

Self-training with Few-shot Rationalization

Meghana Moorthy Bhat[†] Alessandro Sordoni[‡] Subhabrata Mukherjee[‡]

[†]The Ohio State University [‡]Microsoft Research
bhat.89@osu.edu, {alsordon, submukhe}@microsoft.com

Abstract

While pre-trained language models have obtained state-of-the-art performance for several natural language understanding tasks, they are quite opaque in terms of their decision-making process. While some recent works focus on rationalizing neural predictions by highlighting salient concepts in text as justifications or rationales, they rely on thousands of labeled training examples for both task labels as well as annotated rationales for every instance. Such extensive large-scale annotations are infeasible to obtain for many tasks. To this end, we develop a multi-task teacher-student framework based on self-training language models with limited task-specific labels and rationales, and judicious sample selection to learn from informative pseudo-labeled examples¹. We study several characteristics of what constitutes a good rationale and demonstrate that the neural model performance can be significantly improved by making it aware of its rationalized predictions particularly in low-resource settings. Extensive experiments in several benchmark datasets demonstrate the effectiveness of our approach.

1 Introduction

Recent success in several natural language understanding tasks can be attributed to training large-scale and complex neural network models. While these models work very well for specific tasks, they offer limited insights into their inner working and are often used as black-box predictors. To address these shortcomings, recent works (DeYoung et al., 2020; Paranjape et al., 2020; Yu et al., 2019) have focused on designing interpretable NLP systems that can explain the model’s predictions. A typical approach to study this decision-making process has been via annotating *rationales* as a short and sufficient part of the input text leading to the specific

prediction that can be also used as auxiliary supervision for training. Appropriate use of such rationales can improve the downstream task performance as the model learns to focus on the task-relevant parts of the input (Pruthi et al., 2020b).

However, a significant resource challenge is to obtain large-scale annotated rationales to train these models as explored in fully supervised setting in recent work (DeYoung et al., 2020). This requires models to have access to both instance-level task labels as well as token-level binary rationale labels depicting whether a token should be included in the rationale or not. Such extensive annotations are infeasible to obtain for many tasks, hence devising models that can effectively exploit a limited number of annotated rationales is utterly important. Therefore, our objective is two-fold: (1) improve downstream task performance and (2) improve rationale extraction – with few labeled examples for the downstream task and corresponding rationales.

Recent works (Mukherjee and Awadallah, 2020; Wang et al., 2020; Xie et al., 2019) on few-shot learning have explored *self-training* as a mechanism to train neural network models with limited labeled data. These methods usually train a teacher and then a student model to imitate the teacher in turn. They usually assume access to a set of unlabeled instances and use stochastic regularization techniques such as dropout and data augmentation obtained from pseudo-labeled examples. In this work, we leverage self-training as a mechanism to train neural network models with self-generated rationales and task labels over unlabeled data. Since pseudo-labeled rationales from the teacher model can be noisy, we show that judicious sample selection to upweight informative examples and downweight noisy ones is beneficial. Furthermore, we predict task and rationale labels in a multi-task learning (MTL) setup, where we share parameters between the task objective and the rationale prediction objective. We show that the MTL setup for

¹Code available at <https://aka.ms/RationaleST>

joint learning is more effective than the decoupled learning, which consists of first extracting rationales and then using them for classification, as explored in some of the prior works.

Given the paucity of rationale labels, a critical part in the MTL setup is to understand what constitutes a good rationale. We build over insights from prior work (Yu et al., 2019; Lei et al., 2016) focusing on low-resource settings with access to limited labels via multi-task self-training. To this end, we explore several characteristics of a good rationale in terms of (i) *sufficiency* such that the extracted rationale is adequate for the model to make its decision; (ii) *completeness* such that the model is less confident on its predictions if it ignores the rationale text; and (iii) *coherency* such that the model extracts phrases as rationales rather than isolated disconnected words. In practice, we enforce (i) by matching predictions of the student model with the rationale as input and the teacher model with the full input; (ii) by maximizing entropy in the student predictive distribution when it sees the complementary of the rationale as input; and (iii) by recurring to additional regularization methods. We show that our multi-task joint optimization captures all of the above salient aspects for rationale extraction while improving the downstream task performance. In summary, our contributions are:

- (a) We develop a multi-task self-training framework to train neural models with limited labels along with extracting rationales as justifications for its predictions. Furthermore, we show the impact of judicious sample selection for sample- and token-level re-weighting to learn from informative pseudo-labeled examples during self-training.
- (b) We build over prior work on rationale extraction to encode desired rationale characteristics by judiciously designed loss functions in our multi-task self-training algorithm.
- (c) Extensive experiments on five datasets from the ERASER benchmark (DeYoung et al., 2020) demonstrate the effectiveness of our approach. Ablation experiments demonstrate the impact of different components of our framework.

2 Related Work

Rationale Extraction Prior works (Lei et al., 2016; Yu et al., 2019) on rationale extraction explore encoder-generator based models with two components for extracting rationales and then using them to make a prediction. Alternately, Jain

et al. (2020) propose decoupled architectures for extractor (using attention weights) and predictor. Following these works, DeYoung et al. (2020) develop the ERASER benchmark that contains human annotated rationales in extractive format and provide BERT-to-BERT baselines for the tasks. Paranjape et al. (2020) propose a weakly supervised model with user controlled sparsity threshold for rationale extraction and predictions based on the extracted rationale. Similarly, Pruthi et al. (2020a) propose a semi-supervised BERT-CRF architecture with few gold annotations and abundant task labels. In contrast to all these prior work requiring thousands of annotations for either rationales or the task labels, our framework is geared for low-resource settings with access to very few labels for both the tasks. We incorporate insights from prior work in rationale extraction via judiciously designed loss functions in our multi-task self-training framework.

Self-training Self-training (Yarowsky, 1995; Nigam and Ghani, 2000; Lee, 2013) trains a base (teacher) model on limited labeled data and applies them to unlabeled data to generate pseudo-labels. The generated pseudo-labels are used for training the student model in an iterative fashion. Self-training has demonstrated state-of-the-art performance in several tasks including text classification (Mukherjee and Awadallah, 2020; Wang et al., 2020) and image classification (Xie et al., 2019; Zoph et al., 2020). We leverage self-training with re-weighting noisy pseudo-labels for both task and rationale extraction in a multi-task learning framework while encoding the desired characteristics of a good rationale.

3 Rationale Extraction with Few Labels

3.1 Problem Statement

Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be a set of n documents with corresponding associated task labels y_1, \dots, y_n , where each $\mathbf{x}_i = \{x_{ij}\}$ is a sequence of tokens. We consider each document to be associated with a ground-truth rationale sequence $\mathbf{r}_i = \{r_{ij}\}$, where $r_{i,j} = 1$ if the j^{th} token in document \mathbf{x}_i is part of the rationale, and 0 otherwise. We consider a *low-resource* setup with *very few* documents labeled with both the task labels (instance-level) and the rationale labels (token-level) for each task, and additional unlabeled data.

Let us denote $\mathcal{D}_l = \{(\mathbf{x}_i, y_i, \mathbf{r}_i)\}$ as the joint task and rationale labeled training set. We also

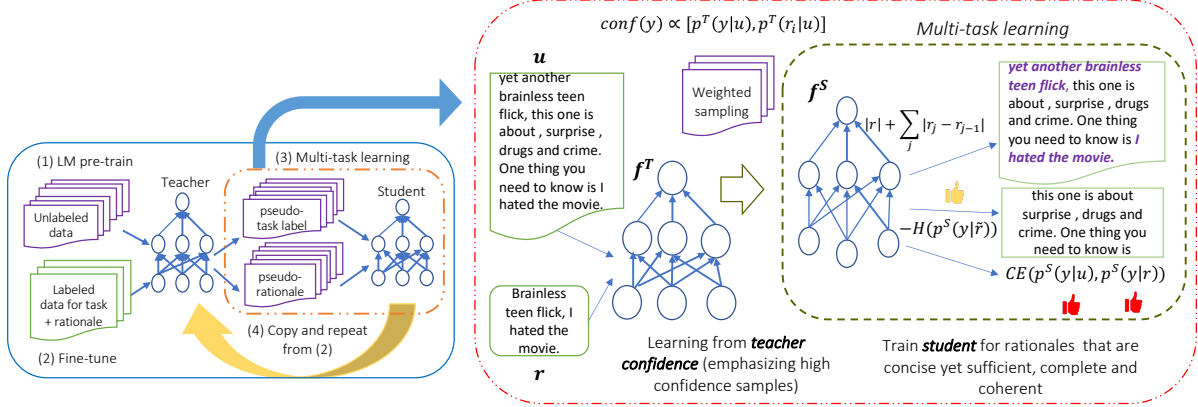


Figure 1: Self-training framework

Algorithm 1: Our self-training algorithm.

Initialize p^T, p^S randomly
while *Not converged* **do**
 $p^T = \min_{p^T} \mathcal{L}_i(p^T)$ # Eq. (1)
 $p^S = \min_{p^S} \mathcal{L}_u(p^S)$ # Eq. (7)
 Copy p^S into p^T .
return p^S

assume to have access to a set of unlabeled documents $\mathcal{D}_u = \{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ for which neither rationales nor task labels are available and $|\mathcal{D}_u| \gg |\mathcal{D}_l|$. Our goal is to learn a model from the few task and token-level rationale labels and additional unlabeled documents to improve its performance for the downstream tasks.

3.2 Self-training

We leverage self-training as the backbone of our framework. The algorithm is composed of two phases that are executed iteratively until convergence. In the first phase, we perform multi-task learning of a teacher model on the few-shot labeled set \mathcal{D}_l by jointly learning to predict instance-level task and token-level rationale labels. Optimizing these losses leads to estimating the parameters of a *teacher* model p^T .

In the second phase, we leverage the teacher model p^T to infer pseudo-labels for the unlabeled set \mathcal{D}_u and train a *student* model p^S to mimic the teacher’s predictions for both the task and associated rationale. Finally, the teacher model is updated with the student model’s parameters and above steps are repeated until convergence.

Due to the noisy nature of pseudo-labels, the above self-training process may result in gradual drifts (Zhang et al., 2016). To address this, we train the student model to explicitly account for

the teacher’s confidence on the generated pseudo-labels with a special weighting scheme. Furthermore, we explore several characteristics of a good rationale and enrich the above framework with additional auxiliary losses.

3.2.1 Multi-Task Teacher

In the first phase, we leverage the small amount of labeled data \mathcal{D}_l to train the teacher model p^T to jointly predict the task labels and the rationale labels. To this end, we leverage a *shared* BERT encoder h^T with two separate softmax classification layers for the two tasks. We denote $p^T(y|\mathbf{x}) = \text{softmax}(h^T(\mathbf{x}); \theta_t^T)$ and $p^T(r_j|\mathbf{x}) = \text{softmax}(h^T(\mathbf{x})_j; \theta_r^T)$ to be the corresponding task and rationale predictions of the teacher model given an instance \mathbf{x} . $h^T(\mathbf{x})_j$ is the BERT hidden state representation corresponding to the j^{th} token and θ_t^T, θ_r^T are the task-specific head parameters. For brevity, we will omit parameter specification in what follows and denote p^T as the shared BERT and task-specific parameters.

We jointly optimize the following losses with respect to (the parameters of) p^T :

$$\mathcal{L}_i(p^T) = \mathbb{E}_{\mathcal{D}_l} \left[\underbrace{-\log p^T(y|\mathbf{x})}_{\text{task prediction}} - \underbrace{\sum_j \log p^T(r_j|\mathbf{x})}_{\text{rationale prediction}} \right], \quad (1)$$

where y is the ground-truth task label for input \mathbf{x} and r_j^T is the ground-truth rationale label for the j^{th} token in input \mathbf{x} .

In contrast to prior work (DeYoung et al., 2020) that leverages two decoupled BERT models in a stage-wise fashion to first extract the rationales, and then use those rationales for task label prediction – our framework learns a single model to predict them jointly in a multi-task learning setup.

This allows the model to capture richer interactions between the two tasks.

There are a few design choices for optimizing the teacher in Eq. 1. For instance, we can optimize the teacher parameters with few-shot labeled data in *each* self-training iteration after the student becomes the new teacher; or *only once* at the beginning of self-training to initialize a good teacher. We observe that executing this phase in each self-training loop is more effective as it diminishes drifting from the ground-truth data distribution.

3.2.2 Multi-Task Student

In the second phase, we self-train a student model p^S on the teacher-generated pseudo-labels with a pseudo-label task loss and rationale loss. In contrast to the teacher model operating on the few labeled examples \mathcal{D}_l , the student model operates on unlabeled data \mathcal{D}_u . Additionally, the student model has the same architecture as the teacher model with a shared encoder h^S and task-specific classification parameters $p^S(y|\mathbf{x}) = \text{softmax}(h^S(\mathbf{x}); \theta_t^S)$ and $p^S(r|\mathbf{x}) = \text{softmax}(h^S(\mathbf{x}); \theta_r^S)$. The pseudo-labeled multi-task loss is formulated as:

$$\mathcal{L}_u(p^S) = \mathbb{E}_{\mathbf{u} \sim \mathcal{D}_u, \mathbf{r}^T \sim p^T(\mathbf{r}|\mathbf{u}), y^T \sim p^T(y|\mathbf{u})} \left[\begin{aligned} & - \underbrace{\log p^S(y^T|\mathbf{u})}_{\text{task pseudo prediction}} - \underbrace{\sum_j \log p^S(r_j^T|\mathbf{u})}_{\text{rationale pseudo prediction}} \end{aligned} \right] \quad (2)$$

where y^T is the teacher-generated task pseudo-label for input \mathbf{u} and r_j^T is the teacher-generated rationale pseudo-label for the j^{th} token in input \mathbf{u} .

3.2.3 Student-Teacher Update

At the end of every self-training iteration, we transfer the knowledge acquired by the student back into the teacher model by setting $h^T = h^S$, $\theta_t^T = \theta_t^S$, $\theta_r^T = \theta_r^S$ and start again by fine-tuning the newly obtained teacher on ground-truth data \mathcal{D}_l .

3.2.4 Re-weighting Pseudo-labeled Samples

Instead of directly imitating the teacher’s predictions as described in Eq. 2, we found it extremely effective to train the student model to explicitly account for the teacher’s confidence for the generated pseudo-labels. This allows us to filter noisy pseudo-labels as the student model can selectively focus more on the pseudo-labeled samples that the teacher is more confident on compared to the less certain ones. Therefore, we optimize a weighted

version of the pseudo-labeled loss in Eq. 2:

$$\mathcal{L}_{wu}(p^S) = \mathbb{E}_{\mathbf{u} \sim \mathcal{D}_u, \mathbf{r}^T \sim p^T(\mathbf{r}|\mathbf{u}), y^T \sim p^T(y|\mathbf{u})} \left[\begin{aligned} & - \underbrace{w_y^T(\mathbf{u})}_{\text{task weight}} \cdot \log p^S(y^T|\mathbf{u}) \\ & - \sum_j \underbrace{w_{r,j}^T(\mathbf{u})}_{\text{rationale weight}} \cdot \log p^S(r_j^T|\mathbf{u}) \end{aligned} \right], \quad (3)$$

where $w_y^T(\mathbf{u}) \propto p^T(y = y^T|\mathbf{u})$ and $w_{r,j}^T(\mathbf{u}) \propto p^T(r_j = r_j^T|\mathbf{u})$. The proportional sign is due to the fact that these weights are normalized across each batch when training by minibatch SGD, so the weights *depend on the batch* and sum to one over the batch. Re-weighting noisy labels with different weighting schemes has been explored with meta-learning (Ren et al., 2018) and uncertain-aware self-training (Mukherjee and Awadallah, 2020).

3.3 Rationale Characteristics

In this section, we encode several characteristics of what constitutes a good rationale from prior work in our self-training framework via several auxiliary loss functions.

3.3.1 Sufficiency

A desired property of a good rationale is *sufficiency*. This imposes the model predictions about the task label considering the entire input text to be similar to the predictions made by looking at only the rationale text. This concept can be promptly translated into an objective function by resorting to a consistency objective:

$$\mathcal{L}_{suff}(p^S) = \mathbb{E} \left[-\log p^S(y = y^T|\mathbf{u} \odot \mathbf{r}^T) \right], \quad (4)$$

where the expectation is taken w.r.t $\mathbf{u} \sim \mathcal{D}_u$, $\mathbf{r}^T \sim p^T(\mathbf{r}|\mathbf{u})$, $y^T \sim p^T(y|\mathbf{u})$ and $\mathbf{u} \odot \mathbf{r}^T$ is the masked version of document \mathbf{u} in which tokens that are not part of the rationale (as predicted by the teacher) are replaced with a special [MASK] token. Here, the teacher model looks at the full input and the student model looks only at the rationale tokens.

The sufficiency loss can be interpreted as an alternative way of integrating rationale information in the model. Current efforts either predict the rationale first and then use it for task prediction sequentially as in BERT-to-BERT (DeYoung et al., 2020); or employ attention regularization such that the BERT attention weights are as close as possible

to uniform on the rationale tokens (Pruthi et al., 2020b). The first approach can be very sensitive to error propagation from the rationale generator since the task label is predicted using only the generated rationales at test-time. The second approach strictly assumes uniform attention on the rationale tokens. In contrast, our sufficiency loss makes very few assumptions on how the model should attend to the rationale tokens, and only requires the student distribution of the task labels to be similar to that of the teacher. This yields more robustness to rationale errors, given that at test-time our model can use the full input to predict the task label.

3.3.2 Completeness

Another desiderata of a rationale is *completeness*. Completeness implies that the rationale should capture all the aspects in the input text that are predictive of the task label. We translate this concept by requiring the student model to be *maximally uncertain* of the task label if it does not look at the rationale i.e. by masking out the teacher-predicted rationale tokens in the input text:

$$\mathcal{L}_{comp}(p^S) = \mathbb{E} \left[-H(p^S(y|\mathbf{u} \odot (1 - \mathbf{r}^T))) \right] \quad (5)$$

where the expectation is w.r.t. $\mathbf{u} \sim \mathcal{D}_u, \mathbf{r}^T \sim p^T(\mathbf{r}|\mathbf{u})$, H is the entropy of the student predictive distribution and $\mathbf{u} \odot (1 - \mathbf{r}^T)$ is the document obtained by masking out the tokens in the rationale.

3.3.3 Coherence Loss

Finally, we desire the rationales to be short and composed of contiguous chunks of text rather than unigrams. To ensure this, we adopt the regularization losses introduced in Lei et al. (2016). This explicitly penalizes the rationale generator for predicting long rationales and encourages rationales to span contiguous chunks of the input text.

$$\mathcal{L}_{co}(p^S) = \mathbb{E}_{\mathbf{u}, \mathbf{r}^S \sim p^S(\mathbf{r}|\mathbf{u})} \left[|\mathbf{r}^S| + \sum_j |r_j^S - r_{j-1}^S| \right] \quad (6)$$

3.4 Training Objectives

Our overall training objective in the teacher learning phase is simply the loss on labeled data $\mathcal{L}_l(p^T)$. For the student, we use a combination of the previously presented loss functions on unlabeled data:

$$\mathcal{L}_u(p^S) = \mathcal{L}_{wu} + \mathcal{L}_{suff} + \mathcal{L}_{comp} + \mathcal{L}_{co} \quad (7)$$

4 Experimental Setup

Datasets We evaluate our framework on five different tasks from the ERASER benchmark (DeYoung et al., 2020). These include Movies for sentiment analysis (Pang and Lee, 2004), e-SNLI (Camburu et al., 2018) for natural language inference, FEVER (Thorne et al., 2018) for fact extraction and verification, BoolQ (Clark et al., 2019) for reading comprehension, and Evidence Inference (Lehman et al., 2019) over scientific articles for medical interventions. All datasets except e-SNLI contain text spans or sentence level annotations as rationale. Following prior works (DeYoung et al., 2020; Pruthi et al., 2020b; Paranjape et al., 2020) we report F1 measures for both the task and token-level rationale extraction performance.

Dataset	#Class	#Full Train	#Validation	#Test
Movies	2	1.6K	200	200
e-SNLI	3	547K	9.8K	9.8K
FEVER	2	98K	6.1K	6.1K
BoolQ	2	6.3K	1.4K	2.8K
Evidence	3	7.9K	972	972

Table 1: Dataset statistics.

Base encoder Following prior work (DeYoung et al., 2020), we adopt BERT-base (Devlin et al., 2019) as the base encoder for our experiments. We use a maximum sequence length of 512 with a batch size of 8. We use the Adam (Kingma and Ba, 2017) optimizer with a learning rate of $3e - 5$. Following (Paranjape et al., 2020), we use domain based checkpoint (BioBERT²) for Evidence Inference. For Movies, we use a BERT checkpoint pre-trained on IMDB reviews. For all other datasets, we initialize the models from publicly available checkpoints³. Due to the question-answering format of tasks in BoolQ, Evidence Inference and FEVER, the query or claim with supporting documents are encoded in the form document [SEP] query.

Few labels setup We consider a resource-constrained setting with only $N = 100$ labeled samples per class and their associated rationales for each task, which are randomly sampled from the training data. We use the rest of the examples in the training set as unlabeled examples by ignoring their labels as in prior work (Mukherjee and

²https://huggingface.co/bionlp/bluebert_pubmed_uncased_L12_H768_A12

³<https://huggingface.co/bert-base-uncased>

Model	Movies		ESNLI		FEVER		BoolQ		Evidence		Avg. Task.
	Task.	Tok.F1	Task.	Tok.F1	Task.	Tok.F1	Task.	Tok.F1	Task.	Tok.F1	
<i>fully supervised</i>											
BERT w/o explanation	91	–	90	–	91	–	62	–	53	–	77
BERT with explanation	93	41	90	30	91	80	65	50	54	51	79
BERT-to-BERT	87	15	73	70	88	81	54	13	71	47	75
Pruthi et al. (2020b)	84	–	–	–	–	–	–	–	–	–	–
Paranjape et al. (2020) (25% rationales)	85	28	–	–	89	64	63	19	46	10	71
<i>few-shot (100 labels per class)</i>											
BERT w/o explanation	82	–	59	–	69	–	52	–	41	–	61
BERT with explanation	82	18	61	5	70	15	56	16	41	30	62
Our model	86	21	67	15	74	51	61	18	43	19	66
<i>few-shot (25% of training data)</i>											
BERT w/o explanation	83	–	86	–	88	–	59	–	43	–	72
BERT with explanation	82	26	86	27	89	32	60	31	44	8	72
Our model	87	31	87	36	89	39	62	51	46	9	74

Table 2: Comparison of task F1 (Task.) and rationale token F1 (Tok. F1) of our model with baselines. Avg. denotes the aggregated task performance averaged across all tasks. Our model includes BERT with explanation and rationale characteristics with self-training. Baseline numbers are reported from corresponding papers.

Awadallah, 2020). We perform early stopping on validation performance for the teacher multi-task training in each loop of self-training. We select the best model based on validation loss from the different self-training iterations. For Evidence, we use a hyper-parameter⁴ to upweight examples from the minority pseudo-labeled class in the self-training loop to combat class imbalance.

Baselines We compare our model performance to the following methods for both *fully supervised* (with access to all training labels) and our *few-shot* setting with 100 labels per class. (1) *BERT w/o explanation* fine-tunes BERT on a set of labeled examples without accounting for rationales. (2) *BERT with explanation* is our multi-task learning setup where the classifiers are trained to predict both task labels and rationales *without* encoding the rationale characteristics. We also compare against the *semi-supervised* setting from prior work Paranjape et al. (2020) that uses 25% annotations for rationales and 100% task labels for each dataset.

Finally, we compare our method against a fully supervised BERT-to-BERT model in the ERASER benchmark (DeYoung et al., 2020) which first performs rationale extraction followed by downstream task prediction using extracted rationale.

5 Results

Overall performance Table 2 summarizes the results of our proposed model and baselines across all the datasets. We observe that our model trained

⁴Set as inversely proportional to the count of pseudo-labeled samples per class

with only $N = 100$ labels per class performs within 10.6% of fully supervised BERT trained with thousands of labels while obtaining an aggregate F1 of 66%. Our self-training framework iteratively improves over the teacher model with a judiciously designed student network with a performance gain of 6.45%.

For the fully supervised models, we observe that BERT with explanation in our multi-task learning setup improves by 2.6% over vanilla BERT, thereby, demonstrating the usefulness of rationales and the effectiveness of multi-task learning.

Our model performance for downstream tasks and rationales further improves with increasing the amount of labeled data. With 25% labeled training data, our model performs at par with Paranjape et al. (2020) (which has been trained with 100% task labels) in 3 out of 4 tasks and has comparable performance with fully supervised BERT-to-BERT baseline. Additionally, our model does not require additional user inputs in the form of desired sparsity threshold as in prior work (Paranjape et al., 2020).

In order to validate the effectiveness of sufficiency and completeness loss as a way of integrating rationales in the model in contrast to attention regularization in Pruthi et al. (2020b), we perform the following experiment. Instead of using teacher-predicted rationales, we use the ground-truth rationales. We observe that our model with sufficiency and completeness loss outperforms the prior method of integrating explanations via attention regularization (results in Figure 3 in Appendix). This may reflect the fact that attention regulariza-

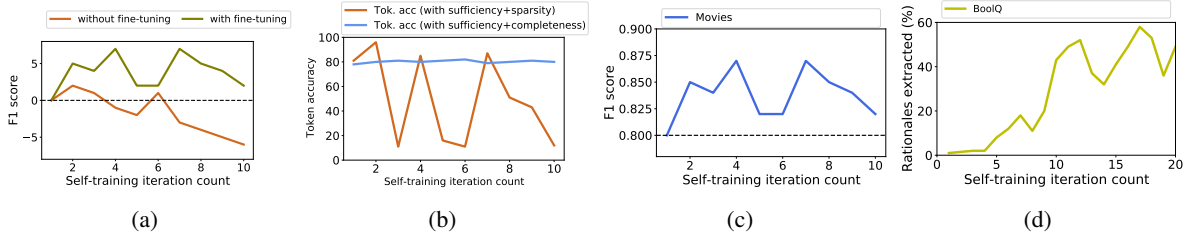


Figure 2: Graphs representing model performance from different ablation study experiments: (a) Impact of fine-tuning with labeled data in self-training loop (b) Ablation study showing impact of adding sparsity loss without completeness, where we observe inconsistent rationale extraction in self-training. (c) Accuracy of Movies dataset across self-training iterations (d) Percentage of rationales extracted for BoolQ.

Model	Movies			e-SNLI			
	BLEU-2	Token P/R/F1	Task.	BLEU-2	Token P/R/F1	Task.	
teacher	13.0	24.1/ 50.1/ 32.5	67	0	0.7/ 5.2/ 1.4	56	
sufficiency	10.1	22.1/ 93.1/ 35.7	69	0.98	10.0/ 12.1/ 10.0	61	
sufficiency + re-weight labels	9.3	9.1/ 96.4/ 16.6	70	2.0	34.4/ 14.6/ 20.5	65	
sufficiency + re-weight labels and rat.	12.2	19.6/ 92.2/ 32.3	71	1.3	33.1/ 10.7/ 16.2	65	
sufficiency + sparsity + re-weight labels and rat.	14.1	19.8/ 61.3/ 29.9	70	1.1	19.3/ 15.7/ 17.3	65	
sufficiency + sparsity + completeness + re-weight labels and rat.	14.0	24.2/ 50.3/ 32.3	71	1.07	12.8/ 7.6/ 5.7	65	
Our model - re-weight labels and rat.	12.7	21.2/ 47.6/ 29.3	71	1.6	15.6/ 11.4/ 13.17	65	
Our model - coherence loss	13.1	25.6/ 44.1/ 31.8	73	6.1	16.4/ 23/ 18.3	66	
Our model	14.1	28.4/ 41.3/ 33.5	75	5.0	17.4/ 13.8/ 15.3	67	

Table 3: Ablation study of various components in our model with corresponding performance in terms of BLEU-2, Token F1 and Task F1 with sequence length 128.

tion imposes the stricter assumption to have uniform attention on the explanation tokens.

We observe the self-training performance to improve on initializing our encoders with pre-trained domain-specific checkpoints. For instance, using BioBERT checkpoint, we observe a gain of 4.8% over that of BERT-base in the Evidence dataset. Table 4 presents few examples from our rationale extractor and the corresponding task label.

5.1 Ablation Study

Table 3 summarizes the impact of different components of our framework with $N = 100$ labels per class as training data for Movies and e-SNLI.

Re-weighting pseudo-labels We found it extremely useful to re-weight noisy pseudo-labeled samples from the teacher model by its confidence. We observe that re-weighting the rationales and task labels work quite well for the Movies dataset. However, this has limited impact for e-SNLI with a low coverage of rationales from the teacher model.

Correspondingly, we did not observe a difference in model performance from re-weighting the task and rationale pseudo-labels.

Impact of different loss functions. We observe that using only the sufficiency loss (Eq. 4) results in the model extracting the entire input as the rationale. This is counteracted by adding penalization via sparsity loss to obtain rationales that are concise yet informative about the task label. From Table 3, we observe sparsity loss to significantly reduce the number of tokens included in the rationale.

However, adding the sparsity loss also caused instability in some cases (Figure 2 (b)). We empirically demonstrate that this instability is mitigated by including the completeness loss that forces the model to be maximally uncertain when it does not look at important tokens constituting the rationale.

Impact of amount of labeled training data We observe our self-training framework to improve both in task and rationale extraction performance

Dataset: Movies	Ground truth: Negative, Prediction: Negative
<p>There're so many things to criticize about I don't know where to start. Recommendation: turn off your brain - don't be like me, decreasing the rating everyday because I think about it too much Firstly, there is nothing outstandingly inferior about the making of the film (nor is there anything outstandingly good about it), but the plot holes make the film corny and stupid.</p>	
Dataset: Movies	Ground truth: Negative, Prediction: Negative
<p>Yet another brainless teen flick, stars Katie Holmes and Sarah Polly couldn't look more bored. One thing you need to know is I really hated this movie. Everything about it annoyed the hell out of me. The acting, and script, the plot, and ending.</p>	
Dataset: e-SNLI	Ground truth: Contradiction, Prediction: Contradiction
<p>A man playing electric guitar on the stage . A man playing banjo on the floor .</p>	

Table 4: Snapshot of correctly classified examples and the corresponding rationales extracted by our model on Movies and e-SNLI.

Dataset: Movies	Ground truth: Positive, Prediction: Negative
<p>Well I'll be damned ... the Canadians can make a good movie. The world is coming to an end. We don't know why or how , but apparently there is no way to stop it most of the rioting and other assorted chaos has passed. director - star Don Mckellar has crafted a highly unique and emotional film. All of the main characters are compelling as they try and do whatever it is.....</p>	
Dataset: Movies	Ground truth: Positive, Prediction: Negative
<p>When I first saw the previews for Ron Howard's latest film, my expectations were discouragingly low. A show about nothing? A guy whose entire life is broadcast 24 hours a day ?.... which is why I was pleasantly surprised by "edtv," which turns out to be a fresh , insightful , and often times hilarious film about the follies of instant celebrity .</p>	
Dataset: e-SNLI	Ground truth: Entailment, Prediction: Neutral
<p>A woman tired from her long day takes a nap on her bed above the sheets and covers. A lady is lying in bed.</p>	

Table 5: Snapshot of mis-classified examples from our model on Movies and e-SNLI dataset.

Legend: Ground-truth rationales not detected by the model Rationales present in both the model and ground-truth Rationales extracted by the model but absent in ground-truth

with increase in number of labeled samples for training (Appendix, Figure 4).

Impact of number of self-training iterations

Figure 2 (c) shows the improvement in the task accuracy of our model over several self-training iterations for Movies. The corresponding plots for other datasets are provided in Figure 6 in Appendix. We observe that self-training gradually improves model performance in the first few iterations for majority of the tasks and converges fast in 12-15 iterations. However, we observe that rationale extraction module drifts after 10 self-training iterations for most of the datasets (Figure 2 (d)) due to error propagation from noisy pseudo-labels, thereby, necessitating early stopping based on rationale validation loss.

Few labeled data fine-tuning At each self-training iteration, the teacher is fine-tuned on labeled data. This is to avoid drifting from the original task description via few annotated labels (Figure 2 (a)). Figure 2 demonstrates the change in the accuracy of our student model with and without this teacher training in every self-training iteration.

5.2 Error Analysis

Table 5 presents a snapshot of the qualitative error analysis of our model. On analyzing the extracted rationales for mis-classified instances, we observe some common failure points where instances have shifts in context; presence of satire or sarcasm; rationales relying on background knowledge; and noisy or incomplete annotated rationales. For instance, the overall polarity of the first example from Movies is positive, although majority of the text

describing the movie plot depicts a negative connotation. We observe a similar trend with reviews involving sarcasm or satire. In the last example from e-SNLI, the annotators marked *{woman, nap}* to entail *{lady, lying in bed}*. In this rationale, annotators do not follow the guidelines for sufficiency and completeness since the spatial qualifier for *nap* is missing in the ground-truth. Surprisingly, our model does not pick up the spatial concept either and marks the sequences to be *neutral* to each other.

6 Conclusion

We develop a multi-task self-training framework for rationale extraction focusing on low-resource settings with access to very few training labels. To this end, we build on insights from prior work on the characteristics of a good rationale to encode them via judiciously designed loss functions in our self-training framework. Extensive experiments on benchmark datasets show our model to outperform other state-of-the-art methods with access to limited labels. We further demonstrate that the performance of pre-trained language model can be improved by making it aware of the rationales for its decision-making process in both high (fully supervised) and low-resource (few label) settings.

References

- 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*, volume 31. Curran Associates, Inc.
- 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.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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.
- Sarthak Jain, Sarah Wiegrefe, 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.
- Diederik P. Kingma and Jimmy Ba. 2017. [Adam: A method for stochastic optimization](#).
- D. Lee. 2013. Pseudo-label : The simple and efficient semi-supervised learning method for deep neural networks.
- 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.
- Subhabrata Mukherjee and Ahmed Awadallah. 2020. [Uncertainty-aware self-training for few-shot text classification](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 21199–21212. Curran Associates, Inc.
- Kamal Nigam and Rayid Ghani. 2000. [Analyzing the effectiveness and applicability of co-training](#). In *Proceedings of the Ninth International Conference on Information and Knowledge Management, CIKM '00*, page 86–93, New York, NY, USA. Association for Computing Machinery.
- Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *ACL*.
- 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.
- Danish Pruthi, Bhuwan Dhingra, Graham Neubig, and Zachary C. Lipton. 2020a. [Weakly- and semi-supervised evidence extraction](#). In *Findings of the Association for Computational Linguistics: EMNLP*

2020, pages 3965–3970, Online. Association for Computational Linguistics.

Danish Pruthi, Bhuwan Dhingra, Livio Baldini Soares, Michael Collins, Zachary C. Lipton, Graham Neubig, and William W. Cohen. 2020b. [Evaluating explanations: How much do explanations from the teacher aid students?](#)

Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. 2018. [Learning to reweight examples for robust deep learning](#). In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 4331–4340. PMLR.

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.

Yaqing Wang, Subhabrata Mukherjee, Haoda Chu, Yuancheng Tu, Ming Wu, Jing Gao, and Ahmed Hassan Awadallah. 2020. [Adaptive self-training for few-shot neural sequence labeling](#). *CoRR*, abs/2010.03680.

Qizhe Xie, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. 2019. [Self-training with noisy student improves imagenet classification](#). *CoRR*, abs/1911.04252.

David Yarowsky. 1995. [Unsupervised word sense disambiguation rivaling supervised methods](#). In *33rd Annual Meeting of the Association for Computational Linguistics*, pages 189–196, Cambridge, Massachusetts, USA. Association for Computational Linguistics.

Mo Yu, Shiyu Chang, Yang Zhang, and Tommi Jaakkola. 2019. [Rethinking cooperative rationalization: Introspective extraction and complement control](#). 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 4094–4103, Hong Kong, China. Association for Computational Linguistics.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2016. [Understanding deep learning requires rethinking generalization](#). Cite arxiv:1611.03530Comment: Published in ICLR 2017.

Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin D. Cubuk, and Quoc V. Le. 2020. [Rethinking pre-training and self-training](#).

A Appendix

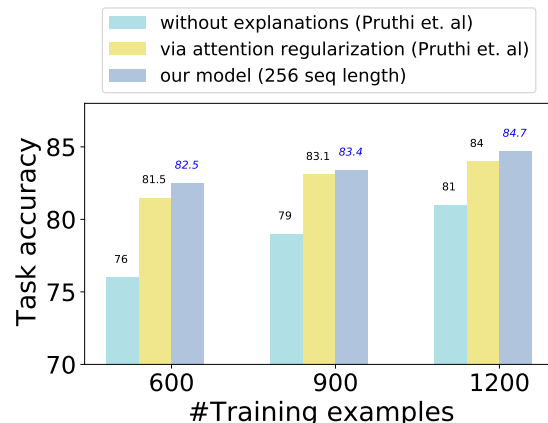


Figure 3: Integrating rationales with sufficiency \mathcal{L}_{suff} and completeness \mathcal{L}_{comp} losses outperforms the recent attention regularization method (Pruthi et al., 2020b). To be comparable, we use only labeled examples, we don’t do any self-training (no unlabeled examples), and we use max sequence length of 256.

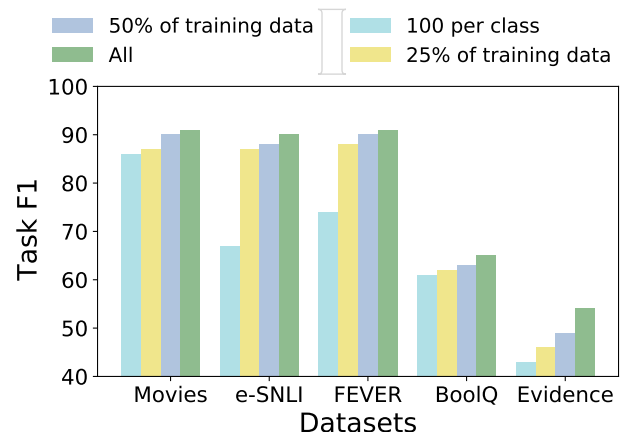


Figure 4: Task F1 across varying percentage of training data across datasets

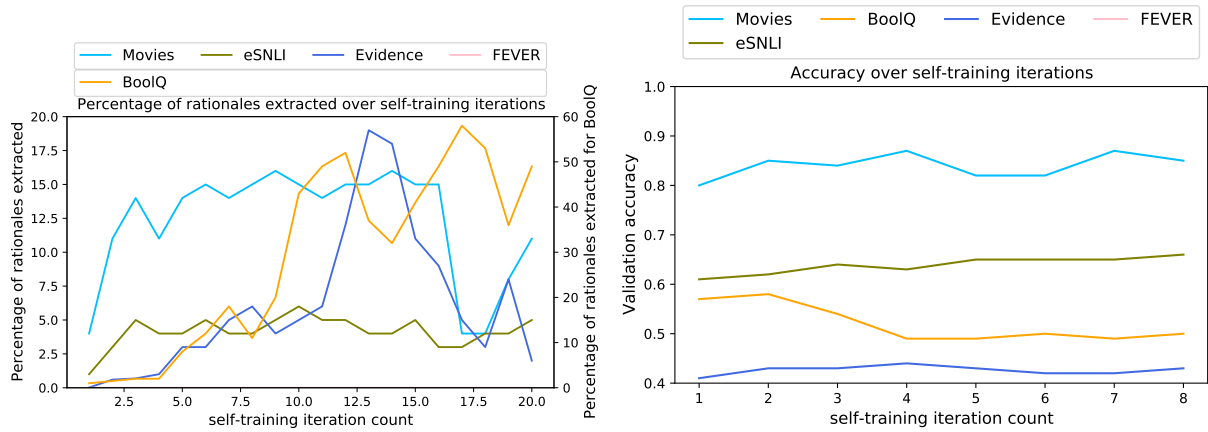


Figure 5: Percentage of rationales extracted over self-training iterations

Figure 6: Accuracy of all the datasets across self-training iterations.

Dataset: Movies	Ground truth: Negative, Prediction: Positive
<p>Though good - looking , its lavish sets , fancy costumes and luscious cinematography can do little to compensate for the emotional wasteland... this is Jodie Foster 's first movie since the jaw - droppingly brilliant contact came out more than two years ago and it isn't the best choice to show off her acting chops .</p>	
Dataset: e-SNLI	Ground truth: Entailment, Prediction: Neutral
<p>A woman tired from her long day takes a nap on her bed above the sheets and covers. A lady is lying in bed.</p>	
Ground truth: Contradiction, Prediction: Entailment	
<p>A mountainous photo is complete with a blue sky. The photo was taken on a cloudy night.</p>	

Table 6: Snapshot of mis-classified examples from Movies and eSNLI.

Legend Ground-truth rationales not detected by the model , Rationales extracted by the model but absent in ground-truth. Rationales present in both the model and ground-truth.