu et al., 2018)	Classify the entailment relation between two sequences	16
ECQA (Aggarwal et al., 2021)	Answer a question, given five answer choices	48
COMVE (Wang et al., 2019)	Select one of two sequences as more nonsensical	24
SBIC (Sap et al., 2020)	Classify a post as offensive or not	24

Table 1: Tasks that we have included in FEB.

have 48 training and 350 evaluation examples.

This gives us 5 datasets, 4 of which are included in FEB and overviewed in Table 1. These datasets span 4 different tasks: natural language inference, multiple-choice commonsense QA, nonsensical sentence selection, and offensiveness classification. We exclude CoS-E (Rajani et al., 2019) as it is too noisy to be useful for modeling and evaluating self-rationalization (Narang et al., 2021).<sup>2</sup>

ECQA contains not only justifications of the correct answer, but also justifications that refute the incorrect answer choices. We use only the former since they answer “why is [input] assigned [label]?”, just as explanations in other datasets that we have included in FEB. The SBIC dataset contains annotations of frames representing the social biases that are implied in language. We format these frames as a self-rationalization task as follows. We allow only two labels: “offensive” and “not offensive”. If a post is not offensive, we assign it the explanation: “*This post does not imply anything offensive.*” A post can be offensive because it targets an individual or a demographic group. In the former cases, a post is assigned the explanation: “*This post is a personal attack.*” Otherwise, we define a set of rules to transform annotations of which identity-based group is targeted and what stereotypes of this group are referenced or implied into a single, coherent sentence; e.g., group: *women*, stereotype: *can’t drive* → “*This post is offensive because it implies that women can’t drive.*”

This is, to the best of our knowledge, the most comprehensive collection of textual self-rationalization tasks that could also be used even when working in a high-resource setting.

**Automatic Evaluation** Evaluating self-rationalization—predicting task labels and generating explanations for the predicted labels—

<sup>2</sup>Since CoS-E is still actively used, we report CoS-E results in Tables 8 and 9 in Appendix.

requires end-task evaluation and assessing the explanation plausibility. We use accuracy as our end-task evaluation metric. Explanation plausibility may be described as a subjective satisfaction with how a given explanation justifies a label/answer (Yang et al., 2019). Kayser et al. (2021) present the largest currently available study on the correlation of NLG metrics with human judgments of free-text explanation plausibility and report that BERTscore (Zhang et al., 2020) is most correlated (although the correlation is still weak). Thus, we use BERTscore to evaluate the similarity between gold and generated explanations. Following Kayser et al., we assign zero BERTscore to explanations of incorrectly predicted instances.

We follow recent recommendations for reliable few-shot evaluation (Bragg et al., 2021). Specifically, we fix hyperparameters (HPs) and use 60 random train-dev splits with 350 examples in each dev set. For classification tasks, the number of shots (examples per label) is chosen such that we construct a balanced training set of size 48.<sup>3</sup> See Table 1 (col. 3) for exact values; for ECQA we sample 48 training examples. For each model, we report the mean and standard error of 60 mean accuracy/BERTscore values calculated on 60 dev sets of 350 examples.<sup>4</sup> Our HPs are reported in Table 7 in Appendix.

**Human Evaluation** For our final models (§4), we conduct a human evaluation of plausibility of generated explanations following prior work (Kayser et al., 2021; Marasović et al., 2020). For each model evaluation, Kayser et al. (2021) take the first 300 dev examples that are correctly predicted by the model. This means that the dev set subsets used for human evaluation differ across models that are evaluated. However, the overlap between the evaluation sets is maximized by fixing the order of dev instances and taking the first 300.

Prior work used a single train-dev split, while FEB has 60 train-dev splits. Multiple splits provides the opportunity to account for the variance caused by changing the random seed to produce a reliable estimate of plausibility of explanations produced with only a few examples. Therefore,

<sup>3</sup>In early studies, we found that 48 gives models that are at least slightly above the random baseline across all four tasks.

<sup>4</sup>To calculate the standard error for accuracy/BERTscore we use  $n = 60$ . The training (and likewise, dev) sets across splits can overlap, so this error reflects the variability expected in average scores when repeating our experiment with 60 new random splits of the same data sets.

we take the first 6 correctly predicted examples per train-dev split, i.e.,  $6 \cdot 60 = 360$  total instances. Moreover, for classification tasks, we propose to take the first  $6/\#\text{labels}$  correctly-predicted examples per label to have a balanced evaluation set.

Following Kayser et al. (2021), we conduct the human evaluation in two steps:

- **Step1:** Select the correct label/answer.
- **Step2:** Assess whether two explanations (gold and generated) justify the label/answer above.

The first step makes sure the annotators understood the task correctly and they are not able to submit their annotations if the answers are wrong. Ground-truth explanations are evaluated to implicitly influence annotators with a gold reference point when they evaluate generated explanations, and to measure the quality of explanation datasets. To evaluate explanations, annotators are asked “Does the explanation justify the answer?” and given the options {“yes”, “weak yes”, “weak no”, “no”}. These options are mapped to plausibility scores of  $\{1, \frac{2}{3}, \frac{1}{3}, 0\}$ , respectively. For each of the 360 examples, we calculate the mean plausibility score of 3 annotators and report the mean and the standard error of 360 mean scores. We also report the inter-annotator agreement calculated with Fleiss’ kappa. Finally, models are evaluated independently to avoid penalizing worse models in the presence of explanations generated by a better model.

### 3 Prompting for Self-Rationalization

We approach few-shot self-rationalization with prompt-based finetuning using natural language (NL) prompts. The key idea behind NL prompts is that a pretrained language model (LM) is already well-positioned to solve the end-task if we format finetuning end-task examples as similar as possible to the format used in pretraining. Following that principle, in this section, we describe our prompting approach with T5 (Raffel et al., 2020) and comprehensively evaluate three distinct prompt types with FEB. Our results show that a unified question-answering (QA) prompt combined with a T5 variant that includes additional supervised multitask QA training (UNIFIEDQA; Khashabi et al., 2020) performs the best overall across tasks, when compared to three different alternative prompts as described below.

Self-rationalization models (Narang et al., 2020; Wiegrefe et al., 2021) are currently based on T5 for at least two reasons. First, T5 has been pretrained

with many supervised tasks including classification and generation tasks, and self-rationalization involves both classification and generation. Second, T5 is one of the largest *open-sourced* and widely studied pretrained models, and higher LM performance is correlated with larger model size (Kaplan et al., 2020). Thus, all of our experiments are based on T5 (and the UNIFIEDQA variant when evaluating prompts based on a QA format). In this section, all results are obtained with the base version of these models and in §4 we scale model size.

When a LM is pretrained with masked language modeling (Devlin et al., 2019) only, an appropriate NL prompt is constructed by adding and infilling masked tokens (Jiang et al., 2020). T5, however, has been pretrained with span infilling and a suite of supervised tasks whose instances were formatted in various ways. One of these supervised tasks includes SQUAD 1.1 (Rajpurkar et al., 2016) which allows us to experiment with prompts based on QA templates. As a result, we were able to design several different types of NL prompts for T5 consistent with different aspects of its pretraining:

1. QA prompts (SQUAD<sub>T5</sub>, QA<sub>SIMPLE</sub>).
2. span-filling prompts (INFILLING),
3. prompts designed by following the formatting of the most similar T5’s pretraining task ( $\approx$ T5; see Table 1),

We illustrate these prompt types for COMVE in Table 11 in Appendix. The following sections describe these formats in detail and compare their performance using FEB.

### 3.1 QA Prompts

Formatting new instances as QA pairs has been shown to be useful for transfer learning from a QA model (Gardner et al., 2019). We first evaluate options for a versatile QA NL prompt for self-rationalization of tasks in FEB before comparing this approach with the other two prompt types (INFILLING and  $\approx$ T5) in §3.3. As alternative QA models, we investigate two models: T5 (which has been pretrained with QA supervision from SQUAD 1.1), and UNIFIEDQA (a T5 variant described in detail below). Since UNIFIEDQA was trained on a multitask mixture of many different QA datasets, these T5 variants allow us to examine the extent to which additional QA supervision can transfer to the few-shot self-rationalization setting.

Prior work (Bragg et al., 2021) introduced UNIFEW, a model based on UNIFIEDQA, that is

finetuned on a few task-specific instances posed as QA. Despite its simplicity, UNIFEW achieves competitive few-shot learning performance with strong baselines for classification tasks. However, Bragg et al.’s prompts do not cover all task types in FEB, and the question structure in their prompts is highly task-specific (see Appendix A.1).

Alternatively, we propose to design QA prompts with a simple principle in mind—given a non-QA task, construct an equivalent QA task in the form of short “Is...?” or “What is...?” questions. Here, “Is...?” questions have yes/no answers (sometimes “maybe”), and task labels verbatim are answers to “What is...?” questions (e.g., “offensive” and “not offensive”). Then, for UNIFIEDQA and non-QA task in FEB, we develop prompts following the formats proposed in UNIFIEDQA (see Appendix A.1). We denote these prompts as QA<sub>SIMPLE</sub>. For T5, we develop prompts following the SQUAD format for the T5’s pretraining (see Appendix A.1). The output takes the form of “*answer* because *explanation*”. There is another factor to consider. We need to decide whether to add *tags*—a single descriptions of different input elements; e.g., “premise:” and “hypothesis:” before the first and second sentence in the E-SNLI input. Without these tags the task seems impossible to understand, but UNIFIEDQA has not been trained with similar tags.

Table 11 in Appendix shows examples of our various QA prompts.

**Results** We present the results of UNIFIEDQA with QA<sub>SIMPLE</sub> in Table 2, and due to space limits, T5’s results with SQUAD<sub>T5</sub> prompts in Table 10 in Appendix. For ECQA with UNIFIEDQA, we use the UNIFIEDQA format for multiple-choice QA (see Appendix A.1).

We observe that for E-SNLI and COMVE it is crucial to add tags (“premise:”/“hypothesis:”; “choice1:”/“choice2:”). This result is intuitive—it should be difficult to pick one of the two sentences, or classify a relation between them, if sentences are not marked.<sup>5</sup> On the other hand, adding label choices is not beneficial and in some cases can even decrease the performance. When tags are included, we see that across all the tasks the “What is...?” question performs the best. This also holds for T5 and SQUAD<sub>T5</sub> prompts (see Table 10). Finally, the prompt with the “What is...?” question

<sup>5</sup>Performance on COMVE with “Is...?” is close to random which suggests that this question form hinders the performance and tags cannot make a difference.

and tags in the input outperforms UNIFEW for both tasks UNIFEW can be applied to. This result shows that this prompt is both versatile and effective.

Finally, we compare the best performing prompts we get with UNIFIEDQA + QA<sub>SIMPLE</sub> and T5+SQUAD<sub>T5</sub>. See prompts “SQUAD<sub>T5</sub> × WHAT IS...? + TAGS” and “QA<sub>SIMPLE</sub> × WHAT IS...? + TAGS” in Table 11. For ECQA and COMVE, we observe notable improvements from using UNIFIEDQA, and minor improvements for SBIC. For E-SNLI, T5 is better, presumably because UNIFIEDQA has lost some useful information from MNLI after extensive continued pretraining for QA. These results suggest that UNIFIEDQA is a better model for prompting self-rationalization with QA prompts.

To recap, the analysis presented in this section suggests that QA prompting for inducing self-rationalization behavior is best done when UNIFIEDQA is combined with the NL prompt below. For true QA tasks, we use UNIFIEDQA formats.

```
explain what is this/more...? \n tag1:
[sequence1] tag2: [sequence2] ...</s>
```

### 3.2 INFILLING Prompts

The simplest way to design an infilling prompt is to prepend the span “<extra\_id\_0> because <extra\_id\_1>” to the input. A model should then replace <extra\_id\_0> with a label/answer and <extra\_id\_1> with an explanation. Besides being similar to T5’s span infilling pretraining task, another benefit of this prompt is that is very flexible—the span above can be added to any task input. This basic infilling prompt could be easily made more natural by prepending phrases such as: “The answer is” (ECQA), “Less common is” (COMVE), or “This is” (E-SNLI, SBIC). We hypothesize that these additional phrases could be beneficial because they suggest which subset of the vocabulary is the right word for filling in <extra\_id\_0>. We test whether it is beneficial to make the infilling prompt more natural sounding.

**Results** T5 results are shown in Table 3. The outcome is mixed—while we observe notable benefits for ECQA/SBIC, for E-SNLI/COMVE there is a minor difference in favor of the basic prompt. A way to explain this is that T5 learned about NLI labels from MNLI during pretraining, so it does not need additional phrase to nudge it in the right

	Prompt	Accuracy	BERTscore
E-SNLI	UNIFEW	61.68 <sub>0.58</sub>	55.85 <sub>0.53</sub>
	+ tags	63.61 <sub>0.44</sub>	57.34 <sub>0.41</sub>
	Is...?	47.47 <sub>0.52</sub>	42.70 <sub>0.47</sub>
	+ tags	66.59 <sub>0.51</sub>	60.05 <sub>0.47</sub>
	+ tags & choices	64.43 <sub>0.53</sub>	58.16 <sub>0.49</sub>
	RANDOM BASELINE	33.33	-
ECQA	UNIFIEDQA	<b>41.37</b> <sub>0.34</sub>	<b>36.72</b> <sub>0.30</sub>
	RANDOM BASELINE	20.00	-
COMVE	Is...?	52.69 <sub>0.35</sub>	47.70 <sub>0.31</sub>
	+ tags	52.47 <sub>0.32</sub>	47.47 <sub>0.30</sub>
	+ tags & choices	52.19 <sub>0.33</sub>	47.27 <sub>0.30</sub>
	What is...?	50.60 <sub>0.22</sub>	45.68 <sub>0.20</sub>
	+ tags	<b>67.33</b> <sub>0.71</sub>	<b>60.97</b> <sub>0.64</sub>
	+ tags & choices	62.56 <sub>0.65</sub>	56.68 <sub>0.59</sub>
RANDOM BASELINE	50.00	-	
SBIC	UNIFEW	66.15 <sub>0.43</sub>	63.84 <sub>0.44</sub>
	Is...?	63.50 <sub>0.44</sub>	61.21 <sub>0.42</sub>
	+ tags	62.64 <sub>0.45</sub>	60.43 <sub>0.45</sub>
	+ tags & choices	63.63 <sub>0.42</sub>	61.31 <sub>0.43</sub>
	What is...?	67.35 <sub>0.38</sub>	65.03 <sub>0.37</sub>
	+ tags	<b>67.55</b> <sub>0.41</sub>	<b>65.29</b> <sub>0.39</sub>
+ tags & choices	65.43 <sub>0.58</sub>	63.07 <sub>0.59</sub>	
RANDOM BASELINE	50.00	-	

Table 2: Prompting UNIFIEDQA with QA<sub>SIMPLE</sub> with “Is...?” and “What is...?” questions, and UNIFEW. See §3.1 for descriptions of these prompts. For ECQA we use the original UNIFIEDQA format for multiple-choice QA. We also inspect the effects of adding label choices and *tags* (defined in §3.1) to the input.

	E-SNLI	ECQA	COMVE	SBIC
B	<b>75.24</b> <sub>0.38</sub>	22.33 <sub>0.29</sub>	<b>50.36</b> <sub>0.31</sub>	61.57 <sub>0.45</sub>
N	75.09 <sub>0.45</sub>	<b>27.60</b> <sub>0.36</sub>	49.02 <sub>0.28</sub>	<b>64.66</b> <sub>0.52</sub>
	E-SNLI	ECQA	COMVE	SBIC
B	<b>67.66</b> <sub>0.35</sub>	19.83 <sub>0.26</sub>	<b>45.51</b> <sub>0.28</sub>	59.18 <sub>0.46</sub>
N	67.52 <sub>0.42</sub>	<b>24.52</b> <sub>0.32</sub>	44.35 <sub>0.26</sub>	<b>62.00</b> <sub>0.54</sub>

Table 3: A comparison of the basic infilling prompt (B) with its more natural sounding version (N). The upper part shows accuracy and the lower part BERTscore.

direction. COMVE results are comparable to the random performance, and the model could not learn

	Task	Accuracy	BERTscore
INFILLING	E-SNLI	75.09 <sub>0.45</sub>	67.52 <sub>0.42</sub>
	ECQA	27.60 <sub>0.36</sub>	24.52 <sub>0.32</sub>
	COMVE	49.02 <sub>0.28</sub>	44.35 <sub>0.26</sub>
	SBIC	64.66 <sub>0.52</sub>	62.00 <sub>0.54</sub>
	Average	54.09	49.57
≈T5	E-SNLI	<b>79.21</b> <sub>0.29</sub>	<b>71.34</b> <sub>0.27</sub>
	ECQA	38.28 <sub>0.33</sub>	33.91 <sub>0.29</sub>
	COMVE	55.88 <sub>0.34</sub>	50.45 <sub>0.30</sub>
	SBIC	65.06 <sub>0.60</sub>	62.77 <sub>0.63</sub>
	Average	59.61	54.62
QA <sub>SIMPLE</sub>	E-SNLI	75.05 <sub>0.34</sub>	67.52 <sub>0.33</sub>
	ECQA	<b>41.37</b> <sub>0.34</sub>	<b>36.72</b> <sub>0.30</sub>
	COMVE	<b>67.33</b> <sub>0.71</sub>	<b>60.97</b> <sub>0.64</sub>
	SBIC	<b>67.55</b> <sub>0.41</sub>	<b>65.29</b> <sub>0.39</sub>
	Average	<b>62.82</b>	<b>57.63</b>

Table 4: A comparison between three prompt types: INFILLING, ≈T5, and QA<sub>SIMPLE</sub> prompts. See §3 for descriptions of these prompts.

the task from the infilling prompt, with or without the additional phrases. Thus, we recommend using the more natural version as it is not detrimental to E-SNLI/COMVE performance while it leads to big improvements for ECQA/SBIC.

### 3.3 INFILLING vs. ≈T5 vs. QA

We have established appropriate QA and INFILLING prompts in §3.1 and §3.2. We now turn to a comparison between all three prompt types: (i) INFILLING (natural), (ii) ≈T5, and (iii) QA<sub>SIMPLE</sub> (“What is...?” with tags). The first two are used to prompt T5 and the last type UNIFIEDQA. To construct ≈T5 prompts, for each task in FEB, we identify the most similar T5’s pretraining task (see Table 6, Appendix) and use that task’s formatting (see, e.g., ≈T5 × COPA in Table 11).

**Results** A comparison of the three prompt types is presented in Table 4. The QA<sub>SIMPLE</sub> prompt outperforms other prompt types for all tasks except E-SNLI for which unsurprisingly ≈T5 is the best. Finally, this brings us to the end of our extensive exploration of natural language prompts for a prompt-based finetuning approach to few-shot self-rationalization. We identify the QA<sub>SIMPLE</sub> prompt as the most effective and we use it to study how few-shot self-rationalization performance scales with the size of the UNIFIEDQA model.

## 4 Improving Self-Rationalization with Increasing Model Size

In §3, we discovered that a QA prompt combined with the base UNIFIEDQA model version is as an effective combination for few-shot self-rationalization through prompt-based finetuning. In this section, we provide two additional evaluations to establish the first approach to few-shot self-rationalization.

First, we assess how plausible the generated explanations are when evaluated by annotators on Amazon MTurk. Details of how we conduct human evaluation of plausibility are given in §2. One HIT contains 10 instances and we pay \$1 per HIT.

Next, we investigate how self-rationalization performance changes with the model size since larger pretrained language models typically give better few-shot performance (Brown et al., 2020). We wonder whether the same trend will hold for a complex generation task of self-rationalization where it is conceivable that an enormous model could overfit on a few examples. To this end, we evaluate three versions of UNIFIEDQA (BASE, LARGE, 3B) and GPT-3 (Brown et al., 2020).

**Experimental Setup** We evaluate GPT-3 using its API and “in-context demonstrations” (Brown et al., 2020). We pack as many demonstrations as we can fit in the input, followed by the input of the test example, then run GPT-3 to generate its output. The number of demonstrations we are able to fit ranges from [28,45] which are randomly selected from the 48 used for UNIFIEDQA. Since evaluation using a single prompt costs us \$1,050, we do not do prompt search for GPT-3. We use the prompts shown in Fig. 1 in Appendix.

A detailed description of evaluation metrics is given in §2. Each dev set size for GPT-3 is 18 instead of 350 (because of the API cost). Ground-truth explanations are evaluated together with explanations generated by 4 models. Therefore, for GOLD explanations, we report the average of 4 plausibility scores, std. errors, and  $\kappa$  values calculated with 4 Mturk batches (corresponding to 4 models).

### 4.1 Results

Results are shown in Table 5. Note that we use T5 with the ≈T5 prompt for E-SNLI, and UNIFIEDQA with QA<sub>SIMPLE</sub> (§3) for other datasets to establish the best possible performance for each dataset. The exact prompts for each task

				Plausibility								
				<i>All</i>		<i>Label<sub>1</sub></i>		<i>Label<sub>2</sub></i>		<i>Label<sub>3</sub></i>		
Model	# Par.	Accuracy	BERTscore	Score	$\kappa$	Score	$\kappa$	Score	$\kappa$	Score	$\kappa$	
E-SNLI	BASE	220M	79.21 <sub>0.29</sub>	71.34 <sub>0.27</sub>	16.75 <sub>1.53</sub>	0.73	15.65 <sub>2.34</sub>	0.67	17.50 <sub>2.88</sub>	0.79	17.13 <sub>2.71</sub>	0.72
	LARGE	770M	84.79 <sub>0.27</sub>	76.56 <sub>0.27</sub>	32.68 <sub>1.92</sub>	0.57	<b>27.31</b> <sub>2.88</sub>	0.43	33.89 <sub>3.44</sub>	0.64	36.85 <sub>3.58</sub>	0.64
	3B	2.8B	<b>87.43</b> <sub>0.23</sub>	<b>79.10</b> <sub>0.23</sub>	41.60 <sub>2.08</sub>	0.62	27.13 <sub>2.85</sub>	0.52	46.76 <sub>3.84</sub>	0.70	<b>50.92</b> <sub>3.63</sub>	0.64
	GPT-3	175B	65.37 <sub>0.53</sub>	59.83 <sub>0.47</sub>	<b>42.44</b> <sub>2.17</sub>	0.54	<b>27.31</b> <sub>2.87</sub>	0.48	<b>66.03</b> <sub>4.37</sub>	0.71	43.80 <sub>3.46</sub>	0.51
	GOLD	-			77.40 <sub>1.59</sub>	0.63	63.50 <sub>3.01</sub>	0.44	87.87 <sub>1.85</sub>	0.74	82.48 <sub>2.42</sub>	0.72
	RAND	-	33.33									
ECQA	BASE	220M	41.37 <sub>0.34</sub>	36.72 <sub>0.30</sub>	25.52 <sub>1.25</sub>	0.32						
	LARGE	770M	57.19 <sub>0.36</sub>	51.00 <sub>0.32</sub>	30.28 <sub>1.53</sub>	0.38						
	3B	2.8B	<b>65.86</b> <sub>0.36</sub>	<b>58.98</b> <sub>0.32</sub>	34.23 <sub>1.56</sub>	0.35						
	GPT-3	175B	60.65 <sub>1.48</sub>	54.42 <sub>1.32</sub>	<b>45.06</b> <sub>1.44</sub>	0.12						
	GOLD	-			70.88 <sub>1.47</sub>	0.45						
	RAND	-	20.00									
COMVE	BASE	220M	67.33 <sub>0.71</sub>	60.97 <sub>0.64</sub>	13.80 <sub>1.26</sub>	0.45						
	LARGE	770M	81.31 <sub>0.39</sub>	73.95 <sub>0.36</sub>	25.59 <sub>1.67</sub>	0.52						
	3B	2.8B	<b>88.96</b> <sub>0.38</sub>	<b>81.02</b> <sub>0.34</sub>	33.40 <sub>1.71</sub>	0.63						
	GPT-3	175B	73.98 <sub>1.40</sub>	67.65 <sub>1.29</sub>	<b>42.16</b> <sub>1.80</sub>	0.73						
	GOLD	-			77.24 <sub>1.30</sub>	0.55						
	RAND	-	50.00									
SBIC	BASE	220M	67.55 <sub>0.41</sub>	65.29 <sub>0.39</sub>	57.96 <sub>2.25</sub>	0.68	21.36 <sub>2.06</sub>	0.54	94.57 <sub>1.08</sub>	0.82		
	LARGE	770M	71.06 <sub>0.39</sub>	68.55 <sub>0.39</sub>	61.82 <sub>2.23</sub>	0.66	27.16 <sub>2.19</sub>	0.43	<b>96.48</b> <sub>0.92</sub>	0.89		
	3B	2.8B	71.66 <sub>0.48</sub>	68.90 <sub>0.49</sub>	64.20 <sub>2.14</sub>	0.68	33.76 <sub>2.65</sub>	0.55	94.63 <sub>1.02</sub>	0.81		
	GPT-3	175B	<b>74.17</b> <sub>1.41</sub>	<b>71.53</b> <sub>1.40</sub>	<b>72.68</b> <sub>1.72</sub>	0.53	<b>52.65</b> <sub>2.51</sub>	0.34	92.72 <sub>1.05</sub>	0.72		
	GOLD	-			79.81 <sub>1.62</sub>	0.67	64.92 <sub>2.66</sub>	0.52	94.69 <sub>1.01</sub>	0.81		
	RAND	-	50.00									

Table 5: The first results on the FEB benchmark using T5/UNIFIEDQA (BASE, LARGE, 3B) and GPT-3. T5+ $\approx$ T5 prompt is used only for E-SNLI, and UNIFIEDQA + QA<sub>SIMPLE</sub> prompt is used for other datasets. The descriptions of these prompts are given in §3 and details of how evaluation metrics are calculated in §2. RAND stands for a random baseline and GOLD for human-authored explanations. *Label<sub>1</sub>/Label<sub>2</sub>/Label<sub>3</sub>* are entailment/neutral/contradiction in E-SNLI and offensive/not offensive in SBIC.

are given in Appendix A.2. We observe that all metrics—accuracy, BERTscore, and plausibility—monotonically increase with the size of UNIFIEDQA for all datasets. That is, larger models learn to predict task labels and generate explanations from a few examples better, despite being overparametrized. UNIFIEDQA-3B has a higher accuracy/BERTscore than GPT-3 for all datasets except SBIC, but GPT-3 generates explanations that are notably more plausible.

The following observations suggest that few-shot self-rationalization is a promising research direction. The difference in plausibility scores between the BASE and 3B model versions ranges from [6.24, 24.85] (on average 14.85). In other words, since it is possible to generate more plausible explanation by only increasing the model size, it is conceivable that further progress could be made with more creative approaches. Next, the plausibility score of the best model (GPT-3) ranges from [42.16, 72.68]

[42.16, 52.65] if we consider only SBIC “offensive” (*Label<sub>1</sub>*) subset. This shows that a moderate plausibility can already be achieved with current models without any task-specific enhancements.

Despite that, the gap between our best models and human-authored explanations remains large. The average plausibility score across datasets is 43.36 (UNIFIEDQA-3B), 50.58 (GPT-3), and 76.33 (GOLD). In other words, the difference in plausibility scores between UNIFIEDQA-3B’s and human explanations is 32.98, and between GPT-3’s and human explanations is 25.75. We expect that the FEB benchmark, our UNIFIEDQA approach, and first results, present a good starting point to tackle this challenge.

**Performance w.r.t. Labels** For E-SNLI and SBIC, we can inspect the metrics with respect to labels. In E-SNLI part of the Table 5, *Label<sub>1</sub>* marks “entailment”, *Label<sub>2</sub>* “neutral”, and *Label<sub>3</sub>* “con-

tradition”. There are notable differences between the plausibility scores for each label. The plausibility score for “entailment” does not scale with the model size and it is much lower than scores for other labels (the best score is 27.31 vs. 66.0/50.92). This issue stems from the difficulty of explaining the entailment label (Camburu et al., 2018). Even people struggle with explaining “entailment” as evident by the lower GOLD score for “entailment” compared to the other two labels. An interesting observation from the other two labels is that UNIFIEDQA-3B explains “contradiction” instances best and GPT-3 “neutral” instances.

In SBIC part of the Table 5,  $Label_1$  marks “offensive” and  $Label_2$  “not offensive” instances. The latter achieve almost perfect plausibility since the models learn to generate “*This post does not imply anything offensive*”. Thus, main plausibility scores for SBIC are those of offensive instances. We can observe that the relative differences between models for offensive instances are much larger than the relative differences when examples of both labels are counted for (column “All / Score”). If we had only looked into a single plausibility score we would not notice these differences. Thus, we recommend breaking down the performance w.r.t. labels whenever possible.

**Annotator Agreement** Finally, we observe challenges in collecting human judgments of plausibility. For all datasets except ECQA, Fleiss’  $\kappa$  is either moderate (between 0.41–0.6) or substantial (between 0.61–0.8). One exception is GPT-3 on SBIC ( $Label_1$ ; offensive) where  $\kappa$  is only 0.34. We also observe that  $\kappa$  for GPT-3’s explanations is lower than  $\kappa$  for UNIFIEDQA’s or GOLD explanations, with the exception of COMVE. The most concerning is ECQA where  $\kappa$  is on average 0.35 for UNIFIEDQA’s explanations, 0.34 for GOLD explanations, and only 0.12 for GPT-3’s. Future work should investigate the reasons behind these differences more carefully.

## 5 Related Work

**Self-Rationalization with Few Human-Written Explanations** Select-then-predict method (Lei et al., 2016) that is standard to creating explanations in the form of *highlights* of the input tokens does not use any human-author highlighting explanations. On the other hand, a standard approach to generating free-text explanations is to use human-written explanations (Liu et al., 2019;

Wu and Mooney, 2019; Narang et al., 2020, among others). To the best of our knowledge, only two prior works generate *free-text* explanations in a weakly-supervised way from the task prediction loss. Laticinnik and Berant (2020) approach commonsense QA in that fashion. Brahman et al. (2021) propose multiple distant supervision approach to explaining a defeasible inference task. In this paper, we introduce the FEB benchmark to unify the evaluation of few-shot self-rationalization and present the first approach and results on FEB.

**Few-Shot Learning** We study natural language prompts (Brown et al., 2020; Schick and Schütze, 2021) to establish the first approach to few-shot self-rationalization. Alternatively, few-shot learning researchers are studying prompts in the form of continuous/soft vectors that do not correspond to real tokens (e.g., Qin and Eisner, 2021). Such methods present a promising research direction for few-shot self-rationalization. Namely, we show that larger models generate notably more plausible explanations, and “prefix tuning” (Li and Liang, 2021) has been show to learn two condition generation tasks using only 0.1% of the parameters, while maintaining comparable performance. In practice, such approaches still require a notable amount of GPU memory. Thus, any efforts to reduce required memory such as compression (Ganesh et al., 2021) may be valuable for few-shot self-rationalization.

## 6 Conclusions

We draw attention to the task of few-shot self-rationalization: predicting task labels and generating *free-text* explanations for the prediction using only a few human-written explanations. We present (i) the FEB benchmark, (ii) the first prompting approach for FEB established through a comprehensive search of natural language prompts, and (iii) results using models with a number of parameters ranging from 220M to 175B. Our human evaluation results show that progress is possible on this task given that just scaling the model size increases both the plausibility of generated explanations and task accuracy by a very large margin. Despite that, few-shot self-rationalization remains very challenging, with plausibility of explanations generated by the best model being 27.75 points behind plausibility of human-authored explanations. We hope that work presented in this paper spurs the community to work on this challenging problem to enable more intuitive interaction with NLP systems.



## References

- Shourya Aggarwal, Divyanshu Mandowara, Vishwa-  
jeet Agrawal, Dinesh Khandelwal, Parag Singla, and  
Dinesh Garg. 2021. [Explanations for Common-  
senseQA: New Dataset and Models](#). In *Proceedings  
of the 59th Annual Meeting of the Association for  
Computational Linguistics and the 11th International  
Joint Conference on Natural Language Processing  
(Volume 1: Long Papers)*, pages 3050–3065, Online.  
Association for Computational Linguistics.
- Jonathan Bragg, Arman Cohan, Kyle Lo, and Iz Beltagy.  
2021. [Flex: Unifying evaluation for few-shot nlp](#). In  
*Proceedings of the Advances in Neural Information  
Processing Systems (NeurIPS)*.
- Faeze Brahman, Vered Shwartz, Rachel Rudinger, and  
Yejin Choi. 2021. [Learning to rationalize for non-  
monotonic reasoning with distant supervision](#). In  
*Proceedings of the AAAI Conference on Artificial  
Intelligence*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
teusz Litwin, Scott Gray, Benjamin Chess, Jack  
Clark, Christopher Berner, Sam McCandlish, Alec  
Radford, Ilya Sutskever, and Dario Amodei. 2020.  
[Language models are few-shot learners](#). In *Ad-  
vances in Neural Information Processing Systems*,  
volume 33, pages 1877–1901. Curran Associates,  
Inc.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas  
Lukasiewicz, and Phil Blunsom. 2018. [e-SNLI: Nat-  
ural language inference with natural language expla-  
nations](#). In *Proceedings of the Advances in Neural  
Information Processing Systems (NeurIPS)*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. [BERT: Pre-training of  
deep bidirectional transformers for language under-  
standing](#). In *Proceedings of the 2019 Conference of  
the North American Chapter of the Association for  
Computational Linguistics: Human Language Tech-  
nologies, Volume 1 (Long and Short Papers)*, pages  
4171–4186, Minneapolis, Minnesota. Association for  
Computational Linguistics.
- Prakhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali  
Khan, Yin Yang, Hassan Sajjad, Preslav Nakov, Dem-  
ing Chen, and Marianne Winslett. 2021. [Compress-  
ing Large-Scale Transformer-Based Models: A Case  
Study on BERT](#). *Transactions of the Association for  
Computational Linguistics*, 9:1061–1080.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021.  
[Making pre-trained language models better few-shot  
learners](#). In *Proceedings of the 59th Annual Meet-  
ing of the Association for Computational Linguistics  
and the 11th International Joint Conference on Natu-  
ral Language Processing (Volume 1: Long Papers)*,  
pages 3816–3830, Online. Association for Computa-  
tional Linguistics.
- Matt Gardner, Jonathan Berant, Hannaneh Hajishirzi,  
Alon Talmor, and Sewon Min. 2019. [Ques-  
tion answering is a format; when is it useful?](#)  
arXiv:1909.11291.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy  
Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-  
hardt. 2021. [Measuring massive multitask language  
understanding](#). In *The International Conference on  
Learning Representations (ICLR)*.
- Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham  
Neubig. 2020. [How can we know what language  
models know?](#) *Transactions of the Association for  
Computational Linguistics*, 8:423–438.
- Jared Kaplan, Sam McCandlish, Tom Henighan,  
Tom B. Brown, Benjamin Chess, Rewon Child, Scott  
Gray, Alec Radford, Jeff Wu, and Dario Amodei.  
2020. [Scaling laws for neural language models](#).  
arXiv:2001.08361.
- 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 explana-  
tions in vision-language tasks](#). In *Proceedings of the  
IEEE/CVF International Conference on Computer  
Vision (ICCV)*.
- Daniel Khashabi, Sewon Min, Tushar Khot, Ashish  
Sabharwal, Oyvind Tafjord, Peter Clark, and Han-  
naneh Hajishirzi. 2020. [UNIFIEDQA: Crossing for-  
mat boundaries with a single QA system](#). In *Find-  
ings of the Association for Computational Linguistics:  
EMNLP 2020*, pages 1896–1907, Online. Association  
for Computational Linguistics.
- Najoung Kim, Ellie Pavlick, Burcu Karagol Ayan, and  
Deepak Ramachandran. 2021. [Which linguist in-  
vented the lightbulb? presupposition verification for  
question-answering](#). arXiv:2101.00391.
- Veronica Latcinnik and Jonathan Berant. 2020. [Explain-  
ing question answering models through text genera-  
tion](#). arXiv:2004.05569.
- Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2016.  
[Rationalizing neural predictions](#). In *Proceedings of  
the 2016 Conference on Empirical Methods in Natu-  
ral Language Processing*, pages 107–117, Austin,  
Texas. Association for Computational Linguistics.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning:  
Optimizing continuous prompts for generation](#). In  
*Proceedings of the 59th Annual Meeting of the Asso-  
ciation for Computational Linguistics and the 11th  
International Joint Conference on Natural Language  
Processing (Volume 1: Long Papers)*, pages 4582–  
4597, Online. Association for Computational Lin-  
guistics.

741	Hui Liu, Qingyu Yin, and William Yang Wang. 2019. <a href="#">Towards explainable NLP: A generative explanation framework for text classification</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5570–5581, Florence, Italy. Association for Computational Linguistics.	797
742		798
743		799
744		
745		
746		
747	Ana Marasović, Chandra Bhagavatula, Jae sung Park, Ronan Le Bras, Noah A. Smith, and Yejin Choi. 2020. <a href="#">Natural language rationales with full-stack visual reasoning: From pixels to semantic frames to commonsense graphs</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2810–2829, Online. Association for Computational Linguistics.	800
748		801
749		802
750		803
751		
752		
753		
754		
755	David Alvarez Melis and Tommi Jaakkola. 2018. <a href="#">Towards robust interpretability with self-explaining neural networks</a> . In <i>Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)</i> .	804
756		805
757		806
758		807
759	Sharan Narang, Hyung Won Chung, Yi Tay, Liam Fedus, Thibault Fevry, Michael Matena, Karishma Malkan, Noah Fiedel, Noam Shazeer, Zhenzhong Lan, Yanqi Zhou, Wei Li, Nan Ding, Jake Marcus, Adam Roberts, and Colin Raffel. 2021. <a href="#">Do transformer modifications transfer across implementations and applications?</a> In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5758–5773, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	808
760		809
761		810
762		811
763		812
764		813
765		814
766		
767		
768		
769		
770	Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan. 2020. <a href="#">WT5?! Training Text-to-Text Models to Explain their Predictions</a> . arXiv:2004.14546.	815
771		816
772		817
773		818
774	Guanghui Qin and Jason Eisner. 2021. <a href="#">Learning how to ask: Querying LMs with mixtures of soft prompts</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5203–5212, Online. Association for Computational Linguistics.	819
775		820
776		821
777		
778		
779		
780		
781	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	822
782		823
783		824
784		825
785		826
786		827
787	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. <a href="#">Explain yourself! leveraging language models for commonsense reasoning</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	828
788		829
789		830
790		831
791		832
792		
793		
794	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">SQuAD: 100,000+ questions for machine comprehension of text</a> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	833
795		834
796		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852

- 853 Jialin Wu and Raymond Mooney. 2019. [Faithful mul-](#)  
854 [timodal explanation for visual question answering.](#)  
855 In *Proceedings of the 2019 ACL Workshop Black-*  
856 *boxNLP: Analyzing and Interpreting Neural Net-*  
857 *works for NLP*, pages 103–112, Florence, Italy. As-  
858 sociation for Computational Linguistics.
- 859 Fan Yang, Mengnan Du, and Xia Hu. 2019. [Evaluat-](#)  
860 [ing explanation without ground truth in interpretable](#)  
861 [machine learning.](#) 1907.06831.
- 862 Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jian-  
863 feng Gao, Kevin Duh, and Benjamin Van Durme.  
864 2018. [Record: Bridging the gap between human](#)  
865 [and machine commonsense reading comprehension.](#)  
866 [arXiv:1810.12885.](#)
- 867 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.  
868 Weinberger, and Yoav Artzi. 2020. [Bertscore: Evalu-](#)  
869 [ating text generation with bert.](#) In *The International*  
870 *Conference on Learning Representations (ICLR).*

## A Appendix

### A.1 Input Formats in Related Work

#### T5 with SQUAD

- question: [question] context: [paragraph]

#### UNIFIEDQA

- **Multiple-choice QA:** [question] \n (A) [choice<sub>1</sub>] (B) [choice<sub>2</sub>]...
- **Extractive QA:** [question] \n [paragraph]

#### UNIFEW

- **Single text classification:** Topic? \n (A) [class<sub>1</sub>] (B) [class<sub>2</sub>] (C) [class<sub>3</sub>] \n [document]
- **Sentence-pair classification:** [sentence<sub>1</sub>] Is [sentence<sub>1</sub>] ? \n (A) Yes (B) No (C) Maybe
- **Relation classification:** [mention<sub>1</sub>] to [mention<sub>1</sub>] ? \n (A) [class<sub>1</sub>] (B) [class<sub>2</sub>] (C) [class<sub>3</sub>]
- **Entity recognition:** What is the type of the entity between the # marks? \n [text... #mention<sub>1</sub># ...text... #mention<sub>2</sub># ...text]

### A.2 Prompts Used in §4

#### UNIFIEDQA/T5 prompts:

- **E-SNLI** explain nli hypothesis: [hypothesis] premise: [premise]
- **ECQA:** explain [question] \n (A) [choice<sub>1</sub>] (B) [choice<sub>2</sub>] (C) [choice<sub>3</sub>] (D) [choice<sub>4</sub>] (E) [choice<sub>5</sub>]
- **COMVE:** explain what is more nonsensical? \n choice1: [sentence<sub>1</sub>] choice2: [sentence<sub>2</sub>]
- **SBIC:** explain what is this? \n post: [post]

Prompts for GPT-3 are given in Figure 1. The output always takes the form of “[answer/label] because [explanation]”.

```
===== SBIC =====
Is the following post offensive? and why?
Post: [Input]
Answer: Yes|No
Reason: [Explanation]

===== E-SNLI =====
Answer the question and provide a reason why
the answer is correct.
[Premise]
Question: Is [Hypothesis]?
Answer: Yes|No|Maybe
Reason: [Explanation]

===== ECQA =====
Answer the question from the provided
choices, and provide a reason why the answer
is correct.
Question: [Question]
Choices: [Choices]
Answer: [one of the choices]
Reason: [Explanation]

===== ComVE =====
Which of the two choices makes more sense?
and why?
Choice1: [Choice1]
Choice2: [Choice2]
Answer: Choice1|Choice2
Reason: [Explanation]
```

Figure 1: GPT-3 prompt templates for all datasets.

FEB Task	Similar T5 Pretraining Tasks
E-SNLI	MNLI (Williams et al., 2018) Classify the entailment relation between two sequences
ECQA	RECORD (Zhang et al., 2018) Answer a cloze-style query about a passage given entities in it
COMVE	COPA (Roemmele et al., 2011) Select one of two sequences as the cause/effect of a premise
SBIC	COLA (Warstadt et al., 2019) Classify a sentence as acceptable or not

Table 6: The first column shows tasks that we have included in FEB. Tasks on the right are included in T5’s pretraining and they are similar to FEB’s tasks. We explore self-rationalization prompts for FEB’s tasks based on the tasks on the right, and compare them to prompts designed as span infilling and QA (§3).

<b>GPUs</b>	NVIDIA A100 on Google Cloud
<b>Implementation</b>	<i>Will be added upon acceptance.</i>

Hyperparameter	Assignment
max step number	300
batch size	4 (1 for UNIFIEDQA-3B)
gradient accumulation steps	1 (4 for UNIFIEDQA-3B)
learning rate	3e-5
learning rate scheduler	linear
warmup steps	0
decoding	greedy

Table 7: Hyperparameters used in our experiments.

		Accuracy	BERTscore
COS-E	INFILLING (b)	34.28 <sub>0.36</sub>	29.60 <sub>0.32</sub>
	INFILLING (n)	40.14 <sub>0.38</sub>	34.70 <sub>0.34</sub>
	≈T5	51.69 <sub>0.41</sub>	44.56 <sub>0.36</sub>
	SQUAD <sub>T5</sub>	51.15 <sub>0.34</sub>	44.13 <sub>0.29</sub>
	QA <sub>SIMPLE</sub>	<b>59.96</b> <sub>0.32</sub>	<b>48.57</b> <sub>0.26</sub>

Table 8: A comparison of all prompt types introduced in §3 on COS-E. We do not support using COS-E in the future given the reported issues with it (Narang et al., 2020; Wiegrefe and Marasović, 2021), especially since ECQA is introduced.

	Size	Accuracy	BERTscore
COS-E	BASE	58.32 <sub>0.28</sub>	50.43 <sub>0.25</sub>
	LARGE	69.44 <sub>0.30</sub>	60.11 <sub>0.26</sub>
	3B	<b>75.37</b> <sub>0.31</sub>	<b>65.34</b> <sub>0.28</sub>
	GPT-3	68.43 <sub>1.35</sub>	59.48 <sub>1.18</sub>

Table 9: The effect of scaling the UNIFIEDQA model size on self-rationalization of COS-E. We do not support using COS-E in the future given the reported issues with it (Narang et al., 2020; Wiegrefe and Marasović, 2021), especially since ECQA is introduced.

	Prompt	Accuracy	BERTscore
E-SNLI	Is...?	38.68 <sub>0.44</sub>	34.74 <sub>0.40</sub>
	+ tags	48.20 <sub>0.62</sub>	43.22 <sub>0.58</sub>
E-SNLI	What is...?	60.76 <sub>0.85</sub>	54.75 <sub>0.77</sub>
	+ tags	<b>77.86</b> <sub>0.34</sub>	<b>70.08</b> <sub>0.32</sub>
ECQA	SQUAD <sub>T5</sub>	<b>36.48</b> <sub>0.34</sub>	<b>32.38</b> <sub>0.30</sub>
	RANDOM BASELINE	20.00	-
COMVE	Is...?	50.38 <sub>0.16</sub>	45.54 <sub>0.14</sub>
	+ tags	50.17 <sub>0.14</sub>	45.35 <sub>0.13</sub>
COMVE	What is...?	50.54 <sub>0.21</sub>	45.67 <sub>0.19</sub>
	+ tags	<b>54.49</b> <sub>0.46</sub>	<b>49.25</b> <sub>0.42</sub>
SBIC	Is...?	63.37 <sub>0.58</sub>	61.15 <sub>0.57</sub>
	+ tags	63.82 <sub>0.54</sub>	61.69 <sub>0.55</sub>
SBIC	What is...?	66.67 <sub>0.49</sub>	64.33 <sub>0.51</sub>
	+ tags	<b>66.99</b> <sub>0.53</sub>	<b>64.60</b> <sub>0.56</sub>

Table 10: A comparison between SQUAD<sub>T5</sub> prompts with “Is...?” and “What is...?” questions. See §3.1 for more info. We also inspect the effects of adding answer choices and *tags* to the input. Tags are a single word descriptions of the input elements; e.g., E-SNLI’s tags are “premise:” / “hypothesis:” before premise / hypothesis.

<p><b>Sentence1:</b> The stove was cleaned with a cleaner. <b>Sentence2:</b> The stove was cleaned with a mop.  <b>Nonsensical Sentence:</b> Sentence2 <b>Explanation:</b> A mop is too large to clean the stove.</p>
<p><b>Prompt:</b> INFILLING × BASIC  <b>Input:</b> explain sensemaking choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i> &lt;extra_id_0&gt; because &lt;extra_id_1&gt;  <b>Output:</b> &lt;extra_id_0&gt; choice2 &lt;extra_id_1&gt; <i>A mop is too large to clean the stove.</i> &lt;extra_id_2&gt;</p>
<p><b>Prompt:</b> INFILLING × NATURAL SOUNDING  <b>Input:</b> explain sensemaking choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i> It is &lt;extra_id_0&gt; that choice2 is less common because &lt;extra_id_1&gt;  <b>Output:</b> &lt;extra_id_0&gt; True &lt;extra_id_1&gt; <i>A mop is too large to clean the stove.</i> &lt;extra_id_2&gt;</p>
<p><b>Prompt:</b> ≈T5 × COPA  <b>Input:</b> explain sensemaking choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i> Less common is choice2  <b>Output:</b> True because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> SQUAD<sub>T5</sub> × YES/NO + TAGS  <b>Input:</b> explain sensemaking question: Is choice2 more nonsensical? context: choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i>  <b>Output:</b> Yes because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> SQUAD<sub>T5</sub> × WHAT IS...? + TAGS  <b>Input:</b> explain sensemaking question: What is more nonsensical? context: choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i>  <b>Output:</b> choice2 because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> QA<sub>SIMPLE</sub> × YES/NO  <b>Input:</b> explain is choice2 more nonsensical? \n <i>The stove was cleaned with a cleaner. The stove was cleaned with a mop.</i>&lt;/s&gt;  <b>Output:</b> yes because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> QA<sub>SIMPLE</sub> × YES/NO + TAGS  <b>Input:</b> explain is choice2 more nonsensical? \n choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i>&lt;/s&gt;  <b>Output:</b> yes because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> QA<sub>SIMPLE</sub> × YES/NO + TAGS + CHOICES  <b>Input:</b> explain is choice2 more nonsensical? \n (A) yes (B) no \n choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i>&lt;/s&gt;  <b>Output:</b> yes because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> QA<sub>SIMPLE</sub> × WHAT IS...?  <b>Input:</b> explain what is more nonsensical? \n <i>The stove was cleaned with a cleaner. The stove was cleaned with a mop.</i>&lt;/s&gt;  <b>Output:</b> choice2 because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> QA<sub>SIMPLE</sub> × WHAT IS...? + TAGS  <b>Input:</b> explain what is more nonsensical? \n choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i>&lt;/s&gt;  <b>Output:</b> choice2 because <i>a mop is too large to clean the stove.</i></p>
<p><b>Prompt:</b> QA<sub>SIMPLE</sub> × WHAT IS...? + TAGS + CHOICES  <b>Input:</b> explain what is more nonsensical? \n (A) choice1 (B) choice2 \n choice1: <i>The stove was cleaned with a cleaner.</i> choice2: <i>The stove was cleaned with a mop.</i>&lt;/s&gt;  <b>Output:</b> choice2 because <i>a mop is too large to clean the stove.</i></p>

Table 11: COMVE self-rationalization prompts that we design and test. INFILLING marks span-filling prompts; ≈T5 prompts made by following the most similar T5 pretraining task (Table 1); SQUAD<sub>T5</sub> prompts designed following SQUAD’s formatting in T5 pretraining; and QA<sub>SIMPLE</sub> prompts made following UNIFIEDQA. This table shows variations of these prompt types. We refer to spans “choice1:”/“choice2:” as TAGS, and to “(A) yes (B) no”/“(A) choice1 (B) choice2” as CHOICES. YES/NO and WHAT IS...? refer to a question type. Following Hendrycks et al. (2021), we add </s> to the end of our QA<sub>SIMPLE</sub> prompts. More info in §3.