

Progressive Generation of Long Text with Pretrained Language Models

Bowen Tan¹, Zichao Yang¹, Maruan Al-Shedivat¹, Eric P. Xing^{1,2,3}, Zhiting Hu^{1,4}

¹Carnegie Mellon University, ²Petuum Inc., ³MBZUAI, ⁴UC San Diego

{btan2, zichao, alshedivat, epxing}@andrew.cmu.edu, zhh019@ucsd.edu

Abstract

Large-scale language models (LMs) pre-trained on massive corpora of text, such as GPT-2, are powerful open-domain text generators. However, as our systematic examination reveals, it is still challenging for such models to generate coherent long passages of text (e.g., 1000 tokens), especially when the models are fine-tuned to the target domain on a small corpus. Previous planning-then-generation methods also fall short of producing such long text in various domains. To overcome the limitations, we propose a simple but effective method of generating text in a progressive manner, inspired by generating images from low to high resolution. Our method first produces domain-specific content keywords and then progressively refines them into complete passages in multiple stages. The simple design allows our approach to take advantage of pretrained LMs at each stage and effectively adapt to any target domain given only a small set of examples. We conduct a comprehensive empirical study with a broad set of evaluation metrics, and show that our approach significantly improves upon the fine-tuned large LMs and various planning-then-generation methods in terms of quality and sample efficiency. Human evaluation also validates that our model generations are more coherent.¹

1 Introduction

Generating coherent long text (e.g., 1000s of tokens) is useful in myriad applications of creating reports, essays, and other long-form content. Yet the problem is particularly challenging as it demands models to capture global context, plan content, and produce local words in a consistent manner. Prior studies on “long” text generation have typically limited to outputs of 50-200 tokens (Shen et al., 2019; Bosselut et al., 2018; Zhao et al., 2020).

¹Code available at <https://github.com/tanyuqian/progressive-generation>

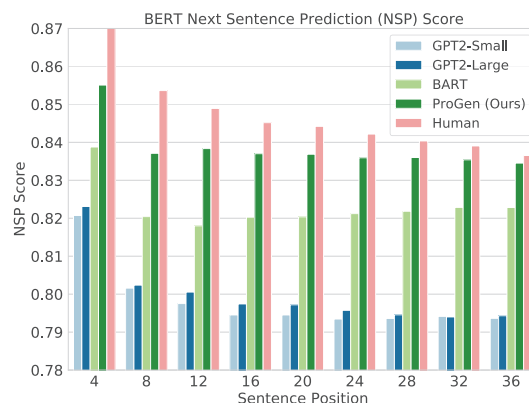


Figure 1: Results of large-scale LMs (GPT-2 and BART) fine-tuned on 10K stories. Coherence of text is evaluated by BERT next sentence prediction (NSP) score, where x-axis is the position of the evaluated sentences in the passage. There is a significant gap in coherence between text by human and text by large-scale LMs. Our proposed ProGen instead generates more coherent samples close to human text.

Recent large-scale pretrained language models (LMs), such as GPT-2 (Radford et al., 2019) and BART (Lewis et al., 2020), emerged as an impressive open-ended text generator capable of producing surprisingly fluent text. The massive LMs are typically pretrained on large corpora of generic text once, and then fine-tuned with small domain-specific data. The latest work has mostly focused on the regime of relatively short text with low hundreds of tokens. For example, Holtzman et al. (2020); See et al. (2019); Hua and Wang (2020) studied GPT-2 and BART generations with a maximum length ranging from 150 to 350 tokens. In this work, we study the problem of generating coherent, much longer passages of text (e.g., 1000 tokens). GPT-3 (Brown et al., 2020) was reported to produce long essays, yet the results seem to need extensive human curations (e.g., MarketMuse; Gardian), and the system is not publicly available to adapt to arbitrary desired domains.

In this work, we examine fine-tuning of large-scale LMs for domain-specific generation of extra-

long text. We find that samples produced by GPT-2 fine-tuned on small domain-specific corpora exhibit various imperfections, including excessive repetitiveness and incoherence between sentences far apart. Figure 1 measures the coherence of text generated by the fine-tuned GPT-2 w.r.t the BERT next sentence prediction (Devlin et al., 2019) score. As the figure shows, GPT-2 models (regardless of the model size) exhibit a significant gap in the score compared with human text, hence falling short in generating coherent text.

We hypothesize that the problem is mainly caused by the sequential generation order of the LMs, which makes global content planning of the passage difficult, especially when the generated text is long and contains thousands of words. One could potentially adopt the recent *planning-then-generation* or *non-monotonic* methods (Sec 2), yet those methods either require specialized neural architectures that need costly retraining for each domain (Gu et al., 2019; Stern et al., 2019; Chan et al., 2019; Fan et al., 2019), or rely on dedicated intermediate content plans (e.g., summaries, SRL labels) (Fan et al., 2019; Yao et al., 2019) with limited flexibility and producing sub-optimal results as shown in our experiments.

To overcome the limitations, we introduce a new method for **Progressive Generation** of Text (ProGen). We observe that generation of some words (e.g., stop words) does not require many contexts, while other words are decisive and have long-term impact on the whole content of the passage. Motivated by this observation, our approach first produces a sequence of most informative words, then progressively refines the sequence by adding finer-grained details in multiple stages, until completing a full passage. The generation at each stage is conditioning on the output of the preceding stage which provides anchors and steers the current generation (Figure 2). The intermediate words produced at each stage are defined based on a simple TF-IDF informativeness metric.

The approach enjoys several core advantages: **(1)** Although the progressive approach implements a conceptually non-monotonic generation process, generation at each stage can still be performed in a left-to-right manner and thus is directly compatible with the powerful pretrained monotonic LMs. The LMs at different stages are easily fine-tuned to accommodate a target domain using only small, independently constructed data. Intuitively, each LM

is addressing a sub-task of mapping a sequence to a finer-resolution one, which is much simpler than the overall task of mapping from conditions to full passages of text. In this work, we use BART (Lewis et al., 2020) for generation at each stage, though one can also plug in other off-the-shelf LMs. As seen from Figure 1, ProGen can generate more much coherent text compared with GPT-2 and nearly match human text in terms of the BERT-NSP score; **(2)** In contrast to the typical 2-stage planning-then-generation in prior work, the simple progressive strategy offers added flexibility for an arbitrary number of intermediate stages, yielding improved results; **(3)** The training data for each stage is extracted from domain corpus using the simple TF-IDF metric, without need of additional resources (e.g., pretrained summarization models) as in prior work, making the method broadly applicable to various domains and languages.

We conduct extensive empirical studies on the CNN News (Hermann et al., 2015) and Writing-Prompts (Fan et al., 2018) corpora, evaluating various systems by a wide-range of automatic metrics as well as human judgement. Results show that ProGen achieves strongly improved performance by decomposing the generation into more progressive stages. Our method produces diverse text passages of higher quality and coherence than a broad set of models, including fine-tuned GPT-2, BART, and other various planning-then-generation strategies.

2 Related Work

Content planning in generation. The idea of separate content planning and surface realization has been studied in early text generation systems (Reiter and Dale, 1997). Recent neural approaches have also adopted similar planning-then-generation strategies for data-to-text (Moryossef et al., 2019; Puduppully et al., 2019), storytelling (Fan et al., 2019; Yao et al., 2019; Xu et al., 2020), machine translation (Ford et al., 2018), and others (Hua and Wang, 2019; Yao et al., 2017). These models often involve customized architectures incompatible with the existing large LMs. Scaling those models for long text generation thus can require expensive training, which restricts systematic studies. On the other hand, it is possible to adopt some of the content planning strategies (e.g., summaries or SRL sequences as the plans (Fan et al., 2019)), and repurpose pretrained LMs for generation in each stage. However, these strategies

with dedicated intermediate plans and a pre-fixed number (typically 2) of stages can have limited flexibility, leading to sub-optimal results as shown in our empirical study. Besides, creating training data for planning requires additional resources (e.g., pretrained summarization models or SRL models) which are not always available (e.g., in certain domains or for low-resource languages). In contrast, we propose a simple way for designing the intermediate stages based on word informativeness, which can flexibly increase the number of stages for improved results, and easily create training data for all stages without additional models.

Non-monotonic generation and refinement.

Another relevant line of research is non-monotonic generation (Welleck et al., 2019; Gu et al., 2019; Stern et al., 2019; Chan et al., 2019; Zhang et al., 2020), infilling (Zhu et al., 2019; Shen et al., 2020; Qin et al., 2020), or refinement (Lee et al., 2018; Novak et al., 2016; Mansimov et al., 2019; Kasai et al., 2020) that differs from the restricted left-to-right generation in conventional LMs. Again, those approaches largely depend on specialized architectures and inference, making them difficult to be integrated with the powerful pretrained LMs. The prior studies have focused on generating short text. Our proposed coarse-to-fine progressive generation conceptually presents a non-monotonic process built upon the pretrained monotonic LMs, which permits fast adaptation to any target domain and generation of much longer text.

Long text generation. Previous work has made attempts to generate text of up to two or three hundred tokens. Those methods often adopt the similar idea of planning-then-generation as above (Shen et al., 2019; Zhao et al., 2020; Bosselut et al., 2018; See et al., 2019; Hua and Wang, 2020; Rashkin et al., 2020). Another line of work instead focuses on extending the transformer architecture (Vaswani et al., 2017) to model longer text sequences (e.g., Dai et al., 2019; Wang et al., 2020; Choromanski et al., 2021, etc). For example, Liu et al. (2018) used a hybrid retrieval-generation architecture for producing long summaries; Dai et al. (2019) showed long text samples qualitatively. Our work systematically examines the pretrained LMs in generating long domain-specific text, and proposes a new approach that empowers pretrained LMs for producing samples of significantly higher-quality.

3 Progressive Generation of Text

One of the main challenges in generating long coherent passages is modeling long-range dependencies across the entire sequences (e.g., 1000 tokens). We propose a progressive generation approach that is conceptually simple yet effective. Intuitively, progressive generation divides the complex problem of generating the full passage into a series of much easier steps of generating coarser-grained intermediate sequences. Contrary to generating everything from left to right from scratch, our progressive generation allows the model to first plan globally and then shift attention to increasingly finer details, which results in more coherent text. Figure 2 illustrates the generation process.

3.1 Generation Process

Let $\mathbf{y} := [y_1, y_2, \dots, y_T]$ be the output text, where each y_i is a token of language (a word or a sub-word). The output sequences are generated either conditionally on any other information \mathbf{x} (e.g., generations of a story given a prompt), or unconditionally (in which case we assume $\mathbf{x} \equiv \emptyset$ while keeping the same notation).

Instead of generating the full passage \mathbf{y} directly, we propose to add multiple intermediate stages: $\mathbf{x} \rightarrow \mathbf{c}_1 \rightarrow \mathbf{c}_2 \dots \rightarrow \mathbf{c}_K \rightarrow \mathbf{y}$, where for each stage $k \in \{1, \dots, K\}$, \mathbf{c}_k is an intermediate sequence containing information of the passage at certain granularity. For instance, at the first stage, \mathbf{c}_1 can be seen as a highest-level content plan consisting of the most informative tokens such as key entities. Then, based on the plan, we gradually refine them into subsequent \mathbf{c}_k , each of which contains finer-grained information than that of the preceding stage. At the final stage, we refine \mathbf{c}_K into the full passage by adding the least informative words (e.g., stop words). The generation process corresponds to a decomposition of the conditional probability as:

$$\mathbb{P}(\mathbf{y}, \{\mathbf{c}_k\} | \mathbf{x}) = \mathbb{P}(\mathbf{c}_1 | \mathbf{x}) \prod_{k=2}^K \mathbb{P}(\mathbf{c}_k | \mathbf{c}_{k-1}, \mathbf{x}) \mathbb{P}(\mathbf{y} | \mathbf{c}_K, \mathbf{x}). \quad (1)$$

As the above intuition, \mathbf{c}_k at early stages as the high-level content plans should contain informative or important words, to serve as skeletons for subsequent enrichment.

We next concretely define the order of generation, namely, which words should each stage generates. Specifically, we propose a simple method

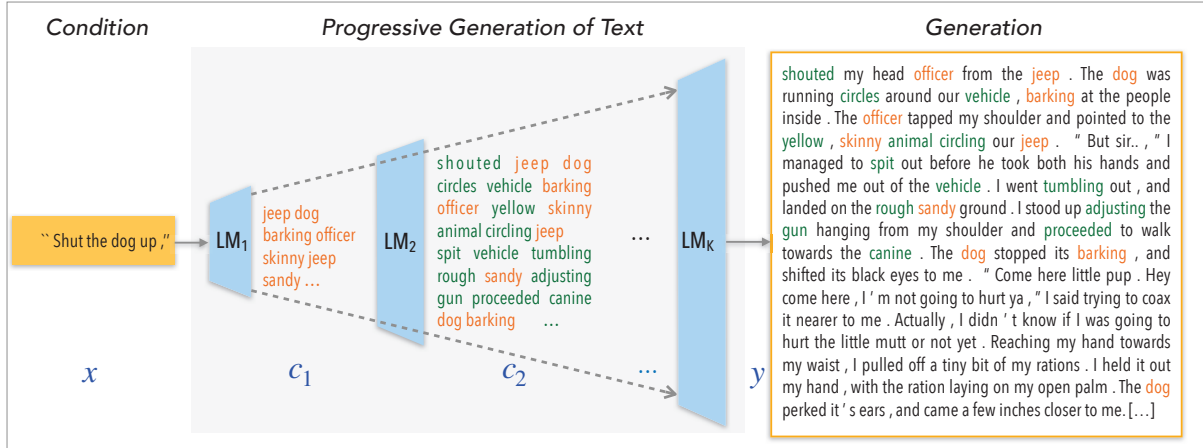


Figure 2: Progressive generation of long text y given any condition x . Each stage refines the results from the previous stage by adding finer-grained details. Added content at each stage is highlighted in different colors.

that constructs a vocabulary \mathcal{V}_k for each stage k , based on the *importance* of words in the target domain. Each particular stage k only produces tokens belonging to its vocabulary \mathcal{V}_k . By the progressive nature of the generation process, we have $\mathcal{V}_1 \subset \dots \subset \mathcal{V}_K \subset \mathcal{V}$. That is, \mathcal{V}_1 contains the smallest core set of words in the domain, and the vocabularies gradually expand at later stages until arriving the full vocabulary \mathcal{V} . Note that vocabularies in later stages are supersets of those in earlier stages. This allows the later stages to remedy and polish potential mistakes made in earlier stages when necessary. We discuss the construction of the vocabularies in the below.

Stage-wise vocabularies based on word importance. Given a text corpus \mathcal{D} of the target domain with the full vocabulary \mathcal{V} , we define the importance scores of words in \mathcal{V} based on the TF-IDF metric. We then rank all the words and assign the top V_k words to the intermediate vocabulary \mathcal{V}_k . Here V_k is a hyper-parameter controlling the size of \mathcal{V}_k .

More concretely, for each word $w \in \mathcal{V}$, we first compute its standard TF-IDF score (Salton and McGill, 1986) in each document $\mathbf{d} \in \mathcal{D}$, which essentially measures how important w is to \mathbf{d} . The importance of the word w in the domain is then defined as the average TF-IDF score across all documents containing w :

$$\text{importance}(w, \mathcal{D}) = \frac{\sum_{\mathbf{d} \in \mathcal{D}} \text{TF_IDF}(w, \mathbf{d})}{\text{DF}(w, \mathcal{D})}, \quad (2)$$

where $\text{TF_IDF}(w, \mathbf{d})$ is the TF-IDF score of word w in document \mathbf{d} ; and $\text{DF}(w, \mathcal{D})$ is the document

Algorithm 1 Training for Progressive Text Generation

Inputs:

Domain corpus \mathcal{D}
 Vocabulary sizes for K stages
 K pretrained LMs (e.g. GPT-2 or BART)

- 1: Construct stage-wise vocabularies $\{\mathcal{V}_k\}$ based on word importance Eq.(2)
- 2: Extract intermediate sequences $\{c_k^*\}$ using $\{\mathcal{V}_k\}$; add data noises (Sec 3.2)
- 3: Fine-tune all LMs independently (Sec 3.2)

Output: Fine-tuned LMs for generation at all stages in a progressive manner

frequency, *i.e.*, the number of documents in the corpus that contain the word w .

Pretrained language models as building blocks.

Compared to many of the previous planning-then-generation and non-monotonic generation methods, one of the key advantages of our progressive generation design is the direct compatibility with the powerful pretrained LMs that perform left-to-right generation. Specifically, although our approach implements a non-monotonic generation process that produces importance words first, we can generate intermediate sequences c_k at each stage still in a left-to-right manner. Thus, we can plug pretrained LM, such as GPT-2 or BART, into each stage to carry out the generation. As described more in section 3.2, for each stage k , we can conveniently construct stage-specific training data from the domain corpus \mathcal{D} using the stage-wise vocabulary \mathcal{V}_k , and fine-tune the stage- k LM in order to generate intermediate sequences at the stage that are pertaining to the target domain.

One can add masks on the pretrained LM's to-

ken distributions to ensure the stage- k LM only produces tokens belonging to \mathcal{V}_k . In practice, we found it is not necessary, as the pretrained LM can usually quickly learn the pattern through fine-tuning and generate appropriate tokens during inference. In our experiments we use BART for all stages, since BART is an encoder-decoder model which can conveniently take as inputs the resulting sequence from the preceding stage and generate new. (For the first stage in an unconditional generation task, we simply set $\mathbf{x} = \emptyset$.) We note that GPT-2, and other relevant pretrained LMs, can indeed also be used as a conditional generator (Radford et al., 2019; Liu et al., 2018) and thus be plugged into any of stages.

3.2 Training

Our approach permits straightforward training/fine-tuning of the (pretrained) LMs at different stages given the domain corpus \mathcal{D} . In particular, we can easily construct independent training data for each stage, and train all LMs in parallel. Note that no additional resources such as pretrained summarization or semantic role labeling models are requested as in previous work, making our approach directly applicable to a potentially broader set of domains and languages. We plan to explore the use of our method in multi-lingual setting in the future.

More concretely, for each stage k , we use the stage vocabularies \mathcal{V}_{k-1} and \mathcal{V}_k to filter all relevant tokens in the documents as training data. That is, given a document, we extract the sub-sequence \mathbf{c}_{k-1}^* of all tokens from the document that are belonging to \mathcal{V}_{k-1} , and similarly extract sub-sequence \mathbf{c}_k^* belonging to \mathcal{V}_k . The \mathbf{c}_{k-1}^* and \mathbf{c}_k^* are then used as the input and the ground-truth output, respectively, for training the LM at stage k with maximum likelihood learning. Therefore, given the stage-wise vocabularies $\{\mathcal{V}_k\}$, we can automatically extract training data from the domain corpus \mathcal{D} for different stages, and train the LMs separately.

In the multi-stage generation, the intermediate sequences are not natural language. Yet we found that fine-tuning pretrained LMs (such as BART and GPT-2) to generate the intermediate sequences is indeed very efficient in terms of data and computation. We tried training other models such as small sequence-to-sequence models and n-gram models from scratch, which we found is much harder, requiring more data, or yielding inferior performance.

This again highlights the importance of using pretrained LMs, as enabled by our simple method design.

Stage-level exposure bias and data noising. In the above training process, the outputs of each LM are conditioning on the ground-truth input sequences extracted from the real corpus. In contrast, at generation time, the LM takes as inputs the imperfect sequences produced at the previous stage, which can result in new mistakes in the outputs since the LM has never be exposed to noisy inputs during training. Thus, the discrepancy between training and generation can lead to mistakes in generation accumulating through the stages. The phenomenon resembles the *exposure bias* issue (Ranzato et al., 2016) of sequential generation models at token level, where the model is trained to predict the next token given the previous ground-truth tokens, while at generation time tokens generated by the model itself are instead used to make the next prediction.

To alleviate the issue and increase the robustness of each intermediate LM, we draw on the rich literature of addressing token-level exposure bias (Xie et al., 2017; Tan et al., 2019). Specifically, during training, we inject noise into the ground-truth inputs at each stage by randomly picking an n -gram ($n \in \{1, 2, 3, 4\}$) and replacing it with another randomly sampled n -gram. The data noising encourages the LMs to learn to recover from the mistakes in inputs, leading to a more robust system during generation.

4 Experiments

4.1 Setup

Domains. We evaluate on two text generation domains including: (1) **CNN News** (Hermann et al., 2015) for unconditional generation. (2) **Writing-Prompts** (Fan et al., 2018) for conditional story generation. The task is to generate a story given a prompt. The two datasets are chosen since they both contain long documents, with CNN’s average and maximum length being 512 and 926, and WritingPrompts’s being 437 and 942, respectively. To demonstrate the data efficiency of our approaches adapting to target domains, we sample 1,000 documents in each dataset for training.

Model configs. We use BARTs for all stages of generation. Due to computation limitations, we experiment models with 2, 3, 4-stages generations. In

our 2-stage model, our first stage covers about 25% of all content; in the 3-stage model, the first and second stages cover 15% and 25% of all content, respectively; and in the 4-stage model, our first three stages cover 15%, 20%, 25% of all content. For model training, we follow the same protocol as (See et al., 2019) to fine-tune all pretrained models until convergence. To combat exposure bias, we add noise to the training data as described in Sec 3.2, with the probability of replacing 1,2,3,4-grams 0.1/0.05/0.025/0.0125. In the generation phase, we use top-p decoding (Holtzman et al., 2020) with $p = 0.95$ to generate 1024 tokens at maximum. Experiments were conducted with RTX6000 GPUs. It took around 4 hours for model fine-tuning and generation with a single GPU.

Comparison methods. We compare with a wide range of baselines, categorized into two groups: (1) The large pretrained LMs including BART (Lewis et al., 2020) and GPT-2 in both small and large sizes (Radford et al., 2019). The LMs generate text in a standard left-to-right manner; (2) Progressive generation with various strategies adopted in the prior planning-then-generation work. Same as our proposed method, each stage adapts a pretrained BART for generation. Specifically, **Summary** first generates a short summary text as the content plan and conditioning on the summary produces the full passage of text (Fan et al., 2019). For training, summaries are obtained using the state-of-the-art pretrained CNN news summarization model based on BART; **Keyword** first generates a series of keywords, based on which the full text is generated in the next stage. Following (Yao et al., 2019), the keywords are extracted with the RAKE algorithm (Rose et al., 2010) for training; **SRL** follows the recent work (Fan et al., 2019) by first generating a sequence of predicates and arguments and then producing the full text conditionally. The same semantic role labeling tool as in the prior work is used here to create training data. **SRL+NER** and **SRL+Coref** further augment the SRL method by an additional stage of generating entity anonymized text conditioning on the predicates sequence prior to the final stage (Fan et al., 2019). SRL+NER uses an NER model to mask all entities, while SRL+Coref applies coreference resolution to mask all clusters of mentions. We use the same NER and coreference tools as in (Fan et al., 2019). Finally, as a reference, we also present the results of **Human**-written text (i.e., the text in the dev set).

4.2 Automatic Evaluation

4.2.1 Evaluation Metrics

To evaluate the generation quality for the domain-specific open-ended generation as studied here, we primarily measure the “closeness” between two sets of text, one generated by the model and the other the real text from the target domain. We evaluate with a broad array of automatic metrics, including *lexical-based quality* metrics and *semantic-based quality* metrics. We also evaluate the generation *diversity*.

MS-Jaccard (MSJ) is a lexical-based metric (Montahaei et al., 2019), where MSJ- n measures the similarity of n -grams frequencies between two sets of text with Jaccard index.

TF-IDF Distance (TID) is defined as the distance between the average TF-IDF features of two text sets. We use it as an additional lexical-based quality measure.

Fréchet BERT Distance (FBD) is a semantic-based metric (Montahaei et al., 2019) that measures the Fréchet Distance in the BERT feature space between the generated and real text. By using the BERT features from shallow (S), medium (M), and deep (D) layers, we can compute FBD-S/M/D, respectively.

Backward BLEU (B-BLEU) is a diversity metric (Shi et al., 2018) measuring how well the generated text covers n -grams occurred in the test set.

Harmonic BLEU (HA-BLEU) (Shi et al., 2018) is an aggregated quality and diversity metric that incorporates both the standard BLEU (i.e., precision) and the Backward BLEU (i.e., recall).

4.2.2 Results

Figures 3 and 4 show the results of the various systems on the news and story domains, respectively, measured with different metrics against test set. We give more complete results in the appendix. We can see that our progressive generation approach consistently outperforms the standard, single-stage LMs (GPT2-Small, GPT2-Large and BART) by a large margin on almost all metrics in both domains. Further, by increasing the number of progression stages, our method steadily achieves even stronger performance. This highlights the benefits of the flexible progressive generation strategy.

The various models using pretrained LMs with previous planning-then-generation strategies show

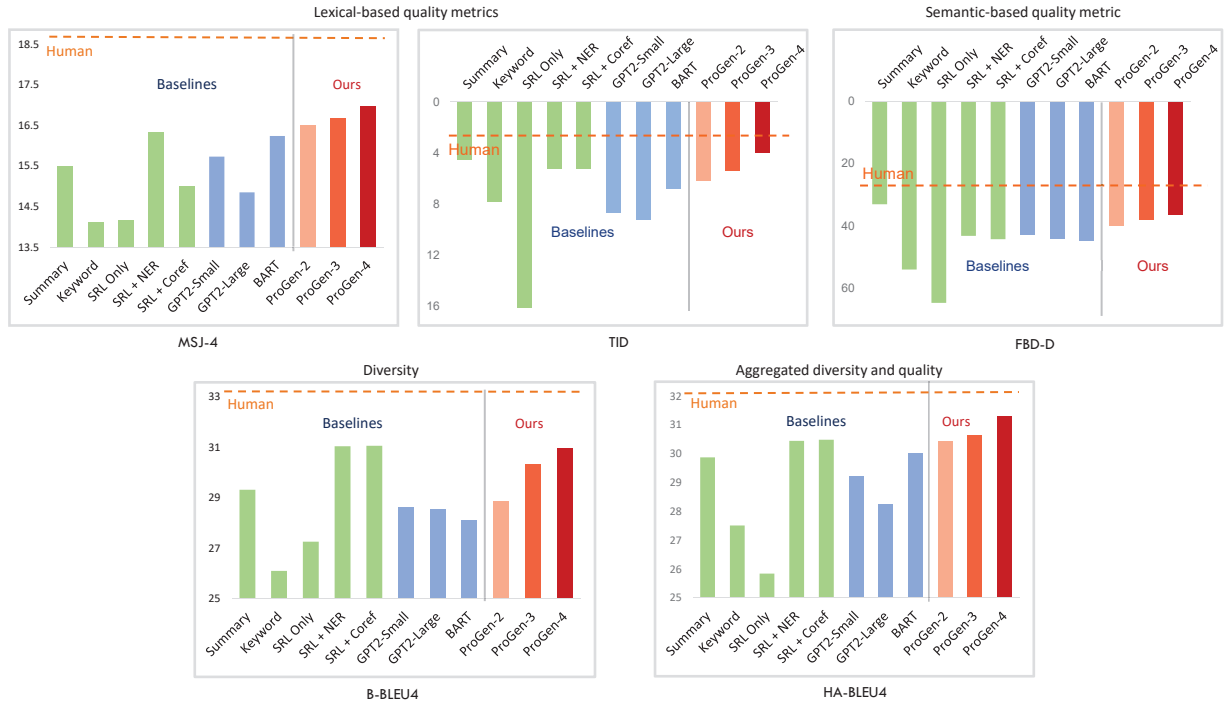


Figure 3: Results on the CNN News domain measured by different metrics. For TID and FBD, the lower value the better. More results (MSJ- n , B-BLEU n and HA-BLEU n with different n values, and FBD-S/M) are included in the appendix. The three sets of comparison methods are shown in different colors, with our ProGen in red, standard large LMs in blue, and the various models with previous planning strategies in green. Human results are shown as dashed lines, often indicating the best potential performance (except for the diversity related metrics).

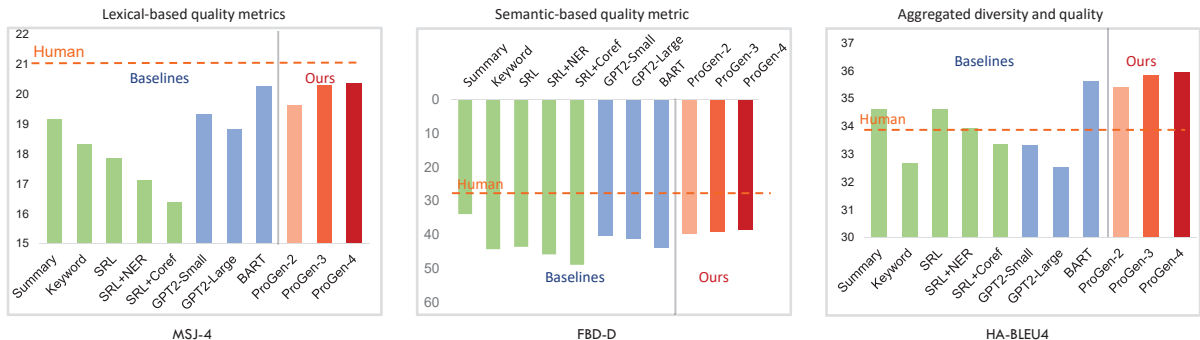


Figure 4: Results on the story domain measured by different metrics. More complete results are in appendix.

mixed results across the different metrics. For example, Summary achieves strong performance in terms of the semantic-based quality metric FBD-D (partially because the summaries are closer to the real text in the BERT feature space), but significantly falls behind other models in terms of diversity (B-BLEU4) and other quality metrics like MSJ and HA-BLEU. Similarly, the SRL-based methods give only mediocre results in terms of the semantic-based FBD-D. In contrast, our approach maintains a relatively consistent performance level. In particular, our 4-stage model, ProGen-4, is steadily among the best across all metrics, further validating

	Fluency	Coherence	
		passage	sentence (%)
BART	4.59	2.95	70.29
GPT2-Small	4.42	3.41	74.69
Summary	4.39	3.37	76.19
ProGen-4 (Ours)	4.46	3.83	86.22

Table 1: Human evaluation results on CNN.

the advantage of the proposed simple yet flexible multi-stage generation.

These results also indicate the necessity of using a large diverse set of automatic metrics for a comprehensive evaluation, and motivate human studies for further assessment.

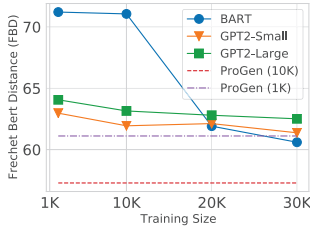


Figure 5: Sample efficiency on the story domain with the FBD metric (the lower, the better).

	FBD-D ↓	MSJ-4 ↑	HA-BL4 ↑
ProGen-2	39.94	16.50	30.45
-Noise	47.18	16.25	31.39
ProGen-3	38.30	16.68	30.64
-Noise	39.64	16.65	30.72
ProGen-4	36.49	16.96	31.32
-Noise	39.78	16.85	30.86

Table 2: Effect of noise on CNN.

	FBD-D ↓	TID ↓
ProGen-2	39.94	6.2
GoldPlan	30.16	3.5
Human	25.63	2.6

Table 3: GoldPlan Results on CNN.

First Stage: beckham \n liverpool bayern chelsea \n beckham chelsea mancini \n liverpool \n club [...]
Second Stage: beckham newcastle \n footballer acquisition transfer liverpool bayern munich chelsea munich \n beckham chelsea season italian siro gianluigi mancini season italy \n season liverpool season \n wembley club \n [...]
Final Stage:
(CNN) – England striker Ashley Beckham has joined Premier League strugglers Newcastle United. \n England Football Association chairman Stephen McClaren confirmed the acquisition of the 20-year-old England international on Wednesday, making the move official following a lucrative transfer agreement with Liverpool, Bayern Munich and Chelsea. \n Beckham spent the first half of this year on loan at Chelsea and did not start a single match during the season at the Italian capital’s San Siro despite being named the player of the season of the year by Italy. \n However, he played 17 times last season and 12 times this season for Liverpool, making 33 league appearances and scoring seven goals. \n He is currently third on the all-time England goalscoring list behind only England manager Alex Ferguson and new club teammate Paul Gascoigne. [...]

Table 4: An excerpt of a 3-stage generated example by ProGen-3 on the CNN News domain.

4.3 Human Evaluation

In our human study, we asked three university students who are proficient English speakers to evaluate the coherence and fluency of the generated text. To better assess the coherence of the long passages of text, we evaluate at both the passage level and the finer-grained sentence level. More concretely, for **passage-level coherence**, human raters assign a coherence score to each full-length text sample, on a 5-point Likert scale. For a more detailed assessment, we further evaluate **sentence-level coherence**, where human raters label each sentence in the text passage with 0 or 1, indicating whether the particular sentence is coherent with the preceding context in the passage. We then calculate the average percentage of coherent sentences in the generated text by each model. Human raters also evaluate the language quality for a **fluency** score on a 5-point Likert scale. We compare our method with the systems that show highest generation quality in automatic evaluation, including BART, GPT2-Small, and Summary. We evaluated 50 examples for each comparison model on the CNN domain. The Pearson correlation coefficient of human scores is 0.52, showing moderate inter-rater agreement.

Table 1 shows the results. All systems receive close fluency scores. Our approach obtained significantly higher coherence scores at both passage and sentence levels. In particular, over 86% sentences

in our model generations are considered as coherent with the context, improving over other models by at least 10 absolute percent.

4.4 Ablation Study and Analysis

Sample efficiency. We study how the progressive generation could improve the sample efficiency of large LMs fine-tuned to target domains. The intuition is that by focusing on the subsets of informative words, the early stages can more efficiently capture the domain-specific characteristics and then steer the subsequent refinement stages. Figure 5 shows the results where we report the FBD score averaged over FBD-S/M/D. We can see our approach can make more efficient use of the training data in learning to generate high quality samples. For example, with only 1K training examples, our method achieves comparable results with large LMs trained on 30K examples.

Generation with gold plans. To investigate the importance of dividing the generation process into stages and what the stages learn separately, we add another set of text into our comparison. It is a 2-stages model whose first stage is the ground truth (gold plan) while the second stage kept the same (a BART model), shown as GoldPlan in Table 3. Note that with gold plan, our model greatly decreases the gap with human text in terms of lexical (TID) and semantic (FBD-D) quality metrics. The results highlight the importance of plans in text

generation. The intermediate plans act as an information bottleneck, and high-quality plans could lead to high-quality text generation.

Effect of data noising. We study the ablation of data noising, to check whether the noising operation really helps reduce stage-wise exposure bias (Sec 3.2) as we expected. Table 2 shows the comparison between models with and without noise in training. The added noise generally brings performance improvement in terms of various metrics.

Example generations. Table 4 shows an example of text generated via three stages. We can see our model first generates the key subject *beckham* and the team name *liverpool* in the very first stage, then adds more fine-grained details like *acquisition*, *transfer* in the second stage and finally expands the keywords into a full document describing Beckham’s joining a new team.

5 Conclusion

We have proposed a new approach for domain-specific generation of long text passages in a progressive manner. Our method is simple and efficient by fine-tuning large-scale off-the-shelf language models. We conduct extensive experiments using a variety of metrics and human studies. We demonstrate that our method outperforms a wide range of large pretrained LMs with single-stage generation or prior planning-then-generation strategies, in terms of quality and coherence of the produced samples. The multi-stage generation also opens up new opportunities to enhance controllability of text generation, which we would love to explore in the future.

References

Antoine Bosselut, Asli Celikyilmaz, Xiaodong He, Jianfeng Gao, Po-Sen Huang, and Yejin Choi. 2018. Discourse-aware neural rewards for coherent text generation. In *NAACL*, pages 173–184.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *NeurIPS*, volume 33, pages 1877–1901.

William Chan, Nikita Kitaev, Kelvin Guu, Mitchell Stern, and Jakob Uszkoreit. 2019. KERMIT: Generative insertion-based modeling for sequences. *arXiv preprint arXiv:1906.01604*.

Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. 2021. Rethinking attention with performers. *ICLR*.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive language models beyond a fixed-length context. In *ACL*, pages 2978–2988.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pages 4171–4186.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *ACL*, pages 889–898.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *ACL*.

Nicolas Ford, Daniel Duckworth, Mohammad Norouzi, and George E Dahl. 2018. The importance of generation order in language modeling. In *EMNLP*.

Gardian. A robot wrote this entire article. are you scared yet, human?

Jiatao Gu, Qi Liu, and Kyunghyun Cho. 2019. Insertion-based decoding with automatically inferred generation order. *TACL*, 7:661–676.

Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *NeurIPS*, pages 1693–1701.

Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *ICLR*.

Xinyu Hua and Lu Wang. 2019. Sentence-level content planning and style specification for neural text generation. In *EMNLP*.

Xinyu Hua and Lu Wang. 2020. PAIR: Planning and iterative refinement in pre-trained transformers for long text generation. In *EMNLP*, pages 781–793.

Jungo Kasai, James Cross, Marjan Ghazvininejad, and Jiatao Gu. 2020. Non-autoregressive machine translation with disentangled context transformer. In *ICML*.

Jason Lee, Elman Mansimov, and Kyunghyun Cho. 2018. Deterministic non-autoregressive neural sequence modeling by iterative refinement. In *EMNLP*, pages 1173–1182.

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 *ACL*, pages 7871–7880.

- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. In *ICLR*.
- Elman Mansimov, Alex Wang, Sean Welleck, and Kyunghyun Cho. 2019. A generalized framework of sequence generation with application to undirected sequence models. *arXiv preprint arXiv:1905.12790*.
- MarketMuse. [Gpt-3 exposed: Behind the smoke and mirrors](#).
- Ehsan Montahaee, Danial Alihosseini, and Mahdiah Soleymani Baghshah. 2019. Jointly measuring diversity and quality in text generation models. *NAACL Workshop*.
- Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019. Step-by-step: Separating planning from realization in neural data-to-text generation. In *NAACL*.
- Roman Novak, Michael Auli, and David Grangier. 2016. Iterative refinement for machine translation. *arXiv preprint arXiv:1610.06602*.
- Ratish Puduppully, Li Dong, and Mirella Lapata. 2019. Data-to-text generation with content selection and planning. In *AAAI*, volume 33, pages 6908–6915.
- Lianhui Qin, Vered Shwartz, Peter West, Chandra Bhagavatula, Jena D. Hwang, Ronan Le Bras, Antoine Bosselut, and Yejin Choi. 2020. [Back to the future: Unsupervised backprop-based decoding for counterfactual and abductive commonsense reasoning](#). In *EMNLP*, pages 794–805. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. *ICLR*.
- Hannah Rashkin, Asli Celikyilmaz, Yejin Choi, and Jianfeng Gao. 2020. [PlotMachines: Outline-conditioned generation with dynamic plot state tracking](#). In *EMNLP*, pages 4274–4295.
- Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87.
- Stuart Rose, Dave Engel, Nick Cramer, and Wendy Cowley. 2010. Automatic keyword extraction from individual documents. *Text mining: applications and theory*, 1:1–20.
- Gerard Salton and Michael J McGill. 1986. Introduction to modern information retrieval.
- Abigail See, Aneesh Pappu, Rohun Saxena, Akhila Yerukola, and Christopher D. Manning. 2019. [Do massively pretrained language models make better storytellers?](#) In *CoNLL*, pages 843–861.
- Dinghan Shen, Asli Celikyilmaz, Yizhe Zhang, Liqun Chen, Xin Wang, Jianfeng Gao, and Lawrence Carin. 2019. Towards generating long and coherent text with multi-level latent variable models. In *ACL*, pages 2079–2089.
- Tianxiao Shen, Victor Quach, Regina Barzilay, and Tommi Jaakkola. 2020. [Blank language models](#). In *EMNLP*, pages 5186–5198.
- Zhan Shi, Xinchu Chen, Xipeng Qiu, and Xuanjing Huang. 2018. Toward diverse text generation with inverse reinforcement learning. *IJCAI*.
- Mitchell Stern, William Chan, Jamie Kiros, and Jakob Uszkoreit. 2019. [Insertion transformer: Flexible sequence generation via insertion operations](#). In *ICML*, volume 97, pages 5976–5985.
- Bowen Tan, Zhiting Hu, Zichao Yang, Ruslan Salakhutdinov, and Eric P Xing. 2019. Connecting the dots between mle and rl for sequence generation. *ICLR Workshop*.
- 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 *NeurIPS*, pages 5998–6008.
- Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. 2020. [Linformer: Self-attention with linear complexity](#). *arXiv preprint arXiv:2006.04768*.
- Sean Welleck, Kianté Brantley, Hal Daumé III, and Kyunghyun Cho. 2019. Non-monotonic sequential text generation. In *ICML*.
- Ziang Xie, Sida I Wang, Jiwei Li, Daniel Lévy, Aiming Nie, Dan Jurafsky, and Andrew Y Ng. 2017. Data noising as smoothing in neural network language models. In *ICLR*.
- Peng Xu, Mostofa Patwary, Mohammad Shoeybi, Raul Puri, Pascale Fung, Anima Anandkumar, and Bryan Catanzaro. 2020. [MEGATRON-CNTRL: Controllable story generation with external knowledge using large-scale language models](#). In *EMNLP*, pages 2831–2845.
- Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. [Plan-and-write: Towards better automatic storytelling](#). In *AAAI*, volume 33, pages 7378–7385.
- Lili Yao, Yaoyuan Zhang, Yansong Feng, Dongyan Zhao, and Rui Yan. 2017. Towards implicit content-introducing for generative short-text conversation systems. In *EMNLP*, pages 2190–2199.
- Yizhe Zhang, Guoyin Wang, Chunyuan Li, Zhe Gan, Chris Brockett, and Bill Dolan. 2020. [POINTER: Constrained progressive text generation via insertion-based generative pre-training](#). In *EMNLP*, pages 8649–8670.

Liang Zhao, Jingjing Xu, Junyang Lin, Yichang Zhang, Hongxia Yang, and Xu Sun. 2020. Graph-based multi-hop reasoning for long text generation. *arXiv preprint arXiv:2009.13282*.

Wanrong Zhu, Zhiting Hu, and Eric Xing. 2019. Text infilling. *arXiv preprint arXiv:1901.00158*.

Appendix: Complete Results

We include complete result numbers of experiments here.

	GPT2-S	GPT2-L	BART	Summ.	RAKE	SRL	SRL-N	SRL-C	ProGen-2	ProGen-3	ProGen-4	Dev
B-BL2	72.84	71.89	71.51	73.28	69.78	70.25	74.50	74.71	72.25	74.10	74.57	75.82
B-BL3	48.53	47.48	47.55	49.26	45.39	46.54	51.19	51.40	48.44	50.38	51.06	52.08
B-BL4	28.64	28.55	28.11	29.31	26.09	27.25	31.04	31.06	28.88	30.32	30.96	32.29
B-BL5	15.87	15.62	15.57	16.35	14.01	14.88	17.58	17.41	16.08	17.09	17.53	19.35
HA-BL2	73.61	71.97	74.56	74.59	71.63	67.47	74.51	75.11	74.64	75.17	75.86	75.72
HA-BL3	49.26	47.83	50.27	50.32	47.34	44.51	50.87	51.18	50.64	51.07	51.88	52.01
HA-BL4	29.21	28.26	30.03	29.88	27.51	25.84	30.45	30.49	30.45	30.64	31.32	32.28
HA-BL5	16.22	15.77	16.77	16.52	14.84	13.91	16.94	16.87	17.09	17.18	17.63	19.40
MSJ-2	49.24	46.94	49.85	46.97	44.19	43.85	49.39	44.37	49.46	50.16	51.00	54.51
MSJ-3	28.79	27.29	29.43	27.99	26.01	25.90	29.58	26.92	29.54	30.04	30.56	32.54
MSJ-4	15.73	14.85	16.24	15.48	14.12	14.15	16.33	14.99	16.50	16.68	16.96	18.60
MSJ-5	8.38	7.91	8.72	8.25	7.36	7.43	8.68	8.02	8.90	8.95	9.10	10.87
TID	8.7	9.2	6.8	4.5	7.8	16.1	5.2	5.2	6.2	5.4	4.0	2.6
FBD-S	16.21	18.50	7.76	2.93	4.17	14.26	11.42	4.66	3.26	3.16	2.64	5.98
FBD-M	24.92	29.61	22.49	15.00	25.92	37.24	22.63	20.28	19.05	18.84	17.38	12.26
FBD-D	43.07	44.15	44.86	33.08	54.12	64.83	43.26	44.34	39.94	38.30	36.49	25.63

Table 5: Complete results on the CNN News domain.

	GPT2-S	GPT2-L	BART	Summ.	RAKE	SRL	SRL-N	SRL-C	ProGet-2	ProGet-3	ProGet-4	Dev
B-BL2	78.38	77.43	76.96	77.19	76.97	77.98	77.90	77.62	78.64	78.73	78.41	79.20
B-BL3	55.51	54.18	54.45	54.45	53.86	55.67	55.49	55.09	56.44	56.50	56.25	56.02
B-BL4	33.41	32.20	33.02	32.88	31.95	33.83	33.75	33.36	34.46	34.62	34.52	34.08
B-BL5	17.59	16.79	17.55	17.53	16.47	17.93	17.98	17.63	18.32	18.49	18.57	18.40
HA-BL2	78.19	76.96	79.99	79.30	77.19	79.24	77.73	77.46	80.57	80.72	80.50	79.51
HA-BL3	55.39	54.33	57.86	56.83	54.71	57.00	55.71	55.14	58.11	58.38	58.35	56.39
HA-BL4	33.32	32.52	35.63	34.63	32.70	34.63	33.93	33.36	35.43	35.84	35.96	34.36
HA-BL5	17.46	16.94	19.16	18.47	16.86	18.26	18.03	17.60	18.72	19.14	19.30	18.55
MSJ-2	55.27	54.21	55.89	52.63	51.88	47.51	45.39	43.36	55.14	56.51	56.18	60.07
MSJ-3	34.48	33.70	35.46	33.46	32.59	30.88	29.51	28.22	34.81	35.80	35.74	37.42
MSJ-4	19.32	18.83	20.27	19.17	18.33	17.87	17.11	16.39	19.63	20.29	20.39	21.22
MSJ-5	10.16	9.90	10.73	10.27	9.57	9.54	9.21	8.82	10.16	10.60	10.76	11.34
TID	4.6	8.3	5.1	4.5	5.8	5.5	5.3	7.0	5.1	5.0	4.8	3.4
FBD-S	3.49	3.43	5.34	5.06	8.28	6.03	7.49	8.63	3.72	3.90	3.81	1.96
FBD-M	19.30	19.41	21.75	18.11	22.97	21.85	23.15	25.01	19.36	19.04	18.62	12.23
FBD-D	40.18	41.22	43.97	33.90	44.32	43.63	45.87	48.92	39.82	39.05	38.68	28.82

Table 6: Complete results on the story domain.