

Discourse-Aware Soft Prompting for Text Generation

Marjan Ghazvininejad

Vladimir Karpukhin
Meta AI

Vera Gor

Asli Celikyilmaz

Abstract

Current efficient fine-tuning methods (e.g., adapters (Houlsby et al., 2019), prefix-tuning (Li and Liang, 2021), etc.) have optimized conditional text generation via training a small set of extra parameters of the neural language model, while freezing the rest for efficiency. While showing strong performance on some generation tasks, they don’t generalize across all generation tasks. We show that soft-prompt based conditional text generation can be improved with simple and efficient methods that simulate modeling the discourse structure of human written text. We investigate two design choices: First, we apply *hierarchical blocking* on the prefix parameters to simulate a higher-level discourse structure of human written text. Second, we apply *attention sparsity* on the prefix parameters at different layers of the network and learn sparse transformations on the softmax-function. We show that structured design of prefix parameters yields more coherent, faithful and relevant generations than the baseline prefix-tuning on all generation tasks.¹

1 Introduction

Recent advances in pre-trained language models (PLMs) (Lewis et al., 2020; Raffel et al., 2020; Radford et al., 2019) have made great impact on text generation research, especially when they are fine-tuned on downstream tasks such as summarization, data-to-text generation, long-question answering, etc. Consequent research have shown that PLMs’ impact can further be improved when trained with more parameters, on more data and with more compute (GPT-3 (Brown et al., 2020), Megatron (Smith et al., 2022)). On the flip side, storing larger LMs or fully fine-tuning them (updating all the parameters) on downstream tasks usually causes resource or over-fitting issues.

To mitigate fine-tuning issues, recent work have proposed *prompt-based learning* (Liu et al., 2021a),

which focus on learning textual prompts to steer PLMs’ continuation towards desired output while keeping the model parameters frozen. While providing strong control of the PLMs, such prompt engineering could be time consuming requiring manual crafting. There is a growing research direction under prompt learning towards lightweight fine-tuning (Houlsby et al., 2019; Lester et al., 2021), which update only a small number of existing or extra parameters while keeping most of the original pre-trained parameters frozen. Among them is *prefix-tuning* (Li and Liang, 2021), which prepends tunable continuous task-specific prompt vectors called *prefixes* to the input and only trains these continuous prompts during fine-tuning. Although prefix-tuning can yield comparable results to full fine-tuning on some generation tasks, it did not generalize to tasks like abstractive summarization.

In this work, we focus on prefix-tuning and investigate ways to improve its generalization on text generation tasks. We start asking the following questions that motivates our design choices: (1) *Do different parts of the transformer network process the prefix parameters more efficiently?*; (2) *Do prefix parameters capture high-level discourse structure of the input text?*; (3) *Can constraining prefix attention distribution to be structurally sparse enable better transfer of the task features?*

To address (1), we conduct empirical analysis on prefix-tuned BART (Lewis et al., 2020), by varying the size of the prefix parameters at the encoder and decoder networks on text generation tasks. We find that the prefix parameters at higher layers impact the performance the most, while sparse prefixes can be sufficient at the lower layers (§ 6.1).

Motivated by this finding and to address (2), we investigate **discourse-aware soft prompting** via **hierarchical blocking** of prefix parameters. Previous text generation work (e.g., abstractive summarization) has shown that abstraction can be better modeled with hierarchically structured architec-

¹All supporting code will be publicly released.

tures (Liu and Lapata, 2019; Fabbri et al., 2019; Xiao et al., 2021). To simulate a hierarchical discourse structure while *only* tuning additional prefix parameters, we first split the input and output text into segments and then assign sets of prefix parameters to each segment at different layers. With this structure, a set of prefixes can only be reached by their designated input or output segments during self-attention. We argue that for conditional generation tasks with hierarchically structured blocking of prefixes, we can simulate the structure of human writing styles: in input text each paragraph is a distinct section of related sentences and in output text (e.g., summary) each output sentence outlines salient concepts. Thus, a set of prefixes designated to each input and output segment at different layers can learn levels of abstractions from each section. We show performance improvements over baseline prefix tuning, yielding comparable results to full fine-tuning in several generation tasks in § 6.2.

Inspired by these findings, we address (3) by applying a suite of known **sparse attention** alternatives to standard full-attention matrix during prefix-tuning. Our goal is to analyze whether sparse prefix-tuned models can encode important features better than dense prefix-tuned models. Prior work have shown that sparsity in self-attention not only improves training efficiency, but also focusing on *salient features* while pushing down unrelated features and relations can impact the model performance. This improves language modeling (Sukhbaatar et al., 2021; Wang et al., 2020), language understanding (Shi et al., 2021; Cui et al., 2019) and text generation (Zaheer et al., 2020; Li et al., 2021; Liu et al., 2021b; Manakul and Gales, 2021). Motivated by previous work we apply **spar-sity into the self-attention** by substituting the softmax function with a sparse alternative under encoder prefix-tuning (without introducing any additional model parameters). With spectral analysis we show that sparse prefix parameters can identify important features compared to dense prefix parameters. Our quantitative analysis yield performance improvements on automatic metrics over best prefix-tuning models, while human judges generally prefer our sparse prefix model generations on factuality and coherence criteria (§ 6.3 and § 6.4).

Efficient tuning of PLMs offers a promising new direction for many NLP tasks including text generation, which we study in this work. Our results suggest that **prompt design with hierarchical**

structure and sparsity in prefix parameters: (i) generate more coherent and faithful text than baseline prefix-tuning across several summarization and data-to-text generation tasks, (ii) trail the performance of fine-tuning on most summarization tasks with a small margin, while at par with fine-tuning on data-to-text generation tasks, (iii) improve all the baselines in low-resource settings.

2 Related Work

Prompt Tuning. Recent years have observed several efficient methods for prompt-based tuning of large-scale PLMs (Liu et al., 2021a). These range from prompt engineering (Petroni et al., 2019; Cui et al., 2021), to more advanced approaches such as prompt ensembling (Mao et al., 2021), composition (Han et al., 2021; Liu et al., 2022; He et al., 2022), or prompt-aware training methods (Lester et al., 2021; Gu et al., 2021). Li and Liang (2021) propose *prefix-tuning* and show strong results on some text generation tasks, leaving room for further generalization. Here, we build directly upon the prefix-tuning from Li and Liang (2021), showing where it falls short and providing several discourse-aware prompt design approaches. We find with human evaluations (§ 6.5) on relevance criteria that the prefix-tuning struggles with encoding of salient concepts that constraint generation models require. This setting bears similarities to discourse modeling, which we discuss below.

Discourse Modeling. Several previous work make architectural design choices to teach models about the overall document discourse structure (Marcu, 1997; Barzilay and Lee, 2004; Barzilay and Lapata, 2008; Li and Hovy, 2014) to improve the summarization task. Recent work investigate different model architectures of discourse structure via: structured attention (Cohan et al., 2018), graph based methods (Dong et al., 2021), or hierarchical encoders (Pasunuru et al., 2021; Cao and Wang, 2022). We simulate the discourse structure of text via hierarchical prefixes and propose discourse-aware prompt-design for efficient PLM tuning.

Sparse Language Models. Most work on sparsity in transformers aim at improving the time and space bottleneck of dense transformers (Tay et al., 2021). Work on text generation imbue sparsity to improve coherence, fluency, n -gram diversity and reduce repetition. These work range from: sparse methods on posterior vocabulary distributions at inference time (Fan et al., 2018; Holtzman

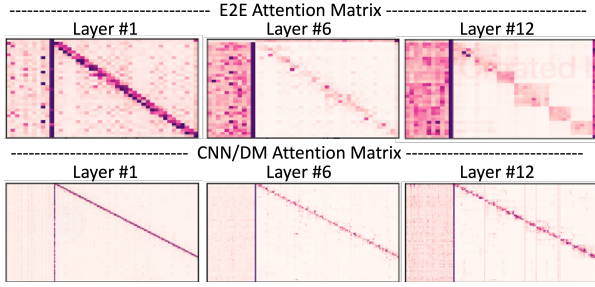


Figure 1: Encoder self-attention matrices A from layers 1, 6 and 12 of prefix-tuned models showing query attention scores (on y-axis) over all prefix+inputs keys (on x-axis). Top row are matrices of models on E2E dataset where the first 10 features on x-axis are prefix features, and bottom row are on CNN/DM dataset where first 100 features are prefix parameters.

et al., 2020), sparse attention mechanisms (Cui et al., 2019; Liu et al., 2021b; Shi et al., 2021; Sukhbaatar et al., 2021), modified softmax Martins et al. (2020), or loss functions (Welleck et al., 2020) to improve LM coherence and generalization. Similarly, we inject sparsity on the attention matrix of prefix+input features to improve the knowledge transferred to downstream text generation tasks and generating more relevant and coherent text (§6).

3 Prefix-Tuning

Following the intuition of the text-based prompt tuning methods (Liu et al., 2021a), prefix-tuning (Li and Liang, 2021) introduces task-specific prompt parameters with the goal of triggering the desired response of the LM without updating any of the original LM parameters. At each layer, it prepends tunable prefix parameters (also called *soft-prompts*) as additional keys and values to the multi-head self-attention. Prefix-tuning defines $h_i^{(l)}$ as the activation at the i -th token ($i=1 \dots T$) of the l -th layer in a L -layer transformer:

$$h_i^{(l)} = \begin{cases} P_\theta[i, :], & \text{if } i \in P_{\text{idx}} \\ LM_\phi(z_i, h_{<i}) & \text{otherwise} \end{cases} \quad (1)$$

$[,]$ indicates concatenation, P_{idx} is the sequence of prefix indices, where the activations of the first $|P_{\text{idx}}|$ positions are directly calculated by P_θ and z_i is the i -th token in the input sequence. During training only the parameters corresponding to the prefix keys and values are updated and the same objective function as finetuning is used².

²For details on prefix-tuning, pls. see (Li and Liang, 2021).

4 Discourse Aware Prefix-Tuning

Visualizing the prompt impact. To motivate the discourse-aware prompt design, we investigate the impact of prefix-parameters on transformer models during prefix-tuning. We first analyze the attention behaviour similar to (Sun and Lu, 2020). We prefix-tune two BART-LARGE models, one on data-to-text generation task with E2E dataset (Dušek et al., 2019), and another on summarization with CNN/DM (Hermann et al., 2015). For E2E we use 10-prefixes (the first 10 keys are from prefix parameters) and 100-prefixes for CNN/DM³ (similar to Li and Liang (2021)). In Figure 1, we plot the encoder self-attention distributions A for different layers averaging over all head vectors. The x -axis represents the keys while y -axis denotes the queries. For attention matrices of all the layers, see Appendix A.4 Figure 7. The attention scores show stronger relations with the prefix-keys in the E2E model compared to CNN/DM, where the prefixes exhibit weaker relations compared to the input keys. We attribute this to a few issues which we investigate in this work:

Modeling hierarchical structure. Firstly, during prefix-tuning, the model should not only focus on learning the task specific semantics, but also the models should learn the corresponding discourse structure of the downstream task datasets. To model the intrinsic structure of the input text, biasing transformer models with a type of hierarchy has been shown to improve the generation performance. For example, previous work (Cohan et al., 2018; Liu and Lapata, 2019; Cao and Wang, 2022) learns the discourse structure of human written text (e.g., the beginning, body, conclusion paragraphs, topic shifts, etc.) with hierarchically structured transformers to capture the salient aspects in the input text necessary for improved performance in summarization. With probing experiments Jawahar et al. (2019) show that BERT (Devlin et al., 2019) captures surface features and phrase-level information in the lower layers, syntactic features in the middle and semantic features and long-distance dependencies at the higher layers. Motivated by these, we apply variations of **hierarchical blocking** on prefix parameters at different layers of the network and investigate their impact on text generation with qualitative and quantitative experiments.

Introducing sparsity. Secondly, the weaker pre-

³The length of per instance prefix+input tokens is 100+512 in CNN/DM and 10+16 in E2E dataset.

fix attention in longer inputs (Figure 1-CNN/DM attention matrices) may imply that the attention neglects important connections, and potentially disturbed by many unrelated words. This issue can be attributed to the softmax function at attention score calculation (Laha et al., 2018; Cui et al., 2019). Softmax produces attention distribution with dense dependencies between words, and fails to assign near/exactly zero probability to less meaningful relations. Thus, the model neglects to put more attention to important connections while also being easily disturbed by many unrelated words (Cui et al., 2019). This issue is more pronounced in tasks like abstractive summarization, since only a handful of salient input aspects is needed to compose a coherent summary. Sparse attention mechanisms (Liu et al., 2021b; Shi et al., 2021) can remedy this issue by learning to put more emphasis on the important features.

Below we describe ways to apply a suite of blocking schemes and sparsity in prefix-tuning models as sketched in Figure 2. Each block represents attention matrix $\mathbf{A} \in \mathbb{R}^{T \times (P+T)}$, while each row vector $a_t \in \mathbb{R}^{(P+T)}$, $t = 1 \dots T$, are the attention weights of P -prefix and T -input-key features.

4.1 Prefix Blocking

As shown in Figure 2-(b) and (e), the two variations of prefix-blocking we apply here are a type of structural bias we imbue the models to simulate high-level discourse structure of documents:

Uniform Blocking (UniBlock): We first split the sequence of input tokens into segments. We allocate different sets of prefix parameters to each segment and apply blocking on the rest of the prefix parameters. In baseline prefix-tuning, a query of a token can bind with all the prefix and input key and value parameters, while in the uniform blocked prefix-tuning, the query of a token in the input or output segment can bind with all input key and values but only with the designated prefix key and value vectors. For example, if 100 prefix parameters are used and we split the input tokens into 2 segments, the first 50 prefix keys and values can only be bound with the query vectors of input tokens from the first input segment and so on. We only apply blocking to the prefix parameters and let all inputs tokens attend to each other, see Figure 2-(b). In uniform blocking, we use the same blocking schema at each layer.

Hierarchical Blocking (HierBlock): To bias the

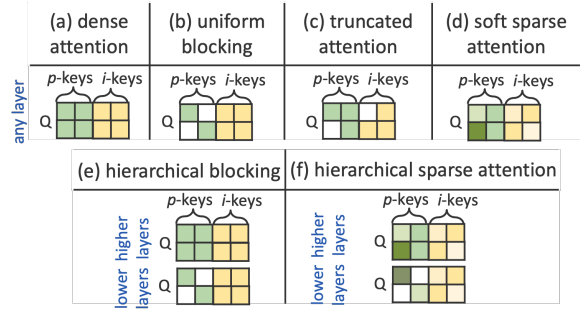


Figure 2: Sketches of attention matrices \mathbf{A} used in prefix-tuning models representing different prefix design patterns. p -keys and i -keys denote P prefix and T input keys. Sparsity of attention scores are indicated by color gradations. White cells in any row represent blocked parameters for the query.

prefix parameters with a form of hierarchy, we use the uniform prefix-blocking on the lower layers of the transformer, while we let all tokens attend to all prefixes at the top layers as shown in Figure 2-(e). The attention matrix of the top layers is same as the standard prefix-tuning of (Li and Liang, 2021) where no blocking on prefixes is applied.

4.2 Sparse Attention Prefix-Tuning

To train a prefix-tuning model that learns to highlight important input content, we apply five sparse attention design options on the encoder.

(1) Truncated Sparse Attention (TruncSA): Dai et al. (2021) used sparse cross-attention using truncation to improve salient feature extraction which showed improvements in downstream tasks performance. To simulate encoding with salient features, we apply top- p truncation on both the prefix and input keys as follows: we first add all the row elements $a_{ti} \in [0, 1]$ of the attention matrix, namely the attention scores contributing from all the queries, then normalize across all key-features, which yields new key-feature row vectors $\tilde{a}_t \in \mathbb{R}^{(P+T)}$, $\tilde{a}_t \in \tilde{\mathbf{A}}$:

$$\bar{a}_t = \sum_i a_{ti} \quad \tilde{a}_t = \bar{a}_t / (\sum_t^{(P+T)} \bar{a}_t) \quad (2)$$

Using top- p truncation (Dai et al., 2021) we truncate the feature key scores and use the top- p portion of the probably mass in each key attention score. We create a binary mask for each key feature via $mask(\tilde{\mathbf{A}}) = \text{top-}p(\tilde{\mathbf{A}}, \tau)$ by assigning 1.0 to the keys that the top- p sampling has selected, 0 otherwise and threshold parameter τ controls sparsity. Lastly, we broadcast point-wise multiplication between the sparse mask and the attention matrix \mathbf{A} to obtain the top- p sparse attention matrix $\tilde{\mathbf{A}} = mask(\tilde{\mathbf{A}}) \odot \mathbf{A}$, as sketched in Figure 2-(c).

The top- p truncation is similar to using dropout

on the features of the network while controlling the dropout rate with a user-defined threshold to compensate for overfitting and performance. Although top- p sparse attention provides automatic control over attention sparsity, truncation completely masks some features. Next, we show how to dynamically learn to apply *soft*-sparsity via sampling from a distribution.

(2) Soft Sparse Attention (SoftSA): Influencing the attention distribution with a stochastic mask to attend to salient tokens can potentially help build higher quality sparse attention for text modeling. Several work investigate novel approaches to learn the sparsity in attention matrix (Li et al., 2021; Roy et al., 2021; Shi et al., 2021) using a sampling method to formulate the right amount of sparsity. They associate the attention scores a_{ti} with each cell (t, i) in \mathbf{A} and define a sampling distribution to learn the attention mask during training as sketched in Figure 2-(d). Similarly, we apply relaxed Bernoulli distribution as a sampler to construct our stochastic mask. Since sampling from Bernoulli distribution is not differentiable, we use the Gumbel Softmax reparameterization trick (Jang et al., 2016) with gumbel-softmax:

$$\tilde{a}_t = \underset{n \in 1 \dots (P+T)}{\text{Softmax}} (a_{tn}, g, \tau) \quad (3)$$

where $g = -\log(-\log(u))$ is an independent Gumbel noise generated from the uniform distribution $u \sim U(0, 1)$ and τ is a temperature. As τ approaches zero, the gumbel output approaches to a discrete distribution in $\{0, 1\}$, becomes identical to those from the Bernoulli distribution. For details on Gumbel-softmax, see Appendix A.1.

(3) & (4) Hierarchical Sparse Attention: To simulate an intrinsic discourse structure of the input text, similar to the hierarchical blocking in § 4.1, we apply sparsity on the parameters only at the lower layers. We train hierarchical models with the dense attention at the higher layers, and apply (c) *truncated* (**HTruncSA**) or (d) *soft* sparse attention (**HSoftSA**) at the lower layers (see Figure 2-(f)).

(5) Hierarchical Blocking with Sparse Attention (HierBlock+SoftSA): The hierarchical blocking models we used in § 4.1 puts restrictions on the prefix parameters that input tokens can bind with at different layers of the network. To analyze the impact of **ensemble of prefix blocking and sparsity**, we apply sparsity on the hierarchically blocked prefix-tuning models. We apply soft sparsity (SoftSA) on the lower layers of the network attention matri-

Dataset	Domain	#Data
Summarization		
XSum (2018)	News	204K/11K/11K
CNN/DM (2015)	News	287K/13K/11K
Wikihow (2018)	DIY	157K/5.6K/5.6k
SAMSum (2019)	Dialog	15.7K/<1K/<1K
Pubmed (2018)	Clinical	203K/6K/6K
Structure to Text (S2T)		
E2E (2017; 2019)	Reviews	33K/4K/4.7K
DART (2021)	Reviews	63K/7K/12.5K

Table 1: Datasets used in the experiments.

ces of HierBlock models and keep the higher layer attention matrices dense.

5 Experiment Setup

Methods. We build all fine/prefix-tuning models on the multi-layered encoder-decoder Transformer architecture using BART-LARGE (Lewis et al., 2020), though our methods can be applied to any transformer architecture with key-value attention. We compare our discourse aware prefix-tuning approaches to full parameter fine-tuning and baseline prefix-tuning (Li and Liang, 2021). Fine-tuning updates all the LM parameters, while all prefix-tuning models freeze LM parameters and only update the prefix parameters. Baseline prefix-tuning models update prefix parameters at each layer of the transformers using dense attention while our proposed models use variations of sparse and blocked attention at different layers of the network. We choose the best models on validation dataset. For setup details see Appendix A.1.

Datasets. We conduct experiments across six datasets on two tasks: abstractive summarization and data-to-text (S2T) generation. We present a summary of the datasets in Table 1 and provide more details about the datasets in Appendix A.2.

Metrics. For all the tasks and datasets we use the n -gram match metrics: ROUGE-1/2/L (Lin, 2004) for summarization. We use BLEU (Papineni et al., 2002), NIST (Belz and Reiter, 2006), METEOR (Lavie and Agarwal, 2007), ROUGE-L, TER (Snover et al., 2006), Movers (Zhao et al., 2019) and BERTScore (Zhang* et al., 2020) for S2T tasks and report human evaluations analysis.

6 Experiment Results

6.1 Are all prefix-parameters useful?

Finding: Prefix-tuning models encode diverse but task specific features at each layer differently, while the top-layer prefixes encode abstract features.

Analysis: Earlier work (Jiang et al., 2021; Elazar

Method	XSum	CNN/DM	PubMed	Wikihow	SAMSum
	R1/R2/R-L	R1/R2/R-L	R1/R2/R-L	R1/R2/R-L	R1/R2/R-L
FineTune (FT)					
FT [*]	45.14/22.27/37.25 ♣	44.16/21.28/40.90 ♣	45.09/19.56/27.42 ♦	42.86/19.71/34.80 ♦	49.30/25.60/47.70 ★
FT (repr.)	45.47/22.40/37.42	44.30/21.14/41.22	44.96/19.00/27.74	43.26/19.38/34.89	53.02/28.30/48.73
PrefixTune (PT)					
PT [*]	43.80/20.93/36.05 ♣	-	-	-	-
PT (repr.)	43.43/20.37/35.47	42.50/19.52/39.08	42.38/16.31/24.57	39.26/15.58/28.27	52.12/26.52/48.05
UniBlock	43.49/20.58/35.80	42.43/19.38/39.15	42.81/16.83/24.87	39.34/15.55/28.75	52.42/27.70/48.12
HierBlock	43.91/20.83/36.38	43.33/20.27/40.12	43.16/16.96/25.73	40.03/15.90/30.15	52.68/27.88/48.56

Table 2: **Summarization: Prefix blocking** experiment results comparing against finetuning and prefix-tuning. Top block are best reported scores from corresponding papers [*]: ♣ (Lewis et al., 2020), ♣ (Li and Liang, 2021), ★ (Chen and Yang, 2020), ♦ (Zhang et al., 2020). (repr.) denotes our replication of finetuned and prefix-tuned BART models. Bottom block are our prefix-tuned models: Uniform (UniBlock) and Hierarchical (HierBlock) prefix blocking represent models which use prefix-blocking at different layers (§ 4.1). All models use BART-large. The best finetune (top block) models are **bolded**, best prefix-tune models (bottom block) are further **underlined**. All hyper-parameters are summarized in Appendix Table 7.

et al., 2021; Yates et al., 2021) suggests that some layers of the transformers are better than others at producing representations that are useful for a given task. To investigate if similar patterns show up in prefix-tuned models, we take XSum dataset and train models with prefix parameters only at the top layers, the bottom layers, and at a single layer.

Layers	Rouge-1/2/L
Top (8-12)	42.3/18.1/33.4
Low (1-7)	35.7/13.1/26.9
All (1-12)	42.6/19.3/34.2

Table 3: Validation Rouge scores of prefix-tuned models on XSum using only top/low layers.

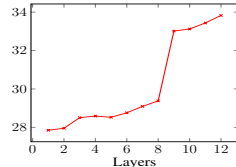


Figure 3: Validation Rouge-L on single-layer prefix-tuning with XSum.

We show layer-specific prefix-tuned models’ validation performance results in Table 3. The ‘Top’ layers model is tuned with only the top-layer prefix parameters (i.e., top 4 layers have additional prefix parameters), the ‘Low’ layers model uses only the lower-layer prefix-parameters (i.e., bottom 7 layers have additional prefix parameters) and ‘All’ layers prefix parameters is same as baseline prefix-tuning. On inspection, we see a moderate/huge performance gap between the models trained with top/lower layers, while we obtain the best performance when we tune all-layer prefix parameters. We see similar patterns on the SAMSum dialog summarization and E2E structure to text generation tasks (in Appendix A.5). We also build models when prefix parameters are used at a single layer of the network. On single layers in Figure 3, all layers contribute to the performance, the top layer prefixes perform best suggesting they might be encoding summary related abstract information.

6.2 Are hierarchical prompts effective?

Finding: Hierarchical design of prefix parameters can yield more robust information transfer in text

generation improving baseline prefix-tuning.

Analysis: To simulate learning the discourse related representations we bias prefix parameters with a structure of input documents (as discussed in §4) and experiment with two hierarchical structures: *uniform* (UniBlock) and *hierarchical* (HierBlock) from § 4.1. In Table 2 we report the performance of our models in comparison to finetuning and baseline prefix-tuning on abstractive summarization tasks. Our results indicate that prefix-blocking models improve over the baseline prefix-tuning on all summarization tasks by up to +1.1 ROUGE-L score overall. Especially for Wikihow, which are considered long document summarization task, we observe up to +2 ROUGE-L score improvement. We further observe that hierarchical blocking on prefixes also helps for data-to-text tasks, though the performance impact of structural bias is more prominent in summarization tasks. We show detailed results of data-to-text tasks and provide samples of generated outputs in Appendix A.6.

6.3 Does sparse attention help prefix-tuning?

Finding: With hierarchically structured sparsity training, prefix tuning show more sparse patterns at the lower layers. Sparse prefix parameters at lower layers, and dense at higher layers enable more efficient tuning of the prefix-parameters.

Spectrum Analysis (Statistical Proof): To investigate if our sparse models do in fact learn sparse representations, we conduct spectrum analysis on the encoder attention matrix \mathbf{A} zooming in on the prefix parameters⁴. To analyze the variation of attention scores we calculate the principal components of the attention scores of prefix parameters⁵ and plot in

⁴A similar spectrum analysis is used to prove the sparsity of the attention matrix in Linformer (Wang et al., 2020).

⁵Eigenvalues capture the variation of the attention scores distribution along different principal components.

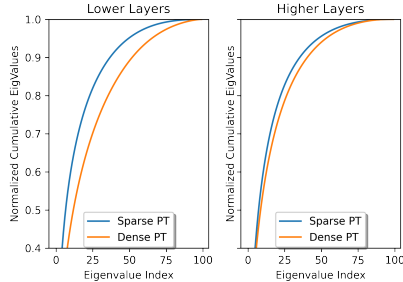


Figure 4: Spectrum analysis of the self-attention matrix comparing the baseline Dense and our Sparse Prefix-Tuned (PT) transformer model zooming in on prefix parameters of size 100. The Y-axis is the normalized cumulative singular value of the self-attention matrix \mathbf{A} , and the X-axis denotes the index of largest eigenvalue. The results are based on BART-Large on XSum dataset. The left plots the averages of all \mathbf{A} on the lower layers, while right plots averages over higher layers.

Method	XSum	CNN	PubMed	Wikihow	SAMSum
<i>Finetune</i>					
Dense	37.42	41.22	27.42	34.89	48.73
(1) TruncSA	37.02	39.96	26.37	35.58	48.12
(2) SoftSA	37.23	39.67	26.26	32.53	48.45
<i>Prefix-tune</i>					
Dense	35.47	39.08	24.57	28.27	48.05
(1) TruncSA	35.39	39.52	25.61	28.94	47.88
(2) SoftSA	35.94	39.24	25.72	28.94	47.35
(3) HTruncSA	36.42	40.00	25.28	30.02	48.00
(4) HSoftSA	36.13	39.83	24.90	30.01	48.33

Table 4: Sparse Attention experiment **ROUGE-L** results on Finetuning, and Prefix-tuning using dense and soft sparse attention designs in §4.2. The best finetuned models (top-block) are **bolded**, the best prefix-tune models (bottom-block) are further **underlined**. Full results are included in Appendix 11.

Figure 4. We observe that the spectrum distribution of prefixes in lower layers is more skewed than in higher layers, meaning that, in lower layers, more information is concentrated in the largest singular values and the rank of \mathbf{A} is lower. In summary, with sparse attention at the lower layers and dense attention at the top layers, the prefix-tuned models enables encoding important features necessary for the factual and consistent summarization. Details on spectrum analysis are provided in Appendix A.7. Next, we empirically support this statistical proof.

Sparsity Analysis: We investigate the impact of sparsity on the performance of the prefix-tuning models. For a fair comparison, we also apply attention sparsity on the finetuned models. We build prefix-tuning models with (1) Truncated Sparse Attention (TruncSA), (2) Soft Sparse Attention (SoftSA), (3) Hierarchical TruncSA (HTruncSA), with top- p sparsity at the lower layers, and dense attention at the top layers, (4) Hierarchical Soft Sparse Attention (HSoftSA), with soft sparse attention at the lower layers but dense at top layers.

We show the ROUGE-L results in Table 4.

Dataset	HierBlock	HierBlock+SoftSA
Summarization	R1/R2/R-L	R1/R2/R-L
XSum	43.91/20.83/36.38	44.00/20.93/36.59
CNN/DM	43.33/20.27/40.12	43.33/20.31/40.10
PubMED	43.16/16.83/24.87	42.96/16.64/25.00
Wikihow	39.34/15.55/28.75	39.30/15.42/28.67
SAMSum	52.58/27.58/48.42	52.83/27.94/48.72
Data-to-Text		
	BLEU/R-L/CiDER	BLEU/R-L/CiDER
E2E	67.2/69.1/2.35	68.0/69.6/3.38
	BLEU/MET/TER↓	BLEU/MET/TER↓
DART	46.6/0.39/0.45	46.0/0.39/0.46

Table 5: What happens when we apply sparsity to hierarchically blocked prompt design? Results comparing dense and sparse prefix-tuning with structurally biased prefix design (via hierarchical blocking) on various text generation tasks. The best results across two models are **bolded**.

We observe that when sparsity is used on the prefix-parameters, the prefix-tuned models outperform baseline all-dense prefix-tuning models on all datasets. The performance improvements are more pronounced on long document summarization tasks such as Wikihow, reaching close to 2.0 ROUGE-L improvements. Comparing all layers sparse models of (1) and (2) to hierarchically biased sparsity models of (3) and (4), we observe improvements with the hierarchically structured sparse prefix-tuning models. More details on quantitative analysis are provided in Appendix A.7 and Table 11.

6.4 Does sparsity on hierarchically blocked prefixes further improve performance?

Finding: The most performance gains are obtained when sparsity constraints are applied on the hierarchically blocked prefixes (Table 5).

Analysis: Recall from the earlier discussions in §6.2 that, if applying blocking on the lower layered prefixes, while letting all tokens attend to all prefixes at the top layers (HierBlock models) can improve performance. On separate set of ablations in §6.3, we also observe that if we apply sparsity at different layers of the network, the sparse parameters influence the performance compared to the dense prefix tuned parameters at all layers. We now apply sparsity on the hierarchically blocked prefix-models, combining the best hierarchically blocked models with the sparse attention.

In Table 5 we show results of our hierarchical prefix blocking (HierBlock) model against hierarchical prefix blocking model with soft sparse attention (HierBlock+SoftSA) from §4.2. To build the HierBlock+SoftSA models, we apply soft sparsity at the lower layers with blocked prefix parameters, while the top layers use dense prefixes with all tokens attending to all prefixes. In Table 5 we repeat the results of the last row from Table 2 for easy

Faithfulness		Wins % matches			
		PT	HSoftSA	HB+SoftSA	HB
Loses %	PT		50.0	64.3	50.0
	HSoftSA	50.0		60.0	60.0
	HB+SoftSA	35.7	40.0		53.9
	HB	50.0	40.0	46.1	
Overall		Wins % matches			
Loses %	PT		67.5	62.0	48.1
	HSoftSA	32.5		63.2	57.9
	HB+SoftSA	38.0	36.8		50.0
	HB	51.9	42.1	50.0	

Table 6: Human evaluation results on *Faithfulness* (top) and *Overall* (bottom) ratings. PT: Prefixtune, HSoftSA: Hierarchical Soft Attention, HB: HierBlock, HB+SoftSA: HierBlock with Soft Sparse Attention. Bold win %s indicate significance ($p < .05$).

comparison. We observe performance improvements on summarization tasks where the output summaries are shorter, (e.g., XSum SAMSUM) and less on the longer summaries (e.g., Pubmed, Wikihow). On the data-to-text generation tasks the sparsity on hierarchical blocking only improves on E2E, though both HierBlock and HierBlock+SoftSA perform better than baseline prefix-tuning models (see App. Table 12). More details are provided in Appendix A.8. Our analysis suggests that discourse aware design can improve prefix-tuning when the output generations are short (<100 tokens).

6.5 Do human evals. support our claims?

Finding: Humans generally prefer generated text from hierarchically blocked prefix-tuned models over all other models, find overall quality of generations indistinguishable from fine-tuning.

Analysis: To evaluate the generated text from our proposed methods against baseline models, we ask human annotators to rate generations on five criteria: *faithfulness* (consistent with the context), *relevance* (captures key-points), *grammaticality*, *coherence* (form a cohesive whole), and *overall quality* (informative). Table 6 shows the results of the study on faithfulness, and overall metrics. The columns show the percentage of wins of the model against its opponent on a given row. Our HierBlock (HB) and Hierarchical Soft Sparse Attention (HSoftSA) models beat prefix-tuning and HierBlock significantly ($p < .05$) beats most of our sparse models on all axes including factuality. On relevance metrics, all our models perform better than prefix-tuning and even improving finetuning. More details about

the evaluation setup as well as results on all the criteria comparing against fine-tuning and prefix-tuning can be found in Appendix A.10. In Table 13 we provide comparisons with fine-tuning and observe that HierBlock models perform as good as finetuning on all criteria.

6.6 Which structural features are harder to transfer in low-resource settings?

Finding: In low-resource settings, hierarchically designed sparse prefix parameters can efficiently transfer knowledge and represent the semantics and structure of input text yielding more accurate output generations.

Analysis: We simulate a low-resource setting by randomly sampling $k\%$ ($k=5,10,25,50$) from the training dataset of two summarization tasks: XSum on news, and Wikihow on DIY domains (see train data sizes in Table 1). We use the same hyperparameter settings as our previous models detailed in § 5. We compare our approach to finetuning and prefix-tuning under low-resource settings.

In Figure 5 on the right, we plot ROUGE-L averaging scores of models trained on XSUM and Wikihow. Our structured prefix-tuned models, HierBlock (blue) and its sparse extension which uses sparse features, HierBlock+SA (red) outperforms fine-tuned (green) and prefix-tuned models (olive), while using the same number of parameters in low resources settings (when <50% training samples are used). Although HierBlock models show consistent performance, on low-resource settings HierBlock-SA performance is more stable. (See Appendix A.11 for more details.)

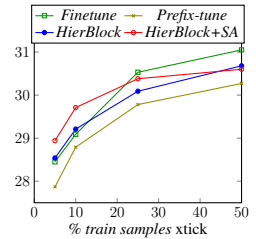


Figure 5: Average ROUGE-L scores on low-resource settings.

7 Conclusion

We have described simple but effective prompt design options for prefix-tuning of text generation tasks. We enrich prefix parameters with structural biases by way of: prefix-blocking at different layers of the network, sparsity on prefix-parameters and an ensemble of both biases. We show with quantitative and human evaluations on metrics such as coherence and faithfulness that discourse aware prefix improves prefix-tuning across all text generation tasks even at low data settings.

8 Limitations

We note a few limitations of our work: (1) our experiments are limited by available datasets, and only evaluated on limited closed domain text generation tasks; (2) we focused on efficient prefix-tuning, while ensemble of different efficient tuning models can boost performance even further; (3) we conduct experiments with ~300M parameter models extending previous work, but it will be valuable for future work to scale to larger models which may exhibit more faithful and consistent generations.

9 Ethics Statement

In this work we apply several changes to the state-of-the-art encoder-decoder modeling architecture and build several models to benchmark our new architecture with baseline architectures on several open source text generation datasets.

Intended use. Our architecture is designed to build models of abstractive document summarization and table summarization. Potentially our architecture could be used to train models for summarizing any type of datasets (e.g., any documents, textual conversational dialogues, blog posts, reports, meetings, legal forms, etc.) to further improve the productivity and efficiency of the users in their daily activities without needing to read/listen to long documents/conversations/meetings.

Failure mode. Even though our models yield factually consistent summaries, as judged by us and raters, they can still generate factually inconsistent summaries or sometimes hallucinate information that the source document does not include. This might be due to the bias or noise in the training data. Model builders wanting to use our architecture to build models on their datasets should build models with consideration of intellectual properties and privacy rights.

Misuse Potential. We note the models to be built with our architecture should be used with careful consideration especially if used to build summarization models. The generated summaries produced by our models are not controlled and use generative approaches, therefore, they could generate unreliable text. Researchers working on abstractive summarization should focus on generating factually correct, ethical and reliable text. If our models are trained on news datasets, a careful consideration should be made on factuality of the generated

text and measures have been taken to prevent model hallucinations.

References

- Regina Barzilay and Mirella Lapata. 2008. [Modeling local coherence: An entity-based approach](#). *Computational Linguistics*, 34(1).
- Regina Barzilay and Lillian Lee. 2004. [Catching the drift: Probabilistic content models, with applications to generation and summarization](#). In *Proceedings of the 2004 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Anja Belz and Ehud Reiter. 2006. [Comparing automatic and human evaluation of NLG systems](#). In *11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 313–320, Trento, Italy. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Shuyang Cao and Lu Wang. 2022. [HIBRIDS: Attention with hierarchical biases for structure-aware long document summarization](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 786–807, Dublin, Ireland. Association for Computational Linguistics.
- Jiaao Chen and Diyi Yang. 2020. [Multi-view sequence-to-sequence models with conversational structure for abstractive dialogue summarization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4106–4118, Online. Association for Computational Linguistics.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. [A discourse-aware attention model for abstractive summarization of long documents](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 615–621, New Orleans, Louisiana. Association for Computational Linguistics.

- Baiyun Cui, Yingming Li, Ming Chen, and Zhongfei Zhang. 2019. [Fine-tune BERT with sparse self-attention mechanism](#). 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 3548–3553, Hong Kong, China. Association for Computational Linguistics.
- Leyang Cui, Yu Wu, Jian Liu, Sen Yang, and Yue Zhang. 2021. [Template-based named entity recognition using BART](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1835–1845, Online. Association for Computational Linguistics.
- Wenliang Dai, Samuel Cahyawijaya, Zihan Liu, and Pascale Fung. 2021. [Multimodal end-to-end sparse model for emotion recognition](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5305–5316, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Yue Dong, Andrei Mircea, and Jackie Chi Kit Cheung. 2021. [Discourse-aware unsupervised summarization for long scientific documents](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1089–1102, Online. Association for Computational Linguistics.
- Ondřej Dušek, David M. Howcroft, and Verena Rieser. 2019. [Semantic noise matters for neural natural language generation](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 421–426, Tokyo, Japan. Association for Computational Linguistics.
- Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2021. [Amnesic probing: Behavioral explanation with amnesic counterfactuals](#). *Transactions of the Association for Computational Linguistics*, 9:160–175.
- Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. [Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. [SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization](#). In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics.
- Xiaodong Gu, Kang Min Yoo, and Sang-Woo Lee. 2021. [Response generation with context-aware prompt learning](#). *ArXiv*.
- Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. 2021. [PTR: prompt tuning with rules for text classification](#). *ArXiv*, abs/2105.11259.
- Pan He, Yuxi Chen, Yan Wang, and Yanru Zhang. 2022. [Protum: A new method for prompt tuning based on "\[mask\]"](#). volume abs/2201.12109.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. [Teaching machines to read and comprehend](#). In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text de-generation](#). In *International Conference on Learning Representations*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for NLP](#). In *International Conference of Machine Learning*.
- Eric Jang, Shixiang Gu, and Ben Poole. 2016. [Categorical reparameterization with gumbel-softmax](#). In *International Conference on Learning Representations*.
- Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy. Association for Computational Linguistics.
- Yichen Jiang, Asli Celikyilmaz, Paul Smolensky, Paul Soulos, Sudha Rao, Hamid Palangi, Roland Fernandez, Caitlin Smith, Mohit Bansal, and Jianfeng Gao. 2021. [Enriching transformers with structured tensor-product representations for abstractive summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Online. Association for Computational Linguistics.
- Mahnaz Koupaei and William Yang Wang. 2018. [Wikihow: A large scale text summarization dataset](#).

- Wojciech Kryscinski, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. [Neural text summarization: A critical evaluation](#). 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 540–551, Hong Kong, China. Association for Computational Linguistics.
- Anirban Laha, Saneem A. Chemmengath, Priyanka Agrawal, Mitesh M. Khapra, Karthik Sankaranarayanan, and Harish G. Ramaswamy. 2018. On controllable sparse alternatives to softmax. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, page 6423–6433, Red Hook, NY, USA. Curran Associates Inc.
- Alon Lavie and Abhaya Agarwal. 2007. Meteor: An automatic metric for mt evaluation with high levels of correlation with human judgments. In *Proceedings of the Second Workshop on Statistical Machine Translation, StatMT '07*, page 228–231, USA. Association for Computational Linguistics.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Haoran Li, Arash Einolghozati, Srinivasan Iyer, Bhargavi Paranjape, Yashar Mehdad, Sonal Gupta, and Marjan Ghazvininejad. 2021. [EASE: Extractive-abstractive summarization end-to-end using the information bottleneck principle](#). In *Proceedings of the Third Workshop on New Frontiers in Summarization*, pages 85–95, Online and in Dominican Republic. Association for Computational Linguistics.
- Jiwei Li and Eduard H Hovy. 2014. A model of coherence based on distributed sentence representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). 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 4582–4597, Online. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021a. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *CoRR*, abs/2107.13586.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. [P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 61–68, Dublin, Ireland. Association for Computational Linguistics.
- Yang Liu and Mirella Lapata. 2019. [Hierarchical transformers for multi-document summarization](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5070–5081, Florence, Italy. Association for Computational Linguistics.
- Ye Liu, Jianguo Zhang, Yao Wan, Congying Xia, Lifang He, and Philip Yu. 2021b. [HETFORMER: Heterogeneous transformer with sparse attention for long-text extractive summarization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 146–154, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Potsawee Manakul and Mark Gales. 2021. [Long-span summarization via local attention and content selection](#). 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 6026–6041, Online. Association for Computational Linguistics.
- Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen tau Yih, and Madian Khabza. 2021. [Unipelt: A unified framework for parameter-efficient language model tuning](#). *ArXiv*, abs/2110.07577.
- Daniel Marcu. 1997. [From discourse structures to text summaries](#). In *Intelligent Scalable Text Summarization*.
- Pedro Henrique Martins, Zita Marinho, and André F. T. Martins. 2020. [Sparse text generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4252–4273, Online. Association for Computational Linguistics.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xianguo Tang, Aadit Vyas, Neha Verma, Pranav Krishna, Yangxiaokang Liu, Nadia Irwanto, Jessica

- Pan, Faiaz Rahman, Ahmad Zaidi, Murori Mutuma, Yasin Tarabar, Ankit Gupta, Tao Yu, Yi Chern Tan, Xi Victoria Lin, Caiming Xiong, Richard Socher, and Nazneen Fatema Rajani. 2021. Dart: Open-domain structured data record to text generation. *arXiv preprint arXiv:2007.02871*.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.
- Jekaterina Novikova, Ondrej Dušek, and Verena Rieser. 2017. The E2E dataset: New challenges for end-to-end generation. In *Proceedings of the 18th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, Saarbrücken, Germany. ArXiv:1706.09254.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, page 311–318, USA. Association for Computational Linguistics.
- Ramakanth Pasunuru, Asli Celikyilmaz, Michel Galley, Chenyan Xiong, Yizhe Zhang, Mohit Bansal, and Jianfeng Gao. 2021. Data augmentation for abstractive query-focused multi-document summarization. In *AAAI*.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- PurdueOWL. 2019. Journalism and journalistic writing: The inverted pyramid structure.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- 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. *Journal of Machine Learning Research*, 21(140):1–67.
- Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. 2021. Efficient Content-Based Sparse Attention with Routing Transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68.
- Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.
- Han Shi, Jiahui Gao, Xiaozhe Ren, Hang Xu, Xiaodan Liang, Zhenguo Li, and James Tin-Yau Kwok. 2021. Sparsebert: Rethinking the importance analysis in self-attention. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 9547–9557. PMLR.
- Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, Elton Zheng, Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi, Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. 2022. Using deepspeed and megatron to train megatron-turing NLG 530b, A large-scale generative language model. volume abs/2201.11990.
- Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.
- Sainbayar Sukhbaatar, Da Ju, Spencer Poff, Stephen Roller, Arthur Szlam, Jason Weston, and Angela Fan. 2021. Not all memories are created equal: Learning to forget by expiring. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*.
- Xiaobing Sun and Wei Lu. 2020. Understanding attention for text classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3418–3428, Online. Association for Computational Linguistics.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2021. Long range arena : A benchmark for efficient transformers. In *International Conference on Learning Representations*.
- Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. 2020. Linformer: Self-attention with linear complexity. *CoRR*, abs/2006.04768.
- Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2020. Neural text generation with unlikelihood training. In *International Conference on Learning Representations*.

Wen Xiao, Iz Beltagy, Giuseppe Carenini, and Arman Cohan. 2021. [Primer: Pyramid-based masked sentence pre-training for multi-document summarization](#).

Andrew Yates, Rodrigo Nogueira, and Jimmy Lin. 2021. [Pretrained transformers for text ranking: Bert and beyond](#). In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining, WSDM '21*, page 1154–1156, New York, NY, USA. Association for Computing Machinery.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. [Big bird: Transformers for longer sequences](#). volume 33.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2020. [Pegasus: Pre-training with extracted gap-sentences for abstractive summarization](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 11328–11339. PMLR.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. [MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance](#). 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 563–578, Hong Kong, China. Association for Computational Linguistics.

A Appendix

Context: How to straighten a curly wig?

It is the best method for both straightening and curling a wig. It is also the best method for synthetic wigs that may be more delicate to heated styling tools.

Allow the hair to fall straight down without resting on anything. Tripods work brilliantly if you fit a mannequin head on them.

Boil about 1.5L for a long wig and 1L or less for a short wig. The heat of the water will soften the plastic fibers and will pull them down. Use a wide soft-bristled brush.

For extremely curly wigs, it will take more than one treatment and a lot of time. It does no damage to the wig and can be done while the wig is still wet. This is a permanent style; you will need to re-curl the wig if you don't want it straight anymore. You're basically remolding plastic into a new shape, so will need to mold it back if you want the curls again.

Output Summaries:

Human : consider trying the hot water method. place the wig on a stand. boil a pot full of water. pour it carefully over the wig, making sure that it completely covers it. avoid combing or brushing a wig when it is wet. wait until it is dry. repeat the method if necessary.

Finetune: use this method **for very curly wigs**. place the wig on a flat surface. **boil water** in a large pot. remove the pot from the heat and **allow the water to cool**. brush the wig with a wide-bristled brush. repeat the process until the entire wig is straight. curl the wig if you don't want it straight anymore.

Prefixtune: understand the benefits of this method. lay the wig flat on a flat surface. **boil a small amount of water**. brush the wig with a soft-bristled brush. repeat the process if necessary. re-curl the wig if desired. finished.

HSoftSA: use this method if you want to straighten the wig. place the wig on a mannequin head. boil a pot of water. brush the wig with a soft-bristled brush. repeat as needed. re-curl the wig if necessary.

HierBlock: heat the water in a large pot over medium heat. put the wig in the pot and allow it to sit for a few minutes. remove the wig from the pot. brush the wig with a soft-bristled brush. repeat the process with the other wig.

HierBlock+SoftSA: wash the wig with **warm** water. put the wig on a mannequin head. rinse the wig. brush the hair with a soft-bristled brush. repeat the process until the hair is completely straight. re-curl the wig if you want.

Figure 6: Model Generated Output Text on Wikihow Dataset. The red colored text indicates factual errors, repetitions, and incoherent text.

A.1 Hyperparameters (Cont. from § 5)

We fit our BART-LARGE models to their respective datasets with the hyperparameters shown in Table 7. Encoder/decoder block sizes indicate the size of the segments we split the input/output tokens. For instance, if the encoder block size is 2, we split the input tokens into two segments. Each segment has designated set of prefixes which can vary at each layer. In hierarchical blocking models

Parameter	Xsum	CNN/DM	PubMed	Wikihow	SAMSum	E2E	DART
learning rate	5e-05	5e-05	5e-05	5e-05	5e-05	5e-05	5e-05
# epochs	30	30	30	30	20	10	10
batch size	8	8	8	8	16	16	16
prefix-length	100	50-100	100-200	100-200	10-20	5-10	5-10
beamsize	5	5	5	5	5	5	5
Hierarchical Blocking							
encoder block size	1,2,3	1,2,3	1,2,3	1,2,3	1,2,3	1,2,3	1,2,3
decoder block size	1,2	1,2	1,2	1,2	1,2	1,2	1,2
Sparse Attention							
top- p	95.%	95.%	95.%	95.%	95.%	95.%	95.%
τ (top- p)	1.0,0.1	1.0,0.1	1.0,0.1	1.0,0.1	1.0,0.1	1.0,0.1	1.0,0.1
τ (soft attn.)	1.0,0.1,0.01	1.0,0.1,0.01	1.0,0.1,0.01	1.0,0.1,0.01	1.0,0.1,0.01	1.0,0.1,0.01	1.0,0.1,0.01

Table 7: Hyperparameters of different prefix-tuned models.

Corpus	Version	License	Citation	Link
XSum	v1	MIT	Narayan et al. (2018)	https://github.com/EdinburghNLP/XSum
CNN/DM	v1	MIT	Hermann et al. (2015)	https://github.com/abisee/cnn-dailymail
PubMed	v1	Creative Commons	Cohan et al. (2018)	https://github.com/armancohan/long-summarization
WikiHow	v1	CC-BY-NC-SA	Koupae and Wang (2018)	https://github.com/mahnazkoupae/WikiHow-Dataset
SAMSum	v1	CC BY-NC-ND 4.0	Gliwa et al. (2019)	https://github.com/giancolu/Samsung-dataset
E2E	v1	CC4.0-BY-SA	Dušek et al. (2019)	https://github.com/tuetschek/e2e-cleaning
DART	v1	MIT	Nan et al. (2021)	https://github.com/Yale-LILY/dart

Table 8: Additional documentation of scientific artifacts used in our paper.

(HierBlock) we segment the lower layers, so the prefixes are blocked for different segments, while at the top layers no segmentation or blocking is applied. We use at most two segments in the output text since the text generations tasks we investigate in this work contain much shorter output tokens compared to the input tokens.

Gumbel Softmax Reparameterization Trick: Sampling introduces discrete valued parameters which are not differential at training time. Thus, we resort to the GumbelSoftmax trick (Jang et al., 2016), which provides a tool for sampling from a continuous approximation of a discrete distribution. The Gumbel-Softmax trick considers a discrete variable with class probabilities π_1, \dots, π_k , and draws samples g_1, \dots, g_k from a Gumbel distribution, Gumbel(0,1), as follows:

$$y = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log \pi_j + g_j)/\tau)} \quad (4)$$

for $i = 1, \dots, k$. The Gumbel(0,1) distribution can be sampled using inverse transform sampling by drawing $u \sim \text{Uniform}(0,1)$ and computing $g = -\log(-\log(u))$. Plugging in the samples g_i and the class probabilities π_i in Eq. 4, we generate a k -dimensional sample vector $y \in \Delta^{k-1}$, that is the continuous approximation of the one-hot-encoded representation of the discrete variable d . In fact, as τ approaches 0, samples from the Gumbel-Softmax distribution become one-hot, making it discrete.

For more details please refer to (Jang et al., 2016).

A.2 Dataset Details (Cont. from §5)

All datasets are in English language. The summarization datasets range from extreme abstractive summarization with XSum (Narayan et al., 2018) to summarize documents into one summary sentence, conversational summarization using SAMSum dataset (Gliwa et al., 2019), long clinical document summarization with PubMed (Cohan et al., 2018)⁶ and DIY domain with Wikihow (Koupae and Wang, 2018), and commonly used CNN/DM (Hermann et al., 2015; See et al., 2017) news article summarization dataset with an "Inverted Pyramid" (PurdueOWL, 2019) document structure (Kryscinski et al., 2019). We also investigate S2T datasets on customer reviewers including E2E (Novikova et al., 2017; Dušek et al., 2019) and DART (Nan et al., 2021) with each input being a semantic RDF triple set derived from data records in tables and sentence descriptions that cover all facts in the triple set.

XSum (Narayan et al., 2018) is a collection of 227k BBC News articles ranging from 2010 to 2017. The dataset covers a wide range of subjects. The single-sentence summaries are written by pro-

⁶We acknowledge that the source of dataset is the NLM Catalog, and the citations used in Pubmed corpus may not reflect the most current/accurate data available from NLM, which is updated regularly.

fessionals.

CNN/DailyMail (Hermann et al., 2015) dataset contains 93k news articles extracted from CNN News, and around 220k articles extracted from the Daily Mail newspapers. The summaries are human written bullet point text which are provided in the same source documents. In our experiments we use the non-anonymized version, which is commonly used in summarization research papers.

PubMed (Cohan et al., 2018) is a long document dataset of 215K scientific publications from PubMed. The task is to generate the abstract from the paper body.

WikiHow (Koupaee and Wang, 2018) is a large-scale dataset of 200K instructions from the online WikiHow.com website. Each instance consists of multiple instruction-step paragraphs and an accompanying summary sentence of each paragraph. The task is to generate the concatenated summary-sentences from the paragraphs.

SAMSum (Gliwa et al., 2019) is a multi-turn dialog corpus of 16K chat dialogues and manually annotated summaries. The task is to generate an abstractive summary of the dialog with coherent discourse structure of the original dialog.

E2E (Dušek et al., 2019) is a structured data to natural language summary dataset that provides information about restaurants. The structured inputs consists of different attributes (slots) such as name, type of food or area and their values. It contains 50K instances of diverse descriptions of the structured input introducing challenges, such as open vocabulary, complex syntactic structures and diverse discourse phenomena.

DART (Nan et al., 2021) is a text generation dataset for open-domain structured data-record to text generation. It consists of 82K examples from variety of domains. The inputs are in semantic RDF triple set form which are derived from data records in tables and tree ontology of the schema. The output generations are human annotated with sentence descriptions that cover all facts in the triple set.

Licence details In our experiments, we use several datasets (as detailed above) from public resources. Table 8 summarizes the licences. All data are solely used for research purposes.

A.3 Compute Infrastructure and Run time

Each experiment runs on a single machine with 8 GPUs. Depending on the training dataset size, summarization models require from 5.5 hours to 18 hours to train. The data-to-text datasets are much smaller which takes less than 4 hours. All fine-tuned models follow the BART-large transformer architecture with a total of 12 layers, 1024 hidden dimensions, and 406M parameters. The prefix-models increase the parameters size of fine-tune models by 0.1% up to 2% depending on the number of prefix parameters. See hyperparameters details in Appendix A.1.

A.4 Visualization of Prefix Parameters (Cont. from § 4)

To analyze the attention behaviour (similar to (Sun and Lu, 2020)) we plot the attention matrix of the prefix-tuned models focusing on the prefix parameters. We use a prefix-tuned BART-LARGE (12-layer stacked transformer) on two tasks: data-to-text generation on E2E (Dušek et al., 2019) and summarization on CNN/DM (Hermann et al., 2015). In Figure 7, we plot the encoder self-attention distributions A for different layers averaging over head vectors. The x -axis represent the keys, while y -axis denote the queries.

A.5 Are All Prefix Parameters Useful? (Cont. from § 6.1)

We investigate the influence of prefix parameters on different layers of the network. For this experiments we trained BART-LARGE and add prefix parameters only at the top layers, lower layers and all layers (this is same as baseline prefix-tuning models). On XSum dataset, we observed a large performance gap between the models trained with top/lower layers, while we obtain the best performance when we tune all-layer prefix parameters (in Table 3 in the main text). Here, we investigate if similar performance gains are observed on dialog summarization (SAMSum) and data to text generation (E2E) tasks.

We show the performance scores of our experiments on validation datasets in Table 10. We observe similar results as the analysis on XSum dataset. Top layers prefix parameters learn salient features related to the task, though using prefixes at all layers yields better performance.

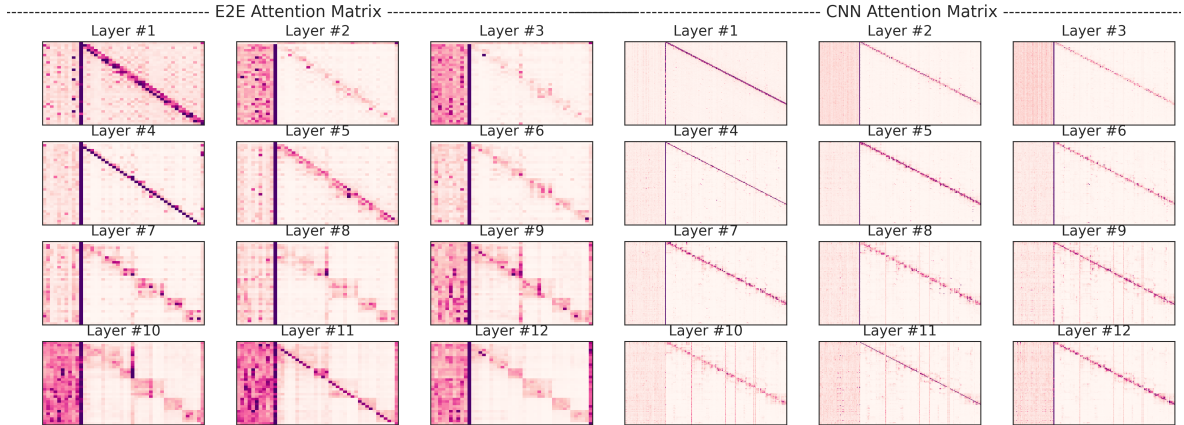


Figure 7: Encoder self-attention matrices A of prefix-tuned models indicating the query attention scores over all keys (prefix+inputs) on the y-axis. The scores are averaged over all heads. The left block is for E2E dataset where the first 10 features represent prefix features, while CNN/DM dataset on the right with first 100 features represent the prefixes.

Method	#Parm.	E2E					DART				
		BLEU	NIST	MET	R-L	CIDEr	BLEU	MET	TER ↓	Mover	BERT
GPT-2-Large											
Finetune (*)	774M	68.5	8.78	46.0	69.9	2.45	47.0	0.39	0.46	0.51	0.94
Prefixtune (*)	774M+%0.1	70.3	8.85	46.2	71.7	2.47	46.7	0.39	0.45	0.51	0.94
BART-Large											
Finetune	406M	67.8	8.76	45.1	69.5	2.38	46.1	0.38	0.46	<u>0.53</u>	0.95
Prefixtune	406M+%0.1	67.3	8.66	44.8	68.6	2.34	45.9	0.39	0.45	<u>0.53</u>	0.95
Uniblock	406M+%0.1	66.1	8.60	45.0	68.3	<u>2.36</u>	46.5	0.39	0.46	<u>0.53</u>	0.95
HierBlock	406M+%0.1	67.2	8.70	45.1	69.1	2.35	46.6	0.39	0.45	0.54	0.95

Table 9: **Data-to-Text: Prefix-Blocking** Models compared against Finetune and Prefixtune models. (*) Top block are best reported numbers in (Li and Liang, 2021) using GPT-2 LARGE model, twice the size of BART-LARGE. Bottom block are our experiment results. Best results on GPT-2 and BART-LARGE models are **bolded** separately and second best BART-LARGE models are further **underlined**. For all the prefix-tuned BART models, the prefix-length is 5 for E2E and 10 for DART dataset. All hyper-parameters are summarized in Appendix Table 7.

Method	SAMSum		E2E	
	R1/R2/RL	BLEU/RL	BLEU/RL	BLEU/RL
Top (8-12)	49.55/23.72/42.16	64.3/65.7	64.3/65.7	64.3/65.7
Low (1-7)	43.16/19.08/38.14	62.4/62.6	62.4/62.6	62.4/62.6
All (1-12)	50.16/25.03/43.16	65.4/66.5	65.4/66.5	65.4/66.5

Table 10: Results of prefix-tuned models on validation datasets of SAMSum (from the summarization task) and E2E (from the structure to text task) using only the top/low layers.

A.6 Investigation of Hierarchical Prompt Design (Cont. from § 6.2)

We investigate if blocking prefixes helps for data-to-text tasks. Table 9 shows the results. Similar to summarization experiments in § 6.2, we observe improvements with hierarchical blocking on E2E dataset, though the improvement is minimal in DART dataset. We include previous best model results reported in (Li and Liang, 2021). Their results are from GPT-2 LARGE, which is twice the size of BART-LARGE models, so our results are slightly lower. We replicated the fine-tuning and pre-

fixtuning results for fair comparison (top two rows of the bottom block in Table 9). We also provide the model sizes in terms of number of parameters. We conclude from these results that the prefix models tuned with structurally biased additional set of parameters can yield more robust information transfer reaching as good as finetuning models. In Figure 6 we show the output summaries generated by some of our best discourse aware prefix-tuned models in comparison to baseline fine-tuned and prefix-tuned models.

A.7 Investigation of the Impact of Sparsity (Cont. from § 6.3)

Spectrum Analysis: We conduct spectrum analysis of the encoder attention matrix A zooming in on the prefix parameters to investigate if our sparse models do in fact learn sparse representations. A similar spectrum analysis has been used to prove the sparsity of the attention matrix in Linformer (Wang et al., 2020), a sparse transformer. Our goal

Method	Xsum	CNN/DM	PubMed	Wikihow	SAMSum
	R1/R2/R-L	R1/R2/R-L	R1/R2/R-L	R1/R2/R-L	R1/R2/R-L
<i>Finetune</i>					
Dense (reproduce)	45.47/22.40/37.42	44.30/21.14/41.22	45.09/19.56/27.42	43.26/19.38/34.89	53.02/28.30/48.73
(1) TruncSA	45.17/21.98/37.02	43.02/20.03/39.96	41.08/17.02/26.37	44.85/20.00/35.58	52.45/27.56/48.12
(2) SoftSA	45.34/22.02/37.23	42.97/20.44/39.67	40.15/16.30/26.26	41.56/17.80/32.53	52.47/27.88/48.45
<i>Prefix-tune</i>					
Dense (reproduce)	43.43/20.37/35.47	42.50/19.52/39.08	42.38/16.31/24.57	39.26/15.58/28.27	52.12/26.52/48.05
(1) TruncSA	43.56/20.62/35.96	42.80/19.81/39.52	42.50/16.85/25.61	39.43/15.69/28.94	52.09/27.49/47.88
(2) SoftSA	43.80/20.82/35.94	42.25/19.50/39.24	42.90/16.95/ 25.72	39.31/15.61/28.94	51.72/25.51/47.35
(3) HTruncSA	44.17/21.11/36.42	43.17/20.17/40.00	43.30/17.00/25.28	40.00/16.11/30.02	52.12/26.94/48.00
(4) HSoftSA	44.05/20.84/36.13	43.10/20.06/39.83	42.30/16.20/24.90	39.83/16.10/30.01	52.37/27.57/48.33

Table 11: **Summarization: Sparse Attention** experiment results on Finetuning, and Prefix-tuning with Truncated (TruncSA) and Bernoulli Sampling soft attention (SoftSA) and Hierarchical Truncated (HTruncSA) and Soft Attention (HSoftSA) for Prefix-Tuning. All models are based on BART-LARGE. Best finetune results (top block) across models are **bolded** while the best prefix-tune models (bottom block) are further **underlined**. All hyper-parameters are summarized in Appendix Table 7.

Method	#Parm.	E2E					DART				
		BLEU	NIST	MET	R-L	CIDEr	BLEU	MET	TER ↓	Mover	BERT
Finetune (repr.)	406M	67.8	8.76	45.1	69.5	2.38	46.1	0.38	0.46	0.53	0.95
Prefix-tune (repr.)	406M+%0.1	67.3	8.66	44.8	68.6	2.34	45.9	0.39	0.45	0.53	0.95
HSoftSA	406M+%0.1	66.2	8.57	45.0	68.7	2.33	46.2	0.39	0.45	0.53	0.95
HierBlock+SoftSA	406M+%0.1	68.0	8.76	45.3	69.6	2.38	46.0	0.39	0.46	0.53	0.95

Table 12: **Data-to-Text: Hierarchical Sparse Attention** Models compared against Finetune and Prefix-tune models (repr=reproduced). All models use BART-Large as backbone and best models are **bolded** and second best are further **underlined**. For all the prefix-tuned models, the prefix-length is 5 for E2E and 10 for DART dataset. All hyper-parameters are summarized in Appendix Table 7.

is to analyze the principal components of the sub-space that captures the variation of the attention scores in prefix parameters. The eigenvalues capture the variation of the attention scores distribution along different principal components. The higher the elbow in the spectrum graph, the less parameters are used and the model learns to represent the inputs with only the salient terms ignoring superfluous details.

For our spectrum analysis, we compare the baseline prefix-tuning, which encodes a *dense* attention matrix everywhere in the network (Dense PT) against one of our sparse prefix-tuned models with truncated attention matrix (Sparse PT), as we explained in § 4.2-(a), using top- p sampling. Both models are a 12-layer stacked transformer (BART-LARGE) trained on XSum extreme summarization task. We apply singular value decomposition into \mathbf{A} across different layers and different heads of the model, and plot the normalized cumulative singular value averaged over 1000 sentences. We compare the models’ sparsity patterns at the top and at the lower layers separately as shown in Figure 4. The two figures exhibit a long-tail spectrum distribution across layers and heads. This implies that most of the information of matrix \mathbf{A} can be recovered from the first few largest singular values. We observe that the spectrum distribution in lower layers

is more skewed than in higher layers, meaning that, in lower layers, more information is concentrated in the largest singular values and the rank of \mathbf{A} is lower. With sparse attention at the lower layers and dense attention at the top layers, the prefix-tuned models can encode salient features controlling the generation.

Sparsity Analysis: In Table 11 we show the ROUGE-1, ROUGE-2 and ROUGE-L scores of fine-tuned and prefix-tuned summarization models comparing dense and sparse attention impact. We observe that when sparsity is used on the prefix-parameters, the prefix-tuned models outperform dense counterparts. The performance improvements are more pronounced on shorter generation tasks such as XSUM but we still see improvements reaching up to 2.0 ROUGE-L score improvements on longer documents such as Wikihow. Similar performance patterns are observed in Table 12 on data-to-text generation tasks using E2E and DART.

A.8 Investigation of the Impact of Sparsity on Hierarchically Blocked Prefixes (Cont. from § 6.4)

In Table 11 we showed ROUGE-L results of our hierarchical prefix blocking (HierBlock) model against hierarchical prefix blocking model with soft sparse attention (HierBlock+SoftSA). We observe

Criteria	PrefixTune HierBlock			PrefixTune HSoftSA			PrefixTune HierBlock+SoftSA		
	wins	wins	same	wins	wins	same	wins	wins	same
factuality	36	36	78	25	45	80	39	38	73
relevance	31	29	90	20	39	91	27	33	90
gramaticality	26	28	96	24	24	102	28	20	102
coherence	32	32	86	28	37	85	32	29	89
overall	40	38	72	30	49	71	41	38	71

Criteria	HierBlock HSoftSA			HierBlock HierBlock+SoftSA			HSoftSA HierBlock+SoftSA		
	wins	wins	same	wins	wins	same	wins	wins	same
factuality	28	42	80	28	42	80	35	41	74
relevance	25	34	91	23	35	92	31	32	87
gramaticality	18	35	97	26	23	101	26	19	105
coherence	18	44	88	34	30	86	35	27	88
overall	28	48	74	32	44	74	38	38	74

Criteria	Finetune HierBlock			Finetune HierBlock+SoftSA		
	wins	wins	same	wins	wins	same
factuality	48	39	63	44	42	64
relevance	37	40	73	31	47	72
gramaticality	38	24	88	43	17	90
coherence	45	32	73	46	30	74
overall	45	49	56	39	49	62

Table 13: Head-to-Head comparison of human evaluations on random subset of Wikihow dataset.

improvements on performance on most summarization tasks including news summarization (XSum and CNN/DM), dialog summarization (SAMSum). We find that HierBlock+SoftSA models show much larger improvements on XSUM and Wikihow summarization, from +0.5 to close to 2 ROUGE-L scores. On the structure to text generation tasks the sparsity on hierarchical blocking helps on some datasets (with E2E), though both HierBlock and HierBlock+SoftSA perform better than the baseline prefix-tuning models (see Table 12).

A.9 Automatic Evaluations (Cont. from § 5)

For model evaluations we use ROUGE-1/2/L using Python rouge-score 0.0.4 version licensed under the Apache 2.0 License. We use the default ROUGE script `rouge.py` from the GEM evaluation [shared task](#). All other metrics are adopted from their linked corresponding papers. We use the official evaluation script for BLEU (Papineni et al., 2002), METEOR (Lavie and Agarwal, 2007), NIST (Belz and Reiter, 2006) TER (Snover et al., 2006), MoverScore (Zhao et al., 2019), BERTScore (Zhang* et al., 2020).

A.10 Human Evaluations (Cont. from § 6.5)

We perform human evaluations to establish that our model’s ROUGE improvements are correlated with human judgments. We compare the generations from four models: baseline prefix-tune (PT), Hierarchically Blocked PT (HierBlock/HB), Hierarchical Soft Sparse Attention PT (HSoftSA) and the ensemble of the blocked sparse model (HierBlock+SoftSA). We use the following as eval-

uation criteria for generated summaries, which we include in the instructions for the annotators.

Faithfulness: Are the details in the summary fully consistent with the details in the source document? The summary must not change any details from the source document. The summary also must not hallucinate any information that is not in the source document.

Relevance: Does the summary capture the key points of the text? Are only the important aspects contained in the summary? Is there any extra/irrelevant information?

Grammaticality: Considers the grammatical quality of each individual sentence in the summary. For each sentence, does it sound natural and grammatically correct?

Coherence: Does the summary form a cohesive, coherent whole? Is it well-written, well-structured and well-organized? Is it easy to follow? It should not be a heap of related information, but should build from sentence to sentence to a coherent body of information about a topic.

Overall Quality: Given the input context, is the summary satisfactory? Does the summary provide good quality information to the user? Is it helpful, informative and detailed enough given the information that’s contained in the text? Which summary of the two do you prefer best overall?

Annotator Details: Human annotation was conducted by 9 professional raters (7 linguist raters, 1

linguist subject-matter-expert and 1 linguist) employed by an outsourcing company handling content moderation. All raters are monolingual native speakers of English; 6 have a minimum of high school degree or equivalent and 3 have a bachelor’s degree. Raters received compensation starting at \$18 per hour (which is close to 2.5 minimum wage in the state where the raters are located) and were also provided with Premium Differential as part of their contracts. Each rater conducted between 44 and 175 pairwise evaluations. Data collection protocol was reviewed by expert reviewers and received expedited approval as the data presented to the raters did not contain any sensitive or integrity-violating content. Participant consent was obtained as part of the non-disclosure agreement signed by each rater employee upon hire. All raters have also signed a sensitive content agreement that outlined the types of content they may encounter as part of their employment, associated potential risks and information and wellness resources provided by the outsourcing company to its employees.

Human Evaluation Procedure: We randomly select 50 samples from the Wikihow test set and ask 9 trained judges to evaluate them on the 5 criteria defined above. We perform head-to-head evaluation (more common in DUC style evaluations), where judges are shown the original document, the ground truth summary and two model summaries in random order. The judges are then asked to compare two model summaries based on each of the five criteria. In each case, a judge either has the option to choose a model summary that ranks higher on a given criterion (i.e., respond by identifying the winning summary), or assert that both summaries are similar given the criterion and rate the comparison as "same". The evaluation of each pair of summaries across all 5 criteria takes on average between 5 and 10 minutes to complete. The raters were shown the data to be rated in a spreadsheet, where each line contained multiple columns in sequence: document, human written summary, model-A generated summary, model-B generation summary, and five additional columns indicating faithfulness, relevance, gramaticality, coherence, overall quality. The headers of the columns were clearly stated. The rates enter a/b/same in each corresponding cell when comparing summaries head-to-head based on each criteria.

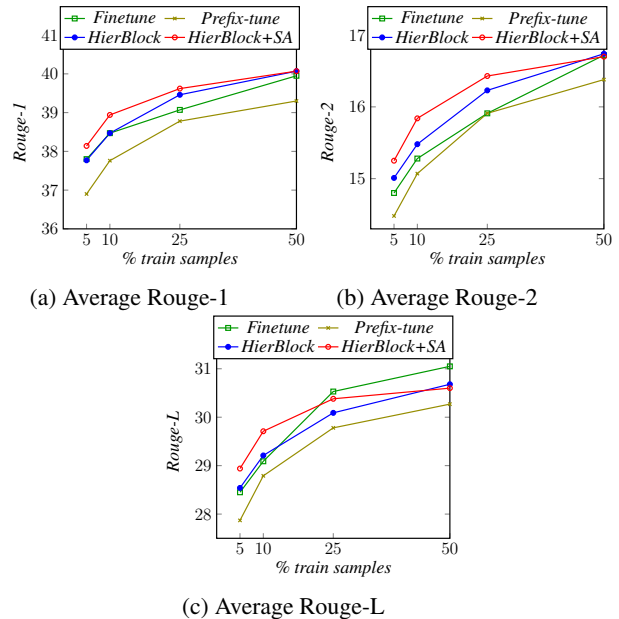


Figure 8: Quantitative analysis on **low-resource settings**. The charts show **average** of ROUGE-1, ROUGE-2, ROUGE-L scores from models trained on two summarization tasks: XSUM and Wikihow. Our structurally biased parameter tuned HierBlock (blue) and HierBlock+SA (red) consistently outperforms the baseline Finetuned (green) and Prefix-tuned models (olive) when <50% training data is used.

Human Evaluation Results: In Table 13 we show head-to-head evaluation scores on all five metrics showing wins from each model as well as when both are selected as equal. Each sub-table compare a different model. Our Hierarchical Blocking (HierBlock) and Hierarchical Soft Sparse Attention (HSoftSA) models beat prefix-tuning and HierBlock significantly ($p < .05$) beats most of our sparse models on all axes including factuality. In

On a small data annotation, we also compare two of our best models HierBlock and HierBlock+SoftSA againsts best finetuning model generations, which are shown in the same Table 13. We observe that in most cases both of our models are preferred as good as finetuning on all criteria, except on overall, the HierBlock summaries are ranked much higher than fine-tuning models.

A.11 Low-data settings (Cont. from § 6.6)

In Figure 8, we plot the ROUGE-1, ROUGE-2 and ROUGE-L scores averaging scores from two summarization tasks (XSUM and Wikihow). Our structured prefix parameter tuned models, HierBlock (blue) and its sparse extension which uses sparse features, HierBlock+SA (red) outperforms Prefix-tuned models (olive), while using the same number of parameters in low resources settings (when

<50% training samples are used). Both models outperform Finetuned models (green) on ROUGE-1 and ROUGE-2 metrics (Figure 8-(a)&(b)). While the HierBlock models show consistent performance, we conclude that on low-resource settings HierBlock-SA performance is more stable.