

# Controlling the Focus of Pretrained Language Generation Models

Anonymous ACL submission

## Abstract

The finetuning of pretrained transformer-based language generation models are typically conducted in an end-to-end manner, where the model learns to attend to relevant parts of the input by itself. However, there does not exist a mechanism to directly control the model’s focus. This work aims to develop a control mechanism by which a user can select spans of context as “highlights” for the model to focus on, and generate relevant output. To achieve this goal, we augment a pretrained model with trainable “focus vectors” that are directly applied to the model’s embeddings, while the model itself is kept fixed. These vectors, trained on automatic annotations derived from attribution methods, act as indicators for context importance. We test our approach on two core generation tasks: dialogue response generation and abstractive summarization. We also collect evaluation data where the highlight-generation pairs are annotated by humans. Our experiments show that the trained focus vectors are effective in steering the model to generate outputs that are relevant to user-selected highlights.

## 1 Introduction

Transformer-based models pretrained on large-scale text data have become the dominant paradigm for natural language generation (NLG) tasks (Roller et al., 2020; Lewis et al., 2019; Raffel et al., 2020). The attention module (Bahdanau et al., 2016; Vaswani et al., 2017), which aggregates information via a weighted average over word-level embeddings, plays a vital role in these models. The attention mechanism serves two major purposes: (1) It captures linguistic phenomena in the input (Clark et al., 2019; Kovaleva et al., 2019; Kobayashi et al., 2020); (2) It helps the model focus on relevant portions of the input (e.g., alignment in machine translation (Bahdanau et al., 2016) and abstractive summarization (Rush et al., 2015)).

The attention module is particularly useful as it does not require any explicit supervision: the

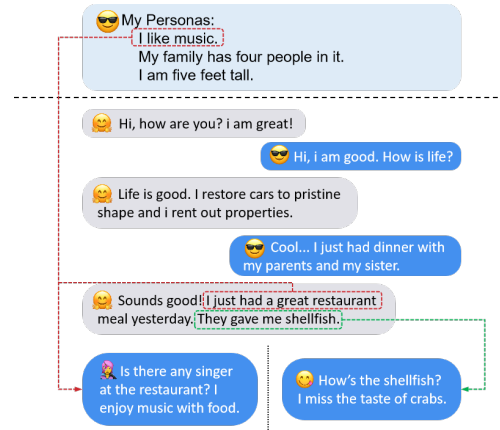


Figure 1: Illustration of our motivation: different highlights in the input (including persona) lead to different generations. This example is from our collected dialogue data for evaluation (Section 3).

model learns to attend to relevant parts of the input **by itself** through end-to-end training. However, this property makes it difficult to explicitly control the model’s focus. **If the model happens to put focus on some span of context that the user thinks is not so important, we do not have a mechanism to correct it.** This is especially sub-optimal in some NLG applications involving a relatively long input such as dialogue or summarization: focusing on different spans of the input could result in completely different generations (illustrated in Figure 1). It would be attractive to give the user an option to control the model’s focus.

In this work, we aim to develop a mechanism to steer the model to generate output relevant to some user-specified input spans (which we term as highlights).<sup>1</sup> This goal, however, brings about significant challenges. For one, many popular NLG datasets are collected in an end-to-end manner, i.e., without annotations of which spans of input are most relevant to the reference target. It would also be ideal for the proposed approach to be compatible

<sup>1</sup>To avoid confusion, our goal is *not* about controlling the attention modules inside the model, instead, we care about the actual generation.

with existing pretrained transformer models, as re-training such models is often costly.

In this work, we propose an *focus vector* framework to address the challenges outlined above. Our contributions are as follows:

- To control the model’s focus, we augment the pretrained model with trainable focus vectors which are directly applied to the encoder embeddings. The model itself is kept fixed, and no further changes to the model architecture is needed.
- The training of focus vectors does not require additional annotations. We utilize attribution methods to derive automatic highlight annotations from existing end-to-end training data.
- For principled evaluation and future work in this direction, we collect and release human evaluation data where the highlight-generation pairs are annotated by humans.
- We test our method on two core NLG tasks: dialogue response generation and abstractive summarization. Experiments show that the trained focus vectors are effective in steering the model to generate a relevant output given the selected highlights.

## 2 Model Formulation

We assume the target model is a standard pretrained transformer encoder-decoder model (Vaswani et al., 2017) that has already been finetuned on end-to-end task-specific data (e.g., dialogue or summarization) with the standard negative log-likelihood (NLL) loss. Our goal is to establish a control mechanism whereby the user can highlight several spans of the input, and the model is supposed to generate outputs relevant to the highlighted text. Crucially, this mechanism should not change the base model, in order to allow the user to default back to the original model if desired.

We begin by establishing notation. We denote the end-to-end training data by  $\{\mathbf{x}, \mathbf{y}\}$ , where  $\mathbf{x} = \{x_1, \dots, x_n\}$  refers to the input token sequence, and  $\mathbf{y}$  refers to the corresponding reference target token sequence. During evaluation, some spans of the input  $\mathbf{x}$  will be highlighted, and we use a binary indicator  $c_i$  to indicate whether the  $i^{\text{th}}$  input token is to be highlighted during generation. In this work we only consider a set of complete sentences as a valid highlight span. This design choice is mainly

for convenience during our human-annotated evaluation data collection, and our framework can readily be generalized to phrase-level highlights.

Suppose the encoder model is composed of  $L$  transformer layers. We denote the  $d$ -dimensional output embedding of the  $i^{\text{th}}$  position on the  $l^{\text{th}}$  encoder layer by  $\mathbf{h}_i^l$ . We use  $\{\mathbf{h}_i^0\}$  to denote the input embeddings. Each decoder layer performs multi-head cross-attention on the outputs of the encoder, where the attention weight computation for the  $h^{\text{th}}$  head on the  $l^{\text{th}}$  decoder layer is formulated as below:

$$\alpha_{i,j}^{h,l} = \operatorname{softmax}_{i \in \{1 \dots n\}} \left( \frac{k(\mathbf{h}_i^L) \cdot \mathbf{q}_j^{h,l}}{\sqrt{d}} \right). \quad (1)$$

Here  $k(\cdot)$  is a linear transform, and  $\alpha_{i,j}$  is the attention weight assigned to encoder output  $\mathbf{h}_i^L$ , for the  $j^{\text{th}}$  position decoder query vector  $\mathbf{q}_j$ . We use  $P_M(\mathbf{y}|\mathbf{x})$  to denote the probability assigned to  $\mathbf{y}$  given input  $\mathbf{x}$  by the original target model. For more details of the transformer encoder-decoder architecture, we refer readers to Vaswani et al. (2017).

Our proposed framework involves two stages. We first obtain automatic highlight annotations using attribution methods. Then, these annotations are used to train the focus vectors. In the next section, we review the attribution methods.

### 2.1 Attribution Methods

Many popular NLG datasets are collected end-to-end, i.e., without annotations of which spans of input are most relevant to the reference target. To obtain these annotations for focus vector training, we make use of existing attribution methods.

Attribution methods (Baehrens et al., 2010; Simonyan et al., 2014; Shrikumar et al., 2017; Adedbayo et al., 2018; Sundararajan et al., 2017), also known as *saliency maps*, attribute the prediction of a (potentially black-box) model to its input features. It thus fits our need to extract relevant spans in the input given the reference target. Most saliency methods are originally designed for image classification, where an importance score is assigned for each dimension of the input feature. Therefore, slight modifications (e.g., dot-product with the word embeddings) are needed to apply them to language data (Ding and Koehn, 2021; Denil et al., 2014).

We implement and compare several popular attribution methods, which compute the attribution

score for a given sentence  $S$  (denoting the set of token indexes in the sentence) in the input  $\mathbf{x}$  for the target  $\mathbf{y}$  and model  $P_M$ .

**Leave-one-out (LOO)** We replace the tokens in  $S$  by the  $\langle \text{pad} \rangle$  token, and compute the difference in NLL:

$$A(S) = \log P_M(\mathbf{y}|\mathbf{x}) - \log P_M(\mathbf{y}|\mathbf{x}_{S\text{-padded}}). \quad (2)$$

LOO is also referred to as an *occlusion-based method* (Zeiler and Fergus, 2014; Li et al., 2016) in the literature.

**Attention-weight** We sum up the attention weights assigned to tokens in  $S$  for all attention heads across all decoder layers:

$$A(S) = \sum_{i \in S} \sum_{j, h, l} \alpha_{i, j}^{h, l}. \quad (3)$$

**Grad-norm** We sum the norm of gradient for the input word embeddings in  $S$ :

$$A(S) = \sum_{i \in S} \|\nabla_{\mathbf{h}_i^0} \log P_M(\mathbf{y}|\mathbf{x})\|_2. \quad (4)$$

**Grad-input-product** Instead of taking vector norm, we compute the dot-product between the input embedding and its gradient:

$$A(S) = \sum_{i \in S} \left( \nabla_{\mathbf{h}_i^0} \log P_M(\mathbf{y}|\mathbf{x}) \right) \cdot \mathbf{h}_i^0. \quad (5)$$

While more sophisticated attribution method have been proposed in the literature (Lei et al., 2016; Sundararajan et al., 2017; Bastings et al., 2019), we mainly experiment with the methods listed above due to their simplicity and popularity. Attribution methods have been used for interpreting black-box models—**applying them to derive labels that can further be used to control the focus of a model has to our knowledge not been explored before.**

Which attribution method best reflects the model’s inner working is still an active research area (Ding and Koehn, 2021; Adebayo et al., 2018). The present work is primarily concerned with how well the attribution scores align with human-annotated highlights. In our experiments, we find that leave-one-out (LOO) has the best correlation on the human-annotated development set (Table 1, details given in Section 3). We therefore adopt LOO to derive the automatic highlight annotations.

More specifically, for the input-output pairs in the training set, we sort the LOO attribution scores of the sentences in the input from large to small,

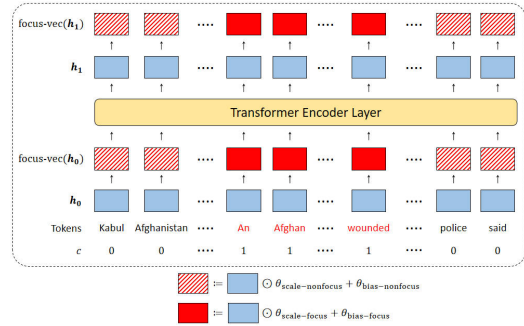


Figure 2: Illustration of our proposed focus vectors applied to a one-layer transformer encoder. The parameters of the transformer model are kept fixed. The highlighted spans are filled by red.

and mark the tokens in the first few sentences (the exact number varies by task) as highlights. We denote the highlight labels obtained from this automatic procedure by a binary indicator variables  $\mathbf{c}^{\text{attr}} = \{c_1^{\text{attr}}, \dots, c_n^{\text{attr}}\}$ , which will be used to train the focus vectors.

## 2.2 Focus Vectors

To control the model’s focus, we introduce a set of  $d$ -dimensional vectors  $\theta$ , named *focus vectors*. They are designed to act as *indicators* for the model, designating which parts of the input to focus on. We now assume the training set contains  $\{\mathbf{x}, \mathbf{c}^{\text{attr}}, \mathbf{y}\}$  triples, where  $\mathbf{c}^{\text{attr}}$  is obtained from the attribution method from the previous section. Focus vectors modify the forward pass of the encoder model by applying a simple transformation  $f$  on the output embeddings of each layer (including the input layer):

$$f(\mathbf{h}_i^l) = \begin{cases} \mathbf{h}_i^l \odot \theta_{\text{scale-focus}}^l + \theta_{\text{bias-focus}}^l, & \text{if } c_i^{\text{attr}} = 1 \\ \mathbf{h}_i^l \odot \theta_{\text{scale-nonfocus}}^l + \theta_{\text{bias-nonfocus}}^l, & \text{if } c_i^{\text{attr}} = 0 \end{cases}. \quad (6)$$

We provide an illustration in Figure 2. The total number of parameters introduced by the focus vectors is therefore  $4 \times (L + 1) \times d$ , which is negligible in comparison to the large number of parameters of the fixed transformer model. We note that as the focus vectors operate directly on the encoder embeddings, it does not require an explicit attention module to exist in the model and is therefore applicable to non-attentional architectures such as LSTMs (Huang et al., 2015).

We train the focus vectors using the standard NLL loss with stochastic gradient descent (SGD):

$$\mathcal{L}(\mathbf{x}, \mathbf{y}, \mathbf{c}^{\text{attr}}; \theta) = -\log P_{\text{focus}}(\mathbf{y}|\mathbf{x}, \mathbf{c}^{\text{attr}}), \quad (7)$$

where  $P_{\text{focus}}(\cdot|\mathbf{x}, \mathbf{c}^{\text{attr}})$  denotes the distribution over the output after the focus vectors are applied. We

Attribution Method	PersonaChat P@1(%)	CNN/Dailymail P@1(%)
attention-weight	29.18	40.31
grad-norm	54.00	43.87
grad-input-product	44.05	32.60
leave-one-out	<b>62.31</b>	<b>64.43</b>

Table 1: Top-1 precision (%) of different attribution methods on the human-labeled development set.

re-iterate that during training of the focus vectors, the transformer model is kept fixed. This allows the user to default back to the pretrained model (i.e., without applying the focus vectors), if the user prefers not to specify any highlights.

Readers may wonder what is the difference between our approach and standard end-to-end training, as both cases use the same  $\mathbf{x}, \mathbf{y}$  pairs. This is related to our key assumption that *different focus of the input lead to different generations*, and the fact that  $\mathbf{c}^{\text{attr}}$  is the relevant span for  $\mathbf{y}$  in the ideal case. Therefore, the focus vectors have the opportunity to give information about which span is more relevant to  $\mathbf{y}$ , before the model observes  $\mathbf{y}$  on the decoder side. To reduce the loss  $-\log P_{\text{focus}}(\mathbf{y}|\mathbf{x}, \mathbf{c}^{\text{attr}})$ , the focus vectors need to steer the model’s focus towards the spans marked by  $\mathbf{c}^{\text{attr}}$ .

At test time, the user will highlight several sentences in the input which we denote by  $\mathbf{c}^{\text{user}}$ . We apply the trained focus vector according to Equation 6, and decode the output from  $P_{\text{focus}}(\cdot|\mathbf{x}, \mathbf{c}^{\text{user}})$ .

### 3 Evaluation Data Collection

We test our method on two NLG tasks: dialogue response generation and abstractive summarization. For the dialogue task, we adopt the *PersonaChat* dataset (Zhang et al., 2018). It is an open domain multi-turn chit-chat dataset, where two participants are required to get to know each other by chatting naturally. Each of them is given a *persona*: several pieces of personal information such as “*I major in Computer Science*”, serving as background information. The participants are required to reflect their assigned persona in the conversation. For summarization, we adopt the *CNN/Dailymail* dataset (Hermann et al., 2015; Nallapati et al., 2016), which is a standard dataset for end-to-end abstractive summarization. To save space, we defer details and statistics of the datasets to Appendix A.

Both PersonaChat and CNN/Dailymail are created end-to-end and do not contain annotated highlight spans. **For principled evaluation, we utilize the Amazon Mechanical Turk (AMT) platform**

**to collect evaluation sets where the highlight-generation pairs are annotated by humans.**

For PersonaChat, each turker<sup>2</sup> is shown a dialogue history and the corresponding persona of the speaker. The dialogue history is randomly selected from the original test set of PersonaChat. Then the turker is required to choose 1-3 sentences as highlights (for example, one sentence in persona, and one sentence in dialogue history), and write down a dialogue response that not only continues the current dialogue, but also is relevant to the chosen highlights. Finally, we ask the turker to repeat the above process, but select a different set of highlights and provide another response. After a few preliminary trials and modifications to our instructions / rewards, we find that turkers comply nicely with our instructions and provide high-quality highlight-response pairs.

For CNN/Dailymail however, we first found that turkers had difficulty writing a high-quality summary for a given news article, with many turkers giving random responses even after we increased the reward. This is perhaps unsurprising given that writing a good summary is challenging and the reference summaries are written by experts. After a few disappointing iterations, we turn to a compromise: we directly provide the turkers with the reference summary, and only ask them to select 2-5 relevant sentences in the article. This simplifies the task, and we are able to collect high-quality labels. This compromise is not ideal, as it reverses the order of highlighting and generation. However, we find that in most cases, the reference summaries in CNN/Dailymail are well covered by several “key” sentences in the article, which are highlighted by the turkers. Therefore, we believe this compromise does not hurt the soundness of our evaluation.

In order to ensure high data quality for both dialogue and summarization, we design a qualification test that turkers need to pass before conducting the actual tasks. Several automatic checks and a minimal time limit are added in the scripts to prevent trivial answers. We also manually monitor the incoming submissions, and ban misbehaving turkers and discard their submissions. More details about our AMT setup are provided in Appendix B.

Our final collected datasets include 3,902 highlight-generation pairs for PersonaChat, and 4,159 pairs for CNN/Dailymail. They are randomly split 50/50 into dev/test sets. **We include a number**

<sup>2</sup>We recruit turkers located in the U.S.

330 **of samples of our collected data in the supplementary materials.** Our code and the collected  
331 dataset will be released in the public version of  
332 this manuscript. We hope that this evaluation data  
333 could facilitate future research in this direction.  
334

335 **Comparison of Attribution Methods** We use  
336 the collected highlight-generation pairs in the dev  
337 set to compare which attribution method aligns  
338 best with human-annotated highlights. In particu-  
339 lar, we compute the top-one precision of the sen-  
340 tence ranked highest by the attribution method. The  
341 result is shown in Table 1. We find that for both Per-  
342 sonaChat and CNN/Dailymail, LOO has the best  
343 alignment. We therefore use LOO to obtain auto-  
344 matic annotations for focus vector training. Inter-  
345 estingly, we observe low alignment between atten-  
346 tion weight-derived attribution scores and human  
347 judgment, which indicates that controlling model  
348 generations via intervening on the attention distri-  
349 butions may not optimal. Finally, we note that this  
350 result does not mean LOO is the “best” attribution  
351 method, as attribution method is supposed to reflect  
352 the model’s inner working, instead of a human’s.

## 353 4 Experiments

### 354 4.1 Experiment Setting and Baselines

355 We use Blenderbot (Roller et al., 2020) as the  
356 base model for PersonaChat and BART (Lewis  
357 et al., 2019) for CNN/Dailymail, both of which are  
358 standard encoder-decoder transformer models. Our  
359 code is based on the *transformers* library (Wolf  
360 et al., 2020). We load the pretrained weights from  
361 facebook/blenderbot-400M-distill  
362 and facebook/bart-base. Blenderbot has  
363 2 encoder layers and 12 decoder layers, while  
364 BART has 6 encoder layers and 6 decoder layers.  
365 To help Blenderbot cope with long dialogue  
366 context in PersonaChat, we extend its maximum  
367 position embedding index from 128 to 256. We use  
368 beam-search for decoding, where we follow the  
369 recommended configuration (Roller et al., 2020;  
370 Lewis et al., 2019) and use a beam size of 10 for  
371 Blenderbot and a beam size of 4 for BART.

372 For both tasks, we first finetune the base model  
373 on the original training set in the standard end-to-  
374 end manner. The model is then fixed and used to  
375 obtain automatic labels  $c^{\text{attr}}$  with the LOO attribu-  
376 tion method on the same training set. For each  
377 training sample, we select the top- $k$  sentences in  
378 the input ranked by LOO. Since we do not know  
379 the best value for  $k$ , we set it to be a random num-

ber from 1 to 3 for PersonaChat, and from 2 to 5  
for CNN/Dailymail.

380 While the highlight labels in the training set  
381 used to train focus vectors are derived automati-  
382 cally, we use the human-labeled dev set for hyper-  
383 parameter tuning. This is to facilitate fair compar-  
384 ison with other baseline approaches which also  
385 utilize the human-labeled dev set. In our abla-  
386 tion study, we will show that this dependence  
387 on human-labeled dev set is not crucial for our  
388 approach to achieve strong performance. We  
389 perform a grid search over learning rate with  
390  $\{1, 3, 5\} \times \{1e^{-4}, 1e^{-3}, 1e^{-2}, 1e^{-1}\}$ . The Adam  
391 optimizer (Kingma and Ba, 2014) is used with  
392  $\beta_1 = 0.9, \beta_2 = 0.999$ , and a L2 decay weight  
393 of 0.01. For both tasks, we set the mini-batch size  
394 to be 16.  
395

396 We compare the proposed focus-vector approach  
397 with several baselines: 398

**Vanilla:** The vanilla model, without any modifi-  
399 cation in both the model and the input. 400

**Padding:** One trivial way to control the model’s  
401 focus is to replace all input by the `<pad>` token,  
402 except the spans highlighted by the user. However,  
403 we find that this direct padding during evaluation  
404 results in drastically worse perplexity. To allevi-  
405 ate this problem, we randomly pad a portion of  
406 sentences in the input during the standard end-to-  
407 end finetuning, to make the model aware that only  
408 partial input would be provided. 409

**Keyword-control:** Keyword-based prompts (Fan  
410 et al., 2017; He et al., 2020) has been a popular  
411 approach for controllable text generation. We adapt  
412 this idea to our focus-control setting. During model  
413 finetuning, we prepend key-phrases extracted from  
414 the reference target sequence to the original input.  
415 We utilize Yake (Campos et al., 2020), which is an  
416 unsupervised keyword extraction method. During  
417 evaluation, we extract and prepend key-phrases  
418 extracted from the highlighted sentences. 419

**Attention-offset:** As a direct way to control the  
420 model’s attention, we add a positive scalar offset  
421  $s^{\text{offset}}$  to the cross-attention heads before the soft-  
422 max operation (Equation 1), for the highlighted  
423 spans. A similar technique has been used in Dong  
424 et al. (2021) to *modulate* the attention distribution  
425 to tackle neural text degeneration problems (Holtz-  
426 man et al., 2019). This approach modifies the at-  
427

Model	PersonaChat			CNN/Dailymail		
	PPL	ROUGE-1/2/L	BERTScore	PPL	ROUGE-1/2/L	BERTScore
vanilla (w.o. highlight)	28.73	17.02/2.73/14.52	85.41	4.51	43.48/21.01/30.98	89.23
padding	38.93	16.69/2.80/13.72	84.42	19.62	39.31/18.44/28.67	88.34
keyword-control	23.64	17.31/3.02/14.81	85.58	4.56	43.81/21.08/31.15	89.26
attention-offset	23.79	<b>21.10</b> /3.77/17.54	86.04	4.49	43.96/20.64/31.26	89.28
<b>focus-vector</b>	<b>22.51</b>	20.81/ <b>3.98</b> / <b>17.58</b>	<b>86.13</b>	<b>4.48</b>	<b>45.92</b> / <b>23.03</b> / <b>32.98</b>	<b>89.78</b>

Table 2: Main evaluation results on the PersonaChat and CNN/Dailymail datasets with annotated highlights. The proposed focus vector approach shows strong performance across different metrics.

tention weights via:

$$\alpha'_{i,j} = \text{softmax}_{i \in \{1 \dots n\}} \left( \frac{k(\mathbf{h}_i^L) \cdot \mathbf{q}_j}{\sqrt{d}} + s^{\text{offset}} \cdot \mathbb{1}_{[c_i=1]} \right), \quad (8)$$

where  $s^{\text{offset}}$  is a hyper-parameter, and is applied to all cross-attention heads in the decoder. We tune  $s^{\text{offset}}$  on the human-annotated development set in a fine-grained manner. More details are given in Appendix C.

Whether the attention distribution faithfully explains a model’s predictions is the subject of much debate (Jain and Wallace, 2019; Wiegrefe and Pinter, 2019; Bastings and Filippova, 2020). Therefore this direct modification of the attention head may not be the optimal solution for focus control. Our proposed focus-vector framework, on the other hand, utilizes attribution methods, and directly operates on the encoder embeddings.

## 4.2 Results and Analysis

During evaluation, human-annotated highlights are fed to the model. In addition to perplexity, we evaluate the generations from different approaches using two popular NLG metrics: ROUGE (Lin, 2004), and BERTScore (Zhang et al., 2019).

We show the main results in Table 2. As expected, the padding baseline has poor performance, as a large portion of input is masked out. Comparing to various baselines, focus-vector obtains significantly improved ROUGE and BERTScore on both tasks. This validates the motivation of this work: focus-vec is effective in steering the model’s focus, which leads towards the desired generation. For CNN/Dailymail, the perplexity of focus-vector is close to the vanilla model even though there is a large difference in ROUGE. We believe this is due to the constrained nature of the summarization task and how perplexity is computed: once the model observes the first few tokens, it is easy to figure out what the current highlight is. The other two metrics, on the other hand, are based on the actual generation, and therefore does not have this issue.

The performance of keyword-control, although better than the vanilla model, is inferior to attention-offset and focus-vector. We surmise this is due to the following two weakness: First, key-phrases can not fully represent the highlighted span. Second, there is a discrepancy of where the key-phrases are extracted between training and evaluation.

The performance gap (in ROUGE/BERTScore) between focus-vector and attention-offset is larger on the CNN/Dailymail dataset. We believe this is because the BART model has a deeper encoder than the Blenderbot model. As the encoder grows deeper, the embeddings become more “contextualized” and its *identifiability* (Brunner et al., 2020) degrades. And since the decoder attends to the last layer of the encoder, this direct manipulation of attention weights could be ineffective with deep encoders.

Table 3 shows generation samples from different focus-control approaches for PersonaChat. Spans of the generation that are relevant to the highlighted persona are marked in red. Comparing to the generation from the vanilla model, the generations from both attention-offset and focus-vector are highly relevant to the respective highlighted persona. One generation from att-offset is a little erratic (“*I am petro, my dog*”), which may be due to the inflexibility of att-offset.

We defer the generation examples for CNN/Dailymail to Table 6 and Table 7 (Appendix D) due to space constraints. We observe that the generation from focus-vector is more focused on the highlighted inputs. On the other hand, attention-offset’s generation still remains similar to the vanilla model.

In Figure 3, we study how the outputs of attribution methods (attention-weight and grad-norm) change with different approaches (vanilla, focus-vector and attention-offset) for the CNN/Dailymail example (Table 6). Note that in this analysis, for the attribution methods we set the target  $\mathbf{y}$  to be the decoded output from the respective modeling,

persona:	<b>I have a dog named pedro.</b> I like to eat muffins. I am five feet tall. I work in the healthcare industry.	I have a dog named pedro. I like to eat muffins. I am five feet tall. <b>I work in the healthcare industry.</b>
context:	[P-1]: Hello, how are you today? [P-2]: [TO ANSWER]	[P-1]: Hello, how are you today? [P-2]: [TO ANSWER]
reference:	Scared, i'm looking for <b>my dog</b> , his name is pedro. have you seen him?	Really tired, it's so busy in <b>healthcare industry</b> . I wish I can have a long vacation.
vanilla:	Hi! I am doing well, how about yourself? What do you do for a living?	Hi! I am doing well, how about yourself? What do you do for a living?
keyword:	I am doing well, how about yourself? What do you like to do for fun?	I am doing well. how about you? What do you do for a living? I am <b>in healthcare</b> .
att-offset:	I am doing well, and you? <b>I am petro, my dog</b> , and you?	I am doing well, and you? I am a <b>healthcare worker in the health industry</b> .
focus-vec:	I am good. Just got back from <b>walking my dog</b> . Do you have any pets?	I am doing well. How about you? What do you do for a living? I am <b>in healthcare</b> .

Table 3: Generation samples from different attention control approaches for PersonaChat. The inputs on both sides are identical, except a different piece of persona is highlighted (marked in bold), which leads to different generations for approaches with control. Another example is provided in Table 5 (Appendix D).

instead of the reference summary. The highlighted sentences are marked by the red rectangles.

We observe that for both attention-weight and grad-norm, the application of focus vector makes the highlighted sentences obtain the highest attribution scores, and the scores differ significantly from the vanilla model. In some of the non-highlighted sentences (marked by the blue rectangles), attention-offset is not strong enough to significantly reduce its attribution. We also tried larger values of  $s^{\text{offset}}$  for attention-offset but found it lead to performance degradation. This analysis shows that despite the small number of parameters associated with the focus vectors, they are able to effectively steer the model’s focus. We provide a simple visualization of the trained focus-vector parameters in Figure 3 (Appendix D).

**Ablation Studies** Table 4 shows several variants of focus vector on CNN/Dailymail. We first tune the hyper-parameters of focus vector only with the original dev set with  $c^{\text{attr}}$ , instead of human-annotated highlights. Despite this discrepancy, focus vector still achieves strong performance on the test set. This result shows that the use of human-annotated dev set is not crucial for our framework. We then conduct an ablation study where we only apply focus vector on the first or last layer of the encoder, which reduces the number of parameters. We find that this results in marginal performance degradation. Finally, we jointly finetune focus vector and the whole model with the same loss function (Equation 7), where a separate and smaller learning rate is used for the model. Interestingly, the gain from model finetuning is very limited, which demonstrates the effectiveness of focus vector.

Model	CNN/Dailymail		
	PPL	ROUGE-1/2/L	BERTScore
all-layer*	4.48	45.92/23.03/32.98	89.78
ori-dev with $c^{\text{attr}}$	4.50	46.41/22.69/32.48	89.62
only first layer	4.48	45.67/22.63/32.45	89.59
only last layer	4.48	46.06/22.84/32.69	89.69
plus model finetune	4.49	46.65/23.54/33.30	89.82

Table 4: Performance of different variants of focus-vector trained on CNN/Dailymail. all-layer\* refers to our proposed modelling (also reported in Table 2).

## 5 Related Work

Our proposed focus-vector framework is closely related to the research topics of controllable text generation, LM adaptation, and attention/attribution analysis, which we review below.

**Controllable Text Generation** Prior work on controllable summarization introduced various types of control mechanisms. Fan et al. (2017); Saito et al. (2020) extract entity, keyword or length, as additional supervision during training. Gehrmann et al. (2018) trains a token-level content selection module, where the supervision is by aligning the summaries to the documents. (Song et al., 2021) proposes a two-staged generation strategy and Goyal et al. (2021) incorporates multiple decoders into a transformer framework. Some recent work (He et al., 2020; Dou et al., 2020) uses prompts to control the generation. Lexically constrained decoding (Post and Vilar, 2018) has also been used to enforce certain key phrases to be included in the summary (Mao et al., 2020).

Existing work on controllable dialogue response generation include using conditional variational au-

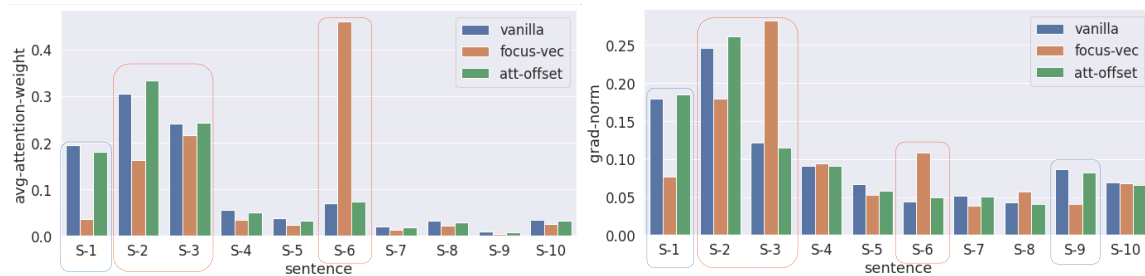


Figure 3: Attribution scores for each sentence in the input, with different focus-control approach applied to BART. The highlighted sentences are marked by red rectangles. The corresponding example is in Table 6 (Appendix D).

toencoders (Zhao et al., 2017; Li et al., 2020), and incorporating external knowledge into the conversational agent using knowledge graphs (Cui et al., 2021; Moon et al., 2019), unstructured documents (Kim et al., 2020), or dialogue context (Zhao et al., 2020).

In open-ended language generation, a series of approaches have been proposed to control for some attribute (e.g., topic) of the generation (Keskar et al., 2019; Dathathri et al., 2020; Krause et al., 2020; Yang and Klein, 2021). Some of these studies utilize a trained classifier to guide the generative model towards the desired attribute.

**LM Adaptation** Our proposed focus vector framework is also inspired by a series of recent works on prompting or light-weight LM adaptation. Li and Liang (2021), followed by Lester et al. (2021) and Zhong et al. (2021), propose *prefix tuning*, where continuous task-specific input vectors are tuned to adapt the pretrained LM to a downstream task with supervised data, and the model is kept fixed.

There is also a line of works on *adapter-tuning*, which insert and finetune task-specific layers (adapters) between each layer of the pretrained LM (Houlsby et al., 2019; Lin et al., 2020; Pfeiffer et al., 2021). More recently, Guo et al. (2021) and Ben-Zaken et al. (2020) propose to finetune only a small subset of a pretrained model’s parameters, and achieves strong performance on the GLUE benchmark (Wang et al., 2018).

### Attention Analysis and Attribution Methods

Due to the ubiquity of the attention module in current NLP models, various work has studied how the module captures linguistic phenomena in the input (Clark et al., 2019; Kovaleva et al., 2019; Kobayashi et al., 2020). It has also been used as a tool to interpret the model’s predictions (Wang et al., 2016; Lee et al., 2017; Ghaeini et al., 2018).

Recently, there have been a series of studies discussing the use of attention weights for inter-

pretability (Jain and Wallace, 2019; Wiegrefe and Pinter, 2019; Bastings and Filippova, 2020; Serano and Smith, 2019), and it has been argued that attribution methods are a better choice to explain the model’s predictions. The poor alignment performance of attention weights that we get in Table 1, on some level, are in agreement with that argument. Our work is also related to the line of work on interpreting black box models through *rationales* (Lei et al., 2016; Bastings et al., 2019), which are typically (discrete) subsets of the input that are used to predict the output. Finally, several recent works (Xu and Durrett, 2021; Ding and Koehn, 2021) have compared different attribution methods for interpreting NLP models.

In comparison to the aforementioned works, our major innovations are two fold: (1) Our goal is to control the *focus* of pretrained models, and thereby steer the model’s generation, and our proposed focus vectors are compatible with the standard transformer architecture; (2) We utilize attribution methods to obtain automatic annotations for focus-vector training. Therefore, our framework can be applied to a wide range of NLG applications.

## 6 Conclusion

In this work we propose the focus vector framework as a light-weight solution to control the focus of pretrained transformer models. It has two major advantages: (1) Focus vectors act as simple transformations to the embeddings in the encoder, and the transformer model is kept fixed; (2) Attribution methods are utilized to get automatic highlight labels for training focus vectors.

We test our approach on two tasks: dialogue response generation, and abstractive summarization. For evaluation, we collect data where the highlight-generation pairs are annotated by humans. Experiments show that the trained focus vectors are effective in steering the model to generate output text that is relevant to the specified highlights.



648  
649  
650  
651  
652  
653  
  
654  
655  
656  
657  
658  
  
659  
660  
661  
  
662  
663  
664  
665  
666  
667  
  
668  
669  
670  
671  
672  
673  
674  
  
675  
676  
677  
  
678  
679  
680  
681  
682  
  
683  
684  
685  
686  
687  
  
688  
689  
690  
691  
692  
693  
694  
  
695  
696  
697  
698  
  
699  
700  
701  
702

## References

Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. 2018. [Sanity checks for saliency maps](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.

David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. 2010. [How to explain individual classification decisions](#). *Journal of Machine Learning Research*, 11(61):1803–1831.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2016. [Neural machine translation by jointly learning to align and translate](#).

Jasmijn Bastings, Wilker Aziz, and Ivan Titov. 2019. [Interpretable neural predictions with differentiable binary variables](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2963–2977, Florence, Italy. Association for Computational Linguistics.

Jasmijn Bastings and Katja Filippova. 2020. [The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?](#) In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 149–155, Online. Association for Computational Linguistics.

Elad Ben-Zaken, Shauli Ravfogel, and Yoav Goldberg. 2020. [Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models](#).

Gino Brunner, Yang Liu, Damian Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger Wattenhofer. 2020. [On identifiability in transformers](#). In *International Conference on Learning Representations*.

Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Jorge, Célia Nunes, and Adam Jatowt. 2020. [Yake! keyword extraction from single documents using multiple local features](#). *Information Sciences*, 509:257–289.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.

Leyang Cui, Yu Wu, Shujie Liu, and Yue Zhang. 2021. [Knowledge enhanced fine-tuning for better handling unseen entities in dialogue generation](#). *arXiv preprint arXiv:2109.05487*.

Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. [Plug and play language models: A simple approach to controlled text generation](#). In

*International Conference on Learning Representations*. 703  
704

Misha Denil, Alban Demiraj, and Nando de Freitas. 2014. [Extraction of salient sentences from labelled documents](#). *CoRR*, abs/1412.6815. 705  
706  
707

Shuoyang Ding and Philipp Koehn. 2021. [Evaluating saliency methods for neural language models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5034–5052, Online. Association for Computational Linguistics. 708  
709  
710  
711  
712  
713  
714

Yue Dong, Chandra Bhagavatula, Ximing Lu, Jena D. Hwang, Antoine Bosselut, Jackie Chi Kit Cheung, and Yejin Choi. 2021. [On-the-fly attention modulation for neural generation](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1261–1274, Online. Association for Computational Linguistics. 715  
716  
717  
718  
719  
720  
721

Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao Jiang, and Graham Neubig. 2020. [Gsum: A general framework for guided neural abstractive summarization](#). *arXiv preprint arXiv:2010.08014*. 722  
723  
724  
725

Angela Fan, David Grangier, and Michael Auli. 2017. [Controllable abstractive summarization](#). *arXiv preprint arXiv:1711.05217*. 726  
727  
728

Sebastian Gehrmann, Yuntian Deng, and Alexander M. Rush. 2018. [Bottom-up abstractive summarization](#). *CoRR*, abs/1808.10792. 729  
730  
731

Reza Ghaeini, Xiaoli Fern, and Prasad Tadepalli. 2018. [Interpreting recurrent and attention-based neural models: a case study on natural language inference](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4952–4957, Brussels, Belgium. Association for Computational Linguistics. 732  
733  
734  
735  
736  
737  
738

Tanya Goyal, Nazneen Fatema Rajani, Wenhao Liu, and Wojciech Kryściński. 2021. [Hydrasum—disentangling stylistic features in text summarization using multi-decoder models](#). *arXiv preprint arXiv:2110.04400*. 739  
740  
741  
742  
743

Demi Guo, Alexander Rush, and Yoon Kim. 2021. [Parameter-efficient transfer learning with diff pruning](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4884–4896, Online. Association for Computational Linguistics. 744  
745  
746  
747  
748  
749  
750  
751

Junxian He, Wojciech Kryściński, Bryan McCann, Nazneen Rajani, and Caiming Xiong. 2020. [Ctrlsum: Towards generic controllable text summarization](#). *arXiv preprint arXiv:2012.04281*. 752  
753  
754  
755

756	Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. <i>Advances in neural information processing systems</i> , 28:1693–1701.	<i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 121–126, Copenhagen, Denmark. Association for Computational Linguistics.	811 812 813 814
761	Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2019. <a href="#">The curious case of neural text degeneration</a> . <i>CoRR</i> , abs/1904.09751.	Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2016. <a href="#">Rationalizing neural predictions</a> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 107–117, Austin, Texas. Association for Computational Linguistics.	815 816 817 818 819
764	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. <a href="#">Parameter-efficient transfer learning for NLP</a> . In <i>Proceedings of the 36th International Conference on Machine Learning</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 2790–2799. PMLR.	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. <a href="#">The power of scale for parameter-efficient prompt tuning</a> . <i>CoRR</i> , abs/2104.08691.	820 821 822
772	Zhiheng Huang, Wei Xu, and Kai Yu. 2015. <a href="#">Bidirectional LSTM-CRF models for sequence tagging</a> . <i>CoRR</i> , abs/1508.01991.	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. <i>arXiv preprint arXiv:1910.13461</i> .	823 824 825 826 827 828
775	Sarthak Jain and Byron C. Wallace. 2019. <a href="#">Attention is not explanation</a> . <i>CoRR</i> , abs/1902.10186.	Chunyuan Li, Xiang Gao, Yuan Li, Baolin Peng, Xiujun Li, Yizhe Zhang, and Jianfeng Gao. 2020. Optimus: Organizing sentences via pre-trained modeling of a latent space. <i>arXiv preprint arXiv:2004.04092</i> .	829 830 831 832
777	Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. <a href="#">CTRL: A conditional transformer language model for controllable generation</a> . <i>CoRR</i> , abs/1909.05858.	Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. <a href="#">Understanding neural networks through representation erasure</a> . <i>CoRR</i> , abs/1612.08220.	833 834 835
781	Byeongchang Kim, Jaewoo Ahn, and Gunhee Kim. 2020. Sequential latent knowledge selection for knowledge-grounded dialogue. <i>arXiv preprint arXiv:2002.07510</i> .	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-tuning: Optimizing continuous prompts for generation</a> . <i>CoRR</i> , abs/2101.00190.	836 837 838
785	Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	839 840 841
788	Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. <a href="#">Attention is not only a weight: Analyzing transformers with vector norms</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7057–7075, Online. Association for Computational Linguistics.	Zhaojiang Lin, Andrea Madotto, and Pascale Fung. 2020. <a href="#">Exploring versatile generative language model via parameter-efficient transfer learning</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 441–459, Online. Association for Computational Linguistics.	842 843 844 845 846 847
795	Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. <a href="#">Revealing the dark secrets of BERT</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4365–4374, Hong Kong, China. Association for Computational Linguistics.	Yuning Mao, Xiang Ren, Heng Ji, and Jiawei Han. 2020. <a href="#">Constrained abstractive summarization: Preserving factual consistency with constrained generation</a> . <i>CoRR</i> , abs/2010.12723.	848 849 850 851
803	Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq R. Joty, Richard Socher, and Nazneen Fatema Rajani. 2020. <a href="#">Gedi: Generative discriminator guided sequence generation</a> . <i>CoRR</i> , abs/2009.06367.	Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. Opendialkg: Explainable conversational reasoning with attention-based walks over knowledge graphs. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 845–854.	852 853 854 855 856 857
808	Jaesong Lee, Joong-Hwi Shin, and Jun-Seok Kim. 2017. <a href="#">Interactive visualization and manipulation of attention-based neural machine translation</a> . In	Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. <i>CoNLL 2016</i> , page 280.	858 859 860 861
810		Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. <a href="#">AdapterFusion: Non-destructive task composition</a>	862 863 864

865	for transfer learning. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 487–503, Online. Association for Computational Linguistics.	
866		
867		
868		
869		
870	Matt Post and David Vilar. 2018. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1314–1324, New Orleans, Louisiana. Association for Computational Linguistics.	
871		
872		
873		
874		
875		
876		
877		
878	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21(140):1–67.	
879		
880		
881		
882		
883		
884	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric M Smith, et al. 2020. Recipes for building an open-domain chatbot. <i>arXiv preprint arXiv:2004.13637</i> .	
885		
886		
887		
888		
889	Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 379–389, Lisbon, Portugal. Association for Computational Linguistics.	
890		
891		
892		
893		
894		
895	Itsumi Saito, Kyosuke Nishida, Kosuke Nishida, and Junji Tomita. 2020. Abstractive summarization with combination of pre-trained sequence-to-sequence and saliency models. <i>arXiv preprint arXiv:2003.13028</i> .	
896		
897		
898		
899	Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? <i>CoRR</i> , abs/1906.03731.	
900		
901	Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 3145–3153. PMLR.	
902		
903		
904		
905		
906		
907	Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Deep inside convolutional networks: Visualising image classification models and saliency maps. In <i>In Workshop at International Conference on Learning Representations</i> .	
908		
909		
910		
911		
912	Kaiqiang Song, Bingqing Wang, Zhe Feng, and Fei Liu. 2021. A new approach to overgenerating and scoring abstractive summaries. <i>arXiv preprint arXiv:2104.01726</i> .	
913		
914		
915		
916	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 3319–3328. PMLR.	
917		
918		
919		
920		
	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 I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems 30</i> , pages 5998–6008. Curran Associates, Inc.	921 922 923 924 925 926 927 928
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	929 930 931 932 933 934 935 936
	Yequan Wang, Minlie Huang, Xiaoyan Zhu, and Li Zhao. 2016. Attention-based LSTM for aspect-level sentiment classification. In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 606–615, Austin, Texas. Association for Computational Linguistics.	937 938 939 940 941 942
	Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 11–20, Hong Kong, China. Association for Computational Linguistics.	943 944 945 946 947 948 949
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	950 951 952 953 954 955 956 957 958 959 960 961
	Jiacheng Xu and Greg Durrett. 2021. Dissecting generation modes for abstractive summarization models via ablation and attribution. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6925–6940, Online. Association for Computational Linguistics.	962 963 964 965 966 967 968 969
	Kevin Yang and Dan Klein. 2021. FUDGE: Controlled text generation with future discriminators. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3511–3535, Online. Association for Computational Linguistics.	970 971 972 973 974 975 976
	Matthew D. Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In <i>Com-</i>	977 978

979 *puter Vision – ECCV 2014*, pages 818–833, Cham.  
980 Springer International Publishing.

981 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur  
982 Szlam, Douwe Kiela, and Jason Weston. 2018. Per-  
983 sonalizing dialogue agents: I have a dog, do you have  
984 pets too? In *Proceedings of the 56th Annual Meet-*  
985 *ing of the Association for Computational Linguistics*  
986 *(Volume 1: Long Papers)*, pages 2204–2213.

987 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q  
988 Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-  
989 uating text generation with bert. *arXiv preprint*  
990 *arXiv:1904.09675*.

991 Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017.  
992 Learning discourse-level diversity for neural dialog  
993 models using conditional variational autoencoders.  
994 *arXiv preprint arXiv:1703.10960*.

995 Xueliang Zhao, Wei Wu, Can Xu, Chongyang Tao,  
996 Dongyan Zhao, and Rui Yan. 2020. Knowledge-  
997 grounded dialogue generation with pre-trained lan-  
998 guage models. *arXiv preprint arXiv:2010.08824*.

999 Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021.  
1000 [Factual probing is \[MASK\]: learning vs. learning to](#)  
1001 [recall](#). *CoRR*, abs/2104.05240.

## Appendices

### A End-to-end Datasets

The PersonaChat dataset contains 8,939 dialogues for training, 1,000 for validation, and 968 for test. For each turn in the dialogue, we concatenate the persona of the speaker and the dialogue history as input, and train the base model to generate the current utterance. In some cases, the dialogue history is long and exceeds the input limit of the model, in which case we truncate the dialogue at the sentence level. The average number of sentences is around 11 after truncation.

The CNN/Dailymail dataset contains 287,113 training examples, and 13,368 / 11,490 examples for validation / test. We apply the same truncation strategy as *PersonaChat* during preprocessing. The processed articles have an average length of 748 tokens, and the reference summaries have an average length of 67 tokens.

### B Human-annotated Evaluation Data Collection

To improve the quality of collected dataset, we design a qualification test, which the turkers need to pass before they can work on real assignments. The test is designed to help turkers understand our task better. For PersonaChat, we give turkers two dialogue samples with pre-selected highlights, and ask them to choose the appropriate response that not only continues the dialogue, but also is relevant to the highlights. For CNN/Dailymail, the turkers are shown two example articles and the corresponding reference summaries. We have already picked some highlights in the article, but there is one highlight missing. And the turker is required to pick the missing highlight. The interface for the PersonaChat qualification test is shown in Figure 5.

We also add multiple checks in our script to prevent trivial answers. We ban trivial copy&paste from the given context. A time check is added that requires turker to spend at least 60 seconds on a single HIT. For the two assignments in PersonaChat, we add a content check that prevents duplicate highlights or response. We show our interface for PersonaChat in Figure 6. Despite these checks and the qualification tests, there still exist a small number of misbehaving turkers who attempt to cheat. Therefore we also manually monitor the incoming submissions, and ban misbehaving turk-

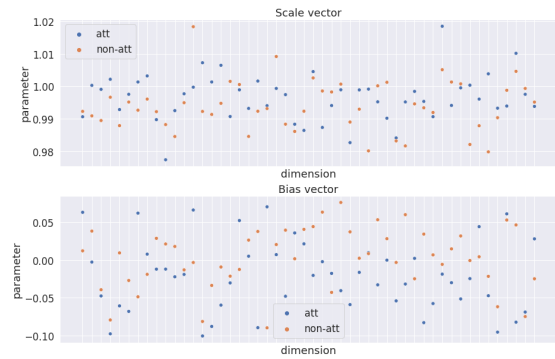


Figure 4: 50 random dimensions of the trained focus vector on first encoder layer of the BART model.

```
Your Persona:
my favorite color is purple.
I have owned two mustangs.
I am currently looking for a job.

Context:
Other: hello there, is not it a great day to adopt a dog?
You: yes it would be a nice thing to do
Other: what are you up to on this fine day?
You: not much just staying home
Other: what do you for fun, i enjoy a good violin piece to play.
You: i like to play video games
Other: what genre? hopefully nothing violent
You: nah, mostly i play rpg ones
Other: nice, enjoy playing while gaming myself, have a favorite food?
You: <Put your response at INPUT area>
```

I'm just staying home and browsing internet to find a new job.

I love spaghetti.

I have owned my mustangs, they're so powerful, really love them.

I love hotpot, a Chinese food. But I haven't had one for a long time, busy looking for a new job.

Sushi.

Figure 5: An example of our AMT qualification test for PersonaChat. We have chosen the highlights in the context, and the turker is supposed to choose a response that not only continues the dialogue, but also is relevant to the highlights.

ers and filter out their submissions. 1051

More examples of our interface and instructions can be found in our uploaded data samples. 1052 1053

### C Implementation Details 1054

For the attention-offset baseline, we tune the offset  $s^{\text{offset}}$  in a fine-grained manner, on the human-annotated dev set. We first set a relatively large max value (100) and get 20 evenly spaced numbers inside the interval (0, 100). Then we calculate model PPL on the dev set with  $s^{\text{offset}}$  set to these different offsets. Then we do another search in the interval that has lowest PPL. We repeat this iteration multiple times, and stops when PPL change is smaller than  $1e^{-3}$ . The final tuned value for Blenderbot is around 3.02, and around 0.17 for BART. 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065

### D Auxiliary Results and Examples 1066

In Figure 4, we provide a simple visualization of the trained focus vectors of BART. To make the 1067 1068

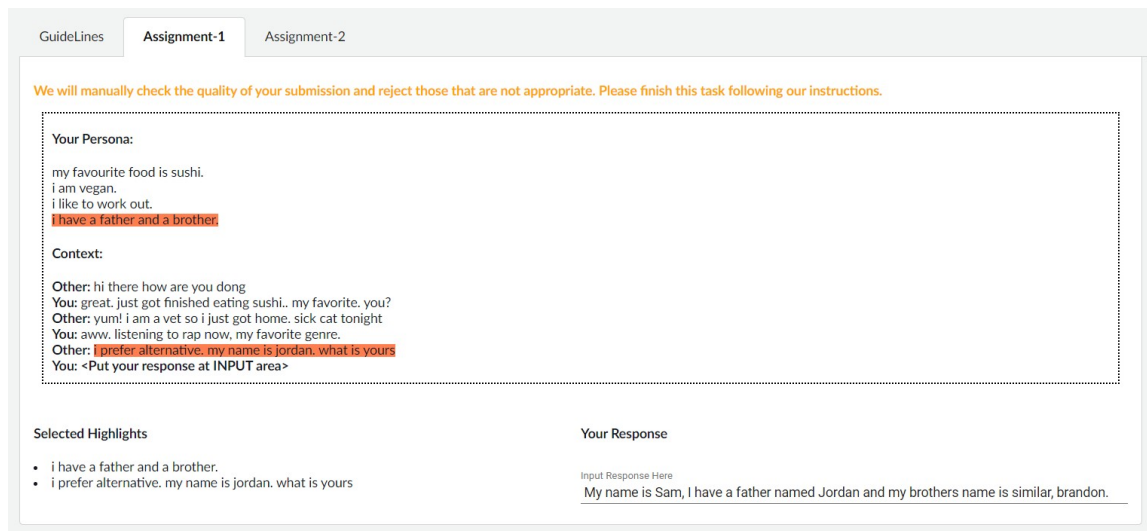


Figure 6: An example of our AMT interface on PersonaChat. The highlights and a response are labeled by a turker.

1069 figure easy to grasp, we randomly sample 50 di-  
 1070 mensions (out of 768) of trained focus-vector pa-  
 1071 rameters. In the figure only the trained parameters  
 1072 for the first encoder layer is shown, and we find  
 1073 that the patterns for the other layers are similar.

1074 We first find that the learned scaling parameters  
 1075 are close to 1 and the bias parameters are close to  
 1076 zero. This implies that the transformation defined  
 1077 by the focus-vec is not drastic, and a relatively  
 1078 small tweak could be enough to steer the model’s  
 1079 focus.

1080 An interesting pattern can be observed for the  
 1081 scaling vectors.  $\theta_{\text{scale-nonfocus}}$  tends to down-scale  
 1082 the embedding, comparing to  $\theta_{\text{scale-focus}}$ . This  
 1083 matches our intuition that the embeddings on the  
 1084 non-highlighted positions should be down-played.

persona:	My parents did not want me. <b>It's a dead end job so i am looking for something different.</b> I was raised by my uncle and aunt.	My parents did not want me. It's a dead end job so i am looking for something different. <b>I was raised by my uncle and aunt.</b>
context:	[P-1]: Hello how are you doing? [P-2]: Pretty good. how are you? [P-1]: I am making dinner getting ready for bed. [P-2]: I just got off work. [P-1]: <b>I am looking for something different to do.</b> [P-2]: [TO ANSWER]	[P-1]: Hello how are you doing? [P-2]: Pretty good. how are you? [P-1]: I am making dinner getting ready for bed. [P-2]: I just got off work. [P-1]: <b>I am looking for something different to do.</b> [P-2]: [TO ANSWER]
reference:	Do you know anywhere that's hiring work from home? I'm <b>looking for a new job.</b>	I grow up with <b>my aunt</b> and we always had scripture study after dinner, you can try it.
vanilla:	My parents wanted me to be a doctor like them, but i do not want to do that.	My parents wanted me to be a doctor like them, but i do not want to do that.
keyword:	What do you do for <b>work</b> ? I am a cashier at a fast food restaurant.	What do you do for work? I am a cashier at a fast food restaurant.
att-offset:	What do you do for a living? I have been <b>looking for a new job</b> myself.	Do you have any relatives you can ask for help with? my family was always there for me.
focus-vec:	What do you do for a living? I am unemployed right now <b>looking for a job.</b>	What do you do for a living? I work with <b>my aunt</b> at her law firm.

Table 5: Generation samples from different attention control approaches for PersonaChat. The inputs on both sides are identical, except a different piece of persona is highlighted, which leads to different generations from att-offset and focus-vec.

article:	England World Cup hopeful Elliot Daly has been named the Aviva Premiership's player of the month for March. <sup>[1]</sup> <b>The uncapped Wasps centre, 22, scored two tries in March as the Coventry-based club maintained on course for a Champions Cup spot next season.</b> <sup>[2]</sup> <b>Equally comfortable at full back, Daly's form at centre has him pushing incumbent England No 13 Jonathan Joseph hard for his place.</b> <sup>[3]</sup> Elliot Daly is handed his player of the month award by Wasps director of rugby Dai Young. <sup>[4]</sup> He looks certain to be included in Stuart Lancaster's England's World Cup training squad named next month. <sup>[5]</sup> <b>Wasps director of rugby Dai Young said: 'If England are going to go out and win the Rugby World Cup, he's the type of player you've got to have in the team so I'd certainly like to think he's pushing for a place in the squad.'</b> <sup>[6]</sup> BT Sport analyst Austin Healey added: 'For a good couple of years now us rugby watchers have known of the potential of Elliot Daly.' <sup>[7]</sup> . . . (LEAVE OUT UNIMPORTANT SENTENCES)
reference:	Elliot Daly was in fine form at outside centre for Wasps in March. Daly, 22, has not yet been capped by England. <b>Wasps head coach Dai Young says England must pick him for World Cup.</b>
vanilla:	England World Cup hopeful Elliot Daly has been named Aviva Premiership's player of the month for March. The 22-year-old scored two tries in March as Coventry-based club maintained on course for a Champions Cup spot next season. Daly's form at centre has him pushing incumbent England No 13 Jonathan Joseph hard for his place.
keyword:	England World Cup hopeful Elliot Daly has been named Aviva Premiership's player of the month. The 22-year-old scored two tries in March as Coventry-based club maintained on course for a Champions Cup spot next season. Daly looks certain to be included in Stuart Lancaster's England's World Cup training squad.
att-offset:	England World Cup hopeful Elliot Daly has been named Aviva Premiership's player of the month for March. The uncapped Wasps centre scored two tries in March as Coventry-based club maintained on course for a Champions Cup spot next season. Daly's form at centre has him pushing incumbent England No 13 Jonathan Joseph hard for his place.
focus-vec:	The uncapped Wasps centre scored two tries in March. Daly's form at centre has him pushing incumbent England No 13 Jonathan Joseph hard for his place. Wasps director of rugby Dai Young said: 'If England are going to go out and win the Rugby World Cup, <b>he's the type of player you've got to have in the team.</b> '

Table 6: Generation samples of different attention control approaches on CNN/Dailymail dataset. The highlighted sentences are marked in bold. The span marked by red in the reference summary is captured by the generation from focus-vec, but not by att-offset. The sentence number marked in the input corresponds to the sentence index in Figure 3.

article:	<b>Kabul, Afghanistan (CNN)A suicide bomber detonated his explosives near a group of protesters in eastern Afghanistan on Thursday, killing 17 people and wounding dozens more, police said.</b> <sup>[1]</sup> <b>"An Afghan lawmaker taking part in the protests in the city of Khost was among the 64 people wounded,"</b> <sup>[2]</sup> <b>said Faizullah Ghairat, the provincial police chief Taliban spokesman Zabiullah Mujahid denied his group was responsible for the attack.</b> <sup>[3]</sup> No other organization has so far claimed responsibility. <sup>[4]</sup> Humayoon Humayoon, an Afghan member of parliament for Khost province, and the other protesters were on their way to join a larger rally against the provincial governor, according to Zahir Jan, an eyewitness. <sup>[5]</sup> The suicide attack hit the group around 10 a.m. local time, police said. <sup>[6]</sup> . . . (LEAVE OUT UNIMPORTANT SENTENCES)
reference:	<b>An Afghan lawmaker is among 64 people wounded in the attack, police say.</b> Taliban spokesman denies his group was responsible for the attack.
vanilla:	A suicide bomber detonates his explosives near a group of protesters, killing 17 people and wounding dozens more. The Taliban spokesman denies his group was responsible for the attack. No other organization has so far claimed responsibility for the attack. The suicide attack hit the group around 10 a.m. local time.
keyword:	Taliban spokesman Zabiullah Mujahid denies his group was responsible. No other organization has so far claimed responsibility. The suicide attack hit the group around 10 a.m. local time.
att-offset:	A suicide bomber detonates his explosives near a group of protesters, killing 17 people and wounding dozens more. The Taliban spokesman denies his group was responsible for the attack. No other organization has so far claimed responsibility for the attack. The suicide attack hit the group around 10 a.m. local time.
focus-vec:	A suicide bomber detonates his explosives near a group of protesters, killing 17 people. <b>An Afghan lawmaker is among the 64 people wounded, police say.</b> Taliban spokesman Zabiullah Mujahid denies his group was responsible for the attack. No other organization has so far claimed responsibility.

Table 7: Generation samples of different attention control approaches on CNN/Dailymail dataset. The span marked by red in the reference summary is captured by the generation from focus-vec, but not by att-offset.