

A Contrastive Framework for Neural Text Generation

Yixuan Su[♣] Tian Lan[◇] Yan Wang[◇] Dani Yogatama[♣]
 Lingpeng Kong[♡] Nigel Collier[♣]

[♣]Language Technology Lab, University of Cambridge

[◇]Tencent AI Lab [♣]DeepMind

[♡]Department of Computer Science, The University of Hong Kong

{ys484, nhc30}@cam.ac.uk

lantiangmftby@gmail.com, brandenwang@tencent.com

dyogatama@deepmind.com, lpk@cs.hku.hk

Abstract

Text generation is of great importance to many natural language processing applications. However, maximization-based decoding methods (e.g. beam search) of neural language models often lead to degenerate solutions—the generated text is unnatural and contains undesirable repetitions. Existing approaches introduce stochasticity via sampling or modify training objectives to decrease probabilities of certain tokens (e.g., unlikelihood training). However, they often lead to solutions that lack coherence. In this work, we show that an underlying reason for model degeneration is the anisotropic distribution of token representations. We present a contrastive solution: (i) *SimCTG*, a contrastive training objective to calibrate the model’s representation space, and (ii) a decoding method—*contrastive search*—to encourage diversity while maintaining coherence in the generated text. Extensive experiments and analyses on three benchmarks from two languages demonstrate that our proposed approach outperforms state-of-the-art text generation methods as evaluated by both human and automatic metrics.¹

1 Introduction

Open-ended neural text generation with Transformer (Vaswani et al., 2017) is an indispensable component in various natural language applications, such as story generation (Fan et al., 2018), contextual text completion (Radford et al., 2019), and dialogue systems (Su et al., 2021g). However, the conventional approach of training a language model with maximum likelihood estimation (MLE) and decoding the most likely sequence is often not sufficient (Holtzman et al., 2020; Welleck et al., 2020). Specifically, this modelling formulation often leads to the problem of *degeneration*, i.e. the generated texts from the language model tend to be dull and contain undesirable repetitions at different levels (e.g., token-, phrase-, and sentence-level) (Dinan et al., 2019). To alleviate this problem, previous solutions modify the decoding strategy by sampling from *less* likely vocabularies (Fan et al., 2018; Holtzman et al., 2020). While reducing the generated repetition, these sampling methods introduce another critical problem (*semantic inconsistency*)—the sampled text tends to diverge from or even contradict to the original semantics defined by the human-written prefix (Basu et al., 2021). Another approach addresses the degeneration problem by modifying the model’s output vocabulary distribution with unlikelihood training (Welleck et al., 2020).

In this work, we argue that the degeneration of neural language models stems from the *anisotropic* distribution of token representations, i.e. their representations reside in a narrow subset of the entire space (Ethayarajh, 2019; Dong et al., 2021; Su et al., 2021c). In Figure 1(a), we showcase a cosine

¹Our code and models are publicly available at <https://github.com/yxuansu/SimCTG>.

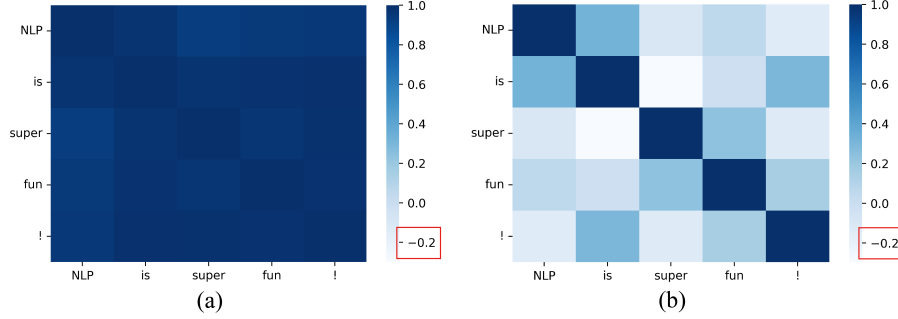


Figure 1: Token cosine similarity matrix of (a) GPT-2 and (b) SimCTG. (best viewed in color)

similarity matrix of token representations (taken from the final layer of the Transformer) produced by GPT-2. We see that the cosine similarities between tokens within a sentence are over **0.95**, meaning that these representations are close to each other. Such high similarity is undesirable as it can naturally cause the model to generate repetitive tokens at different steps, leading to degeneration. In an ideal setting, the token representations of the model should follow an isotropic distribution, i.e. the token similarity matrix should be sparse and the representations of distinct tokens should be discriminative as shown in Figure 1(b). Moreover, during decoding, the sparseness of token similarity matrix of the generated text should be preserved to avoid model degeneration.

Based on the above motivations, we present *SimCTG* (a **simple** **contrastive** framework for neural **text** **generation**) that encourages the model to learn discriminative and isotropic token representations. We also present a novel decoding strategy to complement SimCTG, *contrastive search*. The key intuitions behind contrastive search are: (i) at each decoding step, the output should be selected from the set of most probable candidates predicted by the model to better maintain the semantic coherence between the generated text and the human-written prefix, and (ii) the sparseness of the token similarity matrix of the generated text should be preserved to avoid degeneration.

We comprehensively evaluate our approach on three widely used benchmarks from two languages. The experimental results show that SimCTG improves the intrinsic qualities of the language model, as evaluated by perplexity and token prediction accuracy. Moreover, we demonstrate that the proposed contrastive search significantly outperforms previous state-of-the-art decoding methods in human and automatic evaluations. Furthermore, we provide in-depth analyses to get better insights on the inner-workings of our proposed approach.

2 Background

2.1 Language Modelling

The goal of language modelling is to learn a probability distribution $p_\theta(\mathbf{x})$ over a variable-length text sequence $\mathbf{x} = \{x_1, \dots, x_{|\mathbf{x}|}\}$, where θ denotes model parameters. Typically, the maximum likelihood estimation (MLE) objective is used to train the language model which is defined as

$$\mathcal{L}_{\text{MLE}} = -\frac{1}{|\mathbf{x}|} \sum_{i=1}^{|\mathbf{x}|} \log p_\theta(x_i | \mathbf{x}_{<i}). \quad (1)$$

However, as observed in many recent studies (Ethayarajh, 2019; Dong et al., 2021; Su et al., 2021c), training with likelihood maximization objective often yields an anisotropic distribution of model representations (especially for Transformer-based models) that undermines the model’s capacity.

2.2 Open-ended Text Generation

In this work, we focus on studying the task of open-ended text generation due to its generality in various applications, such as story generation (Fan et al., 2018), contextual text completion (Radford et al., 2019), poetry generation (Li et al., 2020), and dialogue systems (Su et al., 2021g). Formally, conditioned on a human-written prefix (i.e. context) \mathbf{x} , the task is to decode a continuation $\hat{\mathbf{x}}$ from

the language model. The resulting text is $\{x_1, \dots, x_{|\mathbf{x}|}, \hat{x}_{|\mathbf{x}|+1}, \dots, \hat{x}_{|\mathbf{x}|+|\hat{\mathbf{x}}|}\}$. Typically, there are two classes of methods used for decoding, which are (1) deterministic and (2) stochastic methods.

Deterministic Methods. Two widely used deterministic approaches are greedy and beam search which aim to select the text continuation with highest probability based on the model’s probability distribution p_θ . However, solely maximizing the output probability often leads to dullness (Li et al., 2016) and degeneration (Fan et al., 2018; Holtzman et al., 2020) in the generated text.

Stochastic Methods. To remedy the issues of deterministic decoding, several approaches have been proposed to sample from p_θ . To avoid sampling from the unreliable tail of distribution, Fan et al. (2018) proposed top- k sampling which draws sample from the vocabulary subset $V^{(k)}$ that maximizes $\sum_{v \in V^{(k)}} p_\theta(v|\mathbf{x})$. Here, $|V^{(k)}| = k$ and \mathbf{x} is the prefix context. Differently, the current state-of-the-art nucleus sampling (Holtzman et al., 2020) draws sample from the smallest vocabulary subset U with total probability mass above a threshold $p \in [0, 1]$; i.e. U is the smallest vocabulary subset such that $\sum_{v \in U} p_\theta(v|\mathbf{x}) \geq p$. While the sampling approaches help to alleviate model degeneration, the intrinsic stochasticity in these methods could cause the semantic meaning of the sampled text to diverge from or even contradict to the human-written prefix (Basu et al., 2021).

3 Methodology

In this section, we first present how to apply contrastive learning to calibrate the representation space of the language model. Then, we introduce our proposed contrastive search decoding algorithm.

3.1 Contrastive Training

Our goal is to encourage the language model to learn discriminative and isotropic token representations. To this end, we introduce a contrastive objective \mathcal{L}_{CL} into the training of the language model. Specifically, given a variable-length sequence $\mathbf{x} = \{x_1, \dots, x_{|\mathbf{x}|}\}$, the \mathcal{L}_{CL} is defined as

$$\mathcal{L}_{\text{CL}} = \frac{1}{|\mathbf{x}| \times (|\mathbf{x}| - 1)} \sum_{i=1}^{|\mathbf{x}|} \sum_{j=1, j \neq i}^{|\mathbf{x}|} \max\{0, \rho - s(h_{x_i}, h_{x_i}) + s(h_{x_i}, h_{x_j})\}, \quad (2)$$

where $\rho \in [-1, 1]$ is a pre-defined margin and h_{x_i} is the representation of token x_i produced by the model. The similarity function s computes the cosine similarity between token representations as

$$s(h_{x_i}, h_{x_j}) = \frac{h_{x_i}^\top h_{x_j}}{\|h_{x_i}\| \cdot \|h_{x_j}\|}. \quad (3)$$

Intuitively, by training with \mathcal{L}_{CL} , the model learns to pull away the distances between representations of distinct tokens.² Therefore, a discriminative and isotropic model representation space can be obtained. The overall training objective $\mathcal{L}_{\text{SimCTG}}$ is then defined as

$$\mathcal{L}_{\text{SimCTG}} = \mathcal{L}_{\text{MLE}} + \mathcal{L}_{\text{CL}}, \quad (4)$$

where the maximum likelihood estimation (MLE) objective \mathcal{L}_{MLE} is described in Eq. 1. Note that, when the margin ρ in \mathcal{L}_{CL} equals to 0, the $\mathcal{L}_{\text{SimCTG}}$ degenerates to the vanilla MLE objective \mathcal{L}_{MLE} .

3.2 Contrastive Search

We propose a new decoding method, *contrastive search*. At each decoding step, the key ideas of contrastive search are (i) the generated output should be selected from the set of most probable candidates predicted by the model; and (ii) the generated output should be discriminative enough with respect to the previous context. In this way, the generated text can (i) better maintain the semantic coherence with respect to the prefix while (ii) avoiding model degeneration.

Formally, given the context $\mathbf{x}_{<t}$, at time step t , the selection of the output x_t follows

$$x_t = \arg \max_{v \in V^{(k)}} \left\{ (1 - \alpha) \times \underbrace{p_\theta(v|\mathbf{x}_{<t})}_{\text{model confidence}} - \alpha \times \underbrace{(\max\{s(h_v, h_{x_j}) : 1 \leq j \leq t-1\})}_{\text{degeneration penalty}} \right\}, \quad (5)$$

²By definition, the cosine similarity $s(h_{x_i}, h_{x_i})$ of the identical token x_i is 1.0.

where $V^{(k)}$ is the set of k most-probable candidates predicted by the model θ and k is typically set as $3 \sim 10$. In Eq. 5, the first term, *model confidence*, is the probability of candidate v predicted by the model. The second term, *degeneration penalty*, measures how discriminative of candidate v with respect to the previous context $x_{<t}$ and s is defined in Eq. 3. Specifically, it is defined as the maximum cosine similarity between the representation of v and that of all tokens in $x_{<t}$. Here, the candidate representation h_v is computed by the model given the concatenation of $x_{<t}$ and v . Intuitively, a larger degeneration penalty of v means it is more similar to the context, therefore more likely leading to model degeneration. The $\alpha \in [0, 1]$ is a hyperparameter that regulates the importance of the two components. When $\alpha = 0$, contrastive search degenerates to the greedy search method.

Efficiency. The proposed contrastive search can be efficiently implemented. The additional computation required is the calculation of degeneration penalty which can be realized via a simple matrix multiplication. In section §6.4, we show that the decoding speed of contrastive search is better than or comparable with other widely used decoding methods.

4 Document Generation

We first evaluate our approach on the task of open-ended document generation.

Model and Baselines. Our proposed approach is architecture-agnostic and can be applied to any generation model. In this work, we evaluate our method on the representative GPT-2 model (Radford et al., 2019). Specifically, we fine-tune GPT-2 on the evaluated benchmark (detailed below) with the proposed objective $\mathcal{L}_{\text{SimCTG}}$ (Eq. 4) and generates the text continuation with different decoding methods. We perform experiments using the base model (117M parameters) which consists of 12 Transformer layers (Vaswani et al., 2017) with 12 attention heads. We compare our approach with two strong baselines: (1) GPT-2 fine-tuned with the standard MLE objective (Eq. 1); and (2) GPT-2 fine-tuned with unlikelihood objective (Welleck et al., 2020).³ Our implementation is based on the Huggingface library (Wolf et al., 2019).

Evaluation Benchmark. We conduct experiments on the Wikitext-103 dataset (Merity et al., 2017) which contains a large collection of Wikipedia articles with over 100 million words and 260 thousands unique tokens. Wikitext-103 is a document-level dataset and has been widely used for the evaluation of large-scale language modelling (Dai et al., 2019; Khandelwal et al., 2020; Yogatama et al., 2021).

Training. For our SimCTG and the MLE baseline, we fine-tune the models on Wikitext-103 for 40k training steps. For the unlikelihood baseline, following Welleck et al. (2020), we first fine-tune the model with the token-level unlikelihood objective for 38.5k steps and then with the sequence-level unlikelihood objective for 1.5k steps. Therefore, the overall training steps of all compared methods are the same. The batch size is set as 128 and the training samples are truncated to a maximum length of 256. For optimization, we use the Adam optimizer (Kingma and Ba, 2015) and a learning rate of $2e-5$ (with warm-up ratio of 10%).

Decoding. We evaluate the models on document generation by producing text continuations given the prefixes from the test set. In the experiments, the lengths of the prefix and the generated continuation are set as 32 and 128, respectively. We test different models with various decoding methods. For deterministic method, we use greedy search and beam search with beam size of 10. For stochastic method, we use the current state-of-the-art nucleus sampling (Holtzman et al., 2020) with $p = 0.95$. For the proposed contrastive search, the k and α in Eq. 5 are set as 8 and 0.6. These parameters are selected with minimal effort based on a few generated results on the validation set.

4.1 Evaluation Metrics

We perform evaluation from two aspects: (1) *language modelling quality* which measures the intrinsic quality of the model; and (2) *generation quality* which measures the quality of the generated text.

4.1.1 Language Modelling Quality

Following Welleck et al. (2020), we report the results of the model on the metrics below.

³The unlikelihood baseline is implemented with the official code, which can be found here https://github.com/facebookresearch/unlikelihood_training.

Perplexity. The model perplexity (**ppl**) on the test set of Wikitext-103.

Prediction Accuracy. It is defined as: $\mathbf{acc} = \frac{1}{\sum_{x \in \mathcal{D}} |x|} \sum_{x \in \mathcal{D}} \sum_{t=1}^{|x|} \mathbb{1}[\arg \max p_\theta(x|x_{<t}) = x_t]$, where \mathcal{D} is the Wikitext-103 test set, $x_{<t}$ is the prefix, and x_t is the reference token at time step t .

Prediction Repetition. The fraction of next-token (top-1) predictions that occur in the prefix which is defined as: $\mathbf{rep} = \frac{1}{\sum_{x \in \mathcal{D}} |x|} \sum_{x \in \mathcal{D}} \sum_{t=1}^{|x|} \mathbb{1}[\arg \max p_\theta(x|x_{<t}) \in x_{<t}]$.

In addition, the next token repetitions that do not equal to the ground truth token: $\mathbf{wrep} = \frac{1}{\sum_{x \in \mathcal{D}} |x|} \sum_{x \in \mathcal{D}} \sum_{t=1}^{|x|} \mathbb{1}[\arg \max p_\theta(x|x_{<t}) \in x_{<t} \wedge \neq x_t]$ is also reported.

4.1.2 Generation Quality

Generation Repetition. This metric measures the sequence-level repetition as the portion of duplicate n -grams in the generated text (Welleck et al., 2020). For a generated text continuation \hat{x} , the repetition at n -gram level is defined as: $\mathbf{rep-n} = 100 \times (1.0 - \frac{|\text{unique } n\text{-grams}(\hat{x})|}{|\text{total } n\text{-grams}(\hat{x})|})$.

Diversity. This metric takes into account the generation repetition at different n -gram levels and it is defined as: $\mathbf{diversity} = \prod_{n=2}^4 (1.0 - \frac{\mathbf{rep-n}}{100})$. It can be deemed as an overall assessment of model degeneration. A lower diversity means a more severe degeneration of the model.

MAUVE. MAUVE (Pillutla et al., 2021) is a recently proposed metric that automatically measures the token distribution closeness between the generated text and the human-written text. A higher MAUVE score means the model generates more human-like texts.

Semantic Coherence. To automatically measure the semantic coherence (i.e. consistency) between the prefix and the generated text, we employ the advanced sentence embedding method, SimCSE (Gao et al., 2021). Specifically, given the prefix x and the generated text \hat{x} , the coherence score is defined as: $\mathbf{coherence} = \frac{v_x^\top v_{\hat{x}}}{\|v_x\| \cdot \|v_{\hat{x}}\|}$, where $v_x = \text{SimCSE}(x)$ and $v_{\hat{x}} = \text{SimCSE}(\hat{x})$.

Perplexity of Generated Text. Lastly, we evaluate the perplexity of the generated text \hat{x} given the prefix x . The metric is defined as: $\mathbf{gen-ppl} = 2^{f(\mathcal{D}, \theta)}$ and $f(\mathcal{D}, \theta) = \frac{1}{\sum_{x \in \mathcal{D}} |x|} \sum_{x \in \mathcal{D}} \log_2 p_\theta(\hat{x}|x)$. Importantly, the optimal approach should produce text which has a perplexity *close* to that of the human-written text (Holtzman et al., 2020). If gen-ppl is too high, it means the generated text is very *unlikely* given the prefix, and such text tends to have low quality. In contrast, if gen-ppl is too low, it means the generated text has low diversity and gets stuck in repetition loops (Holtzman et al., 2020). In the experiments, we use the model θ trained with $\mathcal{L}_{\text{SimCTG}}$ to measure the gen-ppl of results from different approaches, therefore making sure the numbers are comparable with each other.

4.2 Results

The experimental results on Wikitext-103 are shown in Table 1.

Language Modelling Quality. From the results, we observe that SimCTG achieves the best perplexity and next token accuracy. The reason is that, with more discriminative representations, SimCTG is less confusing when making next token predictions, leading to the improved model performance. On the rep and wrep metrics, the unlikelihood model yields the best result but at the expense of unfavorable performance drops in the perplexity and next token accuracy.

Generation Quality. Firstly, on the rep-n and diversity metrics, SimCTG + contrast search obtains the best performance across the board, suggesting it best addresses the model degeneration problem. Secondly, the MAUVE score demonstrates that SimCTG + contrast search generates texts that are closest to human-written texts in terms of token distribution. Thirdly, among all methods, SimCTG + contrastive search is the only approach that achieves over 0.6 coherence score, showing it produces both high quality and semantically consistent text with respect to the prefix. Lastly, the gen-ppl metric also validates the superiority of SimCTG + contrastive search as it obtains notably better generation perplexity comparing with other approaches.

Moreover, from the results of MLE and Unlikelihood baselines, we see that contrastive search still brings performance boost as compared with greedy and beam search. However, the performance gain still lags behind SimCTG, which demonstrates the necessity of contrastive training. The underlying

Model	Language Modelling Quality				Generation Quality							
	ppl↓	acc↑	rep↓	wrep↓	Method	rep-2↓	rep-3↓	rep-4↓	diversity↑	MAUVE↑	coherence↑	gen-ppl
MLE	24.32	39.63	52.82	29.97	greedy	69.21	65.18	62.05	0.04	0.03	0.587	7.32
					beam	71.94	68.97	66.62	0.03	0.03	0.585	6.42
					nucleus	4.45	0.81	0.43	0.94	0.90	0.577	49.71
					contrastive	44.20	37.07	32.44	0.24	0.18	0.599	9.90
Unlike.	28.57	38.41	51.23	28.57	greedy	24.12	13.35	8.04	0.61	0.69	0.568	37.82
					beam	11.83	5.11	2.86	0.81	0.75	0.524	34.73
					nucleus	4.01	0.80	0.42	0.95	0.87	0.563	72.03
					contrastive	7.48	3.23	1.40	0.88	0.83	0.574	43.61
SimCTG	23.82	40.91	51.66	28.65	greedy	67.36	63.33	60.17	0.05	0.05	0.596	7.16
					beam	70.32	67.17	64.64	0.04	0.06	0.591	6.36
					nucleus	4.05	0.79	0.37	0.94	0.92	0.584	47.19
					contrastive	3.93	0.78	0.31	0.95	0.94	0.610	18.26
Human	-	-	36.19	-	-	3.92	0.88	0.28	0.95	1.00	0.644	24.01

Table 1: Evaluation results on Wikitext-103 test set. “Unlike.” denotes the model trained with unlikelihood objective. ↑ means higher is better and ↓ means lower is better.

reason is that, without using the contrastive objective \mathcal{L}_{CL} (Eq. 2), the token representations obtained by MLE or Unlikelihood are less discriminative (§6.1). Therefore, the degeneration penalty (Eq. 5) of different candidates are less distinguishable and the selection of output is dominated by the model confidence, making contrastive search less effective.

Model	Decoding Method	Coherence	Fluency	Informativeness
Agreement	-	0.51	0.64	0.70
MLE	nucleus	2.92	3.32	3.91
	contrastive	2.78	2.29	2.56
Unlikelihood	nucleus	2.59	3.02	3.58
	contrastive	2.76	2.90	3.35
SimCTG	nucleus	2.96	3.34	3.96
	contrastive	3.25★	3.57★	3.96
SimCTG-large	nucleus	3.01	3.37	3.98
	contrastive	3.33★	3.66★	3.98
Human	-	3.70	3.71	4.21

Table 2: Human evaluation results. ★ results significantly outperforms the results of nucleus sampling with different models (Sign Test with p-value < 0.05).

4.3 Human Evaluation

We also conduct a human evaluation with the help of graders proficient in English from an internal grading platform. We randomly select 200 prefixes with length of 32 from the test set of Wikitext-103. For each prefix, we use different models (MLE, Unlikelihood, and SimCTG) with two decoding methods (nucleus sampling and contrastive search) to generate text continuations with length of 128. To examine the generality of our approach across different model sizes, we include a large size SimCTG which is obtained by fine-tuning GPT-2-large (774M parameters) that consists of 36 Transformer layers with 20 attention heads. All generated results, plus the reference text, are randomly shuffled and evaluated by five graders, which results in 9,000 annotated samples in total. The evaluation follows a 5-point Likert scale (1, 2, 3, 4, or 5) for each of the following features:⁴

- **Coherence:** Whether the generated text is semantically consistent with the prefix.
- **Fluency:** Whether the generated text is fluent and easy to understand.
- **Informativeness:** Whether the generated text is diverse and contains interesting content.

Table 2 presents the human evaluation results, with the first row showing strong inter-annotator agreements as measured by Fleiss’ kappa coefficient Fleiss et al. (1971). Firstly, we see that, directly applying contrastive search with MLE or Unlikelihood model does not yield satisfactory results. This is due to the anisotropic nature of their representation space as discussed in section §4.2. Secondly,

⁴More human evaluation details are provided in Appendix B.

the coherence score of Unlikelihood model is notably lower than MLE and SimCTG, suggesting it generates the most *unlikely* results which is also shown by its generation perplexity (gen-ppl) in Table 1. Furthermore, the results of SimCTG + contrastive search significantly outperforms nucleus sampling with different models in terms of coherence and fluency (Sign Test with p-value < 0.05). Lastly, SimCTG-large + contrastive search achieves the best performance across the board and even performs comparably with human-written text on the fluency metric (Sign Test with p-value > 0.4). This reveals the clear generality of our approach to large size models and future work could focus on extending it to models that contain over billions of parameters such as GPT-3 (Brown et al., 2020).

5 Open-domain Dialogue Generation

To test the generality of our approach across different tasks and languages, we then evaluate our method on the task of open-domain dialogue generation. In this task, given a multi-turn dialogue context (where each turn is an user utterance), the model is asked to generate an adequate response that is semantically consistent with the context. Here, the dialogue context is deemed as the prefix.

Benchmark and Baselines. We conduct experiments on two benchmark datasets from two languages (i.e. Chinese and English). For the Chinese benchmark, we use the LCCC dataset (Wang et al., 2020). For the English Benchmark, we use the DailyDialog dataset (Li et al., 2017).

We compare the GPT-2 models fine-tuned with SimCTG and MLE.⁵ Specifically, for the Chinese benchmark, we use a publicly available Chinese GPT-2 (Zhao et al., 2019).⁶ Same as in section §4, during training, we use a batch size of 128 and truncate the training samples to a maximum length of 256. On the LCCC dataset, we train (i.e. fine-tune) the models for 40k steps. As for the DailyDialog dataset, due to its smaller dataset size, we train the models for 5k steps. For optimization, we use Adam optimizer and a learning rate of 2e-5 (with warm-up ratio of 10%).

For each model, we use four decoding methods, including (1) greedy search; (2) beam search (beam size of 10); (3) nucleus sampling ($p = 0.95$); and (4) contrastive search ($k = 5$, $\alpha = 0.6$).

Evaluation. We rely on human evaluation to assess the model performance. Same as in section §4.3, we randomly select 200 dialogue contexts from the test set and ask five annotators to evaluate the generated responses plus the reference response in three dimensions: (i) coherence, (ii) fluency; and (iii) informativeness. The scores follow a 5-point Likert scale (1, 2, 3, 4, or 5).

Model	Method	LCCC			DailyDialog		
		Coherence	Fluency	Informativeness	Coherence	Fluency	Informativeness
Agreement	-	0.73	0.61	0.57	0.64	0.60	0.55
MLE	greedy	3.01	3.27	1.97	3.28	3.51	2.92
	beam	2.60	2.90	1.55	3.16	3.43	2.78
	nucleus	2.78	3.55	2.64	2.67	3.58	3.42
	contrastive	3.28★	3.84★	3.06★	3.27	3.41	2.82
SimCTG	greedy	3.04	3.32	2.01	3.31	3.50	2.94
	beam	2.57	2.93	1.59	3.19	3.45	2.79
	nucleus	2.84	3.58	2.72	2.75	3.59	3.39
	contrastive	3.32★	3.96★	3.13★	3.73★	3.85★	3.46
Human	-	3.42	3.76	3.20	4.11	3.98	3.74

Table 3: Human evaluation results. ★ results significantly outperforms the results of greedy search, beam search, and nucleus sampling with different models. (Sign Test with p-value < 0.05).

Table 3 shows the evaluation results where the first row shows strong inter-annotator agreements as measured by Fleiss’ kappa coefficient. On both datasets, we see that SimCTG + contrastive search significantly outperforms other methods on various metrics, suggesting that our approach is generalizable to different languages and tasks. It is worth emphasizing that, on the LCCC benchmark, SimCTG + contrastive search surprisingly outperforms the human performance on the fluency metric, while performing comparably on the coherence and informativeness metrics (Sign Test with p-value >

⁵We acknowledge that there are other GPT-like models, such as Zhang et al. (2020) and Thoppilan et al. (2022), that are designed for dialogue generation. We leave the test of our approach on these models to future work.

⁶<https://huggingface.co/uer/gpt2-chinese-cluecorpussmall>

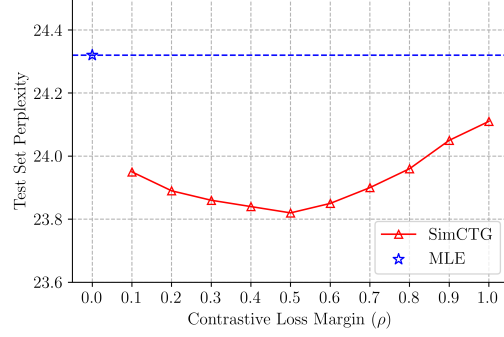
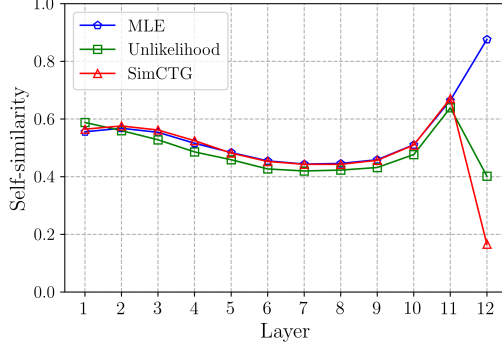


Figure 2: Layer-wise representation self-similarity. Figure 3: The effect of contrastive margin ρ .

0.4). Moreover, even **without** contrastive training, the MLE model performs significantly better when using contrastive search. This is due to the intrinsic property of Chinese language model for which the MLE objective can already yield a representation space that displays a high level of isotropy, making contrastive search directly applicable. This finding is particularly attractive as it reveals the potential applicability of contrastive search on off-the-shelf (i.e. without contrastive training) language models for certain languages such as Chinese.⁷

6 Further Analysis

In this section, we provide and discuss in-depth analyses of our approach.

6.1 Token Representation Self-similarity

To analyze the token representations learned by SimCTG, we follow Ethayarajh (2019) and define the averaged self-similarity of token representations within a text sequence \mathbf{x} as

$$\text{self-similarity}(\mathbf{x}) = \frac{1}{|\mathbf{x}| \times (|\mathbf{x}| - 1)} \sum_{i=1}^{|\mathbf{x}|} \sum_{j=1, j \neq i}^{|\mathbf{x}|} \frac{h_{x_i}^\top h_{x_j}}{\|h_{x_i}\| \cdot \|h_{x_j}\|}, \quad (6)$$

where h_{x_i} and h_{x_j} are the token representations of x_i and x_j produced by the model. Intuitively, a lower $\text{self-similarity}(\mathbf{x})$ indicates the representations of distinct tokens within the sequence \mathbf{x} are less similar to each other, therefore being more discriminative.

We use texts from the test set of Wikitext-103 and compute the self-similarity of token representations over different layers for different models (i.e. MLE, Unlikelihood, and SimCTG). Figure 2 plots the results averaged over all samples. We see that, in the intermediate layers, the self-similarity scores of different models are relatively the same. In contrast, at the output layer (layer 12), SimCTG’s self-similarity becomes notably lower than other baselines. We note that the Unlikelihood model also yields more discriminative representations than MLE, but its language model accuracy is lower than MLE and SimCTG as shown in Table 1. On the other hand, SimCTG obtains the most discriminative and isotropic representations while maintaining the best language model accuracy, which further validates the clear advantage of our proposed approach.

6.2 The Effect of Contrastive Loss Margin

Next, we analyze the effect of contrastive loss margin ρ (Eq. 2) on the performance of the language model. To this end, we fine-tune the GPT-2 by varying the value of ρ from 0.1 to 1.0 and measure the model perplexity on the Wikitext-103 test set. Figure 3 plots the results of SimCTG with different ρ along with the result of the MLE baseline. Note that, as described in section §3.1, SimCTG is equivalent to MLE when $\rho = 0$. From Figure 3, we observe that the contrastive training always helps to improve the model perplexity as compared with MLE. However, when ρ is either too small (e.g.

⁷We provide more analysis for this aspect and several generated examples on the LCCC benchmark in Appendix C and D, respectively.

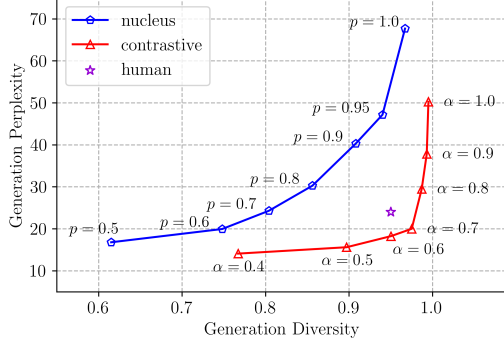


Figure 4: Contrastive search vs nucleus sampling.

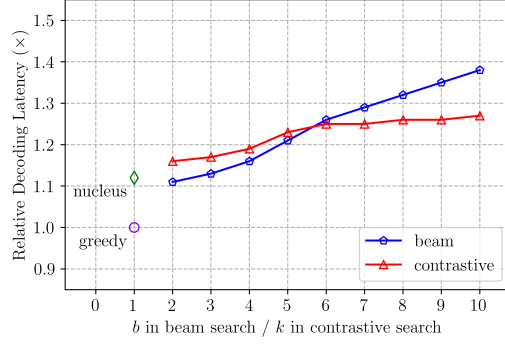


Figure 5: Inference latency comparison.

0.1) or large (e.g. 1.0), the learned representation space of the model would be either less or too isotropic, leading to a sub-optimal perplexity. In our experiments, the most suitable margin $\rho = 0.5$.

6.3 Contrastive Search versus Nucleus Sampling

In this part, we provide a more in-depth comparison between our proposed contrastive search and the current state-of-the-art decoding method, nucleus sampling. To this end, we compare the generated results of SimCTG using these two decoding methods. Specifically, we vary the probability p for nucleus sampling and the α (Eq. 5) for contrastive search to generate results using prefixes from Wikitext-103 test set.⁸ We evaluate the generated results from two aspects: (1) generation diversity and (2) perplexity of the generated text (gen-ppl). Both metrics are described in section §4.1.2. Figure 4 plots the results of different decoding methods along with the human performance. For nucleus sampling, when p is small (i.e. $p \leq 0.7$), its generation perplexity is comparable to that of human. However, the diversity is notably lower than human performance, meaning it stuck in undesirable repetition loops (Holtzman et al., 2020). On the other hand, when p is large (i.e. $p \geq 0.95$), the generation diversity is close to that of human but the generation perplexity is significantly higher. Such high perplexity means the generated text is very *unlikely*, therefore being low quality. As for contrastive search, when $\alpha \in [0.5, 0.8]$, it yields generation diversity and perplexity that are both comparable to human performance. These results demonstrate the superiority of contrastive search as it better balances the trade-off between the generation diversity and perplexity.

6.4 Decoding Latency Comparison

We compare the decoding latency of different decoding methods with SimCTG. For beam search and contrastive search, we vary the beam width b and the k in Eq. 5. The decoding latency is measured by generating fixed length text continuations on Wikitext-103 test cases with a batch size of 1. In Figure 5, we show the averaged relative decoding latency of different methods. From the results, we see that the greedy search is the fastest method and the latency of different methods are generally comparable with each other. Comparing contrastive search with beam search, when b and k are small (i.e. ≤ 6), their latency are nearly identical. When b and k gets larger (i.e. > 6), the contrastive search becomes faster than beam search. These comparisons with other widely used decoding methods illustrate the practical usage of the proposed contrastive search.

6.5 Case Study

In Table 4, we present generated examples of SimCTG with different decoding methods given a specific prefix.⁹ From the results, we see that beam search produces undesirable sequence-level repetitions, resulting in low diversity and low generation perplexity. On the other hand, in the prefix, the person “Buchanan” *criticizes* the game. However, the result from nucleus sampling displays a contradicted semantic, resulting in a low coherence score as well as a high generation perplexity.

⁸For contrastive search, we only vary the value of α and keep k constant to 8 as described in section §4.

⁹More generated examples of SimCTG can be found in Appendix H.

search with stochastic sampling methods. For instance, given the prefix, we could first generate a few tokens (e.g. 1~3 tokens) with nucleus sampling. Then, we switch to contrastive search for the remaining steps.¹⁰ (2) Our approach is architecture agnostic and can be applied to any generation model. Future research could focus on adapting it to other tasks than open-ended text generation (i.e. constrained text generation), such as machine translation and document summarization.

Acknowledgments

The first author would like to thank Jialu Xu and Huayang Li for the insightful discussions and support.

References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1409.0473>.
- Sourya Basu, Govardana Sachitanandam Ramachandran, Nitish Shirish Keskar, and Lav R. Varshney. Mirostat: a neural text decoding algorithm that directly controls perplexity. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL https://openreview.net/forum?id=W1G1JZEIy5_.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 2020. URL <http://proceedings.mlr.press/v119/chen20j.html>.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc Viet Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 2978–2988. Association for Computational Linguistics, 2019. doi: 10.18653/v1/p19-1285. URL <https://doi.org/10.18653/v1/p19-1285>.
- Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca J. Passonneau, and Rui Zhang. Container: Few-shot named entity recognition via contrastive learning. *CoRR*, abs/2109.07589, 2021. URL <https://arxiv.org/abs/2109.07589>.
- Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander H. Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W. Black, Alexander I. Rudnicky, Jason Williams, Joelle Pineau, Mikhail S. Burtsev, and Jason Weston. The second conversational intelligence challenge (convai2). *CoRR*, abs/1902.00098, 2019. URL <http://arxiv.org/abs/1902.00098>.

¹⁰We provide some preliminary experiment results in Appendix G.

- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: pure attention loses rank doubly exponentially with depth. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 2793–2803. PMLR, 2021. URL <http://proceedings.mlr.press/v139/dong21a.html>.
- Kawin Ethayarajh. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and GPT-2 embeddings. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 55–65. Association for Computational Linguistics, 2019. doi: 10.18653/v1/D19-1006. URL <https://doi.org/10.18653/v1/D19-1006>.
- Angela Fan, Mike Lewis, and Yann N. Dauphin. Hierarchical neural story generation. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 889–898. Association for Computational Linguistics, 2018. doi: 10.18653/v1/P18-1082. URL <https://aclanthology.org/P18-1082/>.
- J.L. Fleiss et al. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382, 1971.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6894–6910. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.552. URL <https://doi.org/10.18653/v1/2021.emnlp-main.552>.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=rygGQyrFvH>.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 4904–4916. PMLR, 2021. URL <http://proceedings.mlr.press/v139/jia21b.html>.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=HklBjCEKvH>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush. Open-nmt: Open-source toolkit for neural machine translation. In Mohit Bansal and Heng Ji, editors, *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, System Demonstrations*, pages 67–72. Association for Computational Linguistics, 2017. doi: 10.18653/v1/P17-4012. URL <https://doi.org/10.18653/v1/P17-4012>.
- Wouter Kool, Herke van Hoof, and Max Welling. Stochastic beams and where to find them: The gumbel-top-k trick for sampling sequences without replacement. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine*

- Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 3499–3508. PMLR, 2019. URL <http://proceedings.mlr.press/v97/kool19a.html>.
- Tian Lan, Deng Cai, Yan Wang, Yixuan Su, Xian-Ling Mao, and Heyan Huang. Exploring dense retrieval for dialogue response selection. *CoRR*, abs/2110.06612, 2021. URL <https://arxiv.org/abs/2110.06612>.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting objective function for neural conversation models. In Kevin Knight, Ani Nenkova, and Owen Rambow, editors, *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119. The Association for Computational Linguistics, 2016. doi: 10.18653/v1/n16-1014. URL <https://doi.org/10.18653/v1/n16-1014>.
- Piji Li, Haisong Zhang, Xiaojiang Liu, and Shuming Shi. Rigid formats controlled text generation. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 742–751. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-main.68. URL <https://doi.org/10.18653/v1/2020.acl-main.68>.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. Dailydialog: A manually labelled multi-turn dialogue dataset. In Greg Kondrak and Taro Watanabe, editors, *Proceedings of the Eighth International Conference on Natural Language Processing, IJCNLP 2017, Taipei, Taiwan, November 27 - December 1, 2017 - Volume 1: Long Papers*, pages 986–995. Asian Federation of Natural Language Processing, 2017. URL <https://aclanthology.org/I17-1099/>.
- Fangyu Liu, Ivan Vulic, Anna Korhonen, and Nigel Collier. Fast, effective, and self-supervised: Transforming masked language models into universal lexical and sentence encoders. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 1442–1459. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.109. URL <https://doi.org/10.18653/v1/2021.emnlp-main.109>.
- Yixin Liu and Pengfei Liu. Simcls: A simple framework for contrastive learning of abstractive summarization. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 2: Short Papers), Virtual Event, August 1-6, 2021*, pages 1065–1072. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.acl-short.135. URL <https://doi.org/10.18653/v1/2021.acl-short.135>.
- Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan Le Bras, Lianhui Qin, Youngjae Yu, Rowan Zellers, Noah A. Smith, and Yejin Choi. Neurologic a*esque decoding: Constrained text generation with lookahead heuristics. *CoRR*, abs/2112.08726, 2021a. URL <https://arxiv.org/abs/2112.08726>.
- Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Neurologic decoding: (un)supervised neural text generation with predicate logic constraints. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 4288–4299. Association for Computational Linguistics, 2021b. doi: 10.18653/v1/2021.naacl-main.339. URL <https://doi.org/10.18653/v1/2021.naacl-main.339>.
- Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In Lluís Màrquez, Chris Callison-Burch, Jian Su, Daniele Pighin, and Yuval Marton, editors, *Proceedings of the 2015 Conference on Empirical Methods in Natural*

- Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 1412–1421. The Association for Computational Linguistics, 2015. doi: 10.18653/v1/d15-1166. URL <https://doi.org/10.18653/v1/d15-1166>.
- Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. Typical decoding for natural language generation, 2022.
- Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song. COCO-LM: correcting and contrasting text sequences for language model pretraining. *CoRR*, abs/2102.08473, 2021. URL <https://arxiv.org/abs/2102.08473>.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=Byj72udxe>.
- Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. MAUVE: Measuring the gap between neural text and human text using divergence frontiers. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=Tqx7nJp7PR>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. URL https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 2021. URL <http://proceedings.mlr.press/v139/radford21a.html>.
- Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with pointer-generator networks. In Regina Barzilay and Min-Yen Kan, editors, *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1073–1083. Association for Computational Linguistics, 2017. doi: 10.18653/v1/P17-1099. URL <https://doi.org/10.18653/v1/P17-1099>.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics, 2016. doi: 10.18653/v1/p16-1162. URL <https://doi.org/10.18653/v1/p16-1162>.
- Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, and Sergey Levine. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*, pages 1134–1141. IEEE, 2018. doi: 10.1109/ICRA.2018.8462891. URL <https://doi.org/10.1109/ICRA.2018.8462891>.
- Yixuan Su, Deng Cai, Yan Wang, David Vandyke, Simon Baker, Piji Li, and Nigel Collier. Non-autoregressive text generation with pre-trained language models. In Paola Merlo, Jörg Tiedemann, and Reut Tsarfay, editors, *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pages 234–243. Association for Computational Linguistics, 2021a. doi: 10.18653/v1/2021.eacl-main.18. URL <https://doi.org/10.18653/v1/2021.eacl-main.18>.
- Yixuan Su, Deng Cai, Qingyu Zhou, Zibo Lin, Simon Baker, Yunbo Cao, Shuming Shi, Nigel Collier, and Yan Wang. Dialogue response selection with hierarchical curriculum learning. In Chengqing

- Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 1740–1751. Association for Computational Linguistics, 2021b. doi: 10.18653/v1/2021.acl-long.137. URL <https://doi.org/10.18653/v1/2021.acl-long.137>.
- Yixuan Su, Fangyu Liu, Zaiqiao Meng, Tian Lan, Lei Shu, Ehsan Shareghi, and Nigel Collier. Tacl: Improving BERT pre-training with token-aware contrastive learning. *CoRR*, abs/2111.04198, 2021c. URL <https://arxiv.org/abs/2111.04198>.
- Yixuan Su, Zaiqiao Meng, Simon Baker, and Nigel Collier. Few-shot table-to-text generation with prototype memory. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, pages 910–917. Association for Computational Linguistics, 2021d. doi: 10.18653/v1/2021.findings-emnlp.77. URL <https://doi.org/10.18653/v1/2021.findings-emnlp.77>.
- Yixuan Su, Lei Shu, Elman Mansimov, Arshit Gupta, Deng Cai, Yi-An Lai, and Yi Zhang. Multi-task pre-training for plug-and-play task-oriented dialogue system. *CoRR*, abs/2109.14739, 2021e. URL <https://arxiv.org/abs/2109.14739>.
- Yixuan Su, David Vandyke, Sihui Wang, Yimai Fang, and Nigel Collier. Plan-then-generate: Controlled data-to-text generation via planning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Findings of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 16-20 November, 2021*, pages 895–909. Association for Computational Linguistics, 2021f. doi: 10.18653/v1/2021.findings-emnlp.76. URL <https://doi.org/10.18653/v1/2021.findings-emnlp.76>.
- Yixuan Su, Yan Wang, Deng Cai, Simon Baker, Anna Korhonen, and Nigel Collier. PROTOTYPE-TO-STYLE: dialogue generation with style-aware editing on retrieval memory. *IEEE ACM Trans. Audio Speech Lang. Process.*, 29:2152–2161, 2021g. doi: 10.1109/TASLP.2021.3087948. URL <https://doi.org/10.1109/TASLP.2021.3087948>.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 3104–3112, 2014. URL <https://proceedings.neurips.cc/paper/2014/hash/a14ac55a4f27472c5d894ec1c3c743d2-Abstract.html>.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven 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, Igor Krivokon, Will Rusch, Marc Pickett, Kathleen S. Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. Lamda: Language models for dialog applications. *CoRR*, abs/2201.08239, 2022. URL <https://arxiv.org/abs/2201.08239>.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018. URL <http://arxiv.org/abs/1807.03748>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.

- Yida Wang, Pei Ke, Yinhe Zheng, Kaili Huang, Yong Jiang, Xiaoyan Zhu, and Minlie Huang. A large-scale chinese short-text conversation dataset. In Xiaodan Zhu, Min Zhang, Yu Hong, and Ruifang He, editors, *Natural Language Processing and Chinese Computing - 9th CCF International Conference, NLPCC 2020, Zhengzhou, China, October 14-18, 2020, Proceedings, Part I*, volume 12430 of *Lecture Notes in Computer Science*, pages 91–103. Springer, 2020. doi: 10.1007/978-3-030-60450-9_8. URL https://doi.org/10.1007/978-3-030-60450-9_8.
- Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural text generation with unlikelihood training. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=SJeYe0NtvH>.
- Sam Wiseman, Stuart M. Shieber, and Alexander M. Rush. Challenges in data-to-document generation. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 2253–2263. Association for Computational Linguistics, 2017. doi: 10.18653/v1/d17-1239. URL <https://doi.org/10.18653/v1/d17-1239>.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface’s transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771, 2019. URL <http://arxiv.org/abs/1910.03771>.
- Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga, Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I. Wang, Victor Zhong, Bailin Wang, Chengzu Li, Connor Boyle, Ansong Ni, Ziyu Yao, Dragomir R. Radev, Caiming Xiong, Lingpeng Kong, Rui Zhang, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. Unifiedskg: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. *CoRR*, abs/2201.05966, 2022. URL <https://arxiv.org/abs/2201.05966>.
- Dani Yogatama, Cyprien de Masson d’Autume, and Lingpeng Kong. Adaptive semiparametric language models. *Trans. Assoc. Comput. Linguistics*, 9:362–373, 2021. URL <https://transacl.org/ojs/index.php/tacl/article/view/2693>.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. DIALOGPT : Large-scale generative pre-training for conversational response generation. In Asli Celikyilmaz and Tsung-Hsien Wen, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, ACL 2020, Online, July 5-10, 2020*, pages 270–278. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-demos.30. URL <https://doi.org/10.18653/v1/2020.acl-demos.30>.
- Zhe Zhao, Hui Chen, Jinbin Zhang, Xin Zhao, Tao Liu, Wei Lu, Xi Chen, Haotang Deng, Qi Ju, and Xiaoyong Du. Uer: An open-source toolkit for pre-training models. *EMNLP-IJCNLP 2019*, page 241, 2019.

A Related Work

Neural Text Generation. Neural text generation is a core component in many NLP applications. It can be generally categorized into two classes (1) constrained and (2) open-ended generation.

Constrained generation tasks are always defined over a set of (input, output) pairs, where the output is a transformation of the input following specific constraints. Some typical examples include machine translation (Sutskever et al., 2014; Bahdanau et al., 2015; Luong et al., 2015), text summarization (See et al., 2017; Su et al., 2021a), and data-to-text generation (Wiseman et al., 2017; Su et al., 2021f,d; Xie et al., 2022). As the output is tightly scoped by the input, the generation of repetition and unnaturalness are not that problematic, therefore maximization-based decoding methods such as beam search generally perform well. Still, different variants of beam search have been explored to further improve the model performance in constrained generation tasks (Klein et al., 2017; Kool et al., 2019; Lu et al., 2021b,a).

Open-ended text generation, on the other hand, impose less constrain on the generated text. It aims at producing text that is natural, coherent and informative with respect to the human-written prefix (i.e. context). Several typical applications include story generation (Fan et al., 2018), contextual text completion (Radford et al., 2019), and dialogue systems (Su et al., 2021g,e). However, due to the challenges posed by the increased level of freedom, conventional maximization-based decoding methods (e.g. greedy and beam search) often produce undesirable repetition and unnaturalness in the generated text. To alleviate model degeneration, different sampling approaches (Fan et al., 2018; Holtzman et al., 2020; Meister et al., 2022) have been proposed to generate the text by drawing samples from less likely vocabularies. Welleck et al. (2020) tackled model degeneration from another perspective by introducing unlikelihood objective into the training of the language model.

Contrastive Learning. Generally, contrastive learning methods aim to teach the model to distinguish observed data points from fictitious negative samples. They have been widely applied to various research areas. In the field of computer vision, contrastive learning has been shown to benefit tasks like image (van den Oord et al., 2018) and video (Sermanet et al., 2018) representation learning. Chen et al. (2020) proposed a simple framework, SimCLR, for learning contrastive visual representations. Recently, Radford et al. (2021); Jia et al. (2021) applied contrastive learning for the pre-training of language-image models.

In the field of NLP, contrastive learning has recently gained much more attention. Numerous contrastive approaches have been proposed to learn better token-level (Su et al., 2021c), sentence-level (Meng et al., 2021; Liu et al., 2021; Gao et al., 2021), and discourse-level (Su et al., 2021b; Lan et al., 2021) representations. Beyond representation learning, contrastive learning has also been applied to other NLP applications, such as name entity recognition (NER) (Das et al., 2021) and document summarization (Liu and Liu, 2021).

Our work, to the best of our knowledge, is the first effort on applying contrastive learning to address neural text degeneration. We hope our findings could facilitate future research in this area.

B Human Evaluation Guidelines

Given the human-written prefix, please evaluate the system’s result with respect to the following features: (1) Coherence; (2) Fluency; and (3) Informativeness. In the following, we provide some guidelines regarding how to judge the quality of the system’s result in terms of different features.

B.1 Coherence

This metric measures whether the system’s result is semantically and factually consistent with the human-written prefix. The definitions of different scores are:

- **[5]:** The system’s result is perfectly in line with the semantic meaning defined by the prefix. And all its content is factually supported by or can be logically inferred from the prefix.
- **[4]:** The system’s result is very related to the prefix but with some minor errors that does not affect its overall relevance with respect to the prefix.
- **[3]:** The system’s result is, to some extent, relevant to the prefix with some errors that display minor semantic inconsistency or contradiction.

- [2]: At the first glance, the system’s result seems to be related to the prefix. But with careful inspection, the semantic inconsistency can be easily spotted.
- [1]: The system’s result is obviously off-the-topic or it is semantically contradicted to the content contained in the prefix.

B.2 Fluency

This metric measures the fluency of the system’s result. The definitions of different scores are:

- [5]: The system’s result is human-like, grammatically correct, and very easy to understand.
- [4]: Choose this score when you are hesitant between the score 3 and score 5.
- [3]: The system’s result contains minor errors but they do not affect your understanding.
- [2]: Choose this score when you are hesitant between the score 1 and score 3.
- [1]: The system’s result does not make sense and it is unreadable.

B.3 Informativeness

This metric measures the diversity, informativeness, and interestingness of the system’s result. The definitions of different scores are:

- [5]: The system’s result is very informative and contains novel content. In addition, it displays a high level of diversity and it is enjoyable to read.
- [4]: Choose this score when you are hesitant between the score 3 and score 5.
- [3]: The system’s result contains some new information and it displays a certain level of diversity.
- [2]: Choose this score when you are hesitant between the score 1 and score 3.
- [1]: The system’s result is dull, repetitive, and does not have new information. All its content has already been provided in the prefix.

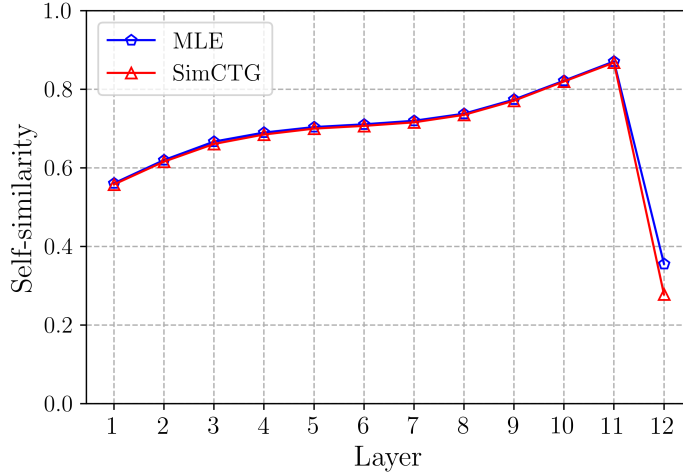


Figure 7: Layer-wise self-similarity of Chinese language models.

C Self-similarity of Chinese Language Models

We follow the same procedure as described in section §6.1 to measure the token self-similarity of Chinese language models. Specifically, we use the test set of LCCC benchmark and compute the self-similarity of the model. Figure 7 plots the layer-wise token self-similarity of the MLE and SimCTG models. We see that in all layers (including the final layer), the MLE model displays a similar self-similarity with respect to SimCTG. This observation is quite different from what we

see from English language models as shown in Figure 2, where the self-similarities of SimCTG and MLE are notably different in the final layer. This discrepancy comes from the intrinsic property of different languages. For English, current state-of-the-art methods always represent the text into subword units, such as BPE (Sennrich et al., 2016)) and the same subword could be shared by many different contexts. Thus, the representations of distinct subwords become less distinguishable which naturally leads to the anisotropy in their representations. On the other hand, languages like Chinese are naturally represented by basic units (i.e. characters). Such natural unit boundary of the text alleviates the share of characters in different contexts. As a result, even the vanilla MLE objective can obtain a representation space that displays a high level of isotropy.

This isotropic property of Chinese language model is particularly attractive as the contrastive search can be directly applied, as shown in Table 3, even **without** contrastive training. In addition, we expect contrastive search could be used on off-the-shelf language models that are trained with MLE in other languages whose texts are naturally tokenized by characters (e.g. Korean and Japanese). This remains to be rigorously tested in future work.

D Generated Examples on Open-domain Dialogue Generation

Dialogue Context	Generated Response	Reference Response
A: 老铁家好吃费三不好吃	哈哈我觉得还可以，就是有点油	我不挑食
A:话说红海真的好看!!! 准备二刷了!!! 但求多排场!!! B:我明天二刷!	哈哈好的! 等我回学校一起吃火锅!!!	必须的! 上班之后可能就没空了呢~
A:你家真有钱 B:从何说起? A:可以买粮	好吧我承认我是土豪	你家的粮自己种的?
A:为何突然伤感?	因为没什么感觉。不知道为什么就觉得自己好孤单	因为我就要回学校了，他舍不得我
A:刺猬很可爱! 以前别人送了只没养，味儿太大!! B:是很可爱但是非常臭 A:是啊，没办法养 B:那个怎么养哦不会扎手吗	我觉得还好，就是有点臭	不会
A:杂糅太多反而摸不到主题了是吗	这就是一个故事，不知道怎么编出来的	不是，故事编不好

Table 5: Case study on the LCCC dataset.

In Table 5, we show some generated responses of our approach (i.e. SimCTG + contrastive search) plus the reference response on examples from the test set of the Chinese LCCC benchmark. We see that, given the dialogue context, our approach is able to generate responses that are both grammatically fluent and semantically consistent with the dialogue context. These results further demonstrate the generality of our approach across different languages and tasks.

E Training Efficiency Comparison

In this part, we compare the training efficiency of different methods (i.e. MLE, Unlikelihood, and SimCTG). To this end, we compute the total floating point operations (FLOPs) required for the training of different models on Wikitext-103. The details of training setup are provided in section §4. Table 6 shows the results, from which we see that SimCTG is more efficient than the unlikelihood method. Comparing with MLE, SimCTG only introduces an negligible 1.48% extra computational overhead, which verifies the practical usage of SimCTG.

	MLE	Unlikelihood	SimCTG
Train FLOPs	8.08e16	8.91e16	8.20e16
Parameters	117M	117M	117M

Table 6: Training efficiency comparison.

F Hyperparameter Analysis on Contrastive Search

Here, we present a more detailed hyperparameter analysis of the proposed contrastive search. Specifically, we simultaneously vary the value of k and α in Eq. 5. k is chosen from $\{5, 8, 10\}$ and α is chosen from $\{0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$. For evaluation, we report the generation diversity and generation perplexity on the test set of Wikitext-103. The results are plotted in Figure 8. We see that, when k is constant, the increase of α generally increases the generation diversity and generation perplexity. When α is constant, a larger k also leads to increased generation diversity and generation perplexity. Nonetheless, for different k , the overall trends are relatively the same and the value of α has more impact on the generated results. In practice, our recommended selection range of k and α are $k \in [5, 10]$ and $\alpha \in [0.5, 0.8]$, as these settings produce results that are more similar to human-written texts as judged by generation diversity and generation perplexity.

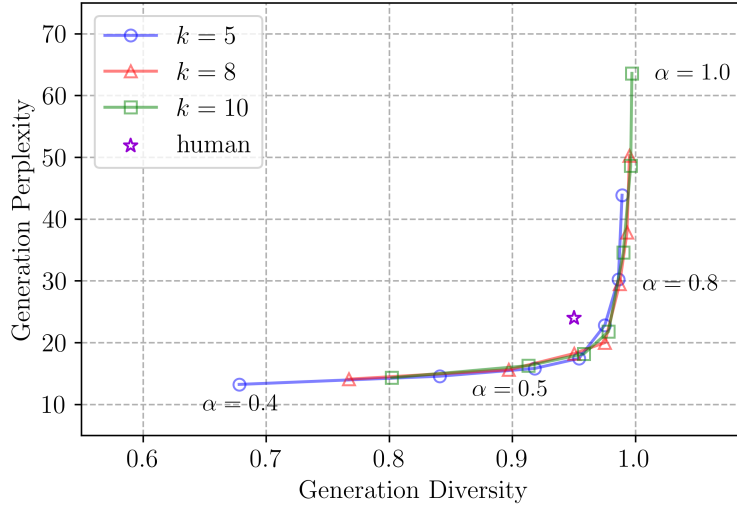


Figure 8: Hyperparameter analysis of contrastive search.

G Diverse Contrastive Search

In this part, we present a diverse version of contrastive search (i.e. diverse contrastive search) which is described in section §7. Specifically, given the prefix with length of 32, we first generate 2 tokens using nucleus sampling with $p = 0.95$, then we use contrastive search to generate the remaining 126 tokens (i.e. 128 generated tokens in total).

Table 7 shows three generated results with diverse contrastive search using the same prefix as in Table 4. We see that only sampling 2 tokens at the start is enough to produce a diverse set of results. In future work, we will investigate other more sophisticated extensions of contrastive search.

H More Generated Examples of SimCTG + Contrastive Search

In Table 8, we provide more generated examples of SimCTG + contrastive search based on prefixes from Wikitext-103. The details of the decoding procedure are described in section §4.

prefix	Butt criticized Donald's controls in certain situations in the game, as well as the difficulty of some levels and puzzles. Buchanan also criticized the controls, calling
<i>Contrastive Search</i>	
result	them "unimpressive" and a "nightmare" of an experience to play with players unfamiliar with Tetris. On the other hand, his opinion was shared by other reviewers, and some were critical of the game's technical design for the Wii version of Tetris. In addition, Tintin's review included a quote from Roger Ebert, who said that Tetris was better than the original game due to its simplicity and ease of play. Ebert's comments were included in the game's DVD commentary, released on March 22, 2010. It is unclear if any of the video commentary was taken from ...
<i>Diverse Contrastive Search</i>	
[1]	the game "unplayable" and adding that, "It's a game in which you have to be careful not to fall into the trap of thinking too much about your opponents' strengths and weaknesses." On the other hand, his review of the game, published in the New York Times, found little to complain about in comparison to some of the other games in the series. In particular, The Escapist's Brian O'Toole wrote that the game was an "absolute joy", and "one of the best-selling games of all time". O'Toole concluded by saying that although ...
[2]	it "a complete waste of time" and "unplayable". On the other hand, his review of Baldur's Gate II was positive, with Buchanan commenting that, "Baldur's Gate II is an adventure game in its own right, full of fun and challenge that makes you want to go back to the first game in your life." Buchanan felt that there were too many elements in the game for players to enjoy without some level-playing to be enjoyable at the same time. He concluded by saying that Baldur's Gate II's controls were well-balanced, and that players ...
[3]	the choice of "a simple jump button to perform a 'jump-and-a-bop' or more complex 'jump-and-a-bop'" an error and a waste of time. On the other hand, Tintin was critical of the game's design, writing that there was "too much going on" at the beginning of the game, and "not enough time" in the final cutscene for the player to make it through the game at all. He felt that the gameplay was lacking in some areas, such as the ...

Table 7: Generated results of SimCTG with diverse contrastive search.

prefix	The Olmec colossal heads are at least seventeen monumental stone representations of human heads sculpted from large basalt boulders. The heads date from at least before
result	the Roman-era, and have been interpreted as representing a group of men and women that lived in the city of Arundel, the same site on which Alexander the Great's tomb was built. In addition, there is evidence to suggest that the city was part of the empire of the Persians and their conquests in Asia Minor during the 7th century BC. This may have been an intentional move to avoid confusion with the Roman city of Acre, which is located in the foothills of the mountain range near the Roman fort of Caesarea. Other evidence points to the construction ...
prefix	Although Rifenburg finished fourth among midwestern Heisman voters in 1948, he did not finish among the top eight. By comparison, Notre Dame end Leon Hart won
result	the Heisman Trophy at the end of the 1948 season after a 4 – 0 victory over Michigan on the final day of the regular season. Hart was named to the Big Ten All-Pac-10 first team, and his Heisman Trophy was awarded to Notre Dame's Bobby Bowden as the team's all-time leader in passing yards (1,000) and completions (2,065). The game ended in a 42 – 21 victory for Notre Dame, which led the nation in total offense (4,000 yards) and touchdowns ...
prefix	Kujira developed from a broad area of disturbed weather as a tropical depression on April 9 well removed from any landmasses. Shortly after development, Kuj
result	ira became extratropical, which would have caused the storm to turn to the west-northwest in the Gulf of Mexico on April 10. However, due to the proximity of the cyclone to an area of low pressure, Kujira's winds were less than 10 mph (15 km / h) at the time of landfall. By April 11, Hurricane Hunters estimated that there were 5,000 people in the area, of whom 1,000 were located in Puerto Rico. The storm weakened to a tropical storm early on April 12, and dissipated the next day. Later that day ...
prefix	The island is composed of a variety of volcanic and sedimentary rock formations from the Devonian period. At that time the Scottish landmass formed part of the Old
result	North Sea, and it was inhabited by an extensive and diverse population of sea-dwelling birds that had lived in the islands for many centuries. There is evidence of human habitation in the Isles of Scilly and St Albans, and some archaeological evidence is in the form of stone mounds dating to the 12th century. The most northerly of these is the Murchison Formation, which is located south of the island's mouth and extends to the northwest as far as Aberdeenshire. Other sites include the Isle of Man, Isle of Wight, Isle of Orkney, The Loch Ness Monster ...
prefix	A report published in the American magazine Time a fortnight later described the incident as " a genuine horror story, calculated to make the most alarming of Rhodesian doomsday
result	scenarios ". In an interview with BBC Radio 4's Tom Pertwee, he said : " There's no evidence that anyone in Rhodesia, other than the police, had done anything to provoke it ". On 20 April 2013, the government announced that it had found no evidence of involvement by Rhodesia in the attack, and would continue its investigation as soon as possible. The Independent Commission Against Corruption (ICAO) was appointed to investigate the incident in December 2013, and is the first body to do so in the UK. ICAO Director-General Mark Rowley said it was " the most ...
prefix	The American poet Ezra Pound was introduced to the group in April 1909 and found that their ideas were close to his own. In particular, Pound's studies of
result	German Expressionism (1891 – 1903) helped him develop a critical attitude towards the subject, and he had the opportunity to speak at the first meeting of the group on 25 April 1909, at New York's Waldorf-Astoria Hotel. This was the first time Pound had been invited to speak at a German Expressionist conference, and it proved to be a success. Pound became an admirer of Expressionism from the outside, and was able to attract many of its members by providing them with information and writing, such as letters to German Expressionists in Germany, as well as photographs of their works ...
prefix	The Olmec civilization developed in the lowlands of southeastern Mexico between 1500 and 400 BC. The Olmec heartland lies on the Gulf Coast of Mexico
result	, about 1,500 miles (2,500 km) west of Acapulco County, where there is a population of 1,500 people. This is the largest city in the state, and it is home to the city's most famous tourist attraction, the El Centro de las Amazonas, which has been named by Mexican historian José Martínez Sánchez for its beauty and architecture. The city, along with its surroundings, were important in the development of the Maya civilization, and in later Mesoamerica and Latin America as well. Today, Mexico's ...
prefix	The discovery of a colossal head at Tres Zapotes in the nineteenth century spurred the first archaeological investigations of Olmec culture by Matthew Stirling in 1938.
result	He showed that the remains of Olmec were in good condition, and suggested that there was something to be said for the fact that they had lived long enough to be recognized as an important part of the culture and its history. Stirling's work was the first to establish what was to become known as the Peruvian-Spanish language, or Peru : the language of the people of Peru and Central America. In doing so, he drew on archaeological evidence and other sources to support his theory of the cultural origins of the Peruvian people. This is reflected in the name " Peruvian ", which stands for Peruvian ...

Table 8: More generated examples of SimCTG + contrastive search.