

---

# Progressive Generation of Long Text

---

Bowen Tan<sup>1</sup>, Zichao Yang<sup>1</sup>, Maruan Al-Shedivat<sup>1</sup>, Eric P. Xing<sup>1,2</sup>, Zhiting Hu<sup>1,2</sup>

<sup>1</sup>Carnegie Mellon University, <sup>2</sup>Petuum Inc.

## Abstract

Large-scale language models pretrained 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 (>1000 tokens), especially when the models are fine-tuned to the target domain on a small corpus. To overcome the limitation, 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 language models 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 GPT-2 in terms of domain-specific quality and sample efficiency. The coarse-to-fine nature of progressive generation also allows for a higher degree of control over the generated content<sup>1</sup>.

## 1 Introduction

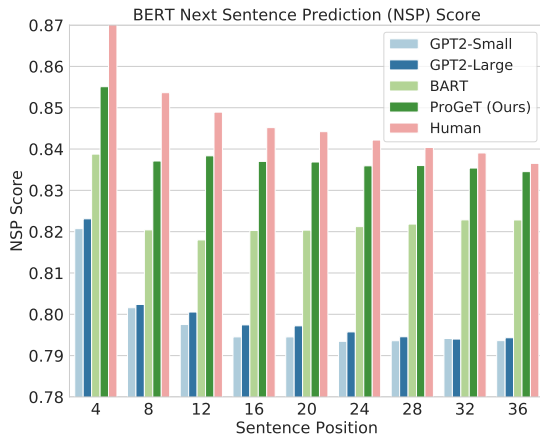
Large-scale pretrained language models (LMs) have enabled groundbreaking performance across a wide range of tasks [20, 5, 22, 34] and have become an essential component of the modern natural language processing toolkit. OpenAI’s GPT-2 [23], in particular, emerged as an impressive open-ended text generator capable of producing surprisingly fluent text.

Typically, models such as GPT-2 or BERT [5] are pretrained on large corpora of unlabeled text once, and then fine-tuned on a small amount of supervised data for downstream tasks, such as summarization, translation, and storytelling [16, 40, 27]. Previous work has mostly focused on the short text regime, where the target text ranges between tens to low hundreds of tokens. For example, Holtzman et al. [11], See et al. [27] studied and improved the quality of GPT-2 generations with a maximum length of 200 and 150, respectively. It remains unclear how the powerful text generators perform for generating coherent *long passages* of text (e.g., 1000s of tokens).

In this work, we examine fine-tuning of large-scale LMs for domain-specific generation of 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, right panel, measures the coherence of text generated by the fine-tuned GPT-2 w.r.t the BERT next sentence prediction [5] 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. An example generation by GPT-2 is shown in the right panel of the figure. We hypothesize that the problem is caused by the sequential generation order of the LM, which makes global planning of the content of the passage difficult, especially when the generated text is long and contains thousands of words.

---

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



**Prompt:** On your way to Mars on the first voyage ever , you discover that you and the countless others on the ship are intended ritual sacrifices by Elon Musk. You now must avoid reaching Mars at all costs .

**GPT-2 Generation:** See? The **sorcery** is a major aspect of human society that brings me great joy. Those behind them usually realise my profound respect for their form, [... 272 words omitted...]

But I couldn't **choose a destination**, like Mars, or anything else at all, and I knew that I had to do something. If I could risk it all I could. It was cold. It was alien. It was violent. [... 1926 words omitted...]

I **lost an arm**, but only after ten hours. If that hadn't tipped me off, what would that be? Would it feel like something bad? [... 287 words omitted ...]

Figure 1: Results of large-scale LMs (GPT-2 and BART) fine-tuned on 10K stories. **Left:** Coherence of text 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 ProGeT instead generates more coherent samples close to human text. **Right:** Example text generated by GPT-2. The text is not coherent by talking about irrelevant topics in different parts. The full example is shown in the appendix.

To overcome this limitation, inspired by progressive image generation [13, 19], we introduce a new method for **Progressive Generation of Text (ProGeT)**. 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).

Importantly, the simple design of progressive generation allows us to still make use of pretrained LMs. In this work, we use GPT-2 [23] or BART [14] for generation at each stage, though one can also plug in other off-the-shelf LMs. The LMs at different stages are easily fine-tuned to accommodate a target domain using independently constructed data. Intuitively, each LM is addressing a sub-task of mapping a sequence to a finer-resolution one (or mapping the condition to a short coarse-grained sequence at the first stage), which is much simpler than the overall task of mapping from conditions to full passages of text. As seen from Figure 1, ProGeT can generate more much coherent text compared with GPT-2 and nearly match human text in terms of the BERT-NSP score.

The coarse-to-fine progressive generation conceptually presents a non-monotonic (non-left-to-right) process. Compared to many of the recent non-monotonic or non-autoregressive generation methods [7, 35, 9, 29, 3] with customized neural architectures or dedicated inference procedures, our approach is simple and permits seamless use of powerful pretrained LMs. The new non-monotonic use of pretrained LMs enables long text generation in any target domain.

Through an extensive empirical evaluation on the CNN News [10] and WritingPrompts [6] corpora, we demonstrate that ProGeT is able to produce diverse text passages of higher quality than fine-tuned GPT-2 and BART models in the target domains, as measured by a wide range of metrics. ProGeT is also much more data efficient when adapting to the target domains.

## 2 Related Work

The idea of planning-then-generation is not new, and has been studied in early text generation systems for specific tasks. For example, conventional data-to-text generation systems [25] separate content planning and surface realization. Recent neural approaches have also adopted similar *two-step* generation strategy for data-to-text [18, 21], storytelling [7, 38], machine translation [8], and others [12, 37]. In a similar spirit, many efforts have been made in developing non-monotonic

or non-autoregressive generation models that differ from the restricted left-to-right generation in conventional LMs [35, 9, 29, 3, 39]. Our work differs in several dimensions: (1) differing from the previous work where the generated text ranges between tens to low hundreds of tokens, we study the generation of long text consisting of 1000 or more tokens; (2) our progressive process can be seen as a generalization of the two-step planning-and-generation approaches by allowing arbitrary multiple intermediate stages. We propose a simple yet effective way for designing the multiple stages that can be trained independently in parallel; (3) Compared to the previous work that usually involves customized model architectures or dedicated inference processes, our approach instead provides a framework for plugging in any common pretrained LMs, and thus combines the power of large-scale pretrained LMs with non-monotonic generation to enable coherent long text generation.

There has been a lot of recent interest in using large-scale pretrained LMs for improving domain-specific generation [16, 40, 27], in which generated text is typically much shorter than the samples considered in this work. Few recent work extended the transformer architecture [32] for generating long text, e.g., Liu et al. [15] used a hybrid retrieval-generation architecture for producing long summaries; Dai et al. [4] 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 Background

We start by reviewing the problem of text generation including language modeling and sequence-to-sequence (seq2seq) methods that serve as building blocks of our progressive generation approach. We also introduce the notations used in the rest of the paper.

**Text generation and language models.** Our discussion centers around methods for text generation using models that can produce outputs of the form  $\mathbf{y} := [y_1, y_2, \dots, y_T]$ , where each  $y_i$  is a token of language (a word or a sub-word). The output sequences are generated either conditionally on another sequence of tokens,  $\mathbf{x} := [x_1, x_2, \dots, x_S]$  (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). To generate sequences, we need to represent conditional (or unconditional) probability distributions of the form  $\mathbb{P}_\theta(\mathbf{y} \mid \mathbf{x})$  with learned  $\theta$ , typically factorized into a product of conditional distributions over tokens:

$$\mathbb{P}_\theta(\mathbf{y} \mid \mathbf{x}) := \prod_t \mathbb{P}_\theta(y_t \mid \mathbf{y}_{<t}, \mathbf{x}) \quad (1)$$

where  $\mathbf{y}_{<t} := [y_1, \dots, y_{t-1}]$  and the probability of each token is computed as  $\mathbb{P}_\theta(y_t \mid \mathbf{y}_{<t}, \mathbf{x}) := \text{softmax}(\mathbf{h}_{t-1} \mathbf{W})$ . Note that the ordering of tokens in Eq. (1) can be arbitrary, but most language models typically use a fixed left-to-right order. To train the model, we can maximize the (conditional) log-likelihood  $\sum_{(\mathbf{y}, \mathbf{x}) \in D} \log \mathbb{P}_\theta(\mathbf{y} \mid \mathbf{x})$  on the training data.

Language models are typically represented with recurrent neural architectures [2]. Specifically, modern large-scale models, including GPT-2 [22] used in our work, are based on the transformers [32]. The models are usually pretrained on large unsupervised corpora crawled from the web [23] and then fine-tuned to specific domains or tasks.

**Decoding.** To generate text, we need to be able to *decode* sequences of tokens from the distribution represented by a language model. The most commonly used approach in non-open-domain generation is *beam search* [e.g., 33] that produces sequences by approximately maximizing the likelihood under the learned model. It has been shown that maximization-based decoding in open-domain generation often leads to degenerate text samples (very unnatural and extremely repetitive text [11]), and thus sequential sampling approaches are typically used.

Due to heavy tail conditional distributions  $\mathbb{P}_\theta(y \mid \mathbf{y}_{<t}, \mathbf{x})$ , sequential sampling can be unstable (i.e., a single unlikely word may derail the rest of the sequence); hence in practice, conditional distributions are truncated before sampling either to the top- $k$  tokens [6, 23] or top- $p$  probability mass [11]. In this work, we use the latter, also called *nucleus sampling*, and sequentially sample  $y_t$ 's from  $\mathbb{P}_\theta(y \mid \mathbf{y}_{<t}, \mathbf{x})$  truncated to the top tokens that contain a total of  $p$  probability mass.

**Sequence-to-sequence models.** A common approach to conditioning language models on input sequences is the encoder-decoder architecture [30] supplemented with the attention mechanism [1]. The architecture represents  $\mathbb{P}_\theta(\mathbf{y} \mid \mathbf{x})$  as follows:

$$\mathbb{P}_\theta(\mathbf{y} \mid \mathbf{x}) := \mathbb{P}_\theta(\mathbf{y} \mid \mathbf{u} = f_\theta(\mathbf{x})), \quad (2)$$

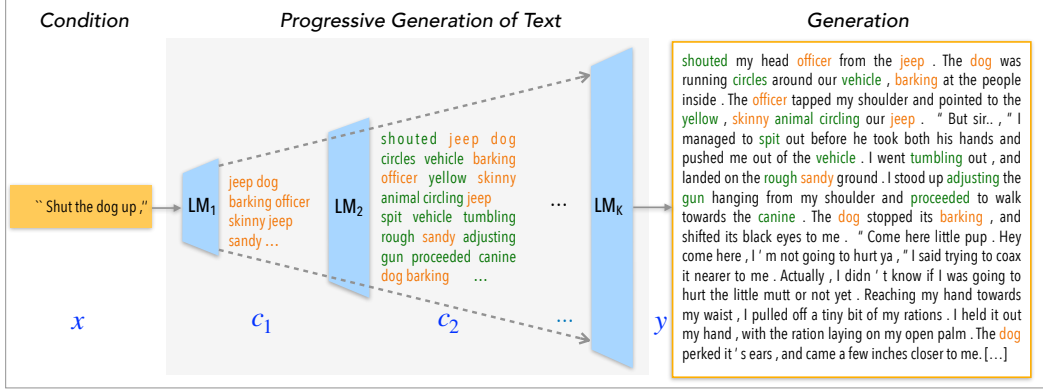


Figure 2: Progressive generation of long text  $y$  given any condition  $x$ . Each stage refines the results from the previous stage by adding finer-grained details. Added content at each stage is highlighted in different colors.

where  $f_\theta(\cdot)$  is a deterministic encoder that transforms input sequences into intermediate representations  $\mathbf{u}$  that are used by the language model  $\mathbb{P}_\theta(y | \mathbf{u})$  as additional inputs at decoding time. Lewis et al. [14] recently proposed BART that uses a BERT-like encoder and a GPT2-like decoder, and is particularly effective when fine-tuned for domain-specific generation. We selected to use BART as a component in ProGeT.

## 4 Progressive Generation of Text

One of the main challenges in generating long coherent passages is modeling long-range dependencies across the entire sequences (1000s of 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.

The key advantages of our approach include (1) enabling seamless use of powerful pretrained monotonic (left-to-right) LMs in the process of progressive generation which is conceptually non-monotonic (Sec 4.1); (2) Straightforward training of models at each stage in parallel with independent data from domain corpus (Sec 4.2).

### 4.1 Generation Process

Instead of generating the full passage  $y$  directly, we propose to add multiple intermediate stages:  $x \rightarrow c_1 \rightarrow c_2 \dots \rightarrow c_K \rightarrow y$ , where for each stage  $k \in \{1, \dots, K\}$ ,  $c_k$  is an intermediate sequence containing information of the passage at certain granularity. For instance, at the first stage,  $c_1$  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  $c_k$ , each of which contains finer-grained information than that of the preceding stage. At the final stage, we refine  $c_K$  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:

$$\mathbb{P}(y, \{c_k\} | x) = \mathbb{P}(c_1 | x) \prod_{k=2}^K \mathbb{P}(c_k | c_{k-1}, x) \mathbb{P}(y | c_K, x). \quad (3)$$

As the above intuition,  $c_k$  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 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  $\mathcal{V}$ . We discuss the construction of the vocabularies in the below.

---

**Algorithm 1** Training for Progressive Generation of Long Text

---

**input** Domain corpus  $\mathcal{D}$

Pretrained LMs (*e.g.* GPT-2 or BART) for  $K$  stages

Vocabulary sizes for  $K$  stages

- 1: Construct stage-wise vocabularies  $\{\mathcal{V}_k\}$  based on word importance Eq.(4)
- 2: Extract intermediate sequences  $\{\mathbf{c}_k^*\}$  using  $\{\mathcal{V}_k\}$ ; add data noises (Sec 4.2)
- 3: Fine-tune all LMs independently (Sec 4.2)

**output** Fine-tuned LMs for generation at all stages in a progressive manner

---

**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 [26] 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  $w$ :

$$\text{importance}(w, \mathcal{D}) = \frac{\sum_{\mathbf{d} \in \mathcal{D}} \text{TF\_IDF}(w, \mathbf{d})}{\text{DF}(w, \mathcal{D})}, \quad (4)$$

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

**Pretrained language models as building blocks.** A key advantage of our design of the progressive generation process is the 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  $\mathbf{c}_k$  at each stage 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 4.2, 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 token distributions to ensure the LM only produces tokens belonging to  $\mathcal{V}_k$ . In practice, we found it is not necessary, as the pretrained LM can usually quickly learns the pattern through fine-tuning and generate appropriate tokens during inference. In our experiments (section 5), we use BART for all  $k > 1$  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, we use GPT-2 for unconditional generation tasks, and BART for conditional generation. We note that GPT-2, and other relevant pretrained LMs, can indeed also be used as a conditional generator [23, 15] and thus be plugged into any of stages.

## 4.2 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.

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  $\mathcal{V}_k$  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.

**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 [24] 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 [36, 31]. 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.

## 5 Experiments

### 5.1 Setup

**Domains.** We evaluate on two text generation domains including: **(1) CNN News** [10] for unconditional generation. It includes about 90K documents and we randomly sample 5K for each of dev and test sets.; **(2) WritingPrompts** [6] for conditional story generation. The task is to generate a story given a prompt. This dataset has 272K/15K/15K examples in train/dev/test set.

**Model configs.** In the CNN News domain, we use GPT2-Large as the first-stage model and use BARTs for the rest of stages, and in the story domain, we use BARTs for all stages. Due to computation limitations, we experiment models with 2-stage and 3-stage generations. In our 2-stage model, our first stage covers about 30% of all content, and in our 3-stage model, the first and second stages cover 20% and 30% of all content respectively. For model training, we follow the same protocol as [27] to fine-tune all pretrained models until convergence. In the generation phase, we use the top-p decoding strategy proposed by [11] with  $p = 0.95$  to generate 1024 tokens at maximum. Refer to our supplementary for more details.

**Evaluation metrics.** We use combinations of automatic metrics and human evaluations to evaluate the domain-specific quality of text generated by GPT-2, BART and our progressive method.

- **MS-Jaccard (MSJ).** Proposed by [17], MS-Jaccard metric measures the similarity of the  $n$ -grams frequencies between two sets of texts with Jaccard index. For 2,3,4,5-grams, we report MSJ2,3,4,5 respectively. Under this metric, the higher, the better.
- **Fréchet BERT Distance (FBD)** Proposed by [17], Fréchet BERT Distance (FBD) is similar to the Fréchet Inception Distance (FID) metric in computer vision tasks, and it uses BERT as a feature provider. We compute the distance with features from every layer of BERT-large, and report the sums of layers 1-8, 9-16, 17-24 and call them FBD-S (shallow), FBD-M (medium), FBD-D (deep) respectively. Under this metric, the lower, the better.
- **TF-IDF Feature Distance (TID)** We propose a simple but intuitive metric. A traditional way to classify texts is based on the TF-IDF feature, so we use the distance between the average TF-IDF features of two passage sets to indicate their distance. Under this metric, the lower, the better.
- **Harmonic BLEU** Defined by [28], harmonic BLEU is the F1 score of Forward and Backward BLEUs. Here we use harmonic BLEU as a metric of domain-specific quality and backward BLEU as a metric of diversity.
- **Human Evaluation.** We also conduct human evaluation. Human annotators are asked to rate each generated sample in terms of coherence with 5-Likert scale.

### 5.2 Results

Table 1 shows the results of our progressive model on both CNN News and story domain measured in BLEU, MSJ, FBD and TID against test set. Our progressive models outperform baseline models on overwhelmingly most of the metrics, which indicates our models generate much more similar text to human. Moreover, the improvement is consistent both on unconditional generation (CNN) and conditional generation (WritingPrompts). It is worth noting that 3-steps model performs better than 2-stages model in the CNN News domain, while in the story generation domain, we don’t observe

	News (CNN)					Story (WritingPrompts)				
Metric	BART	GPT2-S	GPT2-L	ProGeT-2	ProGeT-3	BART	GPT2-S	GPT2-L	ProGeT-2	ProGeT-3
BLEU2 <sub>H</sub> ↑	84.99	84.57	84.92	85.21	<b>85.31</b>	87.00	86.49	86.25	<b>87.85</b>	87.56
BLEU3 <sub>H</sub> ↑	65.28	64.03	64.81	65.39	<b>65.71</b>	70.05	68.45	68.28	<b>70.75</b>	70.66
BLEU4 <sub>H</sub> ↑	44.58	42.92	43.89	44.54	<b>44.94</b>	49.59	47.12	47.15	49.94	<b>50.12</b>
BLEU5 <sub>H</sub> ↑	28.41	26.81	27.76	28.28	<b>28.63</b>	31.43	28.80	29.00	31.43	<b>31.75</b>
MSJ2 ↑	60.81	61.72	60.10	<b>61.73</b>	61.43	64.68	<b>67.03</b>	63.37	66.08	65.07
MSJ3 ↑	40.47	40.38	39.87	41.13	<b>41.13</b>	46.22	46.59	44.15	<b>47.20</b>	46.69
MSJ4 ↑	24.89	24.30	24.36	25.27	<b>25.38</b>	29.93	29.21	27.84	30.52	<b>30.37</b>
MSJ5 ↑	14.77	14.09	14.36	15.00	<b>15.11</b>	18.13	17.05	16.38	18.40	<b>18.42</b>
FBD-S ↓	6.71	6.64	6.55	6.53	<b>6.44</b>	3.80	2.92	3.50	<b>2.35</b>	2.78
FBD-M ↓	17.68	15.92	13.51	14.07	<b>12.97</b>	17.16	14.25	14.88	<b>13.40</b>	13.87
FBD-D ↓	36.16	34.25	26.94	27.69	<b>24.62</b>	36.02	30.59	31.05	<b>27.32</b>	28.31
TID ↓	7.49	8.20	7.05	6.38	<b>5.88</b>	3.09	3.08	3.92	<b>2.33</b>	2.37

Table 1: Comparison of our models with baseline models. ProGeT-2 and ProGeT-3 denote progressive generation with two and three steps respectively. BLEU<sub>H</sub> measures domain specific quality, MSJ measures n-gram similarity, FSD measures semantic similarity and TID measures TD-IDF distance. Our methods outperform all other baselines across the board.

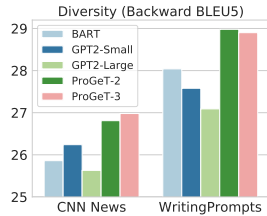


Figure 3: Diversity of text measured by Backward BLEU

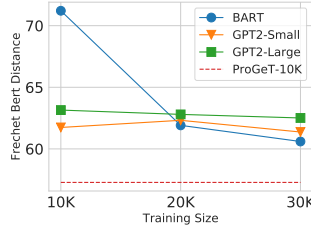


Figure 4: Sample efficiency study on the story domain over the FBD metric (the lower, the better).

	MSJ2	MSJ3	MSJ4	MSJ5
ProGeT-2	<b>66.08</b>	<b>47.20</b>	<b>30.52</b>	<b>18.40</b>
-Noise	65.06	46.67	30.35	18.39
ProGeT-3	<b>65.07</b>	<b>46.69</b>	<b>30.37</b>	<b>18.42</b>
-Noise	64.25	46.15	30.02	18.22

Table 2: Effect of noise on the story domain.

significantly improved quality by using 3-steps model. This suggests that the influence of the number of progressive steps might be different across different domains.

Besides those metrics, we also measure diversity—an aspect in which we usually observe language model suffers—using backward BLEU in Figure 3. Our models can generate much more diverse text compared with GPT, which could be attributed to global planning with progressive generation.

Besides automatic evaluation, we ask human annotators to evaluate the text generated by different models according to our score standard introduced previously. The scores on text samples are averaged and shown in Table 3 with standard deviations. First, we can see that all models get scores higher than 3, which indicates most of the generated texts from baselines are at least moderately coherence. However, all baseline models have scores of around 3.5 while our progressive model gets a significantly higher score of 4.43, suggesting our samples are much harder to distinguish with human text.

In Figure 4, we study the sample efficiency of different models by varying the amount of training data. We can see that the advantage of our model is more prominent when the data size is decreased. Our model can make more efficient use the training data in learning to generate high quality samples.

### 5.3 Ablation Study

**Effect of Data Noising** We try the ablation of data noising, to check whether the noise operation really helps reduce stage-wise exposure bias (Sec 4.2) as we expected. Table 2 shows the comparison between models with and without noise in training. In both 2- and 3-steps generation, noise brings significant performance improvement onto our progressive model.

**Generation with Gold Plan** To investigate the importance of dividing the generation process into stages and what the stages learn separately, we add another set of texts into our comparison. It is a 2-stages model whose first stage is the ground truth (gold plan) while the second stage kept the

	BART	GPT2-S	GPT2-L	Ours
Human	3.17 ± 0.99	3.61 ± 1.26	3.51 ± 1.43	<b>4.43 ± 0.48</b>

Table 3: Human evaluation results on CNN.

	MSJ4	FBD-D
ProGeT-2	25.27	27.69
GoldPlan	27.37	16.40
Human (Dev Set)	28.34	11.64

Table 4: GoldPlan Results on CNN.

<b>First Stage:</b> beckham \n liverpool bayern chelsea \n beckham chelsea mancini \n liverpool \n club [...]
<b>Second Stage:</b> 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 [...]
<b>Final Stage:</b> (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 5: A 3-stage generated example of ProGeT-3 on the CNN News domain.

same (a BART model), shown as GoldPlan in Table 4. Note that with gold plan, our model greatly decreases the gap with human text in MSJ and FBD metrics, which amplifies the importance of plans in text generation. The intermediate plans are information bottleneck and high quality plans indicate high quality text generation.

## 5.4 Example generations

Table 5 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.

Our content-first and coarse-to-fine nature of our progressive approach allows for a higher degree over the generated content. Here we show two easy ways to implement basic content controlling and the examples are provided in supplementary materials due to space limit.

**Content manipulation via editing intermediate sequences.** The first manner is to manipulate the content with an existing intermediate sequence. Specifically, we can replace some words in the input sequences. In this case, the content of generated text can be desired manipulated with other content being similar as before, together with coherence and rationality.

**Controlling the content via prefixes.** The second manner to control content is writing a prefix for the sketch and using the first-step model to complete it, then the content should be closely related to the prefix words you write.

## 6 Conclusion

In this work we have proposed a new approach for domain-specific generation of long text passages in a progressive manner. Our method is simple and efficient and is fully based on fine-tuning large-scale off-the-shelf language and sequence-to-sequence models. We demonstrate that our method outperforms GPT-2—one of the most recent and competitive state-of-the-art story generation methods—in terms domain-specific quality and diversity of the produced samples as measured by a variety of metrics.



## References

- [1] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [2] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.
- [3] W. Chan, N. Kitaev, K. Guu, M. Stern, and J. Uszkoreit. Kermit: Generative insertion-based modeling for sequences. *arXiv preprint arXiv:1906.01604*, 2019.
- [4] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [6] A. Fan, M. Lewis, and Y. Dauphin. Hierarchical neural story generation. *arXiv preprint arXiv:1805.04833*, 2018.
- [7] A. Fan, M. Lewis, and Y. Dauphin. Strategies for structuring story generation. In *ACL*, 2019.
- [8] N. Ford, D. Duckworth, M. Norouzi, and G. E. Dahl. The importance of generation order in language modeling. In *EMNLP*, 2018.
- [9] J. Gu, Q. Liu, and K. Cho. Insertion-based decoding with automatically inferred generation order. *Transactions of the Association for Computational Linguistics*, 7:661–676, 2019.
- [10] K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, pages 1693–1701, 2015.
- [11] A. Holtzman, J. Buys, M. Forbes, and Y. Choi. The curious case of neural text degeneration. In *ICLR*, 2020.
- [12] X. Hua and L. Wang. Sentence-level content planning and style specification for neural text generation. In *EMNLP*, 2019.
- [13] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [14] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [15] P. J. Liu, M. Saleh, E. Pot, B. Goodrich, R. Sepassi, L. Kaiser, and N. Shazeer. Generating wikipedia by summarizing long sequences. In *ICLR*, 2018.
- [16] H. H. Mao, B. P. Majumder, J. McAuley, and G. W. Cottrell. Improving neural story generation by targeted common sense grounding. *arXiv preprint arXiv:1908.09451*, 2019.
- [17] E. Montahaei, D. Alihosseini, and M. S. Baghshah. Jointly measuring diversity and quality in text generation models. *arXiv preprint arXiv:1904.03971*, 2019.
- [18] A. Moryossef, Y. Goldberg, and I. Dagan. Step-by-step: Separating planning from realization in neural data-to-text generation. In *NAACL*, 2019.
- [19] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019.
- [20] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.
- [21] R. Puduppully, L. Dong, and M. Lapata. Data-to-text generation with content selection and planning. In *AAAI*, volume 33, pages 6908–6915, 2019.
- [22] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. Improving language understanding by generative pre-training. URL [https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language\\_understanding\\_paper.pdf](https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf), 2018.
- [23] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9, 2019.

- [24] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732*, 2015.
- [25] E. Reiter and R. Dale. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87, 1997.
- [26] G. Salton and M. J. McGill. Introduction to modern information retrieval. 1986.
- [27] A. See, A. Pappu, R. Saxena, A. Yerukola, and C. D. Manning. Do massively pretrained language models make better storytellers? *arXiv preprint arXiv:1909.10705*, 2019.
- [28] Z. Shi, X. Chen, X. Qiu, and X. Huang. Toward diverse text generation with inverse reinforcement learning. *arXiv preprint arXiv:1804.11258*, 2018.
- [29] M. Stern, W. Chan, J. Kiros, and J. Uszkoreit. Insertion transformer: Flexible sequence generation via insertion operations. *arXiv preprint arXiv:1902.03249*, 2019.
- [30] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- [31] B. Tan, Z. Hu, Z. Yang, R. Salakhutdinov, and E. P. Xing. Connecting the dots between mle and rl for sequence generation. 2018.
- [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [33] A. K. Vijayakumar, M. Cogswell, R. R. Selvaraju, Q. Sun, S. Lee, D. Crandall, and D. Batra. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424*, 2016.
- [34] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- [35] S. Welleck, K. Brantley, H. Daumé III, and K. Cho. Non-monotonic sequential text generation. In *ICML*, 2019.
- [36] Z. Xie, S. I. Wang, J. Li, D. Lévy, A. Nie, D. Jurafsky, and A. Y. Ng. Data noising as smoothing in neural network language models. In *ICLR*, 2017.
- [37] L. Yao, Y. Zhang, Y. Feng, D. Zhao, and R. Yan. Towards implicit content-introducing for generative short-text conversation systems. In *EMNLP*, pages 2190–2199, 2017.
- [38] L. Yao, N. Peng, R. Weischedel, K. Knight, D. Zhao, and R. Yan. Plan-and-write: Towards better automatic storytelling. In *AAAI*, volume 33, pages 7378–7385, 2019.
- [39] Y. Zhang, G. Wang, C. Li, Z. Gan, C. Brockett, and B. Dolan. Pointer: Constrained text generation via insertion-based generative pre-training. *arXiv preprint arXiv:2005.00558*, 2020.
- [40] Z. M. Ziegler, L. Melas-Kyriazi, S. Gehrmann, and A. M. Rush. Encoder-agnostic adaptation for conditional language generation. *arXiv preprint arXiv:1908.06938*, 2019.

## A Effects of Generation Order

To check the effect of vocabulary picking for intermediate stages. We compare our model with 2 different word selection strategies: one is randomly picking a set of words, and another is picking words following the reversed order of our model. To compare the 3 strategies, we keep the proportion of content taken by the intermediate stage the same. As Table 6 shows, our model performs best, while random progressive order gets comparable results with baseline GPT2 and using reversed order would even worsen the performance.

Order	Frechet BERT Distance ↓
GPT2-Small	62.05
Random	61.81
Reversed	69.28
Ours	<b>57.30</b>

Table 6: Performance of different orders, measured by Fréchet BERT Distance (sum over all BERT layers) using 1000 generated examples, the lower, the better

## B Controllable Generation

### B.1 Content manipulation via editing intermediate sequences

Inter. Seq. 1	zimbabwe currency bank confronting zimbabwe currency currency pillar said governor reserve bank zimbabwe
Generation 1	Zimbabwe 's president resigned Friday, just days after his government launched a program to bring end to the hyperinflation rate and introduce a stable currency after the government's central bank warned that the situation was getting worse. "I think the most critical issue confronting Zimbabwe is currency – the real currency that is the pillar of our economy," said Reserve Bank of Zimbabwe President Gideon Gono.
Change the template: Replace "Zimbabwe" → "Britain"	
Inter. Seq. 2	Britain currency bank confronting Britain currency currency pillar said governor reserve bank Britain
Generation 2	Britain 's pound has fallen below \$1 for the first time since 1990, the British government's central bank chief warned Friday. "We are confronting a very serious economic challenge in which British businesses and people depend upon the currency as a fundamental force and a key payment mechanism, and where the currency has to be a pillar of that economy," said Jim O'Neill, Governor of the Bank of England, "Britain as a country is at a crossroad."

Figure 5: Example of Content Manipulation. "Inter. Seq." stands for "Intermediate Sequence".

As Figure 5 shows, we have "Intermediate Sequence 1" which is extracted from an existing document. Then we get another "Intermediate Sequence 2" by replacing the word "Zimbabwe" with "Britain". Comparing the generations from these two intermediate sequences, we can find that the content in generated documents are coherent to the intermediate sequence, e.g., when "Zimbabwe" is in the intermediate sequence, there are coherent content like "hyperinflation" in the generation, and when "Britain" is in the intermediate sequence, our model generates "pound", "England" etc.

### B.2 Controlling the content via prefixes

As Table 7 shows, with a human-written prefix "obama jobs", our model first generates a full intermediate sequence and then generates a whole document coherent to the prefix, progressively.

## C Additional Evaluation Metrics

**BERT NSP Score** We leverage a pretrain task of BERT, Next Sentence Prediction, as a metric to measure the coherence between every sentence and its next sentence.

**Human-Written Prefix:** obama jobs

**Generated Intermediate Sequence:** obama jobs \nobama unemployment jobs marian said \n [...]

**Generated Document:** President Barack Obama is calling for businesses and businesses alike to help create jobs to help Americans. \n " Obama said the unemployment rate is rising. "Jobs have great benefits and contribute to people’s well-being. You could see people walking down the street, smiling, saying how great it was," association President and CEO Marian O’Neill said. \n ...

Table 7: Examples of Controlling sketch prefix. The matching between the sketch and generated text is marked in red. The matching between the controlling prefix and the sketch is underlined.

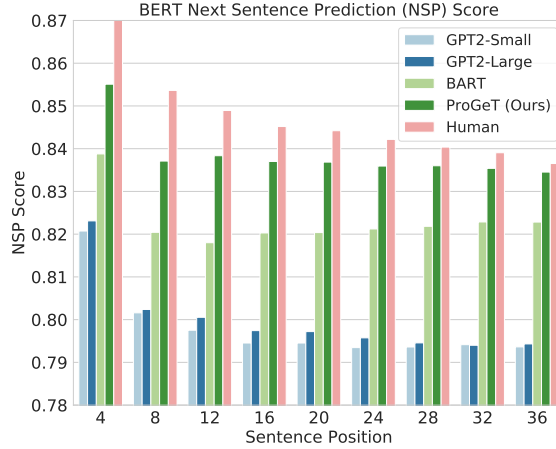


Figure 6: Bert NSP Scores.

**N-gram Similarity** We tried to measure the pattern of n-gram repetition as a metric. Specifically, we calculate the size of n-gram overlap between the first half and the second half for every prefix of a document and make a curve, as Figure 7 shows. This simple metric is not able to distinguish baselines and human texts.

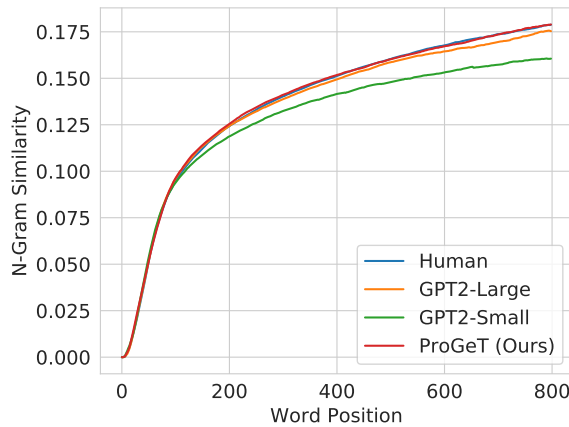


Figure 7: N-Gram Similarity Curve.

**Sentence Repetition** We measure repetition on the sentence level by counting exact same sentences in generated texts. Specifically, for every sentence position, we check whether it is a duplication of a previous sentence. As Figure 8 shows, this simple metric is also not able to distinguish baselines and human texts.

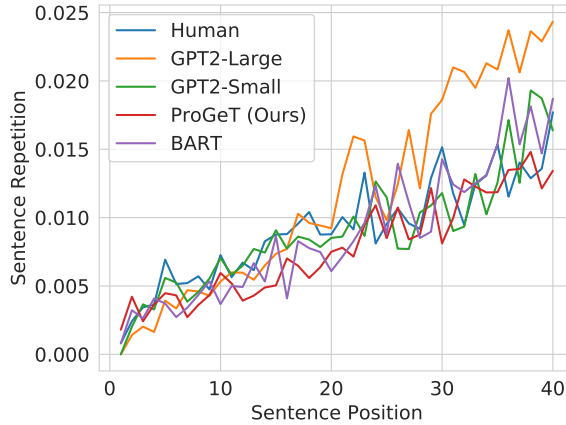


Figure 8: Sentence Repetition

## D Model Configurations

For the size of vocabularies of intermediate steps, we first select a vocabulary size which covers about 30% of content for the 2-step model, and then add another vocabulary size which covers about 20% of content for the 3-step model. In the CNN News domain, we use GPT2-Large as the first-step model and use BARTs for the rest of steps, and in the story domain, we use BARTs for all steps. We add noise in progressive steps by replacing N-grams with same number of random words. In both domains, the probability of replacing 1,2,3,4-grams are 0.1/0.05/0.025/0.0125. For model training, we follow the same setting of finetuning pretrain models as [27] to all of our pretrain models. Specifically, using default optimization hyperparameters provided by HuggingFace repository, we finetune all models a same number of steps to convergence. In the generation phase, we use the top-p decoding strategy with  $p = 0.95$ .

## E Generation Examples of ProGeT

Table 8, 9, and 10 show generated documents of our progressive model. We can see the generated documents are long but their contents are coherent from the beginning to the end.

---

(CNN) – Europe’s largest utility, Uniper, will seek to raise prices again this year by putting solar power units up for bid via a recently unveiled regional electricity utility plan.

Uniper already has received three regional announcements from European governments promising to double its rate structure by 2020, competition for low-cost solar power and the government supporting the project have forced Uniper Energy to raise rates for its second round of new utility rates in the region in three years this month.

Europe’s renewable producers, meanwhile, say U.S. gas prices and electricity prices threaten to continue to drive home costs for energy consumers.

"The number two factor pushing on consumers here is the free fall in costs – it’s to the point where consumers will say ‘OK, I have to hike that again.’ Uniper is taking a tougher stance. You already know that there are a number of factors: high fuel costs on household bills and the weak economy. But in the meantime, U.S. dollar moves higher; there’s not enough demand. And diesel prices are low as well," said Mike Gilbert, Uniper’s operations director.

Uniper Energy officials say the company is seeking to raise – at least – the price of solar power on non-renewable electricity products and to provide them more competitively at home with solar panels.

Uniper has announced new rates in the lower (15%) and upper (20%) price bands for solar panels and solar installations, and expects to combine sales of solar panels and rooftop solar with more coal and natural gas power from its various suppliers this spring.

Uniper’s survey, released last year, was part of a larger report by Uniper CEO Frank Fiume, found that over half of European consumers are not aware of the escalating costs for energy, according to June and July data released by the European Union.

Uniper concluded that energy prices that once seemed like a given had been raised again in March by more than 15

Although Uniper said the government will reappraise its proposals after the year, the result of the report won’t be known until July of this year.

Uniper surveyed 1,583 users in 11 European nations, public utility companies and cooperative electricity producers across the world, a year before it was due to publish the April/May results of a new report this year.

"Uniper’s Europe survey found that people in the three member states that are more exposed to the more expensive gas market are now more likely to use more expensive electricity," said Fiume. "This reflects what consumers have told us, that the more expensive gas prices mean for them, because they are more likely to stick with their current business models. These are their economic economics. This means a lot of change."

According to Uniper, the average European user spent more than \$45,000 (\$48,000) in 2012 on electricity, and an average of over \$60,000 (\$62,800) on electricity in the first two quarters of 2013.

"That’s double what we’ve seen so far this year," Fiume said. "The more the average customer pays for their energy, the more they have to cut costs. It’s getting rid of their unnecessary carbon. It’s affecting them."

Uniper’s survey found that in all countries surveyed, 71 percent of the 3,106 respondents have at least tried to cover the cost of renewable energy. Some 1.1 million users in Europe had both solar panels and coal or natural gas on their homes, and more than 34,000 residents should switch to solar power on their homes before the end of the year.

---

Table 8: A generated example of ProGeT on the CNN News domain.

---

(CNN) – A tornado shower ripped through North Fulton County in Mississippi during a freak thunderstorm Friday afternoon.

In those 20 minutes, a tornado formed in South Fulton County and ripped through the city.

"I heard the tornado. Then I was sitting there and I started to shake. I didn't think it was real, but I thought it was definitely a tornado. And then it came."

The storm hit the southeast side of Fulton County with 90 mph winds at 60 mph," the National Weather Service in Columbia, Mississippi, said in a statement. "One minute, 15 seconds, 15 seconds the size of a baseball field. That's how much hail was falling on one town, even over the course of a country town."

"The winds continued to build during this time," the weather service said. "Little by little, with every pulse of the tornado making a small shake, the storm moved through the county, eliminating damage on three roads."

Gonda Fleck said the storm passed over Hill County (population 1,369). "It was within a mile of the home" she said. "We didn't see any indication."

Fleck had no doubt about the tornado. "It felt like, right now, we're standing before a tornado. I said no way, not even trying to imagine it," she said. "It felt like a tornado coming at us. I couldn't believe it."

"Within a minute or two, it was – the infernos kept coming," spokeswoman Dena Mendenhall said. "Within an hour, the tornado was gone."

"We lost the radar. We lost all the weather cameras, and no one got an inch of radar," she said.

"It did not take long for us to realize, at least, that there was going to be a tornado in here. But, as we continued to make our way through the area, the winds picked up a little bit and made things not seem so bad," she said.

In one stop, the "Dip" method of getting traction worked. To make an impact, the "Dip" tornado tracked in all directions and started to push as it straightened out.

These moments came quickly after the storm progressed.

First the hail.

Fennell said about 12 minutes later, about 2.5 minutes later, there were "well over 3 inches of hail, which is six inches horizontally," in places.

Another minute passed with ice forming on windshields.

After the first 2.5 minutes of hail, the increased hail got so bad that one fellow attended school in Fulton County, which was evacuated, and covered in 3-4 inches of snow.

"You can just tell the ground is moving," Feck said.

"The ice gets a little faster."

Later, the tornado tracked out of the city toward Florence.

Forecasters with the National Weather Service in Fargo, North Dakota, said tornado sightings have been common throughout the weekend.

Irs reporting areas included Champaign, Fontana, McHenry, Lewisville and Carlisle.

"It appears there was a tornado with an embedded ice storm, a tornado with a tornado, and a tornado followed by a tornado," Faulk said.

A tornado warning of a EF0, or enhanced tornado threat, was issued for Taylorsville.

Late Friday evening, a potentially deadly storm developed over a 10 mile radius of Columbia, Missouri, according to the weather service.

---

Table 9: A generated example of ProGeT on the CNN News domain.

---

(CNN)— It wasn't long ago we heard of jet fuel jumbo jets and seen the cry of "TOYOJEE!" as engines descended from the sky. Then, we heard today of the world's biggest jet engine and what it's going to cost. Not just \$6 billion, but \$700 billion. As of June 3, 1958," according to a US trade report.

These days, we hear of even smaller, rotary-rotating rockets and laser-guided missiles that burn through the air at a rate of up to 2,000 miles per hour (4,000 kilometers per hour). But down under, we get slow, narrow engines that take up a great deal of fuel. Airbus, for example, has been manufacturing small turbocharged engines for decades. They have ballooned in size at an annual rate of 5,000 pounds (2,500 kilograms). "This is the largest possible increment that you can have in a human engine using engine type design. Sixteen pounds, a full cubic inch," said Thomas Pilcock, president of Alkera Inc., a company that manufactures compact turbofan engines. "It's the biggest business in the world. Everything is moving. We need more fuel, not less fuel. ... We can't afford to go over our top targets. We have to put together an amazing mini fleet."

It's estimated that around 600,000 airplanes and submarines use jet engines each year, according to figures from USA Aerospace that were released last year.

Nissan, a Japanese company, has been building more powerful engines and available in the United States since 2001, while the United States Aerospace and Defense Administration is working on other solutions.

The \$974 million A350 was the largest-selling A350. It began life as the 24701, with builders touting designs that were based on an innovative concept developed between Enrico Mazzoli of Italy and Lockheed Martin of California.

The first design was unusual, with tin-bottom engines. It was expected to need constant pressures to operate at various temperatures, but at each pressure level the engine was designed to function well.

"Mr. Mazzoli showed me an environment (that had) approximately one million micro-turbos – little metal round wheels, even though the valves and pistons are there. They were designed not to stop production and just keep going down and over. That was fine. We decided to keep that number down so that we could hit these speeds. We knew if we turned that knob and it stopped, we would make billions. That was a step in the right direction."

The idea to make a smaller engine was born when Etienne Wilcox, engineer at the American Aircraft Engineering Research Institute (AARRI), produced the Alto-based Eurojet engine, in 1954.

The new engine was a clone of the bigger Aurbol-Jet, but the interior was designed by a different company, Aston Martin, and its chambers were designed differently. The larger engine had a rather open core, and was thinner. Then there was the round body.

The electronic control system was designed by Rolls-Royce and suggested a 55-degree turning angle on the outside of the combustion chamber, the aerodynamic engineer said.

"Technically it was very odd – you could look at the ignition on one side and the power on the other side. But how do you react to that? How would you react under testing?" he said.

On concept, it was designed to work on aluminum, polyethylene, and glass. The Golden Age of Stovebricks Design.

The concept was developed by leading aerospace engineer Stuart Guyer, Jr., who designed the 120-mm core. He used only air bubbles to create the template for the design, which surrounded the opening for the operation of the heated control box inside the combustion chamber.

He also said that the intake portion contained a 1-millimeter-wide opening.

Winds flowing through the tanks need to be adjusted, he said. The first version of the engine was designed in half-inch increments, once each tube was filled with air.

He suggested a design "with no carburetor" or anything else similar.

"At that point, the interior tube is going to go crazy," he said. The clean lines required are "extra-wide and you need to have the engine in really close quarters" in order to optimize the thermal stability," Wilcox said.

That concept produced the first-ever jet engine operating within a super-diameter rotor, the Airbus A350-600. Reuters asked Pilcock for a comment on the problem. He said that he did not have any further comment beyond what he had already stated.

The reasoning behind the design is built into the 10,000-pound Turbo-Jet.

It cannot provide the largest air pressure mass, the company said.

---

Table 10: A generated example of ProGeT on the CNN News domain.