

AdapLeR: Speeding up Inference by Adaptive Length Reduction

Anonymous ACL submission

Abstract

Pre-trained language models have shown stellar performance in various downstream tasks. But, this usually comes at the cost of high latency and computation, hindering their usage in resource-limited settings. In this work, we propose a novel approach for reducing the computational cost of BERT with minimal loss in downstream performance. Our model dynamically eliminates less contributing tokens through layers, resulting in shorter lengths and consequently lower computational cost. To determine the importance of each token representation, we train a Contribution Predictor for each layer using a gradient-based saliency method. Our experiments on several diverse classification tasks show speedups up to 17x during inference time. We also validate the quality of the selected tokens in our method using human annotations in the ERASER benchmark. In comparison to other widely used strategies for selecting important tokens, such as *saliency* and *attention*, our proposed method has significantly less false positive rate in generating rationales.

1 Introduction

While large-scale pre-trained language models exhibit remarkable performances on various NLP benchmarks, their excessive computational costs and high inference latency have limited their usage in low-resource settings. In this regard, there have been various attempts at improving the efficiency of BERT-based models (Devlin et al., 2019), including knowledge distillation (Hinton et al., 2015; Sanh et al., 2019; Sun et al., 2019, 2020; Jiao et al., 2020), quantization (Gong et al., 2014; Shen et al., 2020; Tambe et al., 2021), weight pruning (Han et al., 2016; He et al., 2017; Michel et al., 2019; Sanh et al., 2020), and progressive module replacing (Xu et al., 2020). Despite providing significant reduction in model size, these techniques are generally static at inference time, i.e., they dedicate the

same amount of computation to all inputs, irrespective of their difficulty.

A number of techniques have been also proposed in order to make efficiency enhancement sensitive to inputs. *Early exit* mechanism (Schwartz et al., 2020; Liao et al., 2021; Xin et al., 2020; Liu et al., 2020; Xin et al., 2021; Sun et al., 2021; Eyzaquirre et al., 2021) is a commonly used method in which each layer in the model is coupled with an intermediate classifier to predict the target label. At inference, a halting condition is used to determine whether the model allows an example to exit without passing through all layers. Various halting conditions have been proposed, including Shannon’s entropy (Xin et al., 2020; Liu et al., 2020), softmax outputs with temperature calibration (Schwartz et al., 2020), trained confidence predictors (Xin et al., 2021), or the number of agreements between predictions of intermediate classifiers (Zhou et al., 2020).

Most of these techniques compress the model from the depth perspective (i.e., reducing the number of involved encoder layers). However, one can view compression from the width perspective (Goyal et al., 2020; Ye et al., 2021), i.e., reducing the length of hidden states. (Ethayarajh, 2019; Klafka and Ettinger, 2020). This is particularly promising as recent analytical studies showed that there are redundant encoded information in token representations (Klafka and Ettinger, 2020; Ethayarajh, 2019). Among these redundancies, some tokens carry more task-specific information than others (Mohebbi et al., 2021), suggesting that only these tokens could be considered through the model. Moreover, in contrast to layer-wise pruning, token-level pruning does not come at the cost of reducing model’s capacity in complex reasoning (Sanh et al., 2019; Sun et al., 2019).

POWER-BERT (Goyal et al., 2020) is one of the first such techniques which reduces inference time by eliminating redundant token representa-

tions through layers based on self-attention weights. Several studies have followed (Kim and Cho, 2021; Wang et al., 2021); However, they usually optimize a single token elimination configuration across the entire dataset, resulting in a static model. In addition, their token selection strategies are based on attention weights which can result in a sub-optimal solution (Ye et al., 2021). In this work, we introduce **Adaptive Length Reduction (AdapLeR)**. Instead of relying on attention weights, our model trains a set of Contribution Predictors (CP) to estimate tokens’ saliency scores at inference. We show that this choice results in more reliable scores than attention weights in measuring tokens’ contributions.

The most related study to ours is TR-BERT (Ye et al., 2021) which leverages reinforcement learning to develop an input-adaptive token selection policy network. However, as pointed out by the authors, the problem has a large search space, making it difficult for RL to solve. To mitigate this, they resorted to extra heuristics such as imitation learning (Hussein et al., 2017) for warming up the training of the policy network, action sampling for limiting the search space, and knowledge distillation for transferring knowledge from the intact backbone fine-tuned model. All of these steps significantly increase the training cost. Hence, they only perform token selection at two layers. In contrast, we propose a simple but effective method to gradually eliminate tokens in each layer throughout the training phase using a soft-removal function which allows the model to be adaptable to various inputs in a batch-wise mode. It is also worth noting above studies are based on top-k operations for identifying the k most important tokens during training or inference, which can be expensive without a specific hardware architecture (Wang et al., 2021).

In summary, our contributions are threefold:

- We couple a simple Contribution Predictor (CP) with each layer of the model to estimate tokens’ contribution scores to eliminate redundant representations.
- Instead of an instant token removal, we gradually mask out less contributing token representations by employing a novel soft-removal function.
- We also show the superiority of our token selection strategy over the other widely used strategies by using human rationales.

2 Background

2.1 Self-attention Weights

Self-attention is a core component of the Transformers (Vaswani et al., 2017) which looks for the relation between different positions of a single sequence of token representations (x_1, \dots, x_n) to build contextualized representations. To this end, each input vector x_i is multiplied by the corresponding trainable matrices Q , K , and V to respectively produce query (q_i) , key (k_i) , and value (v_i) vectors. To construct the output representation z_i , a series of weights is computed by the dot product of q_i with every k_j in all time steps. Before applying a softmax function, these values are divided by a scaling factor and then added to an *attention mask* vector \mathbf{m} , which is zero for positions we wish to attend and $-\infty$ (in practice, -10000) for padded tokens (Vaswani et al., 2017). Mathematically, for a single attention head, the weight attention from token x_i to token x_j in the same input sequence can be written as:

$$\alpha_{i,j} = \text{softmax}_{x_j \in \mathcal{X}} \left(\frac{q_i k_j^\top}{\sqrt{d}} + m_i \right) \in \mathbb{R} \quad (1)$$

The time complexity for this is $O(n^2)$ given the dot product $q_i k_j^\top$, where n is the input sequence length. This impedes the usage of self-attention based models in low-resource settings.

While self-attention is one of the most white-box components in transformer-based models, relying on raw attention weights as an explanation could be misleading given that they are not necessarily responsible for determining the contribution of each token in the final classifier’s decision (Jain and Wallace, 2019; Serrano and Smith, 2019; Abnar and Zuidema, 2020). This is based on the fact that raw attentions are being faithful to the local mixture of information in each layer and are unable to obtain a global perspective of the information flow through the entire model (Pascual et al., 2021).

2.2 Gradient-based Saliency Scores

Gradient-based methods provide alternatives to attention weights to compute the importance of a specific input feature. Despite having been widely utilized in other fields earlier (Ancona et al., 2017; Simonyan et al., 2013; Sundararajan et al., 2017; Smilkov et al., 2017), they have only recently become popular in NLP studies (Bastings and Filippova, 2020; Li et al., 2016; Yuan et al., 2019).

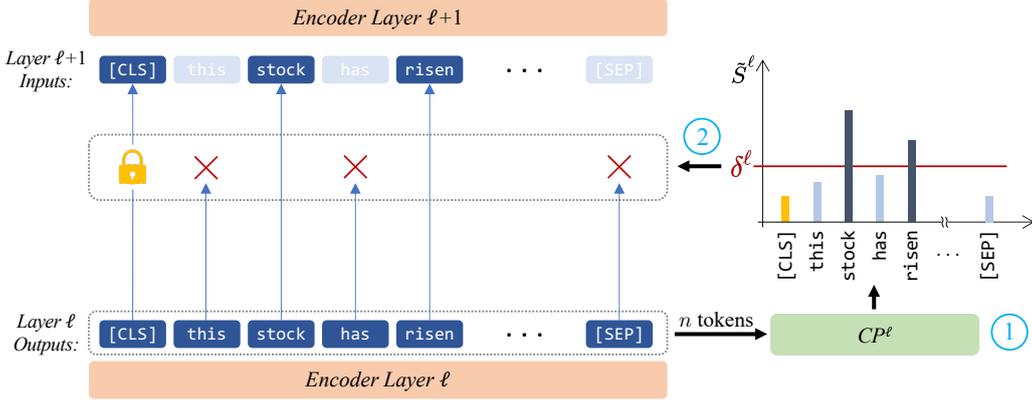


Figure 1: To reduce the inference computation, in each layer (1) the attribution score of the token representation is estimated and (2) based on a reduced uniform-level threshold ($\delta^\ell = \eta^\ell/n$) token representations with low importance score are removed. Since the final layer’s classifier is connected to the [CLS] token and it could act as a pooler within each layer it is the only token that would remain regardless of its score.

180 These methods are based on computing the first-
 181 order derivative of the output logit y_c w.r.t. the
 182 input embedding h_i^0 (initial hidden states), where
 183 c could be true class label to find the most impor-
 184 tant input features or the predicted class to interpret
 185 model’s behavior. After taking the norm of output
 186 derivatives, we get *sensitivity* (Ancona et al., 2017),
 187 which indicates the changes in model’s output with
 188 respect to the changes in specific input dimensions.
 189 Instead, by multiplying gradients with input fea-
 190 tures, we arrive at *gradient \times input* (Bastings and
 191 Filippova, 2020), also known as *saliency*, which
 192 also considers the direction of input vectors to de-
 193 termine the most important tokens. Since these
 194 scores are computed for each dimension of embed-
 195 ding vectors, an aggregation method such as L2
 196 norm or mean is needed to produce one score per
 197 input token (Atanasova et al., 2020a):

$$S_i = \left\| \frac{\partial y_c}{\partial h_i^0} \odot h_i^0 \right\|_2 \quad (2)$$

199 3 Methodology

200 As shown in Figure 1, our approach relies on drop-
 201 ping low contributing tokens in each layer and
 202 passing only the more important ones to the next.
 203 Therefore, one important step is to measure the
 204 importance of each token. To this end, we opted
 205 for saliency scores which is a more reliable crite-
 206 rion in measuring token’s contributions (Bastings
 207 and Filippova, 2020; Pascual et al., 2021). We will
 208 show in Section 5.1 results of a series quantitative
 209 analyses that supports this choice. In what follows,
 210 we first describe how we estimate saliency scores

211 at inference time using a set of Contribution Pre-
 212 dictors (CPs) and then we elaborate on how we
 213 leverage these predictors during inference (Section
 214 3.2) and training (Section 3.3) phase.

215 3.1 Contribution Predictor

216 Computing gradients during inference is problem-
 217 atic as back-propagation computation prolongs in-
 218 ference time, which is contrary to our main goal.
 219 To circumvent this, we simply add a CP after each
 220 layer ℓ in the model to estimate contribution score
 221 for each token representation, i.e., \tilde{S}_i^ℓ . The model
 222 then decides on the tokens that should be passed to
 223 the next layer based on the values of \tilde{S}_i^ℓ . CP com-
 224 putes \tilde{S}_i^ℓ for each token using an MLP followed
 225 by a softmax activation function. We argue that,
 226 despite being limited in learning capacity, the MLP
 227 is sufficient for estimating scores that are more gen-
 228 eralized and relevant than vanilla saliency values.
 229 We will present a quantitative analysis on this topic
 230 in Section 5.

231 3.2 Model Inference

232 Most BERT-based models consist of L encoder
 233 layers. The input sequence of n tokens is usually
 234 passed through an embedding layer to build the
 235 initial hidden states of the model h^0 . Each encoder
 236 layer then produces the next hidden states using the
 237 ones from the previous layer:

$$h^\ell = \text{Encoder}_\ell(h^{\ell-1}) \quad (3)$$

238 In our approach, we eliminate less contribut-
 239 ing token representations before delivering hidden
 240 states to the next encoder. Tokens are selected
 241

based on the contribution scores \tilde{S}_i^ℓ obtained from the CP of the corresponding layer ℓ . As the sum of these scores is equal to one, a uniform level indicates that all tokens contribute equally to the prediction and should be retained. On the other hand, the lower-scoring tokens could be viewed as unnecessary tokens if the contribution scores are concentrated only on a subset of tokens. Given that the final classification head uses the last hidden state of the [CLS] token, we preserve this token’s representation in all layers. Despite preserving this, other tokens might be removed from a layer when [CLS] has a significantly high estimated contribution score than others. Based on this intuition, we define a cutoff threshold based on the uniform as: $\delta^\ell = \eta^\ell \cdot 1/n$ with $0 < \eta^\ell \leq 1$ to distinguish important tokens. Tokens are considered important if their contribution score exceeds δ (which is a equal or smaller value than the uniform score). Intuitively, a larger η provides a higher δ cutoff level, thereby dropping a larger number of tokens, hence, yielding more speedup. The value of η determines the extent to which we can rely on CP’s estimations. In case the estimations of CP are deemed to be inaccurate, its impact can be reduced by lowering η . We train each layer’s η^ℓ using an auxiliary training objective, which allows the model to adjust the cutoff value to control the speedup-performance tradeoff. Also, since each input instance has a different computational path during token removal process, it is obvious that at inference time the batch size should be equal to one (single instance usage), similarly to other dynamic approaches (Zhou et al., 2020; Liu et al., 2020; Ye et al., 2021; Eyzaguirre et al., 2021; Xin et al., 2020).

3.3 Model Training

Training consists of three phases: initial finetuning, saliency extraction, and adaptive length retraining. In the first phase, we simply finetune the backbone model (BERT) on a given target task. We then extract the saliencies of three top-performing checkpoints from the finetuning process and compute the average of them to mitigate potential inconsistencies in saliency scores (cf. Section 2.2). The final step is to train a pre-trained model using an adaptive length reduction procedure. In this phase, a non-linear function gradually fades out the representations throughout the training process. Each CP is jointly trained with the rest of the model using the saliencies extracted in the pre-

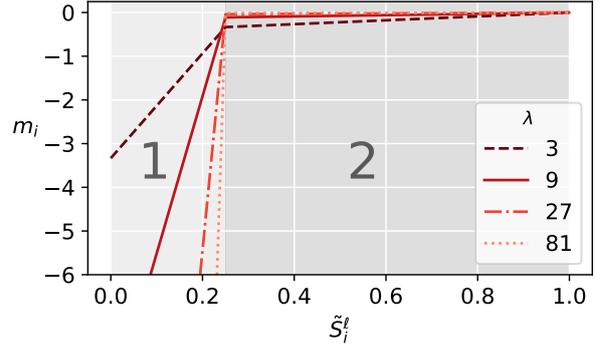


Figure 2: The soft-removal function plotted with $\lambda \in \{3, 9, 27, 81\}$ and $\delta^\ell = 0.25$. As λ increases, the removal region (1) gets steeper while the other zone (2), which is almost horizontal, approaches the zero level.

vious phase alongside with the target task labels. We also define a speedup tuning objective to determine the thresholds (via tuning η) to control the performance-speedup trade-off. In the following, we elaborate on the procedure.

Soft-removal function. During training, if tokens are immediately dropped similarly to the inference mode, the effect of dropping tokens cannot be captured using a gradient back-propagation procedure. Using batch-wise training in this scenario will also be problematic as the structure will vary with each example. Hence, inspired by the padding mechanism of self-attention models (Vaswani et al., 2017) we introduce a new method that gradually masks out less contributing token representations. In each layer, after predicting contribution scores, instead of instantly removing the token representations, we accumulate a negative mask to the attention mask vector M using a soft-removal function:

$$m_i^-(\tilde{S}_i^\ell) = \begin{cases} \lambda_{adj}(\tilde{S}_i^\ell - \delta^\ell) - \frac{\beta}{\lambda} & \tilde{S}_i^\ell < \delta^\ell \\ \frac{(\tilde{S}_i^\ell - 1)\beta}{(1 - \delta^\ell)\lambda} & \tilde{S}_i^\ell \geq \delta^\ell \end{cases} \quad (4)$$

This function consists of two main zones (Figure 2). In the first term, the less important tokens with scores lower than the threshold (δ^ℓ) are assigned higher negative masking as they get more distant from δ . The slope is determined by $\lambda_{adj} = \lambda/\delta$, where λ is a hyperparameter that is increased exponentially after each epoch (e.g., $\lambda \leftarrow 10 \times \lambda$ after finishing each epoch). Increasing λ makes the soft-removal function stronger and more decisive in masking the representations. To avoid under-going zero gradients during training, we define

324 $0 < \beta < 0.1$ to construct a small negative slope
 325 (similar to the well known Leaky-ReLU of [Maas](#)
 326 [et al. 2013](#)) for those tokens with higher contribut-
 327 ing scores than δ^ℓ threshold. Consider a scenario in
 328 which η^ℓ sharply drops, causing most of \hat{S}_i^ℓ get over
 329 the δ^ℓ threshold. In this case, the non-zero value
 330 in the second term of Equation 4, which facilitates
 331 optimizing η^ℓ .

332 **Training the Contribution Predictors.** The CPs
 333 are trained by an additional term which is based
 334 on the KL-divergence¹ of each layer’s CP output
 335 with the extracted saliencies. The main training
 336 objective is a minimization of the following loss:

$$337 \quad \mathcal{L} = \mathcal{L}_{\text{CE}} + \gamma \mathcal{L}_{\text{CP}} \quad (5)$$

338 Where γ is a hyperparameter which that specifies
 339 the amount of emphasis on the CP training loss:

$$340 \quad \begin{aligned} \mathcal{L}_{\text{CP}} &= \sum_{\ell=0}^{L-1} (L - \ell) D_{\text{KL}}(\hat{S}^\ell || \tilde{S}^\ell) \\ &= \sum_{\ell=0}^{L-1} (L - \ell) \sum_{i=1}^N \hat{S}_i^\ell \log\left(\frac{\hat{S}_i^\ell}{\tilde{S}_i^\ell}\right) \end{aligned} \quad (6)$$

341 Since S is based on the input embeddings, the
 342 [CLS] token usually shows a low amount of con-
 343 tribution due to not having any contextualism in
 344 the input. As we leverage the representation of
 345 the [CLS] token in the last layer for classification,
 346 this token acts as a pooler and gathers information
 347 about the context of the input. In other words, the
 348 token can potentially have more contribution as it
 349 passes through the model. To this end, we amplify
 350 the contribution score of [CLS] and renormalize
 351 the distribution (\hat{S}^ℓ) with a trainable parameter θ^ℓ :

$$352 \quad \hat{S}_i^\ell = \frac{\theta^\ell S_1^\ell \mathbf{1}[i = 1] + S_i^\ell \mathbf{1}[i > 1]}{\theta^\ell S_1^\ell + \sum_{i=2}^n S_i^\ell} \quad (7)$$

353 By this procedure, the next objective (discussed
 354 in the next paragraph) will have the capability of
 355 tuning the amount of pooling, consequently con-
 356 trolling the amount of speedup. Larger θ push the
 357 CPs to shift the contribution towards the [CLS] to-
 358 ken to gather most of the task-specific information
 359 and avoids carrying redundant tokens through the
 360 model.

¹Inclusive KL loss. Check Appendix A.

Speedup Tuning. In the speedup tuning process,
 we combine the cross-entropy loss of the target
 classification task with a length loss which is the
 expected number of unmasked token representa-
 tions in all layers. Considering that we have a
 non-positive and continuous attention mask M ,
 the length loss of a single layer would be the
 summation over the exponential of the mask val-
 ues $\exp(m_i)$ to map the masking range $[-\infty, 0]$
 to a $[0$ (fully masked/removed), 1 (fully retained)]
 bound.

$$372 \quad \begin{aligned} \mathcal{L}_{\text{SPD./PERF.}} &= \mathcal{L}_{\text{CE}} + \phi \mathcal{L}_{\text{LENGTH}} \\ \mathcal{L}_{\text{LENGTH}} &= \sum_{\ell=1}^L \sum_{i=1}^n \exp(m_i^\ell) \end{aligned} \quad (8)$$

373 In Equation 8, demonstrates how the length loss is
 374 computed inside the model and how its added to
 375 the main classification loss. During training, we
 376 assign a separate optimization process which tunes
 377 η and θ to adjust the thresholds and the amount of
 378 [CLS] pooling² alongside with the CP training.

379 The reason that this objective is treated as a sepa-
 380 rate problem instead of merging it with the previ-
 381 ous one, is because in the latter case the CPs
 382 could be influenced by the length loss and try to
 383 manipulate the contribution scores for some tokens
 384 regardless of their real influence. So in other words,
 385 the first objective is to solve the task and make it
 386 explainable with the CPs, and the secondary objec-
 387 tive builds the speedup using tuning the threshold
 388 levels and the amount of pooling in each layer.

389 4 Experiments

390 4.1 Datasets

391 To verify the effectiveness of our proposed method
 392 on adaptive length reduction, we selected eight
 393 various text classification datasets. In order to
 394 incorporate a variety of tasks, we utilized SST-
 395 2 ([Socher et al., 2013](#)) and IMDB ([Maas et al.,](#)
 396 [2011](#)) for sentiment, MRPC ([Dolan and Brock-](#)
 397 [ett, 2005](#)) for paraphrase, AG’s News ([Zhang](#)
 398 [et al., 2015](#)) for topic classification, DBpedia
 399 ([Lehmann et al., 2015](#)) for knowledge extraction,
 400 MNLI ([Williams et al., 2018](#)) for NLI, QNLI ([Ra-](#)
 401 [jpurkar et al., 2016](#)) for question answering, and
 402 HateXplain ([Mathew et al., 2021](#)) for hate speech.
 403 Evaluations are based on the test split of each
 404 dataset. For those datasets that are in the GLUE

²Since θ is not in the computational DAG, we employed a dummy variable inside the model. See Appendix B.

Model	SST-2		IMDB		HateXplain		MRPC		MNLI		QNLI		AG’s news		DBpedia	
	Acc.	FLOPs	Acc.	FLOPs	Acc.	FLOPs	F1.	FLOPs	Acc.	FLOPs	Acc.	FLOPs	Acc.	FLOPs	Acc.	FLOPs
BERT	92.7	1.00x	93.8	1.00x	68.3	1.00x	87.5	1.00x	84.2	1.00x	90.3	1.00x	94.4	1.00x	99.3	1.00x
DistilBERT	92.2	2.00x	92.9	2.00x	68.2	2.00x	88.0	2.00x	81.8	2.00x	88.1	2.00x	94.2	2.00x	99.3	2.00x
PoWER-BERT	92.1	1.18x	92.2	1.70x	66.9	2.69x	88.0	1.07x	82.9	1.10x	89.7	1.23x	92.1	12.5x	98.1	14.8x
TR-BERT	93.4	1.09x	93.2	2.90x	67.9	2.23x	81.9	1.16x	84.8	1.00x	89.0	1.09x	93.2	10.2x	98.9	10.01x
AdapLeR	92.3	1.49x	91.7	3.21x	68.6	4.73x	87.6	1.27x	82.9	1.42x	89.3	1.47x	92.5	17.1x	98.9	22.23x

Table 1: Comparison of our method (AdapLeR) with other baselines in eight classification tasks in terms of performance and speedup (FLOPs). For each dataset the corresponding metric has been reported (Accuracy: Acc., F1: F-1 Score). In the MNLI task, the speedup and performance values are the average of the evaluations on the matched and mismatched test sets.

Benchmark (Wang et al., 2018), test results were acquired by submitting the test predictions to the evaluation server. For other tasks results were computed based on the test set provided.

4.2 Experimental Setup

To compare our approach, we set our first baseline to be the pre-trained BERT (base-uncased) (Devlin et al., 2019) which is also the backbone model of our model and the other three baselines: DistilBERT (uncased) (Sanh et al., 2019) as a static model, TR-BERT and PoWER-BERT as dynamic approaches. We used the same implementations and suggested hyperparameters³ to train these baselines. To fine-tune the backbone model we used similar hyperparameters over all tasks that are provided in Section D. The backbone model and our model implementation is based on the HuggingFace’s Transformers library (Wolf et al., 2020). Trainings and evaluations were conducted on a dual 2080Ti 11GB GPU machine with multiple runs.

Hyperparameter Selection. Overall, we introduced four hyperparameters ($\gamma, \phi, \lambda, \beta$)⁴ which are involved in the training process. However, the main two primary terms that are the most influential and have considerable effects on both the output performance and the speedup of the trained model are ϕ and γ . This makes our approach comparable to existing techniques (Goyal et al., 2020; Ye et al., 2021) which usually have two or three hyperparameters adjusted per task. While using grid search for these two terms, we kept other hyperparameters constant over all datasets. The selected hyperparameters and more details are discussed in Section D.

³Since some of the datasets were not used originally, we had to search the hyperparameters based on the given ranges.

⁴Note that θ and η are trainable terms that are tuned by the model during training.

FLOPs Computation. As we wish to determine the computational complexity of models independently of the operating environment (e.g., CPU/GPU), following Ye et al. (2021) and Liu et al. (2020), we computed FLOPs, i.e., the number of floating-point operations (FLOPs) in a single inference procedure. To have a fair comparison, we computed FLOPs for PoWER-BERT in a single instance mode, described in Section C.

4.3 Results

The performance and speedup values of our proposed method and other baselines are presented in Table 1. We can observe that with a low performance gap in all tasks, our approach significantly outperforms others in terms of efficiency. It is noteworthy that the results also reveal some form of dependency on the type of tasks. Some tasks may need less amount of contextualism during inference and could be classified by using a fraction of input tokens. For instance, in AG’s News, the topic of a sentence might be identified with a single token (e.g. Basketball \rightarrow Topic: Sports, see Figure 5 in the Appendix as an example).

We illustrate speed-accuracy curves for HateXplain in Figure 6 in the Appendix to provide a closer look at the efficiency of AdapLeR in comparison with other state-of-the-art methods for length reduction. For each curve, the points were obtained by tuning the most influential hyperparameters of the corresponding model.

5 Analysis

In this section, we first conduct an experiment to support our choice of saliency scores as a supervision in measuring the importance of token representations. Next, we validate the behavior of Contribution Predictors in identifying most important tokens in an AdapLeR model.

Strategy	Movie Reviews		MultiRC	
	Acc.	FLOPs	Acc.	FLOPs
Full input	93.3	1x	67.7	1x
Human rationale	96.7	3.7x	76.6	4.6x
Saliency	92.3	3.7x	66.4	4.4x
Attention ALL	78.5	3.7x	62.9	4.4x
Attention [CLS]	70.3	3.7x	63.7	4.4x

Table 2: Accuracy and speedup when the most important input tokens during fine-tuning are computed based on attention and saliency strategies and human rationale (the upper bound). The bold values indicate the best corresponding strategy for each task (the closest performance to the upper bound).

5.1 Saliency vs. Attention

In dealing with token pruning, a natural question that might arise is what would be the most appropriate criterion for assessing the relative importance of tokens within a sentence? To arrive at an empirical and reliable upper bound in measuring token importance, we resort to human rationale. To this end, we used the ERASER benchmark (DeYoung et al., 2020), which contains multiple tasks for which important spans of the input text have been highlighted as supporting evidence (aka “rationale”) by human. Among the benchmark tasks, we opted for two diverse classification tasks: Movie reviews (Zaidan and Eisner, 2008) and MultiRC (Khashabi et al., 2018) (see Sec. E in the Appendix for task descriptions).

In order to verify the reliability of human rationales, we fine-tuned BERT just on rationales by excluding those tokens that are not highlighted in the input. In Table 2, the first two rows show the performance score of BERT on target tasks with full tokens and only rationales in the input. We see that fine-tuning merely on rationales not only yielded less computation cost, but also resulted in higher performance when compared with the full input setting. Obviously, human annotations are not available for a whole range of downstream NLP tasks; therefore, this criteria is infeasible in practice and can only be viewed as an upper bound for evaluating different strategies in measuring token importance. We investigated the effectiveness of saliency and self-attention weights as two commonly used strategies for measuring the importance of tokens in pre-trained language models.

To compute these, we first fine-tuned BERT with all tokens in the input for a given target task. We

then obtained saliency scores with respect to the tokens in the input embedding layer. This gives us two advantages. First, representations in this layer are non-contextualized, allowing us to measure the importance of each token individually. Second, the fact that the gradient passes from the end to the beginning of the model results in aggregated values for the relative importance of each token based on the entire model. Similarly, we aggregated self-attention weights across all layers of the model using a post-processed variant of attentions called *attention rollout* (Abnar and Zuidema, 2020), a popular technique in which each attention weight matrix in each layer is multiplied by the ones before it to form aggregated attention values. To assign an importance score to each token, we examined two different interpretation of attention weights. The first strategy is the one adopted by PoWER-BERT (Goyal et al., 2020) in which for each token we accumulate attention values from other tokens. Additionally, we measured how much the [CLS] token attends to each token in the input, a strategy which has been widely used in interpretability studies around BERT (Abnar and Zuidema, 2020; Chrysostomou and Aletras, 2021; Jain et al., 2020, *inter alia*). For a fair evaluation, for each sentence in the test set, we selected the top- k salient and attended words, with k being the number of words that are annotated as rationales.

Results in Table 2 show that fine-tuning on the most salient tokens outperforms that based on the most attended tokens. This denotes that saliency is a better indicator for the importance of tokens. Nonetheless, recent length reduction techniques (Goyal et al., 2020; Kim and Cho, 2021; Wang et al., 2021) have mostly adopted attention weights as their criterion for selecting important tokens as these weights are convenient to compute during the inference.

5.2 Contribution Predictor Evaluation

The goal of this section is to validate our Contribution Predictors in selecting the most contributed tokens. Figure 3 shows an input example from SST-2 dataset. As we can see, the CPs can identify and drop the irrelevant tokens gradually through layers, finally focusing mostly on ‘pedestrian’ (Adj.) and [CLS] token representations which is highly aligned with human interpretation.

Next, we attempted to quantify how much our model can preserve rationales without requiring

Layer 0: [CLS] what was once original has been co - opted so frequently that it now seems pedestrian . [SEP]
 Layer 5: [CLS] what was once original has been co - opted so frequently that it now seems pedestrian . [SEP]
 Layer 11: [CLS] what was once original has been co - opted so frequently that it now seems pedestrian . [SEP]

Figure 3: The illustration of contribution scores obtained by CPs in three different layers of the model for an input example from SST-2 (sentiment) task. The color intensity indicates the degree of contribution scores. Only the highlighted token representations are processed in each layer. See more full-layer plots in the appendix 5.

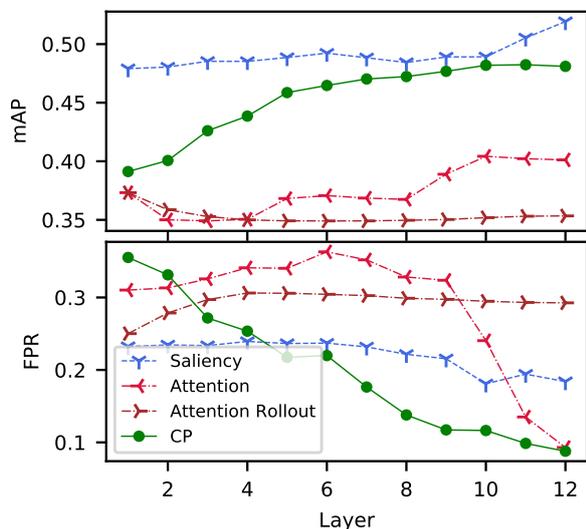


Figure 4: Agreement with human rationales in terms of mean Average Precision and False Positive Rate for CP and three alternative techniques.

direct human annotations. For evaluation, we used two Average Precision (AP) and False Positive Rate (FPR) metrics by comparing the remaining tokens to the human rationale annotations. The first metric measures whether the model assigns higher continuous scores to those tokens that are annotated by humans as rationales. Whereas, the intuition behind the second metric is how many irrelevant tokens are selected by the model to be passed to subsequent layers.

First, we fine-tuned the model on the Movie Review dataset and computed layer-wise raw attention, attention rollout, and saliency scores for each token representation. We also trained a model using our proposed approach and computed the output probability scores of CPs in each layer. Since human rationales are annotated at the word level, we sum the scores across tokens corresponding to each word to arrive at word-level importance scores. In addition to these soft scores, we used the uniform-level threshold to reach a binary score indicating tokens selected in each layer. We used soft scores for computing AP and binary scores for computing FPR.

Figure 4 shows the agreement between human rationales and the selected tokens based on these two metrics. As we can see, in comparison to other widely used strategies for selecting important tokens, such as saliency and attention, our Contribution Predictors have significantly less false positive rate in preserving rationales through the layers. Though attention and CP converge at the same point, note that, CPs can also identify rationales at earlier layers, allowing the model to combine the most relevant token representations to build the final representation and gain better performance results, as we have seen in the main results. There is also a line of research in which practitioners attempt to guide models to perform human-like reasoning by training rationale generation simultaneously with the target task that requires human annotation (Atanasova et al., 2020b; Zhao et al., 2020; Li et al., 2018). As a by-product, our trained CPs are able to generate these rationales at inference without the need for human-generated annotations.

6 Conclusion

In this paper we introduced AdapLeR, a model that dynamically identifies and drops less contributing token representations through layers. Specifically, AdapLeR accomplishes this by training a set of Contribution Predictors based on saliencies extracted from a finetuned model and applying a gradual masking technique to simulate input-adaptive token removal during training. Empirical results on seven diverse text classification tasks show considerable improvements over previous methods. Furthermore, we demonstrated that contribution predictors generate rationales that are highly in line with those manually specified by humans. As future work, we aim to apply our technique to more tasks and see whether it can be adapted to those tasks that rely on all token representations (e.g., question answering). Additionally, combining our width-based strategy with a depth-based one (e.g., early exiting) might potentially yield greater efficiency, something we plan to pursue as future work.

References

Samira Abnar and Willem Zuidema. 2020. [Quantifying attention flow in transformers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, Online. Association for Computational Linguistics.

Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. 2017. Towards better understanding of gradient-based attribution methods for deep neural networks. *arXiv preprint arXiv:1711.06104*.

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020a. [A diagnostic study of explainability techniques for text classification](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3256–3274, Online. Association for Computational Linguistics.

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020b. [Generating fact checking explanations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7352–7364, Online. 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.

George Chrysostomou and Nikolaos Aletras. 2021. Enjoy the salience: Towards better transformer-based faithful explanations with word salience. *ArXiv*, abs/2108.13759.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C. Wallace. 2020. [ERASER: A benchmark to evaluate rationalized NLP models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4443–4458, Online. Association for Computational Linguistics.

William B. Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.

Kawin Ethayarajh. 2019. [How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65, Hong Kong, China. Association for Computational Linguistics.

Cristóbal Eyzaguirre, Felipe del Río, Vladimir Araujo, and Álvaro Soto. 2021. [Dact-bert: Differentiable adaptive computation time for an efficient bert inference](#). *arXiv preprint arXiv:2109.11745*.

Yunchao Gong, L. Liu, Ming Yang, and Lubomir D. Bourdev. 2014. Compressing deep convolutional networks using vector quantization. *ArXiv*, abs/1412.6115.

Saurabh Goyal, Anamitra Roy Choudhury, Saurabh Raje, Venkatesan Chakaravarthy, Yogish Sabharwal, and Ashish Verma. 2020. [Power-bert: Accelerating bert inference via progressive word-vector elimination](#). In *International Conference on Machine Learning*, pages 3690–3699. PMLR.

Song Han, Huizi Mao, and William J. Dally. 2016. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. *arXiv: Computer Vision and Pattern Recognition*.

Yihui He, Xiangyu Zhang, and Jian Sun. 2017. [Channel pruning for accelerating very deep neural networks](#). *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 1398–1406.

Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. [Distilling the knowledge in a neural network](#). *ArXiv*, abs/1503.02531.

Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. 2017. [Imitation learning: A survey of learning methods](#). *ACM Computing Surveys (CSUR)*, 50(2):1–35.

Sarthak Jain and Byron C. Wallace. 2019. [Attention is not Explanation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3543–3556, Minneapolis, Minnesota. Association for Computational Linguistics.

Sarthak Jain, Sarah Wiegrefe, Yuval Pinter, and Byron C. Wallace. 2020. [Learning to faithfully rationalize by construction](#). In *ACL*.

Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. [TinyBERT: Distilling BERT for natural language understanding](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174, Online. Association for Computational Linguistics.

683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737

852		<i>Linguistics</i> , pages 6640–6651, Online. Association for Computational Linguistics.	
853			
854	Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2931–2951, Florence, Italy. Association for Computational Linguistics.		
855			
856			
857			
858			
859	Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. 2020. Q-bert: Hessian based ultra low precision quantization of bert. In <i>AAAI</i> .		
860			
861			
862			
863	Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. <i>arXiv preprint arXiv:1312.6034</i> .		
864			
865			
866			
867	D Smilkov, N Thorat, B Kim, F Viégas, and M Wattenberg. 2017. Smoothgrad: removing noise by adding noise. <i>arxiv. arXiv preprint arxiv:1706.03825</i> .		
868			
869			
870	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank . In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.		
871			
872			
873			
874			
875			
876			
877			
878	Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. 2019. Patient knowledge distillation for BERT model compression . 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 4323–4332, Hong Kong, China. Association for Computational Linguistics.		
879			
880			
881			
882			
883			
884			
885			
886	Tianxiang Sun, Yunhua Zhou, Xiangyang Liu, Xinyu Zhang, Hao Jiang, Zhao Cao, Xuanjing Huang, and Xipeng Qiu. 2021. Early exiting with ensemble internal classifiers. <i>arXiv preprint arXiv:2105.13792</i> .		
887			
888			
889			
890	Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. 2020. MobileBERT: a compact task-agnostic BERT for resource-limited devices . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2158–2170, Online. Association for Computational Linguistics.		
891			
892			
893			
894			
895			
896			
897	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In <i>Proceedings of the 34th International Conference on Machine Learning-Volume 70</i> , pages 3319–3328.		
898			
899			
900			
901	Thierry Tambe, Coleman Hooper, Lillian Pentecost, Tianyu Jia, En-Yu Yang, Marco Donato, Victor Sanh, Paul N. Whatmough, Alexander M. Rush, David Brooks, and Gu-Yeon Wei. 2021. Edgebert: Sentence-level energy optimizations for latency-aware multi-task nlp inference . <i>MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture</i> .		
902			
903			
904			
905			
906			
907			
908			
	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>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.		909
			910
			911
			912
			913
	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.		914
			915
			916
			917
			918
			919
			920
			921
	Hanrui Wang, Zhekai Zhang, and Song Han. 2021. Spaten: Efficient sparse attention architecture with cascade token and head pruning . <i>2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)</i> , pages 97–110.		922
			923
			924
			925
			926
	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference . 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 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.		927
			928
			929
			930
			931
			932
			933
			934
			935
	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.		936
			937
			938
			939
			940
			941
			942
			943
			944
			945
			946
			947
	Ji Xin, Raphael Tang, Jaeyun Lee, Yaoliang Yu, and Jimmy Lin. 2020. DeeBERT: Dynamic early exiting for accelerating BERT inference . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2246–2251, Online. Association for Computational Linguistics.		948
			949
			950
			951
			952
			953
	Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. 2021. BERxiT: Early exiting for BERT with better fine-tuning and extension to regression . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 91–104, Online. Association for Computational Linguistics.		954
			955
			956
			957
			958
			959
			960
	Canwen Xu, Wangchunshu Zhou, Tao Ge, Furu Wei, and Ming Zhou. 2020. BERT-of-theseus: Compressing BERT by progressive module replacing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7859–7869, Online. Association for Computational Linguistics.		961
			962
			963
			964
			965
			966
			967

Deming Ye, Yankai Lin, Yufei Huang, and Maosong Sun. 2021. [TR-BERT: Dynamic token reduction for accelerating BERT inference](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5798–5809, Online. Association for Computational Linguistics.

Hao Yuan, Yongjun Chen, Xia Hu, and Shuiwang Ji. 2019. Interpreting deep models for text analysis via optimization and regularization methods. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5717–5724.

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

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NIPS*.

Chen Zhao, Chenyan Xiong, Corby Rosset, Xia Song, Paul Bennett, and Saurabh Tiwary. 2020. [Transformer-xh: Multi-evidence reasoning with extra hop attention](#). In *International Conference on Learning Representations*.

Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. [Bert loses patience: Fast and robust inference with early exit](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 18330–18341. Curran Associates, Inc.

A Inclusive KL Loss Consideration

We opted for an inclusive KL loss since CPs should be trained to cover all tokens considered important by saliency and not to be mode seeking (i.e., covering a subset of high contributing tokens considered by the saliency scores.). Suppose an exclusive KL is selected. Due to the limited learning capacity of the CP and miscalculation possibility from the saliency, the CP may be trained to maximize its contribution on noninformative tokens. While in an inclusive setting, it trains to extend its coverage over all high-saliency tokens.

Additionally, our initial research indicated that using a symmetric loss (e.g. Jensen-Shannon divergence) would produce similar results but with a significantly longer convergence time.

B Optimization of θ

In Section 3.3, we introduced θ^ℓ as a trainable parameter that increases the saliency score of [CLS].

We can deduce from Equations 6 and 7 that this parameter does not exist in the model’s computational DAG and we need to compute the derivative of \tilde{S}^ℓ w.r.t. θ^ℓ to train this parameter. Hence, first we assume that \tilde{S}^ℓ is a close estimate of \hat{S}^ℓ (due to the CPs’ training objective). Second, using a dummy variable θ_d^ℓ —that is involved in the computational graph and is always equal to 1—we reformulate \tilde{S}^ℓ :

$$\hat{S}_i^\ell \approx \tilde{S}_i^\ell = \frac{\theta_d^\ell \tilde{S}_1^\ell \mathbf{1}[i = 1] + \tilde{S}_i^\ell \mathbf{1}[i > 1]}{\theta_d^\ell \tilde{S}_1^\ell + \sum_{i=2}^n \tilde{S}_i^\ell} \quad (9)$$

This reformulation is valid due to $\theta_d^\ell = 1$ and $\sum_{i=1}^n \tilde{S}_i^\ell = 1$. Now we compute the partial derivative w.r.t. θ_d^ℓ which is the gradient that is computed in the backpropagation:

$$\frac{\partial \tilde{S}_i^\ell}{\partial \theta_d^\ell} = \frac{\tilde{S}_1^\ell (\sum_{i=2}^n \tilde{S}_i^\ell \mathbf{1}[i = 1] - \tilde{S}_i^\ell \mathbf{1}[i > 1])}{(\theta_d^\ell \tilde{S}_1^\ell + \sum_{i=2}^n \tilde{S}_i^\ell)^2} \quad (10)$$

By knowing that $\theta_d^\ell = 1$:

$$\frac{\partial \tilde{S}_i^\ell}{\partial \theta_d^\ell} = \tilde{S}_1^\ell ((1 - \tilde{S}_1^\ell) \mathbf{1}[i = 1] - \tilde{S}_i^\ell \mathbf{1}[i > 1]) \quad (11)$$

Now using our initial assumption ($\hat{S}_i^\ell \approx \tilde{S}_i^\ell$), we can substitute \tilde{S}_i^ℓ with \hat{S}_i^ℓ based on Equation 7:

$$\begin{aligned} \frac{\partial \hat{S}_i^\ell}{\partial \theta_d^\ell} &= \hat{S}_1^\ell ((1 - \hat{S}_1^\ell) \mathbf{1}[i = 1] - \hat{S}_i^\ell \mathbf{1}[i > 1]) \\ &= \frac{\theta^\ell S_1^\ell (\sum_{i=2}^n S_i^\ell \mathbf{1}[i = 1] - S_i^\ell \mathbf{1}[i > 1])}{(\theta^\ell S_1^\ell + \sum_{i=2}^n S_i^\ell)^2} \end{aligned} \quad (12)$$

In addition, the gradient of \hat{S}_i^ℓ w.r.t. θ^ℓ is as follows (cf. Equation 7):

$$\frac{\partial \hat{S}_i^\ell}{\partial \theta^\ell} = \frac{S_1^\ell (\sum_{i=2}^n S_i^\ell \mathbf{1}[i = 1] - S_i^\ell \mathbf{1}[i > 1])}{(\theta^\ell S_1^\ell + \sum_{i=2}^n S_i^\ell)^2} \quad (13)$$

By comparing Equations 12 and 13, these derivatives are related with a term of θ^ℓ :

$$\frac{\partial \hat{S}_i^\ell}{\partial \theta^\ell} \approx \frac{\partial \tilde{S}_i^\ell}{\partial \theta^\ell} = \frac{1}{\theta^\ell} \frac{\partial \tilde{S}_i^\ell}{\partial \theta_d^\ell} \quad (14)$$

Therefore, during training, we can compute the gradient w.r.t. the dummy variable θ_d^ℓ and then divide it by θ^ℓ .

C Evaluating PoWER-BERT in Single Instance Mode

Due to the static structure of PoWER-BERT, the speedup ratios reported in Goyal et al. (2020) are

based on wall time acceleration with batch-wise inference procedure. This means that some inputs might need extra padding to make all inputs with the same token length. However, since our approach and other dynamic approaches are based on single instance inference, in our procedure inputs are fed without being padded. To even out this discrepancy, we apply a single instance flops computation on the PoWER-BERT, which means we compute the computational cost for all input lengths that appear in the test dataset. Some instances may have shorter input length than some values in the resulting retention configuration (number of tokens that are retained in each layer). To overcome this issue, we update the retention configuration by selecting the minimum between the input length and each layers’ number of tokens retained, to build a new retention configuration for each input length. For instance, if the retention configuration trained model on a given task be (153, 125, 111, 105, 85, 80, 72, 48, 35, 27, 22, 5), for an input with 75 tokens length, the new configuration which is used for speedup computation will be: (75, 75, 75, 75, 75, 72, 48, 35, 27, 22, 5).

D AdapLeR Training Hyperparameters

For the initial step of finetuning BERT, we used the hyperparameters in Table 3. For both finetuning and training with length reduction, we employed an AdamW optimizer (Loshchilov and Hutter, 2019) with a weight decay rate of 0.1, warmup proportion 6% of total training steps and a linear learning rate decay which reaches to zero at the end of training.

Dataset	Epoch	LR	MaxLen.	BSZ
SST-2	5	2e-5	64	32
IMDB	5	2e-5	512	16
HateXplain	5	3e-5	72	32
MRPC	5	2e-5	128	32
MNLI	3	2e-5	128	32
QNLI	5	2e-5	128	32
AG’s News	5	2e-5	128	32
DBpedia	3	2e-5	128	32

Table 3: Hyperparameters in each dataset; LR: Learning rate; BSZ: Batch size; MaxLen: Maximum Token Length

For the adaptive length reduction training step, we also used the same hyperparameters in Table 3 with two differences: Since MRPC and CoLA have small training sets, to prolong the gradual soft-

removal process, we increased the training duration to 10 epochs. Moreover, we increase the learning rate to 3e-5. Other hyperparameters are stated in Table 4. To set a trend for λ , it needs to start from a small but effective value ($10 < \lambda < 100$) and grow exponentially per each epoch to reach an extremely high amount at the end of the training to mimic a hard removal function ($1e+5 < \lambda$). Hence, datasets with the same amount of training epochs have similar λ trends.

Dataset	γ	ϕ	λ
SST-2	5e-3	5e-4	10^{Epoch}
IMDB	5e-3	5e-4	10^{Epoch}
HateXplain	5e-2	2e-2	50^{Epoch}
MRPC	3e-2	5e-2	10×3^{Epoch}
MNLI	5e-3	5e-4	50^{Epoch}
QNLI	5e-3	1e-4	10^{Epoch}
AG’s News	1e-1	1e-1	10^{Epoch}
DBpedia	1e-1	1e-1	50^{Epoch}

Table 4: ALR hyperparameters in each dataset; Since λ increases exponentially on each epoch the corresponding formula is written.

E Task Descriptions

In the Movie reviews (Zaidan and Eisner, 2008) task, the model predicts the sentiment based on multiple sentences. The MultiRC (Khashabi et al., 2018) dataset contains a passage, a question, and multiple candidate answers, which is cast as a binary classification task of passage/question/answer triplets in ERASER benchmark.

F Additional Qualitative Examples

G Accuracy-Speedup Trade-off

Layer 0: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 1: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 2: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 3: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 4: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 5: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 6: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 7: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 8: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 9: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 10: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]
 Layer 11: [CLS] gi ##ddy Phelps touches gold for first time michael Phelps won the gold medal in the 400 individual medley and set a world record in a time of 4 minutes 8 . 26 seconds . [SEP]

Figure 5: The illustration of contribution scores obtained by CPs in each layers of the model for an input example from AG’s news (topic classification) task. The color intensity indicates the degree of contribution scores. Only the highlighted token representations are processed in each layer

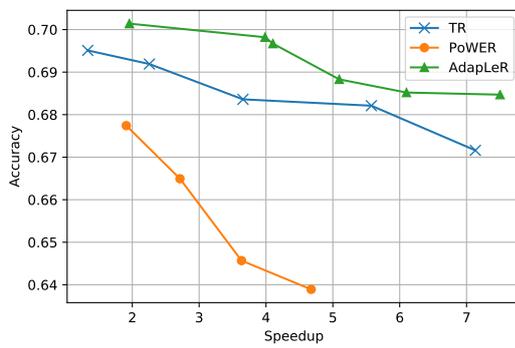


Figure 6: Accuracy-Speedup trade-off curve for AdapLeR and two other state-of-the-art reduction methods; TR: TR-BERT, PoWER:PoWER-BERT on HateXplain task.