Rationales for Sequential Predictions

Keyon Vafa

Columbia University

keyon.vafa@columbia.edu

David M. Blei

Columbia University

david.blei@columbia.edu

Abstract

Sequence models are a critical component of modern NLP systems, but their predictions are difficult to explain. We consider model explanations though rationales, subsets of context that can explain individual model predictions. We find sequential rationales by solving a combinatorial optimization: the best rationale is the smallest subset of input tokens that would predict the same output as the full sequence. Enumerating all subsets is intractable, so we propose an efficient greedy algorithm to approximate this objective. The algorithm, which is called greedy rationalization, applies to any model. For this approach to be effective, the model should form compatible conditional distributions when making predictions on incomplete subsets of the context. This condition can be enforced with a short finetuning step. We study greedy rationalization on language modeling and machine translation. Compared to existing baselines, greedy rationalization is best at optimizing the sequential objective and provides the most faithful rationales. On a new dataset of annotated sequential rationales, greedy rationales are most similar to human rationales.

1 Introduction

Sequence models are a critical component of generation tasks ranging from language modeling to machine translation to summarization. These tasks are dominated by complex neural networks. While these models produce accurate predictions, their decision making processes are hard to explain. Interpreting a model's prediction is important in a variety of settings: a researcher needs to understand a model to debug it; a doctor using a diagnostic model requires justifications to validate a decision; a company deploying a language model relies on model explanations to detect biases appropriated from training data.

Interpretation takes many flavors (Lipton, 2018). We focus on *rationales*, i.e. identifying the most

Yuntian Deng

Harvard University

dengyuntian@seas.harvard.edu

Alexander M. Rush

Cornell Tech

arush@cornell.edu

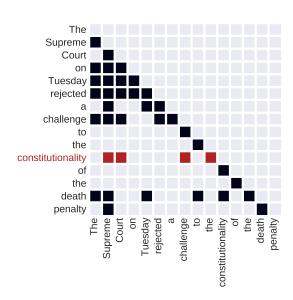


Figure 1. Rationales for sequential prediction on GPT-2. Each row is a predicted word. The dark cells correspond to the context words found by greedy rationalization. To predict "constitutionality", the model only needs "Supreme", "Court", "challenge", and "the".

important subset of input tokens that leads to the model's prediction. For example, consider the sentence: "The Supreme Court on Tuesday rejected a challenge to the constitutionality of the death penalty." Suppose we would like to explain the decision of the model to generate "constitutionality". While the model mathematically conditions on all the previous words, only some are necessary for its predictions. In this case, the rationale produced by our algorithm includes "the", "challenge", and notably "Supreme Court", but not phrases that add no information like "on Tuesday" (Figure 1).

Various rationale methods have been proposed for sequence classification, where each sequence has a single rationale (Lei et al., 2016; Chen et al., 2018; Jain et al., 2020). However, these methods cannot scale to sequence models, where each token in a sequence requires a different rationale.

This work frames the problem of finding se-

quence rationales as a combinatorial optimization: given a model, the best rationale is the smallest subset of input tokens that would predict the same token as the full sequence. Finding the global optimum in this setting is intractable, so we propose **greedy rationalization**, a greedy algorithm that iteratively builds longer rationales. This approach is efficient for many NLP models such as transformers (Vaswani et al., 2017). Moreover, it does not require access to the inner workings of a model, such as gradients.

Underlying this approach is an assumption that the model forms sensible predictions for incomplete subsets of the input. Although we can pass in incomplete subsets to neural models, there is no guarantee that their predictions on these subsets will be compatible with their predictions on full contexts (Arnold and Press, 1989). We show that compatibility can be learned by conditioning on randomly sampled context subsets while training a model. For large pretrained models like GPT-2 (Radford et al., 2019), fine-tuning is sufficient.

In an empirical study, we compare greedy rationalization to various gradient- and attention-based explanation methods on language modeling and machine translation. Greedy rationalization best optimizes the objective, and its rationales are most faithful to the inner workings of the model. We additionally create a new dataset of annotated rationales based on the Lambada corpus (Paperno et al., 2016). We find that greedy rationales are most similar to human annotations, both on our dataset and on a labeled dataset of translation alignments.

Our code and annotated dataset are available.¹

2 Sequential Rationales

Consider a sequence of tokens, $y_{1:T}$, generated by some unknown process $y_{1:T} \sim F$. The goal of sequence modeling is to learn a probabilistic model p_{θ} that approximates F from samples. Maximum-likelihood estimation is an effective way to train these models, where θ is fit according to

$$\underset{\theta}{\operatorname{arg max}} \mathbb{E}_{y_{1:T} \sim F}[\log p_{\theta}(y_{1:T})]. \tag{1}$$

Sequence models are typically factored into conditional distributions:

$$p_{\theta}(y_{1:T}) = f_{\theta}(y_1) \prod_{t=2}^{T} f_{\theta}(y_t|y_{< t}).$$
 (2)

Here, f_{θ} is the specific model parameterizing p_{θ} , such as a transformer (Vaswani et al., 2017), and is trained to take inputs $y_{< t}$. Going forward, we drop the dependence on θ in the notation.

Word-level explanations are a natural way to interpret a sequence model: which words were instrumental for predicting a particular word? Would the same word have been predicted if some of the words had been missing?

Explanations may be straightforward for simpler models; for example, a bigram Markov model uses only the previously generated word to form predictions. However, the most effective sequence models have been based on neural networks, whose predictions are challenging to interpret (Lipton, 2018).

Motivated by this goal, we consider a sequence $y_{1:T}$ generated by a sequence model p. At each position t, the model takes the inputs in the context $y_{< t}$ and uses them to predict y_t . We are interested in forming *rationales*: subsets of the contexts that can explain the model's prediction of y_t .²

What are the properties of a good rationale? Any of the contextual words $y_{< t}$ can contribute to y_t . However, if a model makes the same prediction with only a subset of the context, that subset contains explanatory power on its own. A rationale is *sufficient* if the model would produce the same y_t having seen only the rationale (DeYoung et al., 2020). While rationales consisting of the full context would always be sufficient, they would be ineffective for explaining longer sequences. Intuitively, the smaller the rationale, the easier it is to interpret, so we also prioritize *brevity*.

We combine these desiderata and frame finding rationales as a combinatorial optimization: the best rationale of a word y_t is the smallest subset of inputs that would lead to the same prediction. Each candidate rationale S is an index set, and y_S denotes the subset of tokens indexed by S. Denote by $S = 2^{[t-1]}$ the set of all possible context subsets. An optimal rationale is given by

$$\underset{S \in \mathcal{S}}{\operatorname{arg\,min}} |S| \text{ s.t. } \underset{y'_t}{\operatorname{arg\,max}} p(y'_t|y_S) = y_t. \quad (3)$$

The constraint guarantees sufficiency, and the objective targets brevity. Although the objective may have multiple solutions, we only require one.

Optimizing Equation 3 is hindered by a pair of computational challenges. The first challenge

https://github.com/keyonvafa/
sequential-rationales

²Our paradigm and method extend easily to conditional sequence models, such as those used for machine translation. For full details, refer to Appendix A.

is that solving this combinatorial objective is intractable; framed as a decision problem, it is NP-hard. We discuss this challenge in Section 3. The second challenge is that evaluating distributions conditioned on incomplete context subsets $p(y_t'|y_S)$ involves an intractable marginalization over missing tokens. For now we assume that $f(y_t'|y_S) \approx p(y_t'|y_S)$; we discuss how to enforce this condition in Section 4.

3 Greedy Rationalization

We propose a simple greedy algorithm, **greedy rationalization**, to approximate the solution to Equation 3. The algorithm starts with an empty rationale. At each step, it considers adding each possible token, and it selects the one that most increases the probability of y_t . This process is repeated until the rationale is sufficient for predicting y_t . Figure 2 provides an overview.

Here is the algorithm. Begin with a rationale $S^{(0)} = \emptyset$. Denoting by $[t-1] = \{1, \dots, t-1\}$, the first rationale set is

$$S^{(1)} = \arg\max_{k \in [t-1]} p(y_t|y_k). \tag{4}$$

At each step, we iteratively add a single word to the rationale, choosing the one that maximizes the probability of the word y_t :

$$S^{(n+1)} = S^{(n)} \cup \underset{k \in [t-1] \backslash S^{(n)}}{\arg \max} p(y_t | y_{S^{(n)} \cup k}).$$
 (5)

We continue iterating Equation 5 until $\arg\max_{y_t'} p(y_t'|y_{S^{(n)}}) = y_t$. The procedure will always converge, since in the worst case, $S^{(t-1)}$ contains the full context.

The greedy approach is motivated by approximations to the set cover problem (Chvatal, 1979). In our setting, each set is a single context token, and a rationale "covers" a sequence if it results in predicting the same token.

This procedure is simple to implement, and it is black-box: it does not require access to the inner workings of a model, like gradients or attention.

While greedy rationalization can be applied to any model, greedy rationalization is particularly effective for set-based models such as transformers. If we assume the rationale size m=|S| is significantly shorter than the size of the context t, greedy rationalization requires no extra asymptotic complexity beyond the cost of a single evaluation.

For transformers, the complexity of each evaluation $f(y_t|y_{< t})$ is quadratic in the input set $O(t^2)$.

Each step of greedy rationalization requires evaluating $f(y_t|y_S)$, but y_S can be significantly smaller than $y_{< t}$. A rationale of size m will require m steps of O(t) evaluations to terminate, resulting in a total complexity of $O(m^3t)$. As long as $m = O(t^{1/3})$, greedy rationalization can be performed with the same asymptotic complexity as evaluating a transformer on the full input, $O(t^2)$. In Appendix C, we verify the efficiency of greedy rationalization.

4 Model Compatibility

Greedy rationalization requires computing conditional distributions $p(y_t|y_S)$ for arbitrary subsets S. Using an autoregressive model, this calculation requires marginalizing over unseen positions. For example, rationalizing a sequence $y_{1:3}$ requires evaluating the candidate rationale $p(y_3|y_1)$, which marginalizes over the model's predictions:

$$p(y_3|y_1) = \sum_{k} f(y_3|y_1, y_2 = k) f(y_2 = k|y_1).$$

Given the capacity of modern neural networks, it is tempting to pass in incomplete subsets y_S to f and evaluate this instead as $f(y_t|y_S) \approx p(y_t|y_S)$. However, since f is trained only on complete feature subsets $y_{< t}$, incomplete feature subsets y_S are out-of-distribution (Hooker et al., 2019). Evaluating $f(y_3|y_1)$ may be far from the true conditional $p(y_3|y_1)$. In Figure 4, we show that indeed language models like GPT-2 produce poor predictions on incomplete subsets.

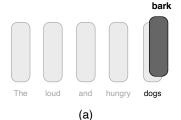
4.1 Fine-tuning for Compatibility

Ideally $f(y_t|y_S)$ approximates $p(y_t|y_S)$, a property known as *compatibility* (Arnold and Press, 1989). Since training with Equation 1 only evaluates f on complete contexts $y_{< t}$, its behavior on incomplete contexts y_S is unspecified. Instead, compatibility can be obtained by training to maximize

$$\mathbb{E}_{y_{1:T} \sim F} \mathbb{E}_{S \sim \text{Unif}(S)} \left[\sum_{t=1}^{T} \log f(y_t | y_{S < t}) \right], \quad (6)$$

where $S \sim \mathrm{Unif}(\mathcal{S})$ indicates sampling word subsets uniformly at random from the power set of all possible word subsets, and $S_{< t}$ denotes the indices in S that are less than t. Jethani et al. (2021) show that the optimum of Equation 6 is the distribution whose conditional distributions are all equal to the ground-truth conditionals.

We approximate Equation 6 with word dropout. In practice, we combine this objective with standard MLE training to learn compatible distributions



context	p (last word is "bark")
"The", "dogs"	0.04
"loud", "dogs"	0.41
"and", "dogs"	0.03
"hungry", "dogs"	0.13
(k	o)

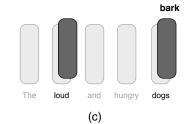


Figure 2. One step of greedy rationalization. In (a), the rationale so far is a single word, "dogs". In (b), each candidate token is considered and "loud" results in the best probability for "bark". In (c), the token "loud" is added to the rationale. This process repeats until the most likely word is the model prediction.

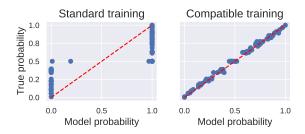


Figure 3. Training with word dropout (right) results in compatible predictions for the majority-class synthetic language. The optimal compatibility is the dashed line.

while maintaining the performance of the original model. The word dropout distribution in Equation 6 is heavily skewed towards contexts containing half the words in the sequence. To alleviate this problem, we modify the word dropout distribution to sample subsets of varying lengths; see Appendix D.

The intuition for Equation 6 is straightforward: if the model sees incomplete contexts while training, it can approximate arbitrary incomplete distributions. Since $f(y_t|y_S)$ approximates $F(y_t|y_S)$ and $f(y_t|y_{< t})$ approximates $F(y_t|y_{< t})$, all the conditional distributions are compatible.

4.2 Compatibility Experiments

To demonstrate the impact of training with the compatibility objective in Equation 6, we consider a synthetic majority-class language over binary strings of 19 tokens. The first 17 are sampled uniformly from $\{0,1\}$, and the 18th token is always '='. The 19th token is 0 if there are more 0's than 1's in the first 17 tokens, and 1 otherwise.

We train two models: one using the standard objective in Equation 1, the other using word dropout to optimize Equation 6. Although both models have the same heldout perplexity on the full context, training with Equation 6 is required to form compatible predictions on incomplete subsets. In Figure 3, we provide both models with random sub-

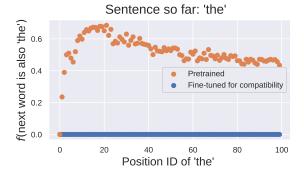


Figure 4. Fine-tuning GPT-2 for compatibility removes pathological repeating on incomplete contexts. For a position t, the vertical axis gives $f(y_{t+1} = \text{"the"})y_t = \text{"the"})$.

sets S and calculate each model's probability that the last token is 1. A model that has only seen a few tokens should be less confident about the prediction of the final majority class, yet models trained without word dropout ignore this uncertainty.

Models do not need to be trained from scratch with Equation 6. A model can be pre-trained with Equation 1, after which it can be fine-tuned for compatibility. As an example, when GPT-2 is not trained with word dropout, it makes insensible predictions for out-of-distribution sequences. For a sequence that contains only the token "the", GPT-2 is trained to give reasonable predictions for $p(y_2|y_1=$ "the"). But when it has only seen the token "the" somewhere besides the first position of the sequence, the top prediction for the word after "the" is also "the". Of course, following "the" with "the" is not grammatical. Fine-tuning for compatibility alleviates this problem (Figure 4).

Finally, we find that that fine-tuning for compatibility does not hurt the heldout performance of the complete conditional distribution of each

³We represent "the" at various positions by changing the positional encoding passed into the transformer.

fine-tuned model (see Appendix D).

5 Connection to Classification Rationales

In this section, we go over related rationalization approaches developed for classification and discuss why they cannot scale to sequence models. We also show that the combinatorial rationale objective in Equation 3 is a global solution to a classification rationale-style objective.

In classification problems, a sequence $x_{1:T}$ is associated with a label y. Rationale methods are commonly used in this setting (Lei et al., 2016; Chen et al., 2018; Yoon et al., 2018; Bastings et al., 2019; Jain et al., 2020; Jethani et al., 2021). The most common approach uses two models: one, a selection model $q(S|x_{1:T})$, provides a distribution over possible rationales; the other, the predictive model $p(y|x_S)$, makes predictions using only samples from the former model. Typically, p and q are both optimized to maximize

$$\mathbb{E}_{x,y\sim F}\mathbb{E}_{S\sim q(S|x,y)}[\log p(y|x_S) - \lambda|S|]. \tag{7}$$

Here, F is the ground truth, unknown data distribution, and λ is a regularizing penalty that encourages smaller rationales.

In practice, it is infeasible to adopt this objective for sequence models. Equation 7 is centered on providing classification models with only the words in a sequence's rationale. In sequential settings, each word has a different rationale. Since sequence models make T predictions per sequence and are trained by sharing all T word representations, each token would be indirectly exposed to words in the rationales of the words it is allowed to use. A remedy would be to train sequence models without sharing representations, but this is computationally infeasible; it requires $O(T^3)$ computations per sequence for transformer architectures.

Most classification rationale methods treat $q(S|x_{1:T})$ as a probability distribution over all possible rationales. However, the q that maximizes Equation 7 is deterministic for any p. To see this, note that q does not appear inside the expectation in Equation 7, so it can place all its mass on a single mode. We provide a formal justification in Appendix B.

Since the optimal selection model q is a pointmass, the optimal rationale can be written as

$$\underset{S \in \mathcal{S}}{\operatorname{arg\,min}} \ \lambda |S| - \log p(y|x_S). \tag{8}$$

This optimization is identical to the combinatorial optimization in Equation 3, albeit with a soft constraint on the rationale's prediction: the true label y is not required to be the maximum of $p(y'|x_S)$. In practice, this soft constraint sometimes results in empty rationales (Jain et al., 2020). Since we view sufficiency as a key component of a good rationale, Equation 3 imposes a hard constraint on the rationale's prediction.

6 Related Work

Finding rationales is similar to feature selection. While global feature selection has been a well-studied problem in statistics (Guyon and Elisseeff, 2003; Hastie et al., 2009), instance-wise feature selection — where the goal is selecting features per-example — is a newer research area (Chen et al., 2018). We review local explanation methods used for NLP.

Gradient saliency. Gradient-based saliency methods have long been used as a measure of feature importance in machine learning (Baehrens et al., 2010; Simonyan et al., 2013; Li et al., 2016a). Some variations involve word embeddings (Denil et al., 2014); integrated gradients, to improve sensitivity (Sundararajan et al., 2017); and relevance-propagation to track each input's contribution through the network (Bach et al., 2015; Voita et al., 2021).

But there are drawbacks to using gradient-based methods as explanatory tools. Sundararajan et al. (2017) show that in practice, gradients are *saturated*: they may all be close to zero for a well-fitted function, and thus not reflect importance. Adversarial methods can also distort gradient-based saliences while keeping a model's prediction the same (Ghorbani et al., 2019; Wang et al., 2020). We compare greedy rationalization to gradient saliency methods in Section 8.

Attention. Recently, NLP practitioners have focused on using attention weights as explanatory tools. The literature has made a distinction between *faithfulness* and *plausibility*. An explanation is faithful if it accurately depicts how a model makes a decision (Jacovi and Goldberg, 2020); an explanation is plausible if it can be understood and interpreted by humans (Wiegreffe and Pinter, 2019). Practitioners have shown that attention-based explanations are generally not faithful (Jain and Wallace, 2019; Serrano and Smith, 2019), but that they may

be plausible (Wiegreffe and Pinter, 2019; Mohankumar et al., 2020; Vashishth et al., 2019). Others show that attention weights should not be interpreted as belonging to single tokens since they mix information across tokens (Brunner et al., 2019; Kobayashi et al., 2020). Bastings and Filippova (2020) argue that general input saliency measures, such as gradients, are better suited for explainability than attention. We compare greedy rationalization to attention-based methods in Section 8.

Local post-hoc interpretability. Another class of methods provides local interpretability for pretrained models. These approaches aim to explain a model's behavior for a single example or for a small subset of inputs. LIME (Ribeiro et al., 2016) trains an interpretable model that locally approximates the pretrained model. Alvarez-Melis and Jaakkola (2017) learn a causal relationship between perturbed inputs and their model outputs. These methods impose no constraints on the pretrained model. However, they are expensive – they require training separate models for each input region. In contrast, the method proposed here, greedy rationalization, can efficiently explain many predictions.

Input perturbation. Practitioners have also measured the importance of inputs by perturbing them (Zeiler and Fergus, 2014; Kádár et al., 2017). Occlusion methods (Li et al., 2016b) replace an input with a baseline (e.g. zeros), while omission methods (Kádár et al., 2017) remove words entirely. Li et al. (2016b) propose a reinforcement learning method that aims to find the minimum number of occluded words that would change a model's prediction. Feng et al. (2018) use gradients to remove unimportant words to see how long it takes for the model's prediction to change. They find that the remaining words are nonsensical and do not comport with other saliency methods. Others have shown that input perturbation performs worse than other saliency methods in practice (Poerner et al., 2018). These methods have mostly focused on subtractive techniques. For this reason, they are inefficient and do not aim to form sufficient explanations. In contrast, greedy rationalization efficiently builds up sufficient explanations.

7 Experimental Setup

There are two goals in our empirical studies. The first is to compare the ability of greedy rationalization to other approaches for optimizing the combi-

natorial objective in Equation 3. The second is to assess the quality of produced rationales.

We measure the quality of rationales using two criteria: faithfulness and plausibility. An explanation is faithful if it accurately depicts how a model makes a decision (Jacovi and Goldberg, 2020); an explanation is plausible if it can be understood and interpreted by humans (Wiegreffe and Pinter, 2019). Although sufficiency is a standard way to measure faithfulness (DeYoung et al., 2020), all the rationales that satisfy the constraint of Equation 3 are sufficient by definition. To measure plausibility, we compare rationales to human annotations. Since there do not exist language modeling datasets with human rationales, we collected annotations based on Lambada (Paperno et al., 2016). The annotated dataset is available online, along with the code used for all experiments.⁴

We compare greedy rationalization to a variety of gradient- and attention-based baselines (see Section 6). To form baseline sequential rationales, we add words by the order prescribed by each approach, stopping when the model prediction is sufficient. The baselines are: l_2 gradient norms of embeddings (Li et al., 2016a), embedding gradients multiplied by the embeddings (Denil et al., 2014), integrated gradients (Sundararajan et al., 2017), attention rollout (Abnar and Zuidema, 2020), the last-layer transformer attention weights averaged-across heads, and all transformer attentions averaged across all layers and heads (Jain et al., 2020).

To compare rationale sets produced by each method to those annotated by humans, we use the set-similarity metrics described in De Young et al. (2020): the intersection-over-union (IOU) of each rationale and the human rationale, along with the token-level F1, treating tokens as binary predictions (either in the human rationale or out of it).

We use transformer-based models for all of the experiments. We fine-tune each model for compatibility using a single GPU. That we can fine-tune GPT-2 Large (Radford et al., 2019) to learn compatible conditional distributions on a single GPU suggests that most practitioners will be able to train compatible models using a reasonable amount of computation. For model and fine-tuning details, refer to Appendix D.

⁴https://github.com/keyonvafa/
sequential-rationales

8 Results and Discussion

The experiments test sequential rationales for language modeling and machine translation. Appendix E contains full details for each experiment.

8.1 Language Modeling

Long-Range Agreement. The first study tests whether rationales for language models can capture long-range agreement. We create a template dataset using the analogies from Mikolov et al. (2013). This dataset includes word pairs that contain either a semantic or syntactic relationship. For each type of relationship, we use a predefined template. It prompts a language model to complete the word pair after it has seen the first word.

For example, one of the fifteen categories is countries and their capitals. We can prompt a language model to generate the capital by first mentioning a country and then alluding to its capital. To test long-range agreement, we also include a distractor sentence that contains no pertinent information about the word pair. For example, our template for this category is,

When my flight landed in **Japan**, I converted my currency and slowly fell asleep. (I had a terrifying dream about my grandmother, but that's a story for another time). I was staying in the capital,

Here, the parenthetical clause is a distractor sentence, since it contains no relevant information about predicting the capital of Japan. The correct capital, "Tokyo", is predicted by GPT-2 both with and without the distractor. We use this template for all of the examples in the country capital category, swapping the antecedent "Japan" for each country provided in Mikolov et al. (2013).

We feed the prompts to GPT-2, which completes each analogy. To measure faithfulness, we calculate the percent of rationales that contain the true antecedent, and the percent of rationales that do not contain any words in the distractor. We only use examples where the prediction is the same both with and without the distractor. We also perform exhaustive rationale search on the objective in Equation 3. This search is highly inefficient, so we only complete it for 40 examples. To measure the approximation ratio, we divide the size of the rationale found by each method by the exhaustive rationale size.

Table 1 contains the results on the compatible model. Although all methods contain the true antecedents in their rationales, greedy rationalization

	Length	Ratio	Ante	No D
Grad norms	22.5	4.1	1.0	0.06
Grad x emb	38.0	7.4	0.99	0.01
Integrated grads	28.1	5.2	0.99	0.00
Attention rollout	36.9	7.1	1.0	0.12
Last attention	16.7	2.9	0.99	0.13
All attentions	14.5	2.6	1.0	0.02
Greedy	7.1	1.2	1.0	0.43

Table 1. Language modeling faithfulness on long-range agreement with templated analogies. "Ratio" refers to the approximation ratio of each method's ratio-nale length to the exhaustive search minimum. "Ante" refers to the percent of rationales that contain the true antecedent. "No D" refers to the percent of rationales that do not contain any tokens from the distractor.

has by far the least distractors in its rationales. The rationales are also universally shorter for greedy rationalization and closer to the optimal rationales, justifying our greedy assumption. To show that fine-tuning GPT-2 for compatibility is not hurting the baselines, we also perform the baseline methods on a pretrained GPT-2 without fine-tuning; see Appendix E.

Annotated Rationales. To test the plausibility of rationales for language models, we collect a dataset of human annotations. We base the collection on Lambada (Paperno et al., 2016), a corpus of narrative passages. Each passage included in Lambada is chosen so that humans need to use both local and global context to reliably predict the final word. By its construction it is guaranteed to have non-trivial rationales.

Our goal is to collect rationales that are both minimal and sufficient for humans. We run an annotation procedure with two roles: a selector and a predictor. First, the selector sees the full passage and ranks the words in order of how informative they are for predicting the final word. Next, the predictor sees one word at a time chosen by the selector, and is asked to predict the final word of the passage. The words the predictor saw before guessing the correct word form a human rationale. This rationale selection method is inspired by Rissanen Data Analysis (Rissanen, 1978; Perez et al., 2021), which uses a minimum description length metric to estimate feature importances. We rely on human annotators to estimate information gains.

Since it could be trivial for humans to predict the final word if it also appears in the context, we only include examples that do not repeat a word. We collect annotations for 107 examples, which we

	Length	IOU	F1
Gradient norms	60.2	0.14	0.22
Gradient x embedding	68.3	0.12	0.21
Integrated gradients	62.8	0.12	0.21
Attention rollout	73.9	0.11	0.19
Last attention layer	54.6	0.15	0.25
All attention layers	48.7	0.20	0.28
Greedy	17.9	0.25	0.35

Table 2. Language modeling plausibility on rationale-annotated Lambada.

Target word: grow

Target word: refuse

It was the kind of smile that I'd seen before. The kind the boxer gave me right before he killed me in that dirty fight.

"I have a proposition for you" he began, pulling his hands down from under his chin and pushing out of the chair. "One that you won't be able to

Figure 5. Examples from our annotated Lambada dataset. Highlighted text denotes greedy rationales, and **bolded text** denotes human-annotated rationales.

also release publicly. We use two sets of annotators for 15% of the examples in order to compute inter-annotator agreement. On this subset, the average token-level Cohen's κ is 0.63 (Cohen, 1960), indicating substantial agreement.

We compare the rationales produced by each method to the annotated rationales. Table 2 shows that the greedy rationales are most similar to the human-annotated rationales. Greedy rationalization is also the most effective at minimizing the combinatorial objective in Equation 3, as its rationales are by far the shortest. Figure 5 contains examples of rationales for this dataset.

It is worth noting that the top few words added by the baselines are quite relevant; after 5 tokens, the all-attention baseline has a better F1 and IOU than greedy rationalization. However, the baselines struggle to form sufficient rationales, which hurts their overall performance.

8.2 Machine Translation

Distractors. To measure faithfulness, we take a transformer trained on IWSLT14 De-En (and fine-

	Mean Crossovers		Crossov	er Rate
	Source	Target	Source	Target
Grad norms	0.40	0.44	0.06	0.06
Grad x emb	6.25	5.57	0.42	0.41
Integrated grads	2.08	1.68	0.23	0.14
Last attention	0.63	2.41	0.09	0.24
All attentions	0.58	0.80	0.08	0.12
Greedy	0.12	0.12	0.09	0.02

Table 3. Translation faithfulness with distractors. "Mean crossovers" refers to the average number of crossovers per rationale, and "Crossover rate" refers to the fraction of rationales that contain at least one crossover.

tuned for compatibility), and generate translations for 1000 source sequences from the test set. We then create a corpus by concatenating random example pairs; for two sampled pairs of source and target sequences, (S_1, T_1) and (S_2, T_2) , we create a new example (S_1S_2, T_1T_2) . Each token in T_1 is generated from S_1 alone, so its rationales shouldn't contain any tokens from S_2 . Similarly, T_2 is generated from S_2 alone, so its rationales shouldn't contain any tokens from S_1 or T_1 .

We evaluate each rationale by counting how many times it has crossed over: a rationale for T_1 crosses over every time it contains a token in S_2 , and a rationale for T_2 crosses over every time it contains a token in S_1 or T_1 (since the model is autoregressive, T_1 's rationales can never contain tokens from T_2).

Table 3 contains the results. Greedy rationalization has by far the fewest average number of crossovers per rationale. Although the percent of source rationales that cross over is slightly higher than the percent using gradient norms, the percentage on the target side is superior.

Annotated Alignments. To test plausibility, we compare the rationales to word alignments (Brown et al., 1993). Using a dataset containing 500 human-labeled alignments for German-English translation,⁵ we compute rationales for each method using the ground truth targets. We measure similarity to the labeled rationales by computing alignment error rate (AER) (Och and Ney, 2000), along with computing the IOU and F1 between sets. To separate the requirement that the rationale be sufficient from each method's global ordering of tokens, we also compare top-1 accuracies, which measure whether

[&]quot;Just who is going to pay for this special feed grain anyway? It must cost a bit if it's that special."
"You're going to pay, obviously," replied Mitch, "since your cows will be eating it. On the other hand, Joe will be planting and irrigating the grain. He'll do all the work to

⁵https://www-i6.informatik.rwth-aachen.
de/goldAlignment/

	Length	AER	IOU	F1	Top1
Grad norms	10.2	0.82	0.30	0.16	0.63
Grad x emb	13.2	0.90	0.16	0.12	0.40
Integrated grads	11.3	0.85	0.24	0.14	0.42
Last attention	10.8	0.84	0.27	0.15	0.59
All attentions	10.7	0.82	0.32	0.15	0.66
Greedy	4.9	0.78	0.40	0.24	0.64

Table 4. Translation plausibility with annotated alignments. The first four columns correspond to using the full source rationale found by each method; the last column "Top1" refers to the accuracy of the first source token added by each method. AER refers to alignment error rate.

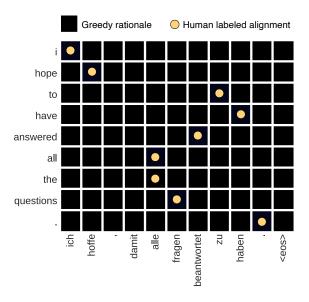


Figure 6. Greedy rationalization for machine translation. Each row depicts the source words contained in a rationale. Although each rationale includes both source and target words, here we only show source-side rationales so they can be compared to annotated alignments.

the top token identified by each baseline is present in the labeled alignment set.

Table 4 contains the results. The rationales learned by greedy rationalization are more similar to human-labeled alignments than those provided by gradient and attention methods. Many methods have similar top-1 accuracies — indeed, the best top-1 accuracy comes from averaging all attention layers. This reinforces the notion that although the baselines may be able to capture first-order information, they struggle to form sufficient rationales. Figure 6 contains an example of greedy rationalization applied to machine translation, along with the human-labeled alignments.

9 Conclusion

We proposed an optimization-based algorithm for rationalizing sequence predictions. Although exact optimization is intractable, we developed a greedy approach that efficiently finds good rationales. Moreover, we showed that models can be fine-tuned to form compatible distributions, thereby circumventing an intractable marginalization step. In experiments, we showed that the greedy algorithm is effective at optimization, and that its rationales are more faithful and plausible than those of gradient- and attention-based methods. We hope that our research, along with the release of an annotated dataset of sequence rationales, catalyzes further research into this area.

Acknowledgments This work is funded by ONR N00014-17-1-2131, ONR N00014-15-1-2209, NSF CCF-1740833, DARPA SD2 FA8750-18-C-0130, Two Sigma, Amazon, and NVIDIA. Keyon Vafa is supported by the Cheung-Kong Innovation Doctoral Fellowship. Alexander Rush and Yuntian Deng are sponsored by NSF 1901030 and NSF CAREER 2037519. We also thank Mark Arildsen, Elizabeth Chen, Justin Chen, Katherine Chen, Nathan Daniel, Alexander Hem, Farzan Vafa, Neekon Vafa, Willy Xiao, and Carolina Zheng.

References

Samira Abnar and Willem Zuidema. 2020. Quantifying attention flow in transformers. In *Association for Computational Linguistics*.

David Alvarez-Melis and Tommi S Jaakkola. 2017. A causal framework for explaining the predictions of black-box sequence-to-sequence models. In *Association for Computational Linguistics*.

Barry C Arnold and S James Press. 1989. Compatible conditional distributions. *Journal of the American Statistical Association*, 84(405):152–156.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140.

David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. 2010. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11:1803–1831.

- Jasmijn Bastings, Wilker Aziz, and Ivan Titov. 2019. Interpretable neural predictions with differentiable binary variables. In *Association for Computational Linguistics*.
- Jasmijn Bastings and Katja Filippova. 2020. The elephant in the interpretability room: Why use attention as explanation when we have saliency methods? In *ACL Workshop on BlackboxNLP*.
- Peter F Brown, Stephen A Della Pietra, Vincent J Della Pietra, and Robert L Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. In *Association for Computational Linguistics*.
- Gino Brunner, Yang Liu, Damian Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger Wattenhofer. 2019. On identifiability in transformers. In *International Conference on Learning Representations*.
- Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. 2014. Report on the 11th IWSLT evaluation campaign. In *International Workshop on Spoken Language Translation*.
- Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. 2018. Learning to explain: An information-theoretic perspective on model interpretation. In *International Conference on Machine Learning*.
- Vasek Chvatal. 1979. A greedy heuristic for the setcovering problem. *Mathematics of Operations Re*search, 4(3):233–235.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.
- Misha Denil, Alban Demiraj, and Nando De Freitas. 2014. Extraction of salient sentences from labelled documents. In *International Conference on Learning Representations*.
- 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 *Association for Computational Linguistics*.
- Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, and Jordan Boyd-Graber. 2018. Pathologies of neural models make interpretations difficult. In *Association for Computational Linguistics*.
- Amirata Ghorbani, Abubakar Abid, and James Zou. 2019. Interpretation of neural networks is fragile. In *Association for the Advancement of Artificial Intelligences*.
- Aaron Gokaslan and Vanya Cohen. 2019. Openwebtext corpus. http://Skylion007.github.io/OpenWebTextCorpus.

- Isabelle Guyon and André Elisseeff. 2003. An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3(Mar):1157–1182.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2009. *The elements of statistical learning: Data mining, inference, and prediction.* Springer Science & Business Media.
- Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. 2019. A benchmark for interpretability methods in deep neural networks. *Neural Information Processing Systems*.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In *Association for Computational Linguistics*.
- Sarthak Jain and Byron C. Wallace. 2019. Attention is not explanation. In *North American Chapter of the Association for Computational Linguistics*.
- Sarthak Jain, Sarah Wiegreffe, Yuval Pinter, and Byron C Wallace. 2020. Learning to faithfully rationalize by construction. In *Association for Computational Linguistics*.
- Neil Jethani, Mukund Sudarshan, Yindalon Aphinyanaphongs, and Rajesh Ranganath. 2021. Have we learned to explain?: How interpretability methods can learn to encode predictions in their interpretations. In *Artificial Intelligence and Statistics*.
- Akos Kádár, Grzegorz Chrupała, and Afra Alishahi. 2017. Representation of linguistic form and function in recurrent neural networks. In *Association for Computational Linguistics*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*.
- Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. Attention is not only a weight: Analyzing transformers with vector norms. In *Association for Computational Linguistics*.
- T. Lei, R. Barzilay, and T. Jaakkola. 2016. Rationalizing neural predictions. In *Association for Computational Linguistics*.
- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016a. Visualizing and understanding neural models in NLP. In *Association for Computational Linguistics*.
- Jiwei Li, Will Monroe, and Dan Jurafsky. 2016b. Understanding neural networks through representation erasure. *arXiv preprint arXiv:1612.08220*.
- Zachary C Lipton. 2018. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. In *Queue*, volume 16, pages 31–57. ACM New York, NY, USA.

- Tomas Mikolov, Kai Chen, Greg S. Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. In *Workshop Track at ICLR*.
- Akash Kumar Mohankumar, Preksha Nema, Sharan Narasimhan, Mitesh M Khapra, Balaji Vasan Srinivasan, and Balaraman Ravindran. 2020. Towards transparent and explainable attention models. In *Association for Computational Linguistics*.
- Franz Josef Och and Hermann Ney. 2000. Improved statistical alignment models. In *Association for computational linguistics*.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. Fairseq: A fast, extensible toolkit for sequence modeling. In *Association for Computational Linguistics*.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. The LAMBADA dataset: Word prediction requiring a broad discourse context. In Association for Computational Linguistics.
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. Rissanen data analysis: Examining dataset characteristics via description length. *arXiv preprint arXiv:2103.03872*.
- Nina Poerner, Benjamin Roth, and Hinrich Schütze. 2018. Evaluating neural network explanation methods using hybrid documents and morphological agreement. In *Association for Computational Linguistics*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why should I trust you?" Explaining the predictions of any classifier. In *Special Interest Group on Knowledge Discovery and Data*.
- Jorma Rissanen. 1978. Modeling by shortest data description. *Automatica*, 14(5):465–471.
- Sofia Serrano and Noah A Smith. 2019. Is attention interpretable? In *Association for Computational Linguistics*.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*.

- Shikhar Vashishth, Shyam Upadhyay, Gaurav Singh Tomar, and Manaal Faruqui. 2019. Attention interpretability across NLP tasks. *arXiv preprint arXiv:1909.11218*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Neural Information Processing Systems*
- Elena Voita, Rico Sennrich, and Ivan Titov. 2021. Analyzing the source and target contributions to predictions in neural machine translation. In *Association for Computational Linguistics*.
- Junlin Wang, Jens Tuyls, Eric Wallace, and Sameer Singh. 2020. Gradient-based analysis of NLP models is manipulable. In *Empirical Methods in Natural Language Processing*.
- Sarah Wiegreffe and Yuval Pinter. 2019. Attention is not not explanation. In *Empirical Methods in Natural Language Processing*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Transformers: State-of-the-art natural language processing. In *Empirical Methods in Natural Language Processing*.
- Jinsung Yoon, James Jordon, and Mihaela van der Schaar. 2018. INVASE: Instance-wise variable selection using neural networks. In *International Con*ference on Learning Representations.
- Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*.

A Algorithm Details

We present greedy rationalization in Algorithm 1.

Algorithm 1: Greedy rationalization

Input: Sequence $y_{1:t}$ generated from p.

Output: Rationale S for y_t .

Initialize: $S = \emptyset$

while
$$\arg\max_{y_t'} p(y_t'|y_S) \neq y_t$$
 do
$$k^* = \arg\max_{k \in [t-1] \setminus S} p(y_t|y_{S \cup k})$$

$$S = S \cup k^*$$

Mose sequence models, including transformers, use the representation of a token y_{t-1} to predict the next token, y_t . As such, a rationale S always needs to contain y_{t-1} . In practice, we initialize $S = \{y_{t-1}\}.$

This method and paradigm extend easily to conditional sequence models, such as those used in machine translation. In this setting, a model uses a source sequence $x_{1:N}$ to generate a target sequence $y_{1:T}$. Thus, a context for a prediction y_t contains both $y_{< t}$ and $x_{1:N}$. The set of all possible rationales is the cross product of power sets $\mathcal{S} = 2^{[N]} \times 2^{[t-1]}$, and the combinatorial objective is

$$\begin{split} S(x_{1:N},y_{1:t}) &= \mathop{\arg\min}_{S_x,S_y \in \mathcal{S}} |S_x| + |S_y| \\ \text{s.t.} & \mathop{\arg\max}_{y_t'} p(y_t'|x_{S_x},y_{S_y}) = y_t. \end{split}$$

To perform greedy rationalization in this setting, we consider adding either a source token or a target token at each step, choosing the one that results in the largest increase in the full model's prediction.

B Optimality of Deterministic Rationales

Here, we show that the selection distribution q(S|x,y) that maximizes the classification rationale objective in Equation 7 is deterministic. We re-write the objective below:

$$\mathbb{E}_{x,y\sim F}\mathbb{E}_{S\sim q(S|x,y)}[\log p(y|x_S) - \lambda |S|]. \tag{9}$$

Theorem 1. For any $p(y|x_S)$, the q(S|x,y) that maximizes Equation 9 is a point-mass.

Proof. Denote by
$$g(x, y, S) = \log p(y|x_S) - \lambda |S|$$
:

$$\max_{q} \mathbb{E}_{x,y \sim F} \mathbb{E}_{S \sim q(S|x,y)}[g(x,y,S)]$$

$$\leq \max_{q} \mathbb{E}_{x,y \sim F} \max_{S}[g(x,y,S)]$$

$$= \mathbb{E}_{x,y \sim F} \max_{S}[g(x,y,S)].$$

Step	Complexity	Evaluations	Total
1	1^2	t	1^2t
2	2^2	t-1	$2^2(t-1)$
:	:	:	:
$O(t^{1/3})$	$O(t^{2/3})$	O(t)	$O(t^{5/3})$
Total			$O(t^2)$

Table 5. For transformers, the asymptotic complexity of greedy rationalization matches the asymptotic complexity of forming a single prediction on the full sequence, as long as the rationale size is $O(t^{1/3})$ for a sequence of length t.

The inequality uses the fact that the expectation of a random variable is bounded by its maximum value. When q(S|x,y) is a point-mass at $\arg\max_S[g(x,y,S)]$, the inequality becomes tight.

The fact that the optimal rationale is deterministic for each example justifies using combinatorial strategies such as our objective in Equation 3.

C Efficiency

In Table 5, we provide a detailed version of our complexity analysis from Section 3: For transformers, greedy rationalization can be performed at no extra asymptotic complexity if the rationale length is $O(t^{1/3})$ for a sequence length t.

We evaluate the computational efficiency of greedy rationalization in Table 6. We compare greedy rationalization to an exhaustive search, which enumerates all possible context context subsets from shortest to longest to optimize Equation 3. To show the efficiency of evaluating transformers on arbitrarily sized inputs, we also compare to a version of greedy rationalization that evaluates a transformer on the full input. To make predictions on sparse subsets, this approach masks tokens that aren't in a candidate rationale during each attention step. In contrast, the efficient version of greedy rationalization only takes as input the tokens in the candidate rationale, so there is no need for masking.

Method	Time (s)
Exhaustive search	>60
Greedy rationalization with full inputs	1.22
Greedy rationalization with sparse inputs	0.30

Table 6. Greedy rationalization is efficient, especially when evaluating transformers on sparse inputs. We report the average wall clock time in seconds for finding rationales on the templated analogies dataset of Mikolov et al. (2013). We cannot complete exhaustive search for the longer examples, so in reality the average runtime is larger than the listed one.

We perform these comparisons on the templated analogies dataset of Mikolov et al. (2013). We use GPT-2 Large as our sequence model (Radford et al., 2019) and perform each method on a single GPU. We compare the two greedy rationalization approaches for all of the examples for which the full model predicts the templated output. Since exhaustive search is intractable, we cannot perform it on every example due to computational constraints. Thus, we only run exhaustive search on examples where the optimal rationale has 6 or less tokens. In reality, the average runtime for exhaustive search is larger than the listed one.

D Training and Fine-Tuning Details

Our experiments consist of three models and datasets: a transformer decoder (Vaswani et al., 2017) trained on a majority-class language, GPT-2 (Radford et al., 2019) fine-tuned on Open WebText (Gokaslan and Cohen, 2019), and a transformer machine translation model trained and fine-tuned with word dropout on IWSLT14 De-En (Cettolo et al., 2014).

For the majority-class language, we generate the dataset as described in Section 4. We include 50,000 examples in the training set, 5,000 in the validation set, and 5,000 in the test set.

We use a 4-layer transformer decoder with 2 attention heads per layer. We use an embedding dimension of 64, and a hidden dimension of 256 for the feedforward layers. This corresponds to 200,000 parameters. We train with 0.1 weight dropout, and optimize using Adam (Kingma and Ba, 2015) with a learning rate of 0.005 and an inverse square root learning rate scheduler. We use a warmup period of 4000 steps and an initial warmup learning rate of 10^{-7} . We include a maximum of 64,000 tokens in each batch. We implement this model in Fairseq (Ott et al., 2019).

To approximate the compatibility objective in Equation 6, we train with varying amounts of word dropout. In practice, this amounts to masking out each token we drop out at each attention layer. We use two levels of word dropout in Figure 3; none (which corresponds to training with the standard maximum likelihood objective in Equation 1) and 0.5. We train each model on a single GPU. Each model takes less than 20,000 steps to converge, less than 90 minutes. Table 7 verifies that fine-tuning with word dropout does not hurt the heldout perplexity.

To fine-tune GPT-2, we use the pretrained GPT-2 Large model available on Hugging Face (Wolf et al., 2019). This model has 774M parameters. We don't change any of the model settings when we fine-tune. Sampling context subsets uniformly at random as stated in the objective in Equation 6 results in a distribution of subsets heavily skewed towards those containing half the words in the sequence. This is fine for the majority-class language, since each sequence contains less than 20 tokens and thus all possible context sizes will be seen during training. However, GPT-2's sequence length is 1,024. 99% of the time, sampling from the objective as stated would result in contexts with size 464-560. Notably, the probability of a context with less than 10 tokens is less than 10^{-284} .

We make two adjustments to make sure the model is trained on both small and large subsets. With probability 0.5, we condition on the full context. With the remaining 0.5, we first randomly sample context sizes uniformly at random from 1 to the sequence length. We then sample a random context subset of this size. This guarantees that all possible sequence lengths will be seen during training.

Since the WebText dataset used to train GPT-2 is not publicly available, we use Open WebText (Gokaslan and Cohen, 2019), an open source reproduction effort. The corpus is in English. Rather than using the entire dataset, we take "Subset 9" and use the first 163M words. Our validation set is also from this subset and contains 160,000 words. We use a test set of 300,000 words from a different subset.

We fine-tune GPT-2 Large using Adam. We use a constant learning rate of 0.0001, using a single batch per training step. We stop training after 62,500 steps. This takes 15 hours on a single GPU. Table 7 shows that fine-tuning with word dropout actually improves the heldout perplexity, although we believe that the improvement is due to our test set bearing more resemblance to the fine-tuning set than to the pretraining set.

We use a standard transformer encoder/decoder to train a machine translation model on IWSLT14 De-En (Cettolo et al., 2014). We follow the preprocessing and model architecture recommended by Fairseq.⁶ The training set has 160,239 translation pairs, the validation set has 7,283, and the test set has 6,750.

As for the model, both the encoder and decoder are transformers with 6 layers, 4 attention heads per layer, 512 embedding dimensions, and 1024 feedforward dimensions. This corresponds to 40M parameters. We train with 0.3 weight dropout and 0.1 label smoothing, using 4,096 tokens for each train step. We train with Adam with a learning rate of 5×10^{-4} and use an inverse square root learning rate scheduler with 4,000 warmup steps.

When we fine-tune for compatibility, we again condition on the full context with probability 0.5. With the remaining probability, we drop out each source and target token independently at each attention head with probability 1 - 1/T, where T is the sequence length (so the dropout probability varies for the source and target sequence). Although we drop out different tokens at each attention head of a layer, we make sure that the same tokens are dropped out at each layer. Our word dropout procedure ensures that our objective will be trained on small contexts since rationales for machine translation are typically very sparse. We fine-tune using Adam with a constant learning rate of 10^{-5} for 410,000 steps. The heldout BLEU scores of both models are equal; see Table 7.

E Experimental Details

E.1 Long-Range Agreement

Table 8 contains the template we used for the set of experiments containing the analogies from Mikolov et al. (2013). To avoid rationales containing partial antecedents, we only include examples where both words in the analogy correspond to single word-pieces using GPT-2's tokenizer. Since it only makes sense to rationalize correct predictions, we also only include the examples where GPT-2 correctly completes the analogy. In total, this results in 175 examples. Of these, we randomly sample 50 to perform exhaustive search, which we use to compute the approximation ratio of each method. Since we cannot run exhaustive search when the minimal sufficient rationale is too large, we use the 40 that converge with rationales of length 6 or less. We use 100 steps to approximate the path integral for the integrated gradients baseline (Sundararajan et al., 2017).

To confirm that the baseline performances are not being hindered by fine-tuning for compatibility, we re-run the experiment for each rationalization method on the pretrained GPT-2 Large, without any fine-tuning. The results are depicted in Table 9. As expected, the baselines perform even worse when GPT-2 is not fine-tuned to form compatible distributions. We do not include comparisons to exhaustive rationales because it is computationally infeasible to run exhaustive search on incompatible models, since optimization takes much longer to converge.

E.2 Machine Translation

For the distractor experiment, we randomly concatenate 500 pairs of source and target sequences generated by our fine-tuned model on the test set. We evaluate rationales by counting how many times they "cross over" and contain words from the distractor sequence. We do not penalize rationales that include special tokens like the beginning of sentence or end of sentence tokens.

⁶https://github.com/pytorch/fairseq/ tree/master/examples/translation

Model	Dataset	Evaluation Metric	Standard Training	Compatible Training
Transformer decoder GPT-2 Transformer encoder/decoder	Open WebText	1 2	1.8 18.3 34.8	1.8 17.1 34.8

Table 7. Fine-tuning for compatibility does not hurt heldout performance. The first two rows are language models and the evaluation metric is heldout perplexity; the last row is machine translation, for which the evaluation metric is BLEU.

Relationship	Example
Capital countries	When my flight landed in Greece , I converted my currency and slowly fell asleep. (I had a terrifying dream about my grandmother, but that's a story for another time). I was staying in the capital, Athens
Currency	As soon as I arrived in Japan , I checked into my hotel and took a long nap. (I had finally finished the book I was reading and it was amazing). I had to figure out the exchange rate to the local currency, which is apparently called the yen
City in state	As soon as I arrived in Florida , I checked into my hotel and watched a movie before falling asleep. (I had a great call with my husband, although I wish it were longer). I was staying in my favorite city, Miami
Family	I initially invited my uncles , who gladly accepted my invitation. (My favorite song just came on, so I was able to relax). When I learned that women were allowed, I went ahead and also invited my aunts
Opposite	I thought it was pleasant . (Just then an ad came on the TV, but that's irrelevant). It was the opposite of that: it was unpleasant
Comparative	I knew it was tall , but that's before I saw it in person. (Just then I thought about my ex-wife, but I had to stop thinking about her). When I did end up seeing it in person, it was even taller
Superlative	I thought it would be the smallest thing I'd ever encounter. (I tried to ignore my phone vibrating in my pocket). But when I did end up encountering it, it turned out it wasn't so small
Present participle	Every other day, it started working in the morning. (I tried to remember the name of the woman at the bar). But today, it did not work
Nationality adjective	I had never been friends with any French people before. (The funniest thing happened to me the other day, but that's a story for another time). In fact, I had never even been to France
Past tense	Although I listened yesterday, I had a million things to do today. (I suddenly felt a pinched nerve, so I made a mental note to get that checked out). So today I wouldn't have time to do any more listen
Plural	I really wanted to buy the computer , more than I ever wanted to buy anything before. (I was also behind on my homework, but that's another story). So I went to the store and asked if they had any computers
Plural verbs	I can usually sing by myself. (I was so behind on work but I tried to distract myself). Although it's so much better when someone else also sings

Table 8. Template using analogies from Mikolov et al. (2013).

	Length	Ante	No D
Gradient norms	24.8	1.0	0.08
Gradient x embedding	41.1	0.99	0.00
Integrated gradients	34.9	1.0	0.00
Attention rollout	38.4	1.0	0.05
Last attention layer	20.1	0.99	0.03
All attention layers	19.5	1.0	0.02
Greedy	13.1	1.0	0.30

Table 9. The performance of each rationalization method on the templated version of the analogies dataset (Mikolov et al., 2013) when we don't fine-tune for compatibility. As expected, fine-tuning for compatibility (Table 1) improves performance across the board.

For the alignment experiment, we use a public corpus of annotated rationales.⁷ Not every word in the dataset has an alignment, and some words have multiple alignments. Although the human annotations are on word-level alignments, our machine translation models are trained on subwords. so the rationales contain subwords in addition to full words. To make these comparable to the human annotations, we define the rationale of a full target word to contain the union of the subword rationales. Since each source word may also be a subword, we also take the union of source words in a rationale. To calculate top-1 accuracy, we define the rationale for a full word to be accurate if the rationales for any of the subwords in the rationale contain any source subwords that are in the annotated alignment.

The alignment dataset contains both "sure" and "possible" alignments. These are used to differentiate between different errors when calculating the alignment error rate (Och and Ney, 2000). For the other metrics, we include both kinds of alignments as part of the annotated alignments.

For both machine translation experiments, we use 50 steps to approximate the path integral for the integrated gradients baseline (Sundararajan et al., 2017).

E.3 Annotated Lambada

We work with volunteers to annotate Lambada (Paperno et al., 2016). Each example requires two annotators: a selector, and a predictor. A selector's goal is to choose the most important words for predicting the final word of a passage, known as the target word. Predictors will only be seeing the words chosen by a selector, and their goal is to predict the final word of the passage.

The selector first takes a passage and ranks 15 words. The top-ranked word always needs to be the word before the final word of the passage. They cannot select the final word of the passage. Each of their selections needs to be a complete word. They cannot select the same word twice, and they need to use all 15 spots. They know that a predictor will be predicting words, one-at-a-time, using the order they create.

When a selector is finished ranking the top 15 words, a predictor begins by seeing the top ranked word. They use this to predict the last word. Words are revealed one-at-a-time in the order chosen by the selector. Selectors can see how much space is between the words that have been revealed. Selectors are not told if they predicted a word correctly; the goal of the exercise is to capture the predictor's true predictions, so if they knew that previous guesses were incorrect, they may use this information to guess a new word at every step. If a predictor is not able to guess the target word at the end of the exercise, we re-assign the example to another predictor.

Figure 7 contains an example given to selectors. Figure 8 contains an example given to predictors.

In total, we annotate 107 examples, and use all of them for the rationalization experiment. For each example, we define a human's rationale to be all the words that were revealed by the selector before the predictor first predicted the true target word or a synonym of it. The average rationale length is 6.0.

To compare human rationales to those found by various methods, we first tokenize the text with GPT-2's tokenizer, and convert an annotated rationale to its set of corresponding subwords. Each method's rationale is also a set of subwords. We use set-comparison metrics like intersection over union (IOU) and F1 to compare the similarity of rationales. We use 100 steps to approximate the path integral for the integrated gradients baseline (Sundararajan et al., 2017).

https://www-i6.informatik.rwth-aachen. de/goldAlignment/

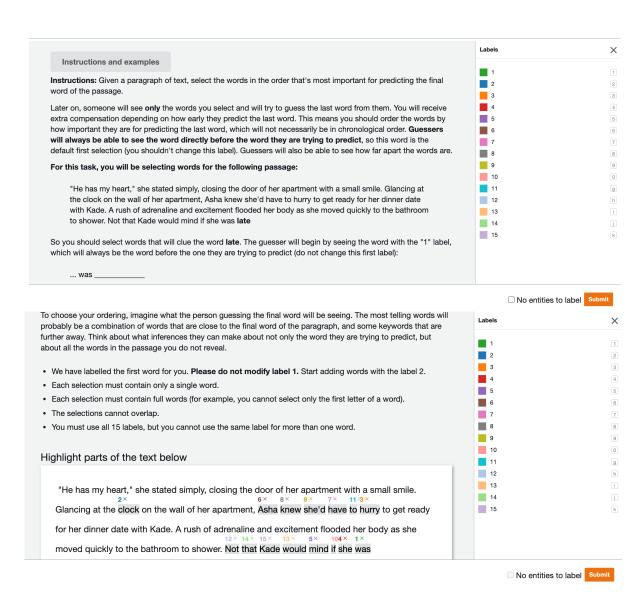


Figure 7. Sample instructions given to selectors to annotate Lambada.

Guess the missing word.

Your job is to guess the last word of the paragraph. You will begin by seeing only the word directly before the last word. At each step, we will reveal a new word to you. You must predict a word at every step. Although you will see new words at each step, you will always be predicting the final word of the paragraph (denoted by "______").

Tips

- The word you are guessing is the last word of the sentence, so make sure your predictions make sense as sentence endings.
- We will never tell you if your prediction is right or wrong. So if you think your prediction is correct as more words are revealed, you should not change it.
- You are predicting one missing word, so make sure your predictions are all single words.
- The spaces between words will give you a sense of how many words are missing between words.
- This may be hard at first, so just use your best judgment. It will get easier as words are revealed, especially because the words will appear in the most helpful order.

Fill in the blank with your best prediction of the final word of the last sentence (denoted by "_____"). The word before the final word is already filled out.

phone	reception able	get a	
Type in your prediction here			

Next

Figure 8. Sample instructions given to predictors to annotate Lambada.

Target word: again I wanted to make sure you were still comfortable with the arrangements. I can always do something different." You're too good to me, Max. But I'm fine. I promise. I'm going to be okay this time. I've learned from my past mistakes. I don't want to make them Target word: fire "We aren't out of danger yet," Horatius said. He headed northwest without thinking about it. It just seemed the right way to go. Chloe squirmed as she became more alert. "We have to go back. Now. We can't leave my house burning with my family in there ' Horatius didn't know what to do about the Target word: ring I joined Mark, Tony, and his son back in the crowd as the event started. I wanted to catch Yegor's match before I got ready for my debut. He was wrestling first; there were only four matches on the card.

Target word: contract

Your services will be required for the period of three months.'

Minutes after I got seated the lights dimmed. That cold music filled with horns played. It was Yegor's time to come

I press my lips together. I was very drunk last night, but I am **sure** he **said** one month. 'Can I speak to him?'

'Of course.' He picks up the phone and speed dials his client's number. 'Mr. Barrington, Miss Bloom would like to have a word about the length of the

Figure 9. Sample rationales from our annotated Lambada dataset. Highlighted text corresponds to greedy rationales, and **bolded text** corresponds to human annotated rationales.

F Qualitative Examples

Figure 9 contains examples of rationales on our annotated Lambada dataset.