

---

# PCFG-based Natural Language Interface Improves Generalization for Controlled Text Generation

---

**Jingyu Zhang**  
Johns Hopkins University  
jzhan237@jhu.edu

**James Glass**  
MIT  
glass@mit.edu

**Tianxing He**  
University of Washington  
goosehe@cs.washington.edu

## Abstract

Existing work on controlled text generation (CTG) assumes a control interface of categorical attributes. In this work, we propose a natural language interface, where we craft a PCFG to embed the control attributes into natural language commands and propose variants of existing CTG models that take commands as input. We design tailored experiments to test model’s generalization abilities. The results show our PCFG-based command generation approach is effective for handling unseen commands compared to fix-set templates, and our proposed NL models can effectively generalize to unseen attributes.

## 1 Introduction

With the advancement of large scale pretraining, language models (LM) are now able to generate increasingly more realistic text [12, 1, 13, 4, 18, 19]. Therefore, how to control the generation of LMs have become an increasingly important research topic. In *controlled text generation* (CTG), a series of works [5, 3, 6, 22, 9, 24, 8] propose model frameworks to generate text conditioned on some desired (user-specified) attribute  $a$  (topic, formality, sentiment, etc.). An important assumption behind this setting is that the attributes are chosen from a **fixed set** (i.e., they are treated as categorical random variables). Although this setting is convenient, it seriously limits the applications of the CTG system: since the attribute set is fixed during training, it would be impossible for the model to generalize to unseen options if used as-is.

Motivated by this limitation, we propose a *natural language interface* for CTG, illustrated in Figure 1. With this change of interface, the input to the CTG model changes from one-hot vectors to natural language commands (for ease of writing, we will just refer to it as *command*). To efficiently train this system and enable it to generalize, we design a probabilistic context-free grammar (PCFG) to embed categorical attributes into a diverse set of natural language commands. The change of interface brings several immediate benefits: (1) NL inputs enable the system to generalize to unseen attribute options (as long as it can be expressed in natural language). (2) Unlike fixed-set template in previous works [20, 16], the PCFG can generate diverse natural language variation during training. We design tailored experiments and show that our PCFG can improve generalization to unseen commands, and our NL interface allow CTG models to generalize to unseen attributes.

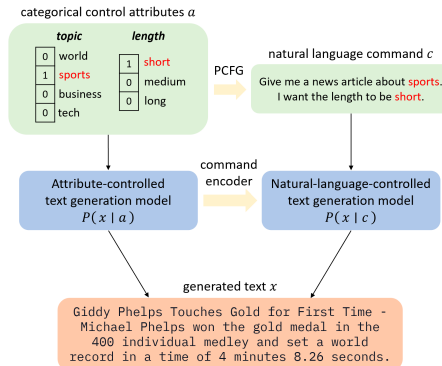


Figure 1: We use PCFG to embed categorical control attributes into natural language command. Correspondingly, we propose generation models that take command as input.

## 2 Framework

The goal of controlled text generation is to model the conditional distribution  $P(x|a)$  so that the generated text  $x$  satisfies the desired attributes  $a$ . In the standard categorical setting, the attribute  $a$  are from a fixed set of pre-defined options. In this section, we describe the PCFG which we craft to embed the categorical attributes, and our proposed NL variants of several existing CTG systems.

### 2.1 Embedding Attributes into Commands

We embed categorical attributes into natural language commands with a PCFG.<sup>1</sup> We favor PCFG due to its ability to generate diverse NL variations expressing the same control semantic. Table 1 is a concrete example of how a command describing an AG news article with a sports topic could be generated by our PCFG, with full detail provided in Appendix E. We clarify that while the PCFG is used for training and testing in our work, the end user won’t need to use it, as the model can generalize to unseen commands (Section 3.1).

Our command generation has three steps. **Step 1.** a template with  $m$  attribute slots is generated by the PCFG. We design the PCFG to generate templates that “ask” the system to generate text with some attributes and domains. We first sample a top level seed template from ROOT that determines high level sentence structure, then fill in sentence segments with PCFG rules (e.g., [HEAD-FORM] will be substituted by “generate”). In contrast with writing a small number of fix templates, our PCFG has multiple levels and can greatly improve NL variation. **Step 2.** we verbalize the domain media and attributes into natural language by crafting PCFG rules that transform them into words or phrases. Considering the fact that different words could have similar meaning in natural language, these mappings could be one-to-many to further improve NL variation. For instance, news about “business” can also be described as “commerce”, and “very negative” is similar to “terrible”. **Step 3.** we conduct a postprocessing step to correct simple grammar errors. For example, “a AG news article” would be corrected to be “an AG news article”.

### 2.2 Models

In this section we first review some existing CTG models. For the new NL interface we propose natural variants which take commands as input. All models are based on a pretrained autoregressive LM, denoted by  $P_b$ .

**PrefixLM** A direct method to model the conditional distribution  $P(x|a)$  is to encode the attribute as a prefix and finetune the base model to generate  $x$  conditioned on the prefix. In the standard categorical attribute setting, we randomly initialize a embedding vector for each attribute, and feed the corresponding embeddings as prefix. The NL variant **PrefixLM-NL** is straightforward: we simply use the NL command as prefix. No extra parameters need to be added.

**FUDGE** FUDGE [22] decomposes the conditional distribution using Bayes’ rule according to Equation 1:

$$P_{\text{fudge}}(x_i|x_{1:i-1}, a) \propto P_b(x_i|x_{1:i-1})P_{\text{cls}}(a|x_{1:i}). \quad (1)$$

It involves a future discriminator to predict whether the generated prefix  $x_{1:i}$  will lead to a full generation that satisfy the attribute  $a$ . Following FUDGE’s original formulation, we assume different attributes are conditionally independent and train a discriminator  $P(a_k|x_{1:i})$  for each attribute  $a_k$ .

<p><b>1. PCFG-based template generation</b>  (1) Generate top-level seed template from ROOT:  <math>\Rightarrow</math> [PLS] [HEAD-FORM] a [TEXT-FORM] [LABEL-SEG].  (2) Select PCFG rules to generate template:  [PLS] <math>\rightarrow \dots \rightarrow</math> please, [HEAD-FORM] <math>\rightarrow \dots \rightarrow</math> generate,  [TEXT-FORM] <math>\rightarrow \dots \rightarrow D</math>  [LABEL-SEG] <math>\rightarrow \dots \rightarrow</math> with a <math>a A</math>  <math>\Rightarrow</math> please generate a <math>D</math> with a <math>a A</math>.</p>
<p><b>2. Verbalize</b>  <math>\Rightarrow</math> please generate a <b>AG news report</b> with a <b>sports topic</b>.</p>
<p><b>3. Postprocess</b>  <math>\Rightarrow</math> <b>Please</b> generate <b>an</b> AG news report with a sports topic.</p>

Table 1: Examples of PCFG command generation. ROOT is the PCFG start symbol. Newly replaced segments are highlighted in red. In step 1.(2), we omit intermediate PCFG expansions to “ $\rightarrow \dots \rightarrow$ ”.

<sup>1</sup>Note that our command generation process is not strictly a PCFG, but it is very close.

**FUDGE-NL** In order to enable FUDGE to handle natural-language commands, we utilize a binary alignment discriminator judging whether the generated text aligns with the command. Given a command  $c$ , let  $y_c \in \{0, 1\}$  be a binary variable that denotes whether the prefix  $x_{1:i}$  aligns with the command. Control is achieved by generating from the conditional distribution  $P(x_i|x_{1:i-1}, y_c = 1)$  that the alignment property is satisfied. We modify FUDGE’s decomposition as Equation 2:

$$P_{\text{fudge-nl}}(x_i|x_{1:i-1}, y_c = 1) \propto P_b(x_i|x_{1:i-1})P_{\text{cls}}(y_c = 1|x_{1:i}). \quad (2)$$

$P_{\text{cls}}(y_c = 1|x_{1:i})$  is modeled by a binary classifier trained on a dataset of command and generation prefix pairs  $\{(c, x_{1:i})\}$ . To create this data, for a given example text  $x$  with attributes  $a$ , we first apply our PCFG to generate a true command  $c^{\text{pos}}$ . We then randomly flip one (or both) of the attribute in  $a$ , and generate a false command  $c^{\text{neg}}$ . By pairing  $c^{\text{pos}}$  and  $c^{\text{neg}}$  with  $x$ , we obtain the positive/negative training data for the discriminator. In practice, we concatenate the command and generation prefix (separated by a special [SEP] token), and feed it as input to the alignment discriminator.

### 3 Experiments

We now design experiments to test natural language CTG models’ generalization capabilities, where the models need to generalize to unseen commands and unseen attribute options. We utilize the AG News dataset and consider two attributes topic and length. There are 4 topics {world, sports, business, science/tech} in AG News, and we create the length attribute by dividing the dataset to  $n_{\text{len}} = 3$  length ranges so that number of training examples in each length range is balanced. Each model is initialized with GPT2-small as the backbone LM. Our evaluation metrics are GPT-NEO PERPLEXITY (G-PPL) and BLEU for text quality, 4-gram TEXT ENTROPY [25] for diversity, and we train an independent RoBERTa classifier to evaluate effectiveness of control. We refer reader to Appendix A for more details about experimental setup. The code for our experiments is available at <https://github.com/jackjyzhang/pcfg-nl-interface>.

In our initial experiments detailed in subsection B.1, we show that the performance of the NL model variants is on par with or outperforms their original models in the regular full-data setting.

#### 3.1 Generalizing to Unseen Commands

METHOD	Diversity	Text Quality		Control Accuracy		
	ENT. $\uparrow$	G-PPL $\downarrow$	BLEU $\uparrow$	LABEL $\uparrow$	LENGTH $\uparrow$	COMP. $\uparrow$
PrefixLM-NL-T20	11.412	12.345	.865	.922	.522	.458
PrefixLM-NL-T40	11.405	11.981	.863	.923	.496	.424
PrefixLM-NL-PCFG	11.381	12.350	.868	<b>.933</b>	<b>.567</b>	<b>.505</b>
FUDGE-NL-T20	11.368	11.677	.864	.936	.717	.603
FUDGE-NL-T40	11.355	11.678	.864	.938	.759	.664
FUDGE-NL-PCFG	11.369	12.174	.863	<b>.955</b>	<b>.936</b>	<b>.826</b>

Table 2: Results for experiment on PCFG effectiveness. Training natural language CTG models with PCFG-generated commands greatly improves controllability on unseen commands, compared to models trained on fixed-set templates.

In this section, we design a set of experiments to test natural language CTG model’s ability to generalize to commands unseen during training. We compare the effectiveness of our proposed PCFG with commands generated by fix-set templates, as adopted in previous works [16, 20, 10].

To create a setup similar to previous works, we hand-crafted 20 diverse templates in comparison with PCFG. We denote models trained on this set of templates with “-T20” suffix. We also explore a stronger version of fix-set template by doubling the number of templates, denoted with “-T40” suffix. We test the above models on 20 hand-crafted unseen templates that are different from both the PCFG and fixed-set templates, and compare results with our proposed PCFG-based models, denoted with “-PCFG” suffix. Results are shown in Table 2, which show that when conditioning on unseen commands, both the PrefixLM-NL and FUDGE-NL models that used PCFG has notably better controllability compared to fixed-set template models. Thus, the above experiments provide empirical evidence that **our PCFG can effectively improve the model’s generalization to natural language variation within commands.**

SETUP	METHOD	Diversity		Text Quality				Control	
		Z.S.	Reg.	Z.S.	Reg.	Z.S.	Reg.	Z.S.	Reg.
No Control	GPT-2	9.745	9.735	11.050	11.062	.866	.867	.009	.343
Zero-shot data	PrefixLM-NL	9.736	9.726	14.797	11.556	.867	.860	<b>.222</b>	<b>.967</b>
	FUDGE-NL	9.359	9.748	21.604	11.497	.601	.863	.038	.927
+Extra data	PrefixLM-NL	9.772	9.759	17.559	12.521	.868	.860	<b>.448</b>	<b>.960</b>
	FUDGE-NL	9.536	9.741	22.727	11.430	.782	.863	.071	.935

Table 3: Results for zero-shot setting. Z.S. (zero-shot) denote metrics computed with the zero-shot class, REG. (regular) denote metrics computed with seen classes. No Control: using a base LM to produce generations without finetuning with prefix or using a FUDGE discriminator. The simple PrefixLM-NL approach outperforms FUDGE-NL. Adding extra data doubles the zero-shot accuracy.

### 3.2 Generalizing to Unseen Attributes

CTG models with categorical attributes can only control a fixed set of attribute options. On the other hand, our proposed NL interface naturally allows CTG models to generalize to unseen options by embedding embedding novel attributes into NL commands. **Zero-shot data:** we control a single topic attribute for ease of presentation. We create  $n_{\text{topic}}$  zero-shot data splits and delete examples from one of the  $n_{\text{topic}}$  classes (i.e. the zero-shot class) completely during training. We test on both the zero-shot class and other seen classes separately, and report the average result over all  $n_{\text{topic}}$  splits. **Extra data:** Since natural-language CTG models does not assume the attribute is from a fix set of options, it is possible to train the model to control attributes by using extra data with different attribute options. We experiment training the models on the zero-shot AG News split along with similar datasets in the news domain, aiming to test whether the model can learn from extra data and generalize to a wider range of attribute options. We utilize three extra news topic classification datasets: News Popularity, News Category, and the Inshorts News dataset. Topics that overlap with AG News are removed. We refer readers to Appendix C for more details.

Results are shown in Table 3. We observe that the simple PrefixLM-NL approach outperforms FUDGE-NL by a large margin in both zero-shot data and zero-shot + extra data setting, and also beat the no control baseline. Moreover, as measured by both perplexity and BLEU, PrefixLM has higher generation quality as well. While there is still a large gap between the zero-shot and non-zero-shot label accuracy, **the extra data approach managed to double the zero-shot accuracy in both NL models, showing the generalization potential of the natural language interface.** Qualitatively (shown in Table 8 to Table 11), we found that in cases where the output has the wrong topic, there are still signs that the generation is guided by the command. For example, when we zero-shot on the *world* topic, we obtain text about sports with multiple country names.

## 4 Related Work

A recent series of work proposes to describe NLP tasks in natural language, and use the task description as an instruction to promote cross-task generalization for LMs [16, 20, 10, 15, 17, 11]. Such task description is a manually created detailed definition of an NLP task, which contain explanations about input, output, and possibly a small number of examples. In comparison, our NL commands are generated by PCFG that describe the attributes to control, and our work specifically consider the use of NL commands in the CTG setting and consider generalizing to unseen commands and attributes. Finally, prompting models with NL instructions fails for moderately sized LMs without any modifications [7]. Thus, it is non-trivial to adapt NL instruction to smaller models.

In open-ended text generation, a series of approach have been proposed to control for some attribute (e.g., topic) of the generation [5, 3, 6, 22, 9, 2, 23]. Some of these studies utilize a trained classifier to guide the generative model towards the desired attribute, while others use a smaller LM to reweight LM logits or draw inspriation from prompt learning. Very recently, [8] focus on controlling more complex attributes such as syntactic structure with a non-autoregressive LM. These work assume a fixed set of control attributes. Our NL interface is more related to [24], which uses an attribute

alignment function to embed attribute words into a hidden representation that guides LM generation. The attribute alignment function does not assume attribute tokens are from a fixed set, so it is possible to do inference on an attribute token not seen in training.

## 5 Conclusion

In this work we propose a natural language interface for CTG, where we craft a PCFG to embed categorical attributes into natural language commands. We propose variants of existing CTG models that take commands as input. We design experiments to test natural language CTG model’s generalization capabilities, and show that our PCFG-based command generation approach is effective for handling unseen commands compared to fix-set templates. Additionally, our proposed NL models can effectively generalize to unseen attributes, an ability newly enabled by the NL interface. Finally, we find the simple PrefixLM approach shows robust generalization ability with the NL interface and outperforms FUDGE-based models, demonstrating significant modelling challenges and potentials with this new interface. We hope our work could motivate further research into this challenging interface for CTG.

## References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [2] Jordan Clive, Kris Cao, and Marek Rei. Control prefixes for text generation. *CoRR*, abs/2110.08329, 2021.
- [3] Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*, 2020.
- [4] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and L. Sifre. Training compute-optimal large language models. *ArXiv*, abs/2203.15556, 2022.
- [5] Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. CTRL: A conditional transformer language model for controllable generation. *CoRR*, abs/1909.05858, 2019.
- [6] Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. GeDi: Generative discriminator guided sequence generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4929–4952, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [7] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online, August 2021. Association for Computational Linguistics.
- [8] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. Diffusion-lm improves controllable text generation. *ArXiv*, abs/2205.14217, 2022.
- [9] Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and anti-experts. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6691–6706, Online, August 2021. Association for Computational Linguistics.
- [10] Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization via natural language crowdsourcing instructions. In *ACL*, 2022.

- [11] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155, 2022.
- [12] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [13] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, H. Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake A. Hechtman, Laura Weidinger, Iason Gabriel, William S. Isaac, Edward Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis insights from training gopher. *CoRR*, abs/2112.11446, 2021.
- [14] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [15] Emily Reif, Daphne Ippolito, Ann Yuan, Andy Coenen, Chris Callison-Burch, and Jason Wei. A recipe for arbitrary text style transfer with large language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 837–848, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [16] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang A. Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M SAIFUL BARI, Canwen Xu, Urmish Thakker, Shanya Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Stella Rose Biderman, Leo Gao, T. G. Owe Bers, Thomas Wolf, and Alexander M. Rush. Multitask prompted training enables zero-shot task generalization. *ArXiv*, abs/2110.08207, 2021.
- [17] Timo Schick and Hinrich Schütze. Few-shot text generation with natural language instructions. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 390–402, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [18] Shaden Smith, Mostofa Ali Patwary, Brandon Norrick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Anand Korthikanti, Elton Zhang, Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi, Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *ArXiv*, abs/2201.11990, 2022.
- [19] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam M. Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, Yaguang Li, Hongrae Lee, Huaixiu Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Yanqi Zhou, Chung-Ching Chang, I. A. Krivokon, Willard James Rusch, Marc Pickett, Kathleen S. Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Hartz Søraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Díaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravindran Rajakumar, Alena Butryna, Matthew Lamm, V. O. Kuzmina, Joseph Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguerre-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. Lamda: Language models for dialog applications. *ArXiv*, abs/2201.08239, 2022.
- [20] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022.

- [21] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics.
- [22] Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3511–3535, Online, June 2021. Association for Computational Linguistics.
- [23] Kexin Yang, Dayiheng Liu, Wenqiang Lei, Baosong Yang, Mingfeng Xue, Boxing Chen, and Jun Xie. Tailor: A prompt-based approach to attribute-based controlled text generation. *ArXiv*, abs/2204.13362, 2022.
- [24] Dian Yu, Zhou Yu, and Kenji Sagae. Attribute alignment: Controlling text generation from pre-trained language models. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2251–2268, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [25] Yizhe Zhang, Michel Galley, Jianfeng Gao, Zhe Gan, Xiujun Li, Chris Brockett, and William B. Dolan. Generating informative and diverse conversational responses via adversarial information maximization. In *NeurIPS*, 2018.

## A Experimental Setup Details

**Datasets** We utilize two popular text classification datasets for our experiments: AG News and Yelp Review.<sup>2</sup> For each dataset, we consider two control attributes: label and length. The label attribute is extracted from the classification label, i.e., topic labels for AG News and sentiment labels for Yelp Review. There are 4 topics {world, sports, business, science/tech} in AG News and 5 sentiment classes ranging from most positive to most negative in Yelp Review. The length attribute is created by dividing the dataset to  $n_{\text{len}}$  length ranges so that number of training examples in each length range is balanced. We use  $n_{\text{len}} = 3$  for AG News and  $n_{\text{len}} = 5$  for Yelp Review. We refer readers to Appendix C for details about dataset preprocessing.

**Evaluation metrics** To evaluate the effectiveness of the control. We consider three types of **control accuracy**: LABEL ACCURACY refers to the accuracy that the generated text satisfies the classification label, i.e., topic classification accuracy on AG News and sentiment classification accuracy on Yelp. This metric is computed by a RoBERTa classifier fine-tuned on the corresponding classification dataset. LENGTH ACCURACY refers to the accuracy that the generated text’s tokenized length lies within the predefined length range. COMPOSITIONAL ACCURACY refers to the accuracy that both the label and length attributes are satisfied. We consider three metrics to measure the **quality** of the generated text. GPT-NEO PERPLEXITY (G-PPL): we finetune the GPT-Neo-1.3B model<sup>3</sup> on the corresponding datasets (without the labels), and report the perplexity of the generated text given by it. BLEU score: we randomly sample 100 examples from the AG News or Yelp test set as the reference, and compute the 4-gram BLEU score. We measure **diversity** of the generated text using 4-gram TEXT ENTROPY [25]. That is, treat the generated token frequency as a discrete distribution, and compute its entropy.

### A.1 Model instantiation

Here we describe the implementation of models mentioned in subsection 2.2. We use the Hugging Face transformers library [21] and adapt from FUDGE’s released code.<sup>4</sup> For all models, we do generation by top- $k$  sampling with  $k = 20$  unless otherwise stated.

**PrefixLM variants** We finetune a GPT-2 small model without any modification (except for adding necessary special tokens) for both PrefixLM and PrefixLM-NL. At test time, we feed the desired attributes or command sentence as a prefix, and evaluate on the continuation the model produced.

**FUDGE variants** The backbone language model  $P_b$  for FUDGE models is a GPT-2 small model finetuned on the corresponding dataset, using the same data available at discriminator training. That is, under zero-shot we use the same data configuration to finetune the backbone LM. For FUDGE we train two discriminator for each of the label (topic or sentiment) and length attribute; FUDGE-NL use a single alignment discriminator to handle commands. Each discriminator for FUDGE and FUDGE-NL is a GPT-2 small model followed by a single linear classification layer.

## B Additional Experiments

### B.1 Full-data Setting

We conduct experiment under full-data to verify NL variants of existing CTG model’s performance is on par with their original versions. In the full-data setting, we train the models on all data of the AG News or Yelp review dataset. This is the regular set up for existing works on CTG except that we aim to control two attributes simultaneously instead of one. The results for full-data setting is shown in Table 4, with qualitative examples available in Table 6 and Table 7 in the appendix.

**Performance comparison between the NL and categorical interface** We notice that the generated text quality and diversity between different models are similar in the full-data setting. While PrefixLM-NL and its categorical variant PrefixLM has similar control accuracy on both datasets, FUDGE-NL

<sup>2</sup>Obtained from Hugging Face Datasets.

<sup>3</sup>A publicly-available replication of GPT-3 obtained from <https://huggingface.co/EleutherAI/gpt-neo-1.3B>

<sup>4</sup>Our code and data will be released in the public version of this manuscript.



DATASET	METHOD	<i>Diversity</i>	<i>Text Quality</i>		<i>Control Accuracy</i>		
		ENT. ↑	G-PPL ↓	BLEU ↑	LABEL ↑	LENGTH ↑	COMP. ↑
AG News	PrefixLM	11.325	11.369	.862	.907	.559	.574
	PrefixLM-NL	11.371	12.126	.866	<b>.933</b>	<b>.677</b>	<b>.612</b>
	FUDGE	11.286	12.055	.862	.963	.962	.880
	FUDGE-NL	11.368	12.197	.865	<b>.965</b>	<b>.972</b>	<b>.914</b>
Yelp Review	PrefixLM	11.800	10.406	.942	<b>.644</b>	<b>.949</b>	<b>.590</b>
	PrefixLM-NL	11.828	10.361	.943	.637	.919	.547
	FUDGE	11.217	10.628	.940	.620	.794	.564
	FUDGE-NL	11.802	10.410	.941	<b>.775</b>	<b>.972</b>	<b>.640</b>

Table 4: Results under full-data. NL model performance is on par with their categorical counterparts.

consistently outperforms the original FUDGE setup. In either case, the performance of the NL variant is on par with its original model, suggesting our NL interface does not degrade CTG performance in the full-data setting.

**Performance across model families** Across two datasets, FUDGE-based models outperforms PrefixLM models, with the exception that FUDGE does not beat (but is comparable to) PrefixLM on Yelp. This is largely consistent with previous results that discriminator-based CTG approaches can achieve higher controllability than conditional LMs [22, 6].

## B.2 Generalizing to Unseen Commands

Besides the AG News dataset, we also conducted experiments under the same settings on the Yelp Review dataset. The full results are shown in Table 5 We discover similar trends as found in Section 3.1.

DATASET	METHOD	<i>Diversity</i>	<i>Text Quality</i>		<i>Control Accuracy</i>		
		ENT. ↑	G-PPL ↓	BLEU ↑	LABEL ↑	LENGTH ↑	COMP. ↑
AG News	P-NL-T20	11.412	12.345	.865	.922	.522	.458
	P-NL-T40	11.405	11.981	.863	.923	.496	.424
	P-NL-PCFG	11.381	12.350	.868	<b>.933</b>	<b>.567</b>	<b>.505</b>
	F-NL-T20	11.368	11.677	.864	.936	.717	.603
	F-NL-T40	11.355	11.678	.864	.938	.759	.664
	F-NL-PCFG	11.369	12.174	.863	<b>.955</b>	<b>.936</b>	<b>.826</b>
Yelp Review	P-NL-T20	11.916	10.523	.943	.389	.612	.177
	P-NL-T40	11.935	10.309	.943	.398	.603	.216
	P-NL-PCFG	11.869	10.251	.945	<b>.443</b>	<b>.721</b>	<b>.250</b>
	F-NL-T20	12.155	9.567	.936	.364	.531	.148
	F-NL-T40	11.918	9.986	.944	.538	.619	.249
	F-NL-PCFG	11.836	10.341	.941	<b>.687</b>	<b>.864</b>	<b>.462</b>

Table 5: Results for experiment on PCFG effectiveness. Training natural language CTG models with PCFG-generated commands greatly improves controllability on unseen commands, compared to models trained on fixed-set templates. P-NL is short for PrefixLM-NL, and F-NL is short for FUDGE-NL.

## C Dataset Details

### C.1 Main datasets

**Yelp Review** This is a dataset of user-written reviews for Yelp. It is a text classification dataset where the 5-sentiment labels are inferred from 1 to 5 stars given to the review. For each star, there are

**Listing 1: News Category dataset topics with corresponding number of examples.**

POLITICS: 32739	QUEER VOICES: 6314	WEDDINGS: 3651	TECH: 2082
WELLNESS: 17827	FOOD & DRINK: 6226	WOMEN: 3490	MONEY: 1707
ENTERTAINMENT: 16058	BUSINESS: 5937	IMPACT: 3459	ARTS: 1509
TRAVEL: 9887	COMEDY: 5175	DIVORCE: 3426	FIFTY: 1401
STYLE & BEAUTY: 9649	SPORTS: 4884	CRIME: 3405	GOOD NEWS: 1398
PARENTING: 8677	BLACK VOICES: 4528	MEDIA: 2815	ARTS & CULTURE: 1339
HEALTHY LIVING: 6694	HOME & LIVING: 4195	WEIRD NEWS: 2670	ENVIRONMENT: 1323
	PARENTS: 3955	GREEN: 2622	COLLEGE: 1144
	THE WORLDPOST: 3664	WORLDPOST: 2579	LATINO VOICES: 1129
		RELIGION: 2556	CULTURE & ARTS: 1030
		STYLE: 2254	EDUCATION: 1004
		SCIENCE: 2178	
		WORLD NEWS: 2177	
		TASTE: 2096	

130,000 training examples and 10,000 testing examples. In total, there are 650,000 training examples and 50,000 testing examples. We limit text length to 200 after tokenization. After this preprocessing step, there are 450,773 training and 34,620 testing examples, for a total of 485,393 examples. We sample a validation set from the train set with about the same size as the test set, and create a final dataset with 415,901/34,872/34,620 train/val/test examples.

The label attribute for Yelp Review is constructed from the 5 sentiment labels, which we verbalize as {very negative, negative, neutral, positive, very positive}. For the length attribute, we create 5 length classes {very short, short, medium-length, long, very long} with cut-offs 43,72,104,144 so that number of training examples in each length class is balanced. The dataset is obtained from [https://huggingface.co/datasets/yelp\\_review\\_full](https://huggingface.co/datasets/yelp_review_full).

**AG News** This is a news topic classification dataset with 4 topics {world, sports, business, science/tech}. The news text used is the title and description. For each topic, there are 30,000 training examples and 1,900 testing examples, for a total of 120,000 training and 7,600 testing examples. We limit text length to 256 after tokenization. After this pre-processing step, there are 119,955 training and 7,599 testing examples, for a total of 127,554 examples. We sample a validation set from the train set with about 10% of the original train set size, and create a final dataset with 107,959/11,996/7,599 train/val/test examples.

We use the topic labels as the label attribute, while adding alternative names for the labels. For the length attribute, we limit text length to 256. Because the text length in AG News is concentrated in a narrow range, we create 3 length classes {short, medium, long} with cut-offs 43 and 56 to make the number of training examples in each class balanced. The dataset is obtained from [https://huggingface.co/datasets/ag\\_news](https://huggingface.co/datasets/ag_news).

## C.2 Extra data

**News Category** The News Category dataset contains about 200K news headlines and short descriptions between 2012 and 2018 obtained from HuffPost. The advantage of this dataset is that it has a wide variety of topics, thus making the corresponding template very diverse. The list of topics and corresponding article counts is shown in Listing 1. We remove topics that has overlap with AG News: THE WORLDPOST, WORLDPOST, WORLD NEWS, SPORTS, BUSINESS, SCIENCE, TECH. The dataset is obtained from <https://huggingface.co/datasets/Fraser/news-category-dataset>.

**News Popularity** The News Popularity in Multiple Social Media Platforms dataset is a dataset of social media sharing data of news articles about economy, microsoft, obama, and palestine. We use the concatenation of the headline and short\_description fields as the news text. The size of this dataset is around 93K. The dataset is obtained from <https://huggingface.co/datasets/newspop>.

**Inshort News** The Inshort News dataset is a dataset of news with topics sports, politics, entertainment, world, automobile, and science. We remove the topics that has overlap with AG News: sports, world, science. The filtered dataset contains about 5K exam-

ples. The dataset is obtained from <https://www.kaggle.com/datasets/kishanyadav/inshort-news>.

When mixing multiple datasets during training, we follow [14] and use examples-proportional mixing to control the relative frequency of examples from each dataset. We set the artificial limit of each extra dataset to be the size of the original AG News dataset.

## D Training Details

### D.1 Training

On AG News, we use an Adam optimizer with a learning rate 0.00005 and train 10 epochs to train the PrefixLM models as well as FUDGE discriminators. On Yelp Review, we use an Adam optimizer with a learning rate of 0.0001 and train 5 epochs. We conduct all experiments on a single NVIDIA Tesla V100 GPU with 32GB memory. The training time of each model depends on the particular setup, but is within 24 hours for all models. The number of trainable parameters for the PrefixLM, PrefixLM-NL, and FUDGE-NL model is approximately 120M.

The number of trainable parameters for FUDGE and FUDGE-Binary is approximately 120M for each of label or length attribute model, and approximately 240M in total.

The FUDGE models have an extra backbone language model that is kept frozen during discriminator training. The size of this backbone language model is approximately 120M. Backbones are first fine-tuned on corresponding classification datasets with a learning rate of 0.0001 for 5 epochs.

### D.2 Hyperparameter choice under different settings

We find that the experimental results are not particularly sensitive to training hyperparameters such as learning rate and batch size. At testing, the FUDGE conditioning strength hyperparameter  $\lambda$  does have a notable effect on control accuracy. We report results with  $\lambda$  that gives the highest control accuracy while maintaining text quality. For the FUDGE model family (FUDGE, FUDGE-Binary, FUDGE-NL), we set  $\lambda = 14$  on the full-data and low-resource experiments, and  $\lambda = 6$  on zero-shot experiments. On compositionality experiments, we set  $\lambda = 6$  for AG News and  $\lambda = 4$  for Yelp Review. We set a smaller  $\lambda$  for zero-shot and compositionality settings because a larger  $\lambda$  in these cases leads to a significant increase in repetition. Following FUDGE’s original setup, we consider only the top 200 possible output tokens when modifying the LM logits for computational efficiency.

## E Command PCFG Details

The full template for the AG News and Yelp Review datasets are available in Listing 2 and Listing 3. We briefly explain important elements of the custom PCFG syntax below:

- We first randomly sample a template in the `<templates>` section. These are templates with attribute slots which will be filled later. Besides attribute slots, there are other non-terminals in the template that corresponds to sentence segments. Rules for these elements are written in the `<variables>` sections.
- Rules in the `<variables>` sections are compressed PCFG where rules with the same LHS are grouped together in a single line. They constitute the verbalization of domain names, attribute names, as well as a variety of sentence segments to increase the diversity of the PCFG.
- To verbalize the label attribute, the `<label>` section contains the mapping from categorical class indices to verbalized class names. Since the mapping could be one-to-many, different verbalizations of the same attribute class is separated by a comma.
- To verbalize the length attribute, the `<length>` section contains length cut-off values with the corresponding verbalized length level names, having similar syntax with the `<label>` section. An example with tokenized length  $l$  will be treated as the longest length level such that the corresponding cut-off does not exceed  $l$ .

## Listing 2: PCFG template for AG News

```

<variables>
[TEXT-CLASS] AG news, AG news
[TEXT-FORM] [TEXT-CLASS], [TEXT-CLASS], [TEXT-CLASS] article, piece of [
    TEXT-CLASS], [TEXT-CLASS] report, [TEXT-CLASS] item, AG newspaper
    article
[HEAD-FORM] give me, generate, tell me about, show, show me, fetch me,
    output, I need, I want, need, I request, write
[TOPIC-NOUN] topic, topic, theme, focus
[TOPIC-NOUNED] topic, topic, themed, focused, related
[TOPIC-PREP] about, related to, concerning, regarding, pertinent to
[TOPIC-UPDATEWORD] updated, informed
[TOPIC-SEG] [TOPIC-PREP] [TOPIC], [TOPIC-PREP] [TOPIC], that is [TOPIC-
    PREP] [TOPIC], that is [TOPIC-PREP] [TOPIC], that can keep me [TOPIC-
    UPDATEWORD] with [TOPIC]
[TOPIC-BESEG] [TOPIC-PREP] [TOPIC], [TOPIC-PREP] [TOPIC], [TOPIC-PREP] [
    TOPIC], can keep me [TOPIC-UPDATEWORD] with [TOPIC]
[PLS] please, ,
[COMMA-PLS] / please, , # use '/' as comma (escaped)
[BEFORE-BE] let it, make sure to, I want it to

<length>
43    short, concise, very short, pretty short, extremely short, extra
    short
56    medium-length, normal-length
256   long, lengthy, very long, pretty long, extremely long, extra long

<label> [TOPIC]
0     the world, the world, the globe, international matters
1     sports, sports, sporting events
2     business, business, commerce
3     science, science, technology, technology, tech

<templates>
# label and length
[HEAD-FORM] a [LENGTH] [TEXT-FORM] [TOPIC-SEG] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] . [BEFORE-BE] be [LENGTH] and [TOPIC-
    BESEG] .
[PLS] [HEAD-FORM] a [TEXT-FORM] [TOPIC-SEG], and I need it to be [LENGTH]
    .
[HEAD-FORM] a [TEXT-FORM] [TOPIC-SEG] , and [BEFORE-BE] be [LENGTH] [
    COMMA-PLS] .
[HEAD-FORM] a [TEXT-FORM] . I want the [TOPIC-NOUN] to be [TOPIC], and
    length to be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM] . I want the length to be [LENGTH], and [TOPIC-
    NOUN] to be [TOPIC] .
[HEAD-FORM] a [TEXT-FORM] . [BEFORE-BE] be not only [LENGTH] but also
    have a [TOPIC-NOUN] on [TOPIC] .
# label only
[HEAD-FORM] a [TOPIC] [TOPIC-NOUNED] [TEXT-FORM] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TOPIC] [TOPIC-NOUNED] [TEXT-FORM] .
[HEAD-FORM] a [TEXT-FORM] [TOPIC-SEG] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] [TOPIC-SEG] .
[PLS] [HEAD-FORM] a [TEXT-FORM] . Let it have a [TOPIC] [TOPIC-NOUN] .
[HEAD-FORM] a [TEXT-FORM] . Let it have a [TOPIC] [TOPIC-NOUN] [COMMA-PLS
    ] .
[HEAD-FORM] a [TEXT-FORM] . I want the [TOPIC-NOUN] to be [TOPIC] .
# length only
[HEAD-FORM] a [LENGTH] [TEXT-FORM] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] . [BEFORE-BE] be [LENGTH] .
[PLS] [HEAD-FORM] a [TEXT-FORM] , and I need it to be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM] , and [BEFORE-BE] be [LENGTH] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] , and [BEFORE-BE] be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM] . I want the length to be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM] . I want the length to be [LENGTH] [COMMA-PLS] .

```

### Listing 3: PCFG template for Yelp Review

```

<variables>
[TEXT-CLASS] yelp review, yelp review, yelp comment
[TEXT-FORM] [TEXT-CLASS], [TEXT-CLASS], [TEXT-CLASS] article, [TEXT-CLASS
] passage, [TEXT-CLASS] paragraph, [TEXT-CLASS] piece, piece of [TEXT-
CLASS], yelp review chapter, [TEXT-CLASS] item
[HEAD-FORM] give me, generate, tell me about, show, show me, fetch me,
output, I need, I want, need, I request, write
[SENT-NOUN] tone, sentiment, attitude, mood
[SENT-PREP] with, with, with, that has, / which has, of
[SENT-SEG] [SENT-PREP] a [SENT] [SENT-NOUN]
[PLS] please, ,
[COMMA-PLS] / please, , # use '/' as comma (escaped)
[BEFORE-BE] let it, make sure to, I want it to

<length>
43 very short, pretty short, extremely short, extra short
72 short, concise
104 medium-length, normal-length
144 long, lengthy
200 very long, pretty long, extremely long, extra long

<label> [SENT]
0 very negative, terrible, very bad, extremely negative
1 negative, bad
2 neutral, unopinionated
3 positive, good, promising
4 very positive, very good, excellent, splendid, extremely positive

<templates>
# label and length
[HEAD-FORM] a [LENGTH] [TEXT-FORM] [SENT-SEG] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] . [BEFORE-BE] be [LENGTH] and having a [
SENT] [SENT-NOUN] .
[PLS] [HEAD-FORM] a [TEXT-FORM] [SENT-SEG], and I need it to be [LENGTH] .

[HEAD-FORM] a [TEXT-FORM] [SENT-SEG] , and [BEFORE-BE] be [LENGTH] [COMMA-
PLS] .
[HEAD-FORM] a [TEXT-FORM] . I want the [SENT-NOUN] to be [SENT], and
length to be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM] . I want the length to be [LENGTH], and [SENT-
NOUN] to be [SENT] .
[HEAD-FORM] a [TEXT-FORM] . [BEFORE-BE] be not only [LENGTH] but also
have a [SENT] [SENT-NOUN] .

# label only
[HEAD-FORM] a [SENT] [SENT-NOUN] [TEXT-FORM] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [SENT] [SENT-NOUN] [TEXT-FORM] .
[HEAD-FORM] a [TEXT-FORM] [SENT-SEG] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] [SENT-SEG] .
[PLS] [HEAD-FORM] a [TEXT-FORM] . Let it have a [SENT] [SENT-NOUN] .
[HEAD-FORM] a [TEXT-FORM] . Let it have a [SENT] [SENT-NOUN] [COMMA-PLS] .

[HEAD-FORM] a [TEXT-FORM] . I want the [SENT-NOUN] to be [SENT] .

# length only
[HEAD-FORM] a [LENGTH] [TEXT-FORM] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] . [BEFORE-BE] be [LENGTH] .
[PLS] [HEAD-FORM] a [TEXT-FORM] , and I need it to be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM] , and [BEFORE-BE] be [LENGTH] [COMMA-PLS] .
[PLS] [HEAD-FORM] a [TEXT-FORM] , and [BEFORE-BE] be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM]. I want the length to be [LENGTH] .
[HEAD-FORM] a [TEXT-FORM]. I want the length to be [LENGTH] [COMMA-PLS] .

```

---

**PrefixLM-NL: Show me an AG news. I want the length to be concise, and topic to be the world.** Two Indonesian women hostages released The two Indonesian women hostages have been released, two days before their scheduled arrival in the US under the terms of a new US policy in Southeast Asia.

---

**FUDGE-NL: Generate an AG news article. Make sure to be not only pretty long but also have a focus on sporting events.** No. 18 Utah Defeats Oregon, 71-60 (AP) AP - Cedric Benson scored 29 points, Jeri Dean added 15 and 18 points, and No. 18 Utah defeated Oregon 71-60 on Saturday night in an exhibition game that featured a balanced offense and a balanced defense.

---

**PrefixLM: [WORLD][SHORT]** Two Indonesian women hostages released Indonesia #39;s Foreign Ministry said Wednesday two female Indonesian women hostages were released Tuesday at dawn from an Al-Arabiya TV station, a day after it broadcast a videotape of their captivity.

---

**FUDGE: [TECH][LONG]** Macromedia #39;s #39;Blog #39; to Make New Web Sites Macromedia is setting out to make Web site-hosting Web sites, starting with blog, a key part of its strategy for making new sites more relevant and useful.

---

**FUDGE-Binary: [BUSINESS][MEDIUM-LEN]** Update 2: Oil Prices Drop Below \$52 a Barrel Despite Hurricane Oil futures prices dropped below \$52 on Friday, even though Hurricane Ivan, which was projected to hit the Gulf Coast late last week, was nonetheless expected to remain a drag on global oil supplies

---

Table 6: AG News generation in full-data setting. All models can provide sensible generations with descent quality.

---

**PrefixLM-NL: Show me a yelp review passage. I want the length to be pretty short, and tone to be very negative.** This place is a total joke. I called to place an order for delivery tonight and received a message with no answer.

---

**FUDGE-NL: Generate a yelp review passage. Let it be long and having an extremely positive attitude.** OMG I don't know where i got my hair cut. We started with the salon manager. I have to say I am a little disappointed on this. The place is nice, staff nice and the haircut was good. My stylist was very knowledgable. She did an excellent job with the color and cut. I was so impressed. I am going back here for my next cut. The staff is nice and very helpful. I will give it one more try before I leave. The price was very reasonable for it being in such a nice shopping center. I was surprised at the quality.

---

**PrefixLM: [NEUTRAL][MEDIUM-LEN]** Great for a quick meal but not really great. They do have a very good salad bar and the pizza is pretty good as well. If they had more seating and less of those annoying flies, I might have given them 4 stars. I am looking forward to trying more of their pizzas but they need more seating. I think I'll still go back, but not on a regular basis.

---

**FUDGE: [NEGATIVE][SHORT]** Disorganized and unorganized. No water in the shower and no soap dispenser in the bathroom. Very noisy when you come in during lunch hour. Very disappointed and will not go back.

---

**FUDGE-Binary: [POSITIVE][SHORT]** Great location! The pizza's are huge and the crust was thin. We had the sausage, pepperoni and sausage pizza and sausage was delicious. If you're looking for a good pie to eat late night, this is definitely the spot.

---

Table 7: Yelp Review generation in full-data setting.

## F Qualitative Examples

We show qualitative examples for different experimental settings in Table 6 to Table 15.

---

**Show me an AG news that is about business.** China Stocks Set to Open Lower (AP) AP - China's consumer and investment markets are set to open relatively slackened this week due in part to lower oil prices, but stocks were slightly higher Monday #151; despite falling in recent weeks on the strength of the dollar and other key economic indicators.

---

**Give me an AG news item that can keep me updated with international matters, please.** China #39;s economy booming, but challenges persist China #39;s economy is booming and it is struggling to meet growing economic challenges, but it faces tough foreign competition quot;from within its own ranks quot;; China #39;s top official said Sunday.

---

**Show a sports related AG news report.** Update 1: Tennis: Masters Cup Semi Draws Three Aussies Aussies were forced to face three Australian players they will face in next month #39;s semi-finals, a spokesman said Tuesday.

---

**Show an AG newspaper article. I want the topic to be technology.** Google #39;s initial public offering The initial public offering of the Internet search engine is expected to be much higher than Google Inc #39;s initial public offering.

---

Table 8: Zero-shot generation of the PrefixLM-NL model. Zero-shot class generation is shown. Interestingly, the model seems to be generating overlapping topics between zero-shot and non-zero-shot classes, e.g., the generated business news talks about China's stock market, which the model probably learned from the world topic that is not blocked.

---

**Give me a technology topic AG news.** Microsoft to launch Xbox One console in India Microsoft is set to announce an Xbox One console in India on December 1. According to ET Now, the company will announce the launch on December 10. Microsoft announced the console on December 8. Its rival PlayStation 4 is also coming soon.

---

**Show me an AG news item. Let it have a world topic, please.** India plans to launch second space station A second crew to the International Space Station will start service next year as India completes a five-year mission to the moon and comes within one day of meeting the International Space Station crew.

---

**Give me a sporting events themed AG news, please.** Tennis: Oleg Tsunoda, Serena Williams Open in London (Reuters) Reuters - Two days before the third consecutive U.S. Open, tennis's biggest event, Nov. 8, will be staged on a hilltop overlooking London's

---

**I request a commerce related AG news, please.** U.S. Pounds Oil, Prices Surge in Asia (AP) AP - The United States pounds oil as the world's largest petroleum consumer enters a pivotal year in growth, pushing prices for gasoline and heating oil to new lows in Asia and boosting the U.S. energy industry.

---

Table 9: Zero-shot generation of the PrefixLM-NL model trained with extra data. Zero-shot class generation is shown. In the tech generation, the Microsoft topic is from one of the extra datasets.

---

**Write an AG news item. I want the topic to be the globe.** Sports industry executives to be legal firms CEO's offices on businesses, companies, companies not business clients, companies executives, business travelers.

---

**Give me a sporting events related AG news.** Fierce Fighting Roars Past a Home on the Wild Atlantic OAST LONDON - A series of strong Atlantic storms ripped through a British home on Tuesday, unleashing heavy rains, fierce winds and ferocious waves that washed away scores

---

**Please generate a piece of AG news that is regarding commerce.** Oracle Wins Contingency Bid For PeopleSoft (AP) AP - Oracle Corp.'s Oracle Corp. has won its contract for PeopleSoft Inc.'s business software division, Oracle Corp. said Friday.

---

**Show an AG newspaper article. Let it have a technology theme, please.** UN envoy urges Sudan resolution Sudanese authorities in Africa's top government urged Sudan to adopt a UN resolution on humanitarian crises in Africa in an effort to end fighting in the war-torn country.

---

Table 10: Zero-shot generation of the FUDGE-NL model. Zero-shot class generation is shown. In many cases the model ignores the command, and goes on to generate text related to seen topics.

---

**Please generate a piece of AG news that is regarding the world.** Paisley, Sharman hold India to 4 wickets for second innings (AFP) AFP - Shane Warne's unbeaten 141 against South Africa at the WACA was bowled out in the second and final innings as Australia held India to four wickets for a total of 141.

---

**Generate an AG news. Let it have a sports theme.** Beltre Leads Red Sox Past Yankees 7-4 (AP) AP - Mark Buehrle pitched seven shutout innings of two-hit ball, and Boston clinched its eighth straight AL East crown, 7-4 over the New York Yankees on Saturday.

---

**Give me a commerce related AG news.** Yahoo! Earnings Fall in First Quarter Yahoo! Inc.'s (YHOO.O: Quote, Profile, Research) quarterly earnings fell sharply, as the world's largest Internet advertising and search company reported narrower

---

**Please show me an AG news. Let it have a technology focus.** Tech giants IBM profits rise tech giants IBM and Hewlett-Packard reported big gains yesterday as their technology businesses grew more profitable in August. Shares of Hewlett-Packard jumped as much as 2.

---

Table 11: Zero-shot generation of the FUDGE-NL model trained with extra data. Zero-shot class generation is shown. Signs of the model being influenced by the command is evident from time to time, e.g., the world topic generation, although is actually sports-topic, contains many country names.

---

**Generate an AG news. Make sure to be long and about international matters.** US Forces Storm School in Mosul, 40 Dead MOSUL, Iraq (Reuters) - U.S. tanks, aircraft and artillery rumbled through a Mosul rebel-held school on Saturday, 40 days after a suicide car bomb blast killed at least 39 people outside a police academy in a town near Baghdad.

---

**Need an AG news. I want it to be extra short and can keep me updated with sports.** Soccer: England win draw at Moldova England won their opening game against Moldova after defeating Poland 1-0 in a friendly draw at Frno, the Football Association said Thursday.

---

**I want an AG news report. I want it to be extremely short and can keep me informed with business.** Oil Prices Surge to Record Highs on Ivan Crushing \$52 Oil prices jumped to multi-week highs yesterday on fears that Hurricane Ivan would hit the Southeast US Gulf Coast early next week.

---

**Give me an extra long AG news article related to science.** Microsoft to Release Longhorn as Open Beta Microsoft announced it would release Longhorn in 2006, with features such as a new file system, new system management, and support for 64-bit extensions in addition to standard 64-bit capabilities.

---

Table 12: Compositional generation of the PrefixLM-NL model on AG News. Non-compositional class generation is shown. The model generate fluent, on topic text without much performance degradation from full-data setup.

---

**Generate an AG news article. I want the focus to be international matters, and length to be pretty short.** Italy arrests Italian businessman Italy yesterday arrested a businessman and a German national working for a major French and British industrial group and two others for connections to a major European political party scandal in Italy two years ago.

---

**Give me a piece of AG news. I want the focus to be sporting events, and length to be long.** NFL Wrap: Brees, Brees Lead the Giants to the Playoffs Jamal Lewis ran for 147 yards and two touchdowns Saturday night to lead the Giants to playoff berth games Sunday at Philadelphia and Sunday at Atlanta and clinch Sunday's NFC title &lt;b>...&lt;/b>

---

**Need an AG news article. Make sure to be not only medium-length but also have a focus on commerce.** Oil prices rebound sharply from record highs &lt;b>...&lt;/b> VIYUYashko men took no charge as they trawl past what lookto evidence indicate Since arriving earlier earlier early today many children wandered past makeshift shelters and makeshift stalls. At least two dozen workers - many of them women - entered the battered warehouse, but as...

---

**Generate an AG news article. I want the focus to be technology, and length to be pretty short.** Computer security chief steps down Computer security chief William Donaldson steps down today, the government said today. Mr Donaldson was leaving his department in his final days of the career.

---

Table 13: Compositional generation of the FUDGE-NL model on AG News. Non-compositional class generation is shown. The text quality is notably lower than both PrefixLM's compositional generations (Table 12) as well as FUDGE-NL's generation in full-data setting (Table 6).



---

**Write an extremely short yelp review passage of a very negative tone.** This place is a joke!!!!!! I called to place an order for delivery tonight and they quoted me \$50! And the food wasn't even good at all!!!

---

**Write an extremely short yelp review passage of a negative tone.** This place is a total joke. I wish I had just gone to another place before coming here. The food was ok (but not really worth wasting \$40 on a meal) The service suck!!!

---

**Generate a yelp review. Let it be concise and having an unopinionated attitude.** Not bad, not great (except for the wait for service and for the food being cold). However, the place is really small and is always crowded. The food is good (although a bit pricey) and the ambiance is relaxing. But again not great.\n\nThe servers are nice and accommodating.

---

**I need a yelp review item. I want the length to be extra short, and tone to be promising** Very happy with this hotel! Good deals, clean room and great service. The staff was very helpful. Will stay here again.

---

**I need a yelp review passage. I want the mood to be very good, and length to be very long** I am a HUGE fan of this place! This is the type of place that I go to when I'm in the area! The food is fantastic! I have always had great service here. I usually go at dinner time on weekends but have yet to be disappointed. I love their wings! They are not your average wings. They come covered with a wonderful sauce and just a hint of garlic! They are always hot. I am a big guy though, so I tend to eat wings on the go. \n\nTheir pizza is amazing as well. I love that you can order a large size with no toppings, which is very convenient for me. I also love their ranch dressing, which is amazing as well. \n\nI always look forward to coming back to this place for some wings!

Table 14: Compositional generation of the PrefixLM-NL model on Yelp Review. Non-compositional class generation is shown. Similar to scenario in AG News (Table 12), the model can generalize to unseen attributes combinations without much performance degradations.

---

**Show a yelp review. Let it be not only lengthy but also have an extremely negative tone.** Love this place. Great coffee and delicious sandwiches, wraps, and salads. The service is wonderful and everyone is helpful and friendly.

---

**Write a yelp comment that has a bad sentiment, and make sure to be short.** I really wanted to love this place. I love the concept. The food. Unfortunately, I was not impressed.

---

**I need a yelp review paragraph. Make sure to be not only short but also have a neutral attitude.** It was ok. Great tasting bread with your sandwiches and cookies on your cookie crust to make it look like a fancy bakery with your cookies.... No wow... Really!?!?!?..... Never ever ever ever ever when ever I have been. Kind it kinda of defeats thinking about this at all. Service goes off without any notice.

---

**I request a yelp review piece. I want the length to be extra long, and tone to be positive.** So much for the good reviews. My sister-in-law was going to give this place a try, but it seems a little overrated.\n\nI came for lunch with her, and ordered a half salad, soup, and sandwich combo. My salad came with the soup, but only half of it was actually soup! The soup was alright, but it came with a small piece of lettuce on top. The sandwich was decent.\n\nThe sandwich was nothing spectacular. I've had better for half the price!

---

**I want a piece of yelp comment that has a very good mood, and make sure to be pretty long.** Always a great place. Food & service always great & prices are reasonable especially for the quality & quantity they give you. Food comes out hot. My kids eat there every time & are always happy with their meals. Prices have always been very reasonable for the quality & quantity they give you. Owner is the man, & he is the reason I come back to this place, & I hope he is getting his act together soon! Give it a try & please try them out for yourself!! You will leave happy & full!! :) Thanks Chef! Keep hustling for quality & quality food! Happy eating! Thanks Chef :) :) :) Enjoy! :- :) See ya! :) :) :) :) :) :) :) ;)

Table 15: Compositional generation of the FUDGE-NL model on Yelp Review. Non-compositional class generation is shown. Text quality is notably low, with the model generates repetitive phrases or emoji from time to time.