

Multi-Objective Linguistic Control of Large Language Models

Dang Nguyen, Jiuhai Chen, Tianyi Zhou

University of Maryland, College Park

{dangmn, jchen169, tianyi}@umd.edu

Project: <https://github.com/tianyi-lab/mctune>

Abstract

Large language models (LLMs), despite their breakthroughs on many challenging benchmark tasks, lean to generate verbose responses and lack the controllability of output complexity, which is usually preferred by human users in practice. In this paper, we study how to precisely control multiple linguistic complexities of LLM output by finetuning using off-the-shelf data. To this end, we propose multi-control tuning (MCTune), which includes multiple linguistic complexity values of ground-truth responses as controls in the input for instruction tuning. We finetune LLaMA2-7B on Alpaca-GPT4 and WizardLM datasets. Evaluations on widely used benchmarks demonstrate that our method does not only improve LLMs' multi-complexity controllability substantially but also retains or even enhances the quality of the responses as a side benefit.

1 Introduction

Large language models have achieved remarkable success in generating free-from texts for different downstream tasks or human instructions (Xu et al., 2023a; Chowdhery et al., 2022; Touvron et al., 2023b). However, existing LLMs still lack precise control over the linguistic complexity of their outputs (Soatto et al., 2023; Chen et al., 2024b; Li et al., 2024a), e.g., the total number of nouns, the variation of verbs, etc. Linguistic controllability is crucial to creating personalized outputs since those complexity indices directly reflect human reading complexity in multiple aspects. For example, a short yes/no answer is required by some users while a detailed explanation is preferred by others. Moreover, recent studies (Singhal et al., 2023; Zhao et al., 2024; Lambert et al., 2024) have discovered a spurious correlation between the quality reward used in LLM alignment and the output length (i.e., a specific linguistic complexity). Consequently, LLMs favor generating verbose responses due to

the length bias, which may increase unnecessary reading complexity. It is still an open problem to mitigate the bias without hurting the output quality.

While existing LLM finetuning techniques such as instruction-tuning (Xu et al., 2023a; Stiennon et al., 2020; Li et al., 2024b) and reinforcement learning from human feedback (RLHF) (Christiano et al., 2017; Bai et al., 2022) have been demonstrated to be effective in aligning the output with human intent or preference, they only focus on maximizing a single objective. Instead, achieving the controllability of multiple complexity indices requires a non-trivial multi-objective optimization that has not been thoroughly studied on LLMs. Rather than solely maximizing or minimizing the complexities, it aims to reach different target complexity values on the Pareto frontier. This requires LLMs to adjust the trade-off among objectives and capture their potential correlations or constraints in the text generation process. In addition, due to the huge space of possible combinations of complexity indices, it could be expensive to collect training data for multi-objective control.

In this paper, we take the first step towards multi-objective control of the linguistic complexity of LLM outputs. Instead of collecting new data, our strategy allows the reuse of existing instruction-tuning data. In particular, we annotate the ground-truth responses in a dataset by their complexity metrics evaluated using tools developed in computational linguistics. The multiple complexity indices and their values are appended as tags to the input so finetuning an LLM on the linguistic-label augmented data helps build a strong connection between the input tags and the linguistic complexity of the output, hence enforcing the LLM to adhere to the complexity requirements during sequential decoding. Surprisingly, we observe that randomly sampling a small subset of tags for each training example suffice to obtain controllability over all the complexities of test examples, thereby reducing

the required amount of training data.

We examine our approach by finetuning LLaMA2-7B (Touvron et al., 2023a) using a linguistic-complexity labeled Alpaca-GPT4 dataset (Peng et al., 2023) (i.e., the prompts are from the original Alpaca dataset while the responses are composed by GPT-4) and a WizardLM dataset (Xu et al., 2023b). We do not only observe an expected substantial improvement in linguistic controllability but also a side benefit of enhanced response quality, indicating that finetuning under multiple linguistic constraints can improve the LLM for general purposes. Compared to unconstrained RLHF and IFT methods, which suffer from the length bias (Chen et al., 2024a), our approach does not introduce the bias but can improve both the controllability and quality simultaneously.

2 Related Work

Linguistic Features and Complexity. The exploration of linguistic features and complexity in language models encompasses a diverse range of research. Seminal studies have investigated the syntactic abilities of LSTMs (Linzen et al., 2016). Fidler and Goldberg (2017) has introduced methods for manipulating the stylistics and syntactic output of text generation models. Additionally, linguistic style transfer (Shen et al., 2017) has also showcased the adaptability of language models to capture and replicate varied linguistic features. Building on the previous efforts to understand the multifaceted nature of linguistics complexity, our work concentrates on producing responses endowed with specific linguistic characteristics, which is less explored in the previous work.

Controllability of LLMs. The topic of personalized language modeling has attracted significant attention across various research papers. Techniques such as user embedding have become common for customizing language models to individual needs (Welch et al., 2020; Rocca and Yarkoni, 2022). More recently, Mireshghallah et al. (2022); Oba et al. (2023) propose prompt based personalized fine-tuning for specific users, and producing personalized responses. Our research shifts the focus from personalization for specific users to the broader goal of controlling large language models (LLMs) to produce outputs with linguistic diversity and complexity, addressing a gap not explored by the aforementioned works.

Another prevalent technique for guiding the out-

put of large language models involves Tagging (Korbak et al., 2023; Prabhume et al., 2023; Lu et al., 2022). This method incorporates appending human-readable text during the training of LLMs. Contrary to previous studies that concentrated on managing aspects like toxicity (Korbak et al., 2023; Prabhume et al., 2023; Li et al., 2024a) and controlling repetition (Lu et al., 2022), our approach employs tagging techniques to control linguistic features across multiple attributes. Additionally, we utilize multiple tags to enable simultaneous consideration of various attributes, a difference from earlier work that primarily uses of a single tag.

Finetuning and Alignment of LLMs. With the emergence of large-scale language models, such as those in the GPT series, aligning language models has become prevalent. Studies like those by (Zhou et al., 2023; Xu et al., 2023b; Li et al., 2023, 2024c; Chen and Mueller, 2024) have concentrated on the process of data curation for instruction fine-tuning to enhance models’ instructions following capabilities. Unlike these data-centric approaches, we keep the original instruction dataset (Chen and Mueller, 2023) but augment instructions with various tags to introduce a richer array of linguistic features, thereby elevating the instruction-following capabilities of the models.

3 Finetuning LLMs for Linguistic Controllability

In this section, we delineate our approach to multi-control tuning. Section 3.1 outlines the linguistic features of interest and describes how we extract them from a given text segment. In Section 3.2, we explain how the extracted features are incorporated into the multi-control tuning process.

3.1 Handcrafted Linguistic Features

We are specifically interested in controlling the *handcrafted* linguistic properties of the model’s generation. This type of feature has been used throughout the NLP field (Bogdanova et al., 2017; Anshika Choudhary, 2021; Lee et al., 2021) and is loosely defined in Lee and Lee (2023) as “*a single numerical value produced by a uniquely identifiable method on any natural language.*” An example of a linguistic feature not considered handcrafted is text embeddings produced by deep neural networks, which usually take the form of a vector. We extract such features from a text segment using the LFTK package proposed in Lee and Lee (2023).

ID	Name	Description
1	t_word	number of words
2	n_noun	number of nouns
3	n_verb	number of verbs
4	n_adj	number of adjectives
5	t_uword	number of unique words
6	n_unoun	number of unique nouns
7	n_uverb	number of unique verbs
8	n_uadj	number of unique adjectives
9	ttr	type-token ratio
10	noun_var	noun variation
11	verb_var	verb variation
12	adj_var	adjective variation
13	fkre	Flesch-Kincaid reading ease
14	rt_average	average reading time

Table 1: The list of linguistic features for controllability tuning. The second column shows the name of each feature with the corresponding descriptions in the third column. We include a detailed explanation of how these features are computed in Appendix A.

It encompasses a diverse set of 220 features that are grouped into different linguistic families. Within the scope of this paper, we sample a reasonably-sized set of 14 features for multi-control tuning, which are presented in Table 1. These features are selected to cover most of the feature families while being simple to understand and verify by human users.

Formally, given a text segment $\mathbf{x} = [x_1, x_2, \dots, x_l]$, with x_i being the i -th token of \mathbf{x} , we denote the feature extractor as a function $f: \mathcal{X} \rightarrow \mathbb{R}^d$ that maps \mathbf{x} to a d -dimensional vector $f(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_d(\mathbf{x})]$, where $f_j(\mathbf{x})$ represents the j -th linguistic feature of interest and \mathcal{X} represents the space of all texts. In this paper, the function f refers to the LFTK feature extractor.

3.2 Multi-Objective Control Tuning

Consider the standard instruction tuning setting where an LLM, denoted as $p_\theta(\mathbf{x}) = \prod_{i=1}^l p_\theta(x_i | x_{<i})$, is trained on a dataset of N instruction-output pairs, $D_{\text{train}} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$. The training objective is to generate \mathbf{y}_i given \mathbf{x}_i ; thus, the loss is

$$\mathcal{L} = - \sum_{i=1}^N \sum_{k=1}^{|\mathbf{y}_i|} \log p_\theta(\mathbf{y}_{i,k} | \mathbf{x}_i, \mathbf{y}_{i,<k})$$

We address multi-control tuning by framing it as a conditional instruction tuning problem, where we utilize the target response \mathbf{y}_i from D_{train} to

generate a linguistic control vector $f(\mathbf{y}_i)$ and append it to \mathbf{x}_i . The model is then trained to generate \mathbf{y}_i condition on both \mathbf{x}_i and $f(\mathbf{y}_i)$. However, to enhance data diversity and better simulate real-world scenarios, where a user may wish to control only a few features, we do not utilize all features for every data example. For each pair $(\mathbf{x}_i, \mathbf{y}_i) \in D_{\text{train}}$, we randomly sample an integer $n_i \sim \text{Uniform}\{1, \dots, m\}$, $m \leq d$. The set of n_i feature indices, denoted by C_i , is then randomly sampled from the pool of all feature indices, i.e., $C_i \sim \text{Uniform}(\{A \subseteq \{1, \dots, d\} : |A| = n_i\})$. The linguistic control vector used for the i -th example is denoted as $f_{C_i}(\mathbf{y}_i)$, where $f_{C_i}(\mathbf{y}_i) = [f_{C_{i,1}}(\mathbf{y}_i), \dots, f_{C_{i,n_i}}(\mathbf{y}_i)]$. The resulting training loss becomes

$$\mathcal{L} = - \sum_{i=1}^N \sum_{k=1}^{|\mathbf{y}_i|} \log p_\theta(\mathbf{y}_{i,k} | \mathbf{x}_i, f_{C_i}(\mathbf{y}_i), \mathbf{y}_{i,<k})$$

The benefits of our training strategy are twofold: (1) it improves the controllability and instruction-following capability simultaneously, thus avoiding the catastrophic forgetting problem that may degrade the model’s generation quality. In fact, we will later show in Section 5.6.2 that LLMs trained with our approach achieve even stronger instruction-following ability compared to those trained with vanilla instruction tuning; (2) it allows us to utilize off-the-shelf datasets without the need to collect new ones.

3.3 Prompt Template

In practice, we format both the instruction \mathbf{x} and the output \mathbf{y} into a predefined template before they are input into the LLMs. This paper adopts the template outlined in (Taori et al., 2023), wherein \mathbf{x} is decomposed into an instruction and an input component. Regarding linguistic controls, we format them into a sequence of the form [name_1: value_1] ... [name_n: value_n], aiming for conciseness by utilizing the features’ abbreviated names listed in Table 1. To assist the LLM in better understanding these abbreviations, we provide a comprehensive list of feature descriptions within the system prompt, which has been confirmed to enhance the effectiveness of our approach in preliminary experiments. This sequence of controls is subsequently attached to the input component. Figure 1 illustrates a detailed example of how an input prompt is constructed.

Below is an instruction that describes a task, paired with an input that provides further context. At the end of the input, there will be a list of tags specifying the desired properties of the response. The following tags are available: [t_word] for the total number of words; [n_noun] for the total number of nouns; [n_verb] for the total number of verbs; [n_adj] for the total number of adjectives; [t_uword] for the total number of unique words; [n_unoun] for the total number of unique nouns; [n_uverb] for the total number of unique verbs; [n_uadj] for the total number of unique adjectives; [simp_ttr] for the simple type-token ratio; [simp_noun_var] for simple noun variation; [simp_verb_var] for simple verb variation; [simp_adj_var] for simple adjective variation; [fkre] for the Flesch-Kincaid Reading Ease; [rt_average] for the average reading time. Write a response that appropriately completes the request and satisfies the tags.

Instruction:
 Arrange the words in the given sentence to form a grammatically correct sentence.

Input:
 quickly the brown fox jumped [t_word: 6] [n_noun: 1] [fkre: 102.05]

Response:
 The brown fox jumped quickly.

Figure 1: An example of how data is formatted before being fed into LLMs in this paper. The first paragraph presents a system prompt containing a complete list of feature descriptions. Our preliminary results indicate that including descriptions enhances the effectiveness of our approach.

4 Evaluation of Linguistic Controllability

During evaluation on a reserved test dataset D_{test} , we aim to measure how the model performs as the linguistic control vector f changes. To conduct such an evaluation, we develop a sampling strategy to sample different control vectors for an instruction \mathbf{x}_i . Specifically, we need a sampling strategy that: (1) is specific to an instruction \mathbf{x}_i , i.e., the sampled control vectors should not stray too far from the reasonable range for a specific \mathbf{x}_i . For example, an instruction to "Generate a short story" should not have a large value for `t_word`; (2) ensures the sampled control vectors are always valid, meaning that no linguistic controls conflict with each other (e.g., `t_uword` should always be less than or equal to `t_word`), and no control is out-of-bound (e.g., `fkre` should always be less than or equal to 121.22).

To achieve the first goal, we utilize \mathbf{y}_i 's linguistic feature vector as a reference point and sample new control vectors f' from the Gaussian distribution centered at $f(\mathbf{y}_i)$, i.e., $f' \sim \mathcal{N}(f(\mathbf{y}_i), I\sigma^2)$. However, since each feature $f_i(\mathbf{y})$ has a different range, a small σ for some features may be large for others. To avoid this inconsistency, we standardize all features to unit variance before sampling. More formally, the new control vector f' is computed by

$$f' = z^{-1}(z(f(\mathbf{y}_i)) + \sigma\epsilon), \quad \epsilon \sim \mathcal{N}(0, I)$$

where $z : \mathbb{R}^d \rightarrow \mathbb{R}^d$ standardizes each feature to unit variance, and z^{-1} is the inverse operation.

While sampling a valid control vector can be a challenging problem, verifying its validity is

straightforward. This can be done using a simple rule-based method, which we outline in Appendix B. We utilize this observation to achieve the second goal by performing rejection-based sampling, i.e., keep resampling a new control vector f' until we find a valid one. We will show in Section 5.7.3 that σ can serve as a hyperparameter to control the evaluation difficulty.

Lastly, for each example in D_{test} , we randomly sample a number n , a set of n feature indices C , and K new control vectors f'_C for controllability and generation quality evaluation.

5 Experiments

5.1 Implementation Details

By default, we set $K = 5$ and $\sigma = 0.1$ unless specified otherwise. We set the maximum number of linguistic controls per example to $m = 5$ and will show in Section 5.7.2 that limiting m to 5 does not affect the model's controllability, even when more than 5 controls are used in the evaluation. For all datasets, we fine-tune LLaMA2-7B using the AdamW optimizer with a linear learning rate schedule. We set the learning rate to 2×10^{-5} , the batch size to 128, and the number of warmup steps to 100. All models are trained on 8 RTX A6000 GPUs for 5 epochs.

5.2 Datasets.

Data Preprocessing. Because LFTK cannot extract linguistic features with perfect accuracy, it sometimes generates invalid control vectors $f(\mathbf{y}_i)$ (e.g., producing an `fkre` greater than 121.22). This

situation complicates the process of sampling new control vectors, as $f(y_i)$ may deviate significantly from the feasible region, drastically reducing the probability of sampling a valid one. Therefore, we exclude data examples with invalid $f(y_i)$.

Alpaca-GPT4. We utilize the Alpaca-GPT4 dataset (Peng et al., 2023) for instruction tuning and evaluation. This dataset shares the same set of instructions as the Alpaca dataset (Taori et al., 2023), but it employs OpenAI’s GPT-4 to generate high-quality responses. Our preliminary experiments demonstrate that training with this dataset yields better results in terms of both controllability and generation quality compared to training with the original Alpaca dataset. After preprocessing, the dataset is divided into a training set of 45,000 examples and a test set of 2,000 examples.

WizardLM. In addition to Alpaca-GPT4, we also evaluate our method on the WizardLM dataset (Xu et al., 2023b). The original dataset contains 70,000 examples of instruction-output pairs that were automatically generated by ChatGPT. We subsampled 50,000 examples from the original dataset. Similar to Alpaca-GPT4, after preprocessing, we split the data into a training set of 40,000 examples and a test set of 2,000 examples.

5.3 Models.

We use LLaMA2-7B (Touvron et al., 2023b) as the base model for multi-control tuning in all experiments.

5.4 Baselines

To evaluate the impact of our method on controllability and generation quality, we conduct comparisons with the same LLaMA2-7B base model but trained with regular instruction fine-tuning. We also include OpenAI’s GPT-3.5 Turbo (gpt-3.5-turbo-0125) as a baseline to compare our model against state-of-the-art proprietary LLMs in terms of controllability. For each baseline, we control the linguistic complexity of their responses using the prompt template shown in Figure 1 in a zero-shot manner.

5.5 Evaluation Metrics

Controllability Error. Given an instruction \mathbf{x} , a linguistic control vector f_C , and a generated response $\hat{\mathbf{y}} \sim p_\theta(\mathbf{y} \mid \mathbf{x}, f_C)$, we measure controllability error by computing the L_1 error between the

specified control vector f_C and the response’s linguistic feature vector $f_C(\hat{\mathbf{y}})$. We denote this error as $\mathbf{e} = |f_C(\hat{\mathbf{y}}) - f_C| \in \mathbb{R}^{|C|}$.

Quality Score. To evaluate generation quality, we follow Zheng et al. (2023) in using a powerful LLM (e.g., ChatGPT-4 Turbo) as a judge to assign quality scores ranging from 1 to 10.

5.6 Main Results and Analysis

5.6.1 Linguistic Controllability Evaluation

This section presents the evaluation of linguistic controllability between our method and various baselines. We consider three evaluation settings: Easy ($\sigma = 0.1$), Medium ($\sigma = 0.2$), and Hard ($\sigma = 0.3$). Our goal is to measure the controllability error of each baseline on each linguistic control and then visualize all of them on the same radar plot for comparison. Formally, for each linguistic control i , and each baseline j , we maintain a list $\mathbf{e}_{i,j} = [e_{i,j}^1, \dots, e_{i,j}^{|\mathbf{e}_{i,j}|}]$ of L_1 errors made by j , where the length of this list depends on the number of text examples that contain i . The matrix of all L_1 errors for feature i is denoted as

$$\mathbf{E}_i = \begin{pmatrix} e_{i,1}^1 & \dots & e_{i,1}^{|\mathbf{e}_{i,1}|} \\ \vdots & \ddots & \vdots \\ e_{i,J}^1 & \dots & e_{i,J}^{|\mathbf{e}_{i,J}|} \end{pmatrix} \quad (1)$$

It is challenging to directly visualize \mathbf{E}_i for all i on the same radar plot because each linguistic control i has a different range of values. Therefore, we normalize each element in \mathbf{E}_i using min-max normalization, with the minimum being $\min E_i$ and the maximum being the 95th percentile of \mathbf{E}_i , or $P_{95}(\mathbf{E}_i)$. For each baseline j , its average normalized L_1 error on control i is calculated as $\sum_k \text{norm}(e_{i,j}^k) / |\mathbf{e}_{i,j}|$, where $\text{norm}(\cdot)$ denotes the normalization described above.

As shown in Figure 2, our method consistently outperforms the other baselines across all settings. Surprisingly, a state-of-the-art model like ChatGPT underperforms in controllability compared to the instruction-finetuned LLaMA2-7B. Upon manual inspection, we observed that ChatGPT tends to produce more verbose responses regardless of the linguistic controls applied, which may explain the observed poor performance. Another possible reason for this discrepancy is that LLaMA2-7B is finetuned on data sharing the same distribution as the test set, giving it an advantage over ChatGPT. Note that for noun_var, verb_var, and adj_var,

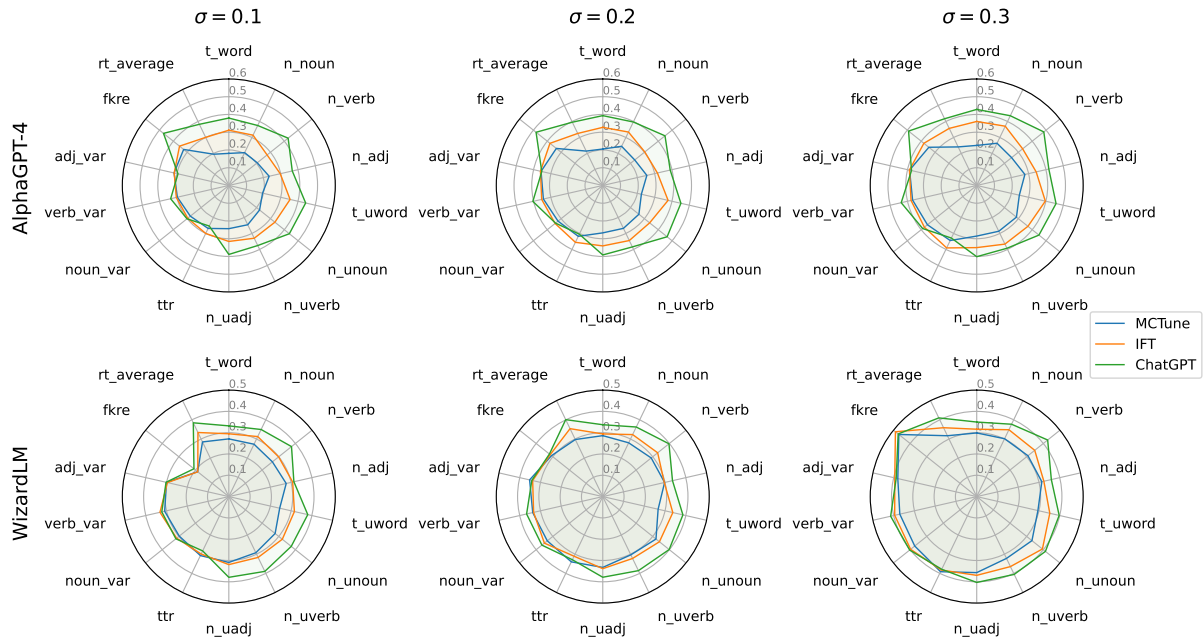


Figure 2: **Comparison of controllability error** (the average normalized L_1 error) of ChatGPT (gpt-3.5-turbo-0125), IFT (finetuning without controls), and MCTune (ours) on AlpacaGPT-4 and WizardLM datasets and three test settings of target linguistic complexity with increasing σ (difficulty). To visualize each linguistic feature’s average error on the same radar plot, we apply min-max normalization to normalize them to a similar scale. To reduce the effect of outliers, the minimum L_1 error of the i -th feature is the minimum among all baselines, and the maximum L_1 error refers to the 95th percentile of errors among all baselines.

the differences between each baseline’s errors are not significant. We hypothesize that, in contrast to other linguistic complexities, the descriptions of these features are somewhat more ambiguous, and we also did not provide a clear description in the system prompt, which may not be helpful for the model to understand and follow these features.

5.6.2 Generation Quality Evaluation

A natural concern when fine-tuning for controllability is how it affects the model’s generation quality. In this section, we answer this question by comparing our method with standard instruction fine-tuning on MT-Bench (Zheng et al., 2023), a widely used benchmark for evaluating LLMs on multi-turn open-ended questions. Specifically, we use GPT-4 Turbo (gpt-4-0125-preview) as a judge and compare the two methods in both single-answer and pairwise settings. In the single-answer setting, the judge assigns a single numeric score from 1 to 10 for each model’s answer. In the pairwise setting, the judge receives two answers from both baselines and returns either a win, a loss, or a tie. The result for the single-answer setting is shown in Figure 4. As shown in the figure, our method does not degrade the model’s general language ability but

even improves it in most categories. Note that the performance of both baselines on Coding and Math questions is poor because the AlpacaGPT-4 dataset does not have a strong specialty in questions of this category. When evaluating under the pairwise setting, MCTune achieves a 51.25% win rate over instruction fine-tuning. These results suggest that training for controllability is beneficial for improving LLMs’ natural language performance.

5.7 Analyses

5.7.1 Relationship Between Linguistic Controllability and Generation Quality

This experiment focuses on studying the relationship between generation quality and the controllability error of LLM responses. To start off, we randomly select a set of 5 linguistic controls C and fix it throughout the experiment for simplicity. Given an LLM p_θ and the Alpaca-GPT4 test dataset $D_{\text{test}} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, we select 50 examples from D_{test} then sample 50 responses $\hat{\mathbf{y}}_i$ from $p_\theta(\mathbf{y}_i | \mathbf{x}_i, f_C(\mathbf{y}_i))$. Each $\hat{\mathbf{y}}_i$ is then evaluated by GPT-4 Turbo on how well it satisfies \mathbf{x}_i and a controllability error of $\hat{\mathbf{y}}_i$ is computed by taking the average of the normalized L_1 errors across all controls in C . We also consider three different settings

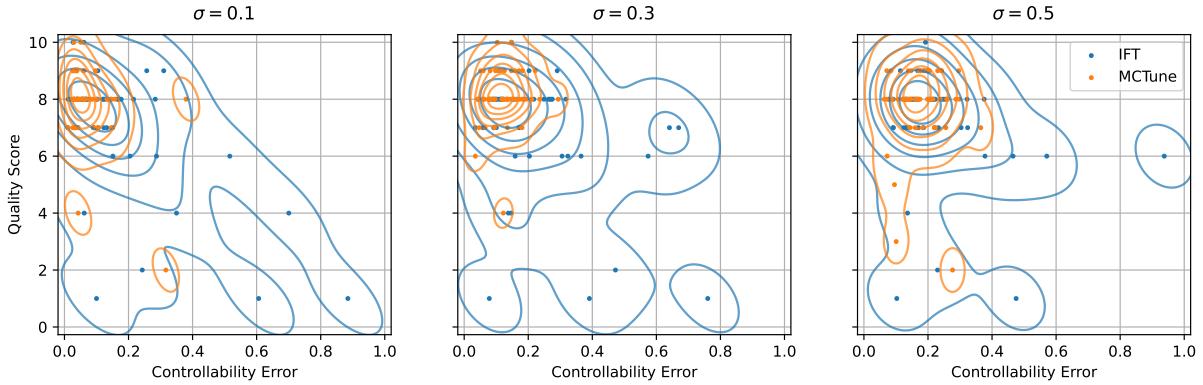


Figure 3: **Trade-off between linguistic controllability and generation quality** in three increasingly difficult test settings. Each dot represents a model’s response \hat{y} to a specific query $[x, f_C]$. The response is given a quality score from 1 to 10 by a judge LLM (GPT-4 Turbo) based on how well \hat{y} addresses x . A controllability error is measured for \hat{y} , which is computed by taking the average of normalized L_1 errors across all linguistic controls in f_C . Blue and orange dots respectively represent responses from models trained by IFT and MCTune using Alpaca-GPT4 dataset.

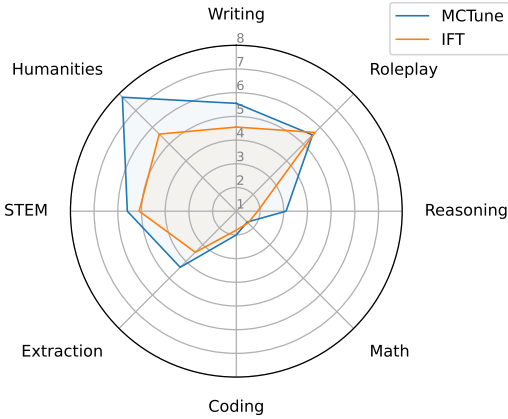


Figure 4: **Comparison between MCTune and IFT-trained models on MT-Bench.** We finetune LLaMA2-7B on Alpaca-GPT4 dataset and GPT-4 Turbo is the judge in the test. The average score per axis ranges from 1 to 10 and are given by the judge.

where $\sigma = 0.1$, $\sigma = 0.3$, and $\sigma = 0.5$ to see how the pattern shifts as σ changes. The results are shown in Figure 3. We can see that, compared to IFT, our method achieves better quality and controllability simultaneously. As σ increases, the controllability error increases, which is expected and consistent with the results in Section 6. We then start to notice a slight degradation in quality at $\sigma = 0.5$, where there are no responses with a score of 10, and more responses with low scores begin to appear. This observation suggests that there might be a positive correlation between controllability and generation quality.

5.7.2 Analysis of Model Controllability with Varying Linguistic Controls

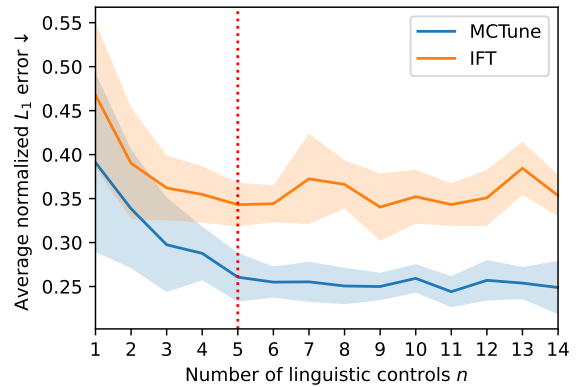


Figure 5: Linguistic controllability error on the test set as the number of linguistic controls n per sample in MCTune’s training increases. The solid curves represent the linguistic controllability error averaged over all linguistic complexities (lower is better) with the shaded areas represent the 95%-confidence interval. The dotted vertical line indicates the maximum number of controls (5) used during MCTune’s training.

In this experiment, we are interested in studying how the model’s controllability is affected as the number of linguistic controls increases. We conduct this experiment using the Alpaca-GPT4 dataset, from which we randomly sample 100 examples from the test set. For each example, we sample five new control vectors f' of length n , compute the average normalized L_1 error per control, and then average these over all controls to

obtain a single numeric score. We repeat this process for $n = 1$ to $n = 14$. The results are shown in Figure 5. Surprisingly, controllability is worse when n is small. This is true for both MCTune and the IFT baseline, suggesting that this is not an inherent limitation of our approach. Notice that, as the number of linguistic controls increases beyond five, the maximum number of controls used during training, the controllability error of IFT fluctuates, while MCTune remains approximately stable. This demonstrates that our model, trained with MCTune, is able to generalize to a larger number of controls than it has seen in training.

5.7.3 Analyzing the Impact of σ on Linguistic Controllability Evaluation

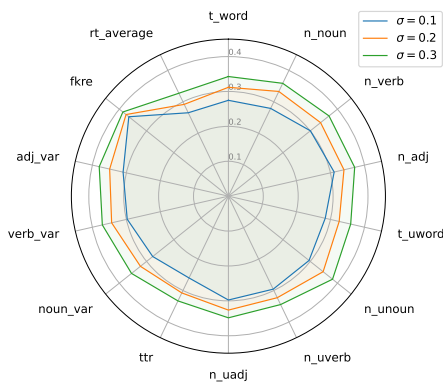


Figure 6: Analysis of how the average normalized L_1 error varies when **changing the test difficulty level** via σ . We use the same model and only vary σ .

This section examines the effect of σ on controllability evaluation. Given a model fine-tuned with MCTune and a test dataset, we follow the process described in Section 4 to sample new linguistic control vectors for evaluating controllability. We consider three settings: $\sigma = 0.1$, $\sigma = 0.2$, and $\sigma = 0.3$. The results are depicted in Figure 6. As illustrated by the results, as σ increases, the average normalized L_1 error uniformly increases across all linguistic controls. This indicates that σ can serve as a hyperparameter to adjust the difficulty of the evaluation, which is intuitive since a larger σ enables the sampling of control vectors that are more distant from the reference control vector $f(y)$.

5.8 Case Studies

This section presents examples of how different values of controls affect the model’s responses. As shown in Figure 7, increasing the t_word con-

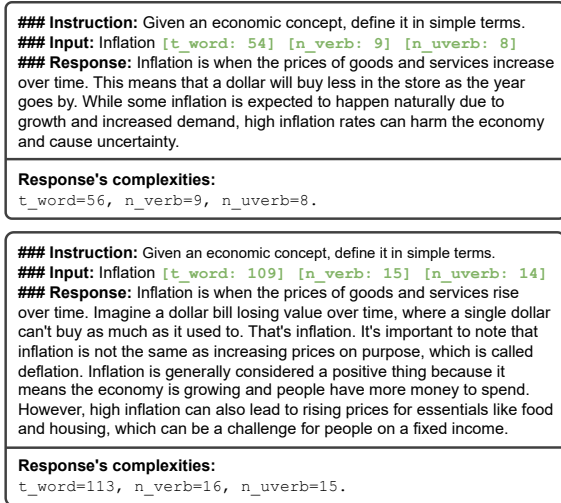


Figure 7: Examples of MCTune model responding to the same instruction under different linguistic controls. The linguistic complexities of the responses are extracted using LFTK with exact values shown at the bottom of each example.

trol leads to a longer response and larger values of n_verb and n_uverb increase the variation in verb choice. Although it is challenging to strictly control the complexities of the response, the linguistic controls can serve as soft constraints on the model’s generation.

6 Discussion and Conclusions

In this paper, we advance the precise control of linguistic complexities in LLMs through multi-objective control tuning. Our method is straightforward and can be seamlessly integrated with existing instruction-tuning datasets without the need to gather new ones. Our training objective concurrently optimizes the LLM’s ability to follow instructions and its controllability. Through our experiments, we find that incorporating this dual-focus strategy significantly improves the LLM’s generative quality, surpassing the results of instruction fine-tuning alone. This finding suggests that controllability and instruction-following ability may have a complementary effect on each other. Additionally, we observe that while state-of-the-art LLMs achieve impressive natural language performance, they are not easy to control. This emphasizes the need for studying methods that improve controllability in LLMs.

7 Limitations

While our method significantly improves controllability compared to regular instruction fine-tuning, we observe that there is still room for improvement. The model trained with MCTune is able to loosely follow the linguistic controls but struggles to produce responses with the exact complexities. An interesting next step would be to improve controllability in a strict setting. Another limitation we observe is the fact that LFTK sometimes extracts incorrect linguistic complexities for a given text. This leads to noisy controls that may confuse and reduce the controllability of the model during training.

References

- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. 2024. [Back to basics: Revisiting REINFORCE style optimization for learning from human feedback in llms](#). *CoRR*, abs/2402.14740.
- Anuja Arora Anshika Choudhary. 2021. Linguistic feature based learning model for fake news detection and classification. *Expert Systems with Applications*.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#). *CoRR*, abs/2204.05862.
- Dasha Bogdanova, Jennifer Foster, Daria Dziedzic, and Qun Liu. 2017. [If you can't beat them join them: Handcrafted features complement neural nets for non-factoid answer reranking](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 121–131, Valencia, Spain. Association for Computational Linguistics.
- Jiuhai Chen and Jonas Mueller. 2023. [Quantifying uncertainty in answers from any language model via intrinsic and extrinsic confidence assessment](#). *CoRR*, abs/2308.16175.
- Jiuhai Chen and Jonas Mueller. 2024. [Automated data curation for robust language model fine-tuning](#). *CoRR*, abs/2403.12776.
- Lichang Chen, Chen Zhu, Davit Soselia, Jiuhai Chen, Tianyi Zhou, Tom Goldstein, Heng Huang, Mohammad Shoeybi, and Bryan Catanzaro. 2024a. [ODIN: disentangled reward mitigates hacking in RLHF](#). *CoRR*, abs/2402.07319.
- Yihan Chen, Benfeng Xu, Quan Wang, Yi Liu, and Zhendong Mao. 2024b. [Benchmarking large language models on controllable generation under diversified instructions](#).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#). *CoRR*, abs/2204.02311.
- Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. [Deep reinforcement learning from human preferences](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4299–4307.
- Jessica Fidler and Yoav Goldberg. 2017. [Controlling linguistic style aspects in neural language generation](#). In *Proceedings of the Workshop on Stylistic Variation*, pages 94–104, Copenhagen, Denmark. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *CoRR*, abs/2310.06825.
- Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Vinayak Bhalerao, Christopher Buckley, Jason Phang, Samuel R. Bowman, and Ethan Perez. 2023. [Pretraining language models with human preferences](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 17506–17533. PMLR.

- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. **Rewardbench: Evaluating reward models for language modeling.** *CoRR*, abs/2403.13787.
- Bruce W. Lee, Yoo Sung Jang, and Jason Hyung-Jong Lee. 2021. **Pushing on text readability assessment: A transformer meets handcrafted linguistic features.**
- Bruce W. Lee and Jason Lee. 2023. **LFTK: Handcrafted features in computational linguistics.** In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 1–19, Toronto, Canada. Association for Computational Linguistics.
- Ming Li, Jiu Hai Chen, Lichang Chen, and Tianyi Zhou. 2024a. **Can llms speak for diverse people? tuning llms via debate to generate controllable controversial statements.** *CoRR*, abs/2402.10614.
- Ming Li, Lichang Chen, Jiu Hai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. 2024b. **Selective reflection-tuning: Student-selected data recycling for LLM instruction-tuning.** *CoRR*, abs/2402.10110.
- Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. 2024c. **Superfiltering: Weak-to-strong data filtering for fast instruction-tuning.** *CoRR*, abs/2402.00530.
- Ming Li, Yong Zhang, Zhitao Li, Jiu Hai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2023. **From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning.** *ArXiv*, abs/2308.12032.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of lstms to learn syntax-sensitive dependencies. *Transactions of the Association for Computational Linguistics*, 4:521–535.
- Ximing Lu, Sean Welleck, Jack Hessel, Liwei Jiang, Lianhui Qin, Peter West, Prithviraj Ammanabrolu, and Yejin Choi. 2022. **Quark: Controllable text generation with reinforced unlearning.** *Advances in neural information processing systems*, 35:27591–27609.
- Fatemehsadat Mireshghallah, Vaishnavi Shrivastava, Milad Shokouhi, Taylor Berg-Kirkpatrick, Robert Sim, and Dimitrios Dimitriadis. 2022. **UserIdentifier: Implicit user representations for simple and effective personalized sentiment analysis.** In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3449–3456, Seattle, United States. Association for Computational Linguistics.
- Daisuke Oba, Naoki Yoshinaga, and Masashi Toyoda. 2023. **Perplm: Personalized fine-tuning of pretrained language models via writer-specific intermediate learning and prompts.** *arXiv preprint arXiv:2309.07727*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. **Instruction tuning with gpt-4.** *arXiv preprint arXiv:2304.03277*.
- Shrimai Prabhunoye, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. 2023. **Adding instructions during pretraining: Effective way of controlling toxicity in language models.** *arXiv preprint arXiv:2302.07388*.
- Roberta Rocca and Tal Yarkoni. 2022. **Language as a fingerprint: Self-supervised learning of user encodings using transformers.** In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1701–1714, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. **Style transfer from non-parallel text by cross-alignment.** *Advances in neural information processing systems*, 30.
- Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. 2023. **A long way to go: Investigating length correlations in RLHF.** *CoRR*, abs/2310.03716.
- Stefano Soatto, Paulo Tabuada, Pratik Chaudhari, and Tian Yu Liu. 2023. **Taming AI bots: Controllability of neural states in large language models.** *CoRR*, abs/2305.18449.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. 2020. **Learning to summarize with human feedback.** In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. **Stanford alpaca: An instruction-following llama model.** https://github.com/tatsu-lab/stanford_alpaca.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,

Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023a. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#).

Charles Welch, Jonathan K. Kummerfeld, Verónica Pérez-Rosas, and Rada Mihalcea. 2020. [Exploring the value of personalized word embeddings](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6856–6862, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhao Feng, Chongyang Tao, and Daxin Jiang. 2023a. [Wizardlm: Empowering large language models to follow complex instructions](#). *CoRR*, abs/2304.12244.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhao Feng, Chongyang Tao, and Daxin Jiang. 2023b. [Wizardlm: Empowering large language models to follow complex instructions](#). *arXiv preprint arXiv:2304.12244*.

Hao Zhao, Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. 2024. [Long is more for alignment: A simple but tough-to-beat baseline for instruction fine-tuning](#). *CoRR*, abs/2402.04833.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#).

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. [Lima: Less is more for alignment](#). *arXiv preprint arXiv:2305.11206*.

A Computing Linguistic Features

In this section, we describe in greater detail how linguistic features are computed. We use LFTK as a feature extractor, which is based on spaCy. Given a string, spaCy tokenizes it into a sequence of tokens along with their annotations (e.g., part-of-speech). The number of words t_word is the number of tokens. The n_noun , n_verb , and n_adj are computed based on the POS provided by spaCy. The t_uword , n_unoun , n_uverb , and n_uadj are computed accordingly. The type-token ratio, noun variation, verb variation, and adjective variation are $\frac{t_uword}{t_word}$, $\frac{n_unoun}{n_noun}$, $\frac{n_uverb}{n_verb}$, and $\frac{n_uadj}{n_adj}$, respectively. Let t_sent be the number of sentences and t_syll is the number of syllables, the formula to compute Flesch-Kincaid reading ease is

$$206.835 - 1.015 \left(\frac{t_word}{t_sent} \right) - 84.6 \left(\frac{t_syll}{t_word} \right)$$

Lastly, the average reading time $rt_average$ is $\frac{t_word}{240}$.

B Verifying Validity of Control Vectors

Given a linguistic control vector, it is straightforward to check its validity by iterating through a list of pre-defined rules.

- $t_word > 0$
- $t_word \geq n_noun + n_verb + n_adj$
- $t_word \geq t_uword$
- $n_noun \geq 0$
- $n_noun \geq n_unoun$
- $n_verb \geq 0$
- $n_verb \geq n_uverb$
- $n_adj \geq 0$
- $n_adj \geq n_uadj$
- $t_uword > 0$
- $t_uword \geq n_unoun + n_uverb + n_uadj$
- $n_unoun \geq 0$
- $n_uverb \geq 0$
- $n_uadj \geq 0$
- $fkre \leq 121.22$

If any of the rules above is violated, we conclude that the control vector is invalid.

Method	t_word	n_noun	n_verb	n_adj	t_uword	n_unoun	n_uverb	n_uadj	ttr	noun_var	verb_var	adj_var	fkre	rt_average
IFT (Zero-shot)	0.31	0.33	0.34	0.32	0.38	0.35	0.33	0.35	0.30	0.30	0.31	0.32	0.34	0.31
IFT (Few-shot)	0.33	0.32	0.33	0.26	0.39	0.41	0.31	0.28	0.38	0.30	0.33	0.32	0.39	0.31
GPT-3.5 (Zero-shot)	0.37	0.50	0.36	0.44	0.55	0.34	0.33	0.54	0.26	0.33	0.48	0.29	0.39	0.47
GPT-3.5 (Few-shot)	0.48	0.53	0.65	0.37	0.67	0.60	0.39	0.44	0.40	0.29	0.35	0.28	0.40	0.48
GPT-4 (Zero-shot)	0.64	0.52	0.34	0.51	0.63	0.73	0.46	0.39	0.21	0.32	0.38	0.28	0.49	0.51
GPT-4 (Few-shot)	0.40	0.40	0.47	0.31	0.46	0.51	0.35	0.36	0.36	0.28	0.30	0.27	0.40	0.32
RLHF (REINFORCE)	0.25	0.24	0.27	0.40	0.27	0.25	0.25	0.27	0.29	0.29	0.38	0.30	0.83	0.36
MCTune	0.18	0.22	0.25	0.25	0.21	0.23	0.22	0.25	0.27	0.26	0.29	0.26	0.31	0.20

Table 2: Additional comparison of controllability error with GPT-4, few-shot and RL-based baselines.

Method	t_word	n_noun	n_verb	n_adj	t_uword	n_unoun	n_uverb	n_uadj	ttr	noun_var	verb_var	adj_var	fkre	rt_average
IFT	0.37	0.37	0.33	0.35	0.41	0.38	0.33	0.35	0.44	0.40	0.43	0.47	0.40	0.36
MCTune	0.34	0.32	0.32	0.29	0.37	0.35	0.32	0.34	0.47	0.41	0.42	0.44	0.33	0.31

Table 3: Comparison of controllability error when training Mistral-7B-v0.1 with IFT and MCTune (ours).

C Additional Experiments With More Baselines

In this section, we conduct experiments with other baselines that include:

GPT-4. Similar to other baselines, we query GPT-4 (gpt-4-0125-preview) using the prompt template in Figure 1 in a zero-shot manner. Since GPT-4 is very expensive, we only evaluate it on 500 test examples.

Few-shot baselines. To introduce stronger baselines, we consider the few-shot setting where a few exemplars are included in the prompt to boost the performance of the underlying model. We carefully select 5 exemplars from the training set that can cover all the linguistic complexities.

RL-based methods. We compare with a variant of RLHF proposed in Ahmadian et al. (2024), where instead of training the LLM policy with PPO, they use REINFORCE. To adapt this method to our task, we use the negative controllability error as the rewards. Since they can be computed deterministically, we do not need to use neural networks for reward modeling. More specifically, given $(\mathbf{x}, \mathbf{y}, f_C) \sim D$ and a response $\hat{\mathbf{y}}$, the reward is

$$R_{f_C}(\hat{\mathbf{y}}, \mathbf{y}) = -\frac{1}{|C|} \sum_{i=1}^{|C|} \frac{|f_{C_i}(\hat{\mathbf{y}}) - f_{C_i}(\mathbf{y})|}{M_i - m_i}$$

where M_i and m_i are the maximum and minimum values of the i -th linguistic complexity in the training set. The REINFORCE training objective is:

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}, f_C) \sim D} [(R_{f_C}(\hat{\mathbf{y}}, \mathbf{y}) - b) \nabla_{\theta} \log p_{\theta}(\hat{\mathbf{y}}|\mathbf{x})]$$

where $\hat{\mathbf{y}} \sim p_{\theta}(\cdot|\mathbf{x}, f_C(\mathbf{y}))$ and b is the average of all achieved rewards in previous steps. For the

Fine-grained RLHF, we are actively adapting their codebase to our task and will report the results shortly. As shown in Table 2, MCTune exhibits the lowest controllability errors on most linguistic complexities.

D Additional Experiments With Other LLM Architectures

In this section, we evaluate our method on Mistral-7B-v0.1 (Jiang et al., 2023) to see how it generalizes to LLM architectures other than LLaMA. As shown in Table 3, MCTune keeps improving the controllability of Mistral on most attributes.