

RULER : A Model-Agnostic Method to Control Generated Length for Large Language Models

Anonymous ACL submission

Abstract

The instruction-following ability of large language models enables humans to interact with AI agents in a natural way. However, when required to generate responses of a specific length, large language models often struggle to meet users' needs due to their inherent difficulty in accurately perceiving numerical constraints. To explore the ability of large language models to control the length of generated responses, we propose the Target Length Generation task (TLG) and design two metrics, Precise Match (PM) and Flexible Match (FM) to evaluate the model's performance in adhering to specified response lengths. Furthermore, we introduce a novel, model-agnostic approach called RULER, which employs Meta Length Tokens (MLTs) to enhance the instruction-following ability of large language models under length-constrained instructions. Specifically, RULER equips LLMs with the ability to generate responses of a specified length based on length constraints within the instructions. Moreover, RULER can automatically generate appropriate MLT when length constraints are not explicitly provided, demonstrating excellent versatility and generalization. Comprehensive experiments show the effectiveness of RULER across different LLMs on Target Length Generation Task, e.g., 28.25 average gain on FM, 18.40 average gain on PM. In addition, we conduct extensive ablation experiments to further substantiate the efficacy and generalization of RULER. Our code and data will be made publicly available upon publication.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across a variety of natural language tasks and are increasingly being utilized in various fields (Vaswani et al., 2017; Devlin et al., 2019; Brown et al., 2020). A primary area of interest is the instruction following ability,

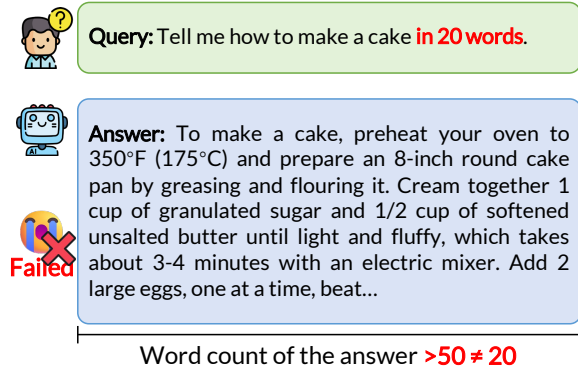


Figure 1: Existing LLMs lack the capability to follow instructions for generating texts of a specified length.

referring to their capability to execute tasks or generate outputs based on instructions (Ouyang et al., 2022; Wei et al., 2022a). It reflects the model's effectiveness in understanding and responding to instructions.

The practical challenges highlight the complexity of achieving precise instruction following, particularly when users require control over the output's length. Users frequently give LLMs various instructions, such as "Tell me how to make a cake in 20 words", "Write a blog post using 50 words", "Compose a 300-word story for me" and so on. These instructions challenge the instruction following capability of LLMs. To explore how well LLMs handle such challenges, we focus on the scenario where users specify the target length of the responses. We pose the question, "Can LLMs accurately generate with target length?" and introduce the *Target Length Generation Task (TLG)*. We create a test dataset with various target lengths and introduce two evaluation metrics: Precise Match (PM) and Flexible Match (FM). Our findings reveal that current LLMs generally perform poorly in this task, indicating considerable room for improvement. Potential reasons for this include the inherent complexity of the lan-

069 guage, limitations in training data, and insufficient
070 understanding of context, among other factors.

071 To address aforementioned issues, we introduce
072 RULER, a model-agnostic approach designed to en-
073 hance the instruction-following capability of LLMs
074 through *Meta Length Tokens (MLTs)*. *MLTs* are de-
075 signed to control model’s responses. By utilizing
076 RULER, LLMs can generate responses that meet tar-
077 get lengths. We create a dataset with *MLTs* \mathcal{D}_{MLT}
078 for end-to-end training of LLMs. LLMs learn to
079 generate *MLT* and the corresponding length re-
080 sponse after training. During inference, if a target
081 length is provided, RULER can transform it into
082 a *MLT* and generate responses that meet the re-
083 quirement. If no target length is specified, it first
084 generates a *MLT*, then the response, ensuring its
085 length aligns with the generated *MLT*

086 We apply RULER to various large language mod-
087 els and test them on *TLG*. Each model demonstrates
088 significant improvements. Furthermore, to rigor-
089 ously test the capabilities of RULER, we randomly
090 sample the dataset provided by Li et al. (2024a).
091 We provide nine target lengths for each question
092 and test the performance. RULER shows a mini-
093 mum accuracy of 52.72, marking an improvement
094 of 25.89 compared to the original models. Ad-
095 ditionally, to test the ability in scenarios without
096 target lengths, we assess whether the automatically
097 generated *MLT* and the corresponding response
098 lengths match. The lowest accuracy is 76.00. Ad-
099 ditionally, we test RULER on three other benchmarks
100 to observe whether the models’ performance is af-
101 fected.

102 Our contributions can be summarized as follows:

- 103 • We introduce the *Target Length Generation*
104 *Task (TLG)*, which designed to assess the in-
105 struction following capability of LLMs. It
106 evaluates how well models generate responses
107 of target lengths as directed by instructions.
- 108 • We propose RULER, a novel and model-
109 agnostic approach which employs the *Meta*
110 *Length Tokens (MLTs)*. Through end-to-end
111 training, it enables models to generate re-
112 sponse matching the target lengths indicated
113 by *MLTs*.
- 114 • We demonstrate that RULER significantly en-
115 hances the performance of various models on
116 the *TLG*. Further experiments have also vali-
117 dated the effectiveness and generalizability of
118 RULER.

2 Related Work 119

2.1 Large Language Model 120

121 The advent of LLMs has revolutionized the field
122 of natural language processing and become a mile-
123 stone (Vaswani et al., 2017; Devlin et al., 2019;
124 Brown et al., 2020; Zhang et al., 2023a). Large lan-
125 guage models have achieved success across various
126 NLP tasks. Models such as GPT-4(Achiam et al.,
127 2023), Llama-3(AI@Meta, 2024), and Qwen(Bai
128 et al., 2023), known for their powerful capabilities,
129 are increasingly serving as the foundation for var-
130 ious applications and making significant inroads
131 into diverse fields, exerting a substantial impact.
132 In-context learning enables LLMs to infer and gen-
133 erate responses solely based on the contextual in-
134 formation provided within a prompt(Dong et al.,
135 2022; Wei et al., 2022b). This capability allows
136 the models to exhibit a high degree of flexibility
137 and adaptability across a variety of tasks(Levine
138 et al., 2022; Chen et al., 2022; Zhao et al., 2021).
139 CoT further excavates and demonstrates the pow-
140 erful logical reasoning capabilities of LLMs(Wei
141 et al., 2022c; Huang and Chang, 2023; Zhang et al.,
142 2023b).

2.2 Instruction Following 143

144 Instruction following refers to the ability of large
145 language models to comprehend and execute given
146 natural language instructions (Brown et al., 2020;
147 Ouyang et al., 2022; Wei et al., 2022a; Zhou et al.,
148 2023a). This capability enables the models to per-
149 form a broad spectrum of tasks, from simple query
150 responses to complex problem-solving and content
151 generation, tailored to specific user requests.

152 In practical deployments, models may not adhere
153 to comply with user instructions, exhibiting behav-
154 iors that deviate from anticipated outcomes. This
155 includes generating responses unrelated to explicit
156 instructions, emitting redundant or erroneous in-
157 formation, or entirely ignoring specified directives
158 (Gehman et al., 2020; Kenton et al., 2021; Wei
159 et al., 2024). To enhance the instruction following
160 capability of LLMs, open-domain instruction fol-
161 lowing data is frequently used for training. Several
162 prominent studies have explored the construction
163 of instruct-tuning data, to achieve efficient and cost-
164 effective results(Li et al., 2024b; Cao et al., 2024;
165 Liu et al., 2024; Xu et al., 2024).

Level	Target Length	Precise Match (PM)	Flexible Match (FM)
Level:0	10	± 10	(0, 20]
	30	± 10	(20, 40]
	50	± 10	(40, 60]
	80	± 10	(60, 100]
Level:1	150	± 20	(100, 200]
	300	± 20	(200, 400]
	500	± 50	(400, 600]
Level:2	700	± 70	(600, 800]
	>800	(800, ∞)	(800, ∞)

Table 1: Nine target lengths and their corresponding match ranges categorized as Precise Match (PM) and Flexible Match (FM). Target lengths are classified into three categories, *Level:0*, *Level:1*, and *Level:2*.

2.3 Meta Token

Recently, an increasing number of studies have employed custom tokens within language models to execute specific functions or enhance performance. Todd et al. (2024) report findings that the hidden states of language models capture representations of these functions, which can be condensed into a Function Vector (FV). Furthermore, their research demonstrates that FV can effectively guide language models in performing specific tasks.

Numerous studies have utilized meta tokens to compress prompts, thereby enhancing the the inference capability of models (Li et al., 2023; Liu et al., 2023; Zhang et al., 2024). Mu et al. (2023) introduce the concept of "gist tokens", which can be cached and reused for compute efficiency. Further Jiang et al. (2024) utilize a hierarchical and dynamic approach to extend the concept, proposing "HD-Gist tokens" to improve model performance.

3 Can LLMs Accurately Generate with Target Length?

In this section, we examine the capability of LLMs to generate responses of a target length. Initially, we introduce *Target Length Generation Task (TLG)*. Subsequently, we establish various target lengths and two evaluation metrics (§3.1). We then detail the experimental setup and assess the ability of LLMs to generate responses at target lengths (§3.2). Finally, we present the outcomes of the experiments (§3.3).

3.1 Target Length Generation Task

To assess the ability of existing LLMs to control the length of generated response, we develop the *TLG*. This task assesses the models' ability in producing responses that match target lengths as directed The

designed target lengths are detailed in Table 1. Additionally, we divide these nine target lengths into three *levels*: *Level:0*, *Level:1*, and *Level:2*.

Given that generating responses with target lengths is challenging for existing LLMs, we develop two metrics to evaluate the accuracy of response lengths.

- **Precise Match (PM):** This metric requires that the length of the generated response be very close to the target length. For different *Level*, a precise tolerance range is set (± 10 , ± 20 , ...) necessitating that the response length stringently conforms to these defined limits.
- **Flexible Match (FM):** This metric requires a broader tolerance interval for target length. For longer texts, the range incrementally widens to meet response generation requirements.

For the N responses, we assess whether each response meets the target length, then calculating the PM and FM scores of the model.

$$PM = \frac{\sum_{i=1}^N \mathbb{1}(\text{lb}_{TL_i}^P < L(c_i) \leq \text{ub}_{TL_i}^P)}{N} \quad (1)$$

$$FM = \frac{\sum_{i=1}^N \mathbb{1}(\text{lb}_{TL_i}^F < L(c_i) \leq \text{ub}_{TL_i}^F)}{N} \quad (2)$$

where: c_i denotes the i -th response generated by LLM. The function $L(\cdot)$ calculates the word count of the input string. $\text{lb}_{TL_i}^P$ and $\text{ub}_{TL_i}^P$ denote the lower and upper bounds of the precise match range associated with the target length of i -th response. $\text{lb}_{TL_i}^F$ and $\text{ub}_{TL_i}^F$ denote the lower and upper bounds of the flexible match range associated with the target length of i -th response.

Model	Params	Target Length Generation Task (TLG)							
		Level:0		Level:1		Level:2		All Level	
		PM	FM	PM	FM	PM	FM	PM	FM
Mistral	7B	20.29	23.50	16.77	48.32	3.62	5.66	15.45	27.70
Gemma	2B	20.95	23.17	8.69	24.24	0.23	0.23	12.35	18.45
	7B	15.52	18.85	11.74	35.82	0.45	0.45	10.95	20.35
Llama3	8B	34.59	<u>40.02</u>	<u>29.73</u>	<u>65.70</u>	18.10	21.04	<u>29.35</u>	<u>44.25</u>
	70B	58.76	64.52	36.59	77.90	36.43	41.18	46.55	63.75
InternLM2	7B	6.65	7.21	8.69	27.44	19.68	22.40	10.20	17.20
	20B	8.98	9.87	10.98	34.45	17.42	20.14	11.50	20.20
DeepSeek-LLM	7B	28.16	31.37	17.68	44.36	10.86	13.12	20.90	31.60
	67B	26.94	30.27	17.07	49.54	9.50	11.99	19.85	32.55
Yi-1.5	6B	23.50	25.83	16.46	48.78	18.10	20.36	20.00	32.15
	9B	25.28	29.16	17.38	44.36	<u>24.43</u>	<u>29.41</u>	22.50	34.20
	34B	28.82	33.59	26.07	65.40	21.27	25.79	26.25	42.30
Qwen1.5	7B	24.28	27.38	14.33	46.19	9.05	11.99	17.65	30.15
	14B	28.27	31.49	18.45	43.90	11.09	14.25	21.25	31.75
	32B	32.59	36.25	22.26	49.39	21.49	25.34	26.75	38.15
	72B	<u>35.59</u>	39.69	18.29	49.70	3.85	6.11	22.90	35.55

Table 2: Overall results of different LLMs of *TLG*. All models used are either chat or instruct models. In models belonging to the same series but varying in parameter sizes, those with larger parameters typically exhibit superior performance. The best-performing model in each *Level* is **in-bold**, and the second best is underlined.

3.2 Experimental Setup

Dataset. We employ a two-stage data construction method for this study. Initially, we randomly sample 2,000 data from OpenHermes2.5 (Teknium, 2023). To enhance the complexity of the task and prevent data leakage, the second stage involved uses only the questions from these samples. Additionally, we randomly assign one of nine target lengths for the responses. The distribution of target length in the *TLG* dataset is shown in Figure 3. Further details regarding the format of the *TLG* dataset are provided in Appendix A.1.

Models & Prompt Templates. We conduct extensive experiments with open-source LLMs, specifically the chat or instruct version. The specific models used are listed in Table 7. We evaluate each model using its own prompt template, as detailed in Table 8.

To integrate the target length into the prompt, we modify the sentence The response should have a word count of {Target Length} words into each question. For target length is >800, we replace this with more than 800.

Hardware & Hyperparameters. All experiments are conducted on NVIDIA A100 GPUs. Inference is performed using the vllm (Kwon et al., 2023), with temperature set to 0 and

max_tokens set to 2,048 in the SamplingParams, thereby employing greedy decoding for inference. The model_max_length for all models is consistent with their respective configurations, as shown in Table 7.

3.3 Results and Analysis

Table 2 displays the PM and FM scores of various models at different *Levels*. Generally, models with advanced capabilities achieve higher PM and FM scores, indicating stronger adherence to instructions. This observation aligns with human expectations. The Meta-Llama-3-70B-Instruct (AI@Meta, 2024) achieves a FM score exceeding 60 at *All Level*. Within models from the same series but with different parameter sizes, larger models, as indicated by parameter size, generally demonstrate improved performance. Notably, the Qwen1.5 (Bai et al., 2023) with 72B parameters underperforms compared to its 32B variant.

For most models, scores are lowest at *Level:2*, suggesting significant potential for enhancement in producing longer responses. In contrast, scores at *Level:1* are typically the highest, suggesting a preference for generating shorter responses, which are more common at this level. This trend may be attributed to the prevalence of shorter responses in the training datasets utilized for model fine-tuning,

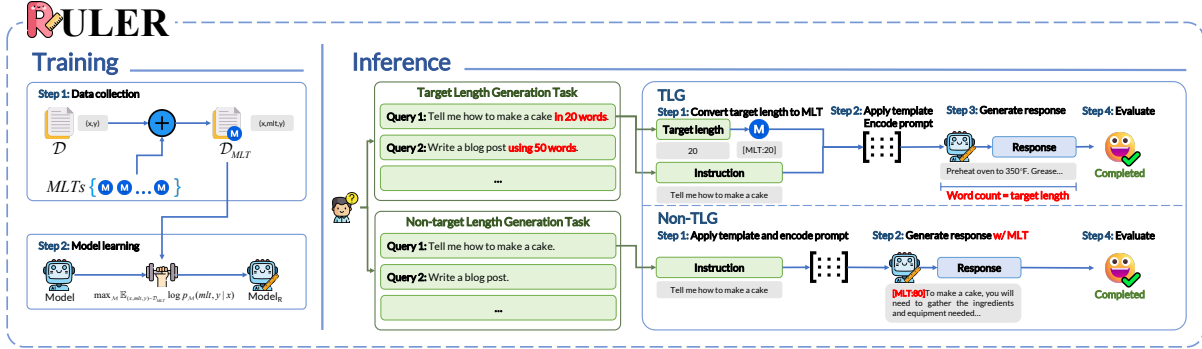


Figure 2: Overview of RULER. The method is divided into two parts: training and inference. The figure illustrates the main content of both sections. Additionally, in the inference section, we show two scenarios: *TLG* and *non-TLG* to show the difference.

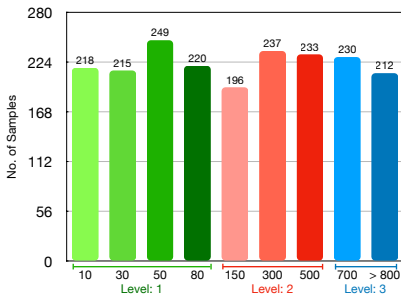


Figure 3: Target length distribution in *TLG* dataset. The count of each target length is approximately 200.

<i>MLT</i>	Range of Variation	No. in \mathcal{D}_{MLT}
[<i>MLT</i> : 10]	[5, 15)	20,000
[<i>MLT</i> : 30]	[25, 35)	20,000
[<i>MLT</i> : 50]	[45, 55)	20,000
[<i>MLT</i> : 80]	[75, 85)	20,000
[<i>MLT</i> : 150]	[145, 155)	20,000
[<i>MLT</i> : 300]	[295, 305)	10,333
[<i>MLT</i> : 500]	[495, 505)	2,317
[<i>MLT</i> : 700]	[695, 705)	497
[<i>MLT</i> : >800]	(800, ∞)	8,082

Table 3: Meta length tokens in RULER showing their range of variation in data collection and counts in \mathcal{D}_{MLT} .

287 which influences their generative biases. Further-
 288 more, the PM and FM scores for each model across
 289 various target lengths are detailed in Appendix A.3.

4 RULER: Meta Length Token Controlled Generation

292 In this section, we first introduce RULER, encom-
 293 passing the design of the *Meta Length Tokens*
 294 (*MLTs*), the data collection and the learning process
 295 associated with the models (§4.1). Subsequently,
 296 we detail the difference in the generation of RULER
 297 under two scenarios: *TLG* and *non-TLG* (§4.2).

4.1 Method

298 **RULER.** We introduce RULER, as illustrated in
 299 Figure 2, to effectively control the response length
 300 of LLMs using *MLTs*. The *MLTs* represent the
 301 model’s response length range and aim to enhance
 302 its capability on the *TLG* task. Our end-to-end
 303 training enables the LLMs to automatically gener-
 304 ate *MLTs* in various scenarios, regardless of target
 305 length requirements. *MLTs* (Table 3) offer more
 306 precise control than traditional text prompt meth-
 307 ods, which often prove insufficiently constraining.
 308

309 **Data collection for RULER.** For common fine-
 310 tuning training datasets, the format typically consist
 311 of input-output pairs (x, y) . Following Zhou
 312 et al. (2023b), we calculate the word count of y
 313 for each entry. Based on the predefined *MLTs* in Table
 314 3 and their range of variation, we aim to match
 315 each y to a corresponding *mlt* based on its word
 316 count. If a match is found, the data is reformatted
 317 as (x, mlt, y) . This method aids in the construction
 318 of the fine-tuning training dataset \mathcal{D}_{MLT} , detailed
 319 in Algorithm B.

320 **RULER learning.** To minimize changes to the
 321 model’s generation pattern and ensure stability in
 322 non-*TLG* scenario, we position the *MLT* immedi-
 323 ately before the original response during the con-
 324 struction of fine-tuning data. This strategy main-
 325 tains the model chat template. Consequently, the
 326 combination of *mlt* and the original response y
 327 forms a new complete response y' .

328 We conduct the training of the RULER \mathcal{M} on
 329 the curated corpus \mathcal{D}_{MLT} , which is augmented
 330 with *Meta Length Tokens* \mathcal{D}_{MLT} , employing the

Model	Target Length Generation Task (TLG)							
	Level:0		Level:1		Level:2		All Level	
	PM	FM	PM	FM	PM	FM	PM	FM
Mistral-7B-Instruct	20.29	23.50	16.77	48.32	3.62	5.66	15.45	27.70
+RULER	68.29 \uparrow 48.00	73.84 \uparrow 50.34	41.92 \uparrow 25.15	72.41 \uparrow 24.09	33.71 \uparrow 30.09	37.56 \uparrow 31.90	52.00 \uparrow 36.55	65.35 \uparrow 37.65
gemma-7b-it	15.52	18.85	11.74	35.82	0.45	0.45	10.95	20.35
+RULER	71.84 \uparrow 56.32	77.16 \uparrow 58.31	47.26 \uparrow 35.52	78.66 \uparrow 42.84	33.26 \uparrow 32.81	36.20 \uparrow 35.75	55.25 \uparrow 44.30	68.60 \uparrow 48.25
Llama-3-8B-Instruct	34.59	40.02	29.73	65.70	18.10	21.04	29.35	44.25
+RULER	76.27 \uparrow 41.68	80.49 \uparrow 40.47	45.58 \uparrow 15.85	78.05 \uparrow 12.35	18.33 \uparrow 0.23	20.59 \downarrow 0.45	53.40 \uparrow 24.05	66.45 \uparrow 22.20
deepseek-llm-7b-chat	28.16	31.37	17.68	44.36	10.86	13.12	20.90	31.60
+RULER	61.75 \uparrow 33.59	66.30 \uparrow 34.93	25.76 \uparrow 8.08	62.35 \uparrow 17.99	9.50 \downarrow 1.36	9.95 \downarrow 3.17	38.40 \uparrow 17.50	52.55 \uparrow 20.95
Yi-1.5-6B-Chat	23.50	25.83	16.46	48.78	18.10	20.36	20.00	32.15
+RULER	58.54 \uparrow 35.04	64.08 \uparrow 38.25	32.47 \uparrow 16.01	69.36 \uparrow 20.58	11.76 \downarrow 6.34	13.35 \downarrow 7.01	39.65 \uparrow 19.65	54.60 \uparrow 22.45
Qwen1.5-7B-Chat	24.28	27.38	14.33	46.19	9.05	11.99	17.65	30.15
+RULER	58.98 \uparrow 34.70	64.75 \uparrow 37.37	29.27 \uparrow 14.94	63.72 \uparrow 17.53	14.71 \uparrow 5.66	17.65 \uparrow 5.66	39.45 \uparrow 21.80	54.00 \uparrow 23.85

Table 4: Overall results of various LLMs with RULER are presented. Additionally, we also annotate the table with the score changes compared to the original model. Consistent improvements in both PM and FM scores are observed across all Levels.

standard next token objective:

$$\max_{\mathcal{M}} \mathbb{E}_{(x, mlt, y) \sim \mathcal{D}_{MLT}} \log p_{\mathcal{M}}(mlt, y|x) \quad (3)$$

We concatenate the *MLT* directly to the beginning of *y* to compute the loss and use the *MLTs* to expand the original vocabulary \mathcal{V} .

4.2 RULER Inference

TLG scenario. In the *Target Length Generation (TLG)* scenario, the user’s instruction specifies a target length, decomposed into a question and a target length. The RULER converts this target length into the corresponding *MLT* and appends it to the model chat template. Subsequent to the *MLT*, RULER generates response that aligns with the target length, ensuring compliance with both the user’s question and the target length, as illustrated in Figure 2. This approach yields superior results compared to controlling outputs solely through prompts.

non-TLG scenario. In the non-*TLG* scenario, users provide straightforward instructions consisting solely of a question. RULER integrates these instructions directly into the model’s chat template for generation. Owing to its innovative design and the use of a standard next-token objective in training (Equation 3), RULER autonomously generates a *MLT* prior to producing the textual response. This *MLT* is designed to match the length of the content generated, thereby ensuring normal generation of the model in non-*TLG* scenarios, as illustrated in Figure 2.

5 Experiments

5.1 Experimental Setup

Dataset \mathcal{D}_{MLT} . To ensure balanced frequency distribution of each *Meta Length Token (MLT)* in \mathcal{D}_{MLT} , we set a maximum occurrence limit of 20,000 for each *MLT*. We construct \mathcal{D}_{MLT} from three datasets: OpenHermes2.5 (excluding data previously used in *TLG*) (Teknium, 2023), LongForm (Köksal et al., 2023), and ELI5 (Fan et al., 2019), in accordance with Algorithm 1. This approach aims to create a diverse dataset, particularly effective for generating longer content that is relatively rare. In total, \mathcal{D}_{MLT} comprises 121,229 entries, with the frequency of each *MLT* in Table 3. Moreover, we calculate the word count for each response in every dataset, allowing us to statistically analyze the *MLT* distribution, as detailed in Table 12.

LLMs. To comprehensively evaluate the performance of RULER across different models, we consider factors such as model size, open-source availability, and overall model performance. We select six representative LLMs are selected: Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), gemma-7b-it (Team et al., 2024), Llama-3-8B-Instruct (AI@Meta, 2024), deepseek-llm-7b-chat (DeepSeek-AI, 2024), Yi-1.5-6B-Chat (AI et al., 2024), and Qwen1.5-7B-Chat (Bai et al., 2023).

Evaluation Metric. Consistent with the *TLG* and compared to previous results, we also calculate

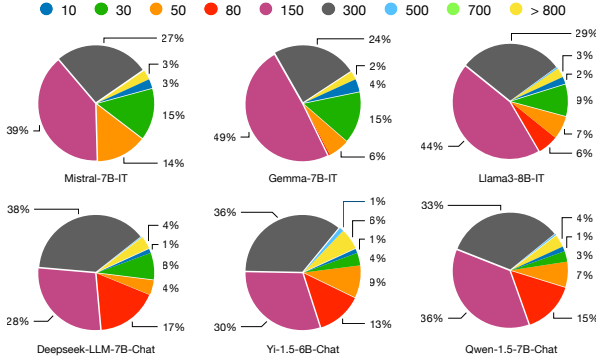


Figure 4: Distribution of *MLTs* generated by RULER in self-generated *MLT* experiment. The models demonstrate a preference for generating responses with target lengths of 150 and 300.

Model	FM	Avg WC
Mistral-7B-Instruct-v0.3 _R	76.00	247
gemma-7b-it _R	79.20	256
Llama-3-8B-Instruct _R	87.00	248
deepseek-llm-7b-chat _R	80.40	238
Yi-1.5-6B-Chat _R	80.40	257
Qwen1.5-7B-Chat _R	80.20	263

Table 5: The FM score and average word count of RULER with different models in self-generated *MLT* experiment. FM scores are notably high. Specifically, Mistral-7B-Instruct-v0.3_R recorded the lowest at 76.00, while Llama-3-8B-Instruct_R achieved the highest at 87.00.

PM and FM scores to assess the effectiveness of RULER.

5.2 Main Results

Table 4 presents a detailed comparison of PM and FM scores across various LLMs using RULER across different *Levels*. For information on model training see Appendix C.2.

Overall Performance Enhancement. Across all evaluated models, we observe a consistent improvement in both PM and FM scores at all *Levels*. The most significant improvement is observed in gemma-7b-it_R¹, with PM and FM scores increasing by 44.30 and 48.25, respectively. In contrast, the least improvement is noted with PM and FM rising by 17.50 and 20.95, respectively. These improvements indicate that RULER effectively enhances the model’s ability to generate content of target lengths. This suggests that using *MLT* to control output length is more effective than using prompts, as the model learns to generate content of corresponding lengths during fine-tuning. Additionally, RULER’s ability to enhance various models demonstrates its generalizability and scalability.

Despite these positive trends, some models, such as deepseek-llm-7b-chat_R, show a slight decrease in scores at *Level:2*. This is attributed to the insufficient data for *Level:2* in \mathcal{D}_{MLT} . The uneven distribution of data likely contributes to the slight decrease in scores at *Level:2*.

Different Level Analysis. At *Level:0*, all models show significant improvements in both PM and FM scores. Compared to other *Level*, each

¹Model name with _R means model with RULER

model achieves the highest PM and FM score improvements at *Level:0*. This enhancement occurs because the models are capable of generating responses of this length; however, their coarse length control impedes precise adherence to target length requirements. Our method significantly improves the models’ capacity to accurately control content length at *Level:0* more accurately, better meeting the target length requirements.

Moving to *Level:1*, while the improvements are not as pronounced as at *Level:0*, the models still exhibit significant gains in both PM and FM scores. At *Level:2*, the extent of score improvements varies across models. For instance, Mistral-7B-Instruct-v0.3_R and gemma-7b-it_R continue to show substantial score increases. In contrast, some models, such as Yi-1.5-6B-Chat_R, exhibit slight decreases, with reductions of 6.34 and 7.01 in PM and FM scores, respectively. These declines can be attributed to the relatively small number of *MLT* at *Level:2* in \mathcal{D}_{MLT} , which might differ from the original training data distribution of these models, leading to slight score reductions.

5.3 Do *MLTs* actually influence the length of the generated content?

To further investigate the effectiveness and scalability of *MLTs*, we designed two additional experiments: multi *MLT* generation experiment and self-generated *MLT* experiment.

Multi *MLT* Generation Experiment. To further validate the efficacy and robustness of RULER, we assess its ability to control response length. We randomly sample 200 entries from Arena-Hard-Auto (Li et al., 2024a) and subject each to all target

Model	Avg	HellaSwag	MMLU	TruthfulQA	Winogrande
Llama-3-8B-Instruct	66.95	78.84	65.77	51.66	71.51
+ RULER	62.85	77.40	52.86	52.55	68.59
deepseek-llm-7b-chat	62.29	79.65	51.35	47.92	70.24
+ RULER	61.80	80.50	50.12	47.34	69.22
Yi-1.5-6B-Chat	66.43	79.12	62.69	52.49	71.43
+ RULER	63.47	76.89	57.86	51.72	67.40
Qwen1.5-7B-Chat	64.52	78.67	60.56	53.58	65.27
+ RULER	61.27	74.79	54.34	50.19	65.75

Table 6: The results of 4 models with RULER in HellaSwag, MMLU, TruthfulQA and Winogrande benchmarks.

lengths (Table 1), culminating in 1,800 entries at last. Subsequently, we calculate the FM scores for each target length, using the original model as a baseline.

The results presented in Table 13 highlight the enhancements in model performance due to RULER. The FM scores achieved by RULER generally surpass those of the baseline models. Notably, even the well-performing Llama-3-8B-Instruct shows significant improvements. However, when the target length is 700, RULER shows a decline in FM if the baseline model already achieves a certain score. In contrast, RULER enhances performance if the baseline model is underperforming. This phenomenon is likely due to an imbalance in the \mathcal{D}_{MLT} , where responses of 700 words are infrequent and differ from the fine-tuning data of the baseline, potentially undermining performance. Overall, RULER significantly improves model performance.

Self-generated *MLT* Experiment. To validate RULER in generating *MLT* and responses under a non-*TLG* scenario, we use the Arena-Hard-Auto dataset without providing *MLTs*, thereby necessitating autonomous response generation by the model. We evaluate performance by cataloging the types and proportions of generated *MLTs* (Figure 4) and evaluating response length using FM score at the target lengths corresponding to the *MLTs* (Table 5).

Models show a preference for producing responses with target lengths of 150 and 300. This inclination is likely attributable to the complex nature of the queries in the Arena-Hard-Auto, which require longer responses for problem resolution. In the non-*TLG* scenario, the FM scores are notably high, with the Mistral-7B-Instruct-v0.3_R recording the lowest at 76.00 and Llama-3-8B-Instruct_R achieving the highest at 87.00. The average word count across all models approximates 250 words.

5.4 Evaluation RULER on Other Tasks

To evaluate the impact of RULER on other tasks, we conduct experiments utilizing four benchmark datasets: HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), TruthfulQA (Lin et al., 2022), and Winogrande (Sakaguchi et al., 2019). These benchmarks provide a comprehensive assessment across different task types. Further details about the experiments on the experiment can be found in Appendix C.4.

Table 6 illustrates that RULER marginally reduces performance on several tasks. Specifically, the MMLU dataset scores decline by 12.91 for Llama-3-8B-Instruct and 6.22 for Qwen1.5-7B-Chat, while other score changes remain within five points. These variations are considered acceptable because dataset \mathcal{D}_{MLT} primarily focuses on response length, without stringent criteria for data quality and distribution, leading to score fluctuations. We contend these fluctuations could be minimized or eliminated with consideration of data quality.

6 Conclusion

This study initially investigate the instruction following abilities of LLMs and introduces *Target Length Generation Task (TLG)*. Additionally, we propose RULER, a novel and model-agnostic method that controls generated length for LLMs. RULER utilizes the *MLT* and end-to-end training to enhance model performance. Experimental results demonstrate that substantial improvements in PM and FM scores across various models. Moreover, two additional experiments are conducted to further validate the efficacy of the proposed method. Finally, we assess performance across four different benchmarks to demonstrate its superiority.

531 Limitations

532 With the emergence of large language models
533 (LLMs), an increasing number of applications are
534 now utilizing LLMs. A particularly interesting
535 aspect is the instruction-following capabilities of
536 LLMs. In this paper, we analyze the capabilities
537 of LLMs solely from the perspective of controlling
538 generated length and propose a solution through
539 RULER. Instructions, which vary widely and repre-
540 sent a real-life scenario or application. We believe
541 addressing the challenges or solving widespread
542 issues across various instructions is crucial. We em-
543 ploy meta token to construct RULER and argue that
544 meta tokens offer more robust control over models
545 than prompts do. Exploring how to develop and
546 utilize models effectively with the help of tokens is
547 a profoundly important question.

548 Ethical Statements

549 This study concentrates on managing the output
550 length of Large Language Models (LLMs). While
551 our primary focus is on the length of generated
552 content, we have not assessed the potential for pro-
553 ducing toxic content. The research does not involve
554 human participants, nor does it handle personal or
555 sensitive information. We have used only open-
556 source or suitably licensed resources, thereby com-
557 plying with relevant standards. Additionally, all
558 training data employed are open-source, ensuring
559 the exclusion of any private or sensitive informa-
560 tion.

561 References

562 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
563 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
564 Diogo Almeida, Janko Altenschmidt, Sam Altman,
565 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
566 *arXiv preprint arXiv:2303.08774*.

567 01. AI, :, Alex Young, Bei Chen, Chao Li, Chen-
568 gen Huang, Ge Zhang, Guanwei Zhang, Heng Li,
569 Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong
570 Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang,
571 Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang,
572 Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng
573 Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai,
574 Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024.
575 *Yi: Open foundation models by 01.ai*.

576 AI@Meta. 2024. *Llama 3 model card*.

577 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
578 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
579 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,

Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, 580
Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, 581
Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong 582
Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng- 583
guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, 584
Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, 585
Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx- 586
uan Zhang, Yichang Zhang, Zhenru Zhang, Chang 587
Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang 588
Zhu. 2023. Qwen technical report. *arXiv preprint* 589
arXiv:2309.16609. 590

Tom Brown, Benjamin Mann, Nick Ryder, Melanie 591
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 592
Neelakantan, Pranav Shyam, Girish Sastry, Amanda 593
Askell, et al. 2020. Language models are few-shot 594
learners. *Advances in neural information processing* 595
systems, 33:1877–1901. 596

Yihan Cao, Yanbin Kang, Chi Wang, and Lichao Sun. 597
2024. *Instruction mining: Instruction data selection* 598
for tuning large language models. 599

Yanda Chen, Ruiqi Zhong, Sheng Zha, George Karypis, 600
and He He. 2022. *Meta-learning via language model* 601
in-context tuning. In *Proceedings of the 60th Annual* 602
Meeting of the Association for Computational Lin- 603
guistics (Volume 1: Long Papers), pages 719–730, 604
Dublin, Ireland. Association for Computational Lin- 605
guistics. 606

DeepSeek-AI. 2024. *Deepseek llm: Scaling open-* 607
source language models with longtermism. *arXiv* 608
preprint arXiv:2401.02954. 609

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 610
Kristina Toutanova. 2019. *BERT: Pre-training of* 611
deep bidirectional transformers for language under- 612
standing. In *Proceedings of the 2019 Conference of* 613
the North American Chapter of the Association for 614
Computational Linguistics: Human Language Tech- 615
nologies, Volume 1 (Long and Short Papers), pages 616
4171–4186, Minneapolis, Minnesota. Association for 617
Computational Linguistics. 618

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiy- 619
ong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and 620
Zhifang Sui. 2022. A survey on in-context learning.
621 *arXiv preprint arXiv:2301.00234*. 622

Angela Fan, Yacine Jernite, Ethan Perez, David Grang- 623
ier, Jason Weston, and Michael Auli. 2019. *Eli5:* 624
Long form question answering. In *Proceedings of* 625
ACL 2019. 626

Samuel Gehman, Suchin Gururangan, Maarten Sap, 627
Yejin Choi, and Noah A. Smith. 2020. *RealToxi-* 628
cityPrompts: Evaluating neural toxic degeneration 629
in language models. In *Findings of the Association* 630
for Computational Linguistics: EMNLP 2020, pages 631
3356–3369, Online. Association for Computational 632
Linguistics. 633

Dan Hendrycks, Collin Burns, Steven Basart, Andy 634
Zou, Mantas Mazeika, Dawn Song, and Jacob Stein- 635
hardt. 2021. Measuring massive multitask language 636

637	understanding. <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> .	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	691
638			692
639	Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards reasoning in large language models: A survey . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.		693
640			694
641			695
642			696
643			
644	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. Mistral 7b .	Junyi Liu, Liangzhi Li, Tong Xiang, Bowen Wang, and Yiming Qian. 2023. TCRA-LLM: Token compression retrieval augmented large language model for inference cost reduction . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9796–9810, Singapore. Association for Computational Linguistics.	697
645			698
646			699
647			700
648			701
649			702
650			703
651	Yichen Jiang, Marco Vecchio, Mohit Bansal, and Anders Johannsen. 2024. Hierarchical and dynamic prompt compression for efficient zero-shot API usage . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 2162–2174, St. Julian’s, Malta. Association for Computational Linguistics.	Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2024. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning . In <i>The Twelfth International Conference on Learning Representations</i> .	704
652			705
653			706
654			707
655			708
656		Jesse Mu, Xiang Lisa Li, and Noah Goodman. 2023. Learning to compress prompts with gist tokens . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	709
657			710
658			711
659	Zachary Kenton, Tom Everitt, Laura Weidinger, Iason Gabriel, Vladimir Mikulik, and Geoffrey Irving. 2021. Alignment of language agents .		712
660			
661	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	713
662			714
663			715
664			716
665			717
666		Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale. <i>arXiv preprint arXiv:1907.10641</i> .	719
667			720
668	Abdullatif K��ksal, Timo Schick, Anna Korhonen, and Hinrich Sch��tze. 2023. Longform: Effective instruction tuning with reverse instructions .		721
669			722
670		Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Riviere, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, L��onard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Am��lie H��liou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Cl��ment Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Miku��a, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas,	723
671	Yoav Levine, Noam Wies, Daniel Jannai, Dan Navon, Yedid Hoshen, and Amnon Shashua. 2022. The inductive bias of in-context learning: Rethinking pre-training example design . In <i>International Conference on Learning Representations</i> .		724
672			725
673			726
674			727
675			728
676	Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. 2024a. From live data to high-quality benchmarks: The arena-hard pipeline .		729
677			730
678			731
679			732
680	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. 2023. Compressing context to enhance inference efficiency of large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6342–6353, Singapore. Association for Computational Linguistics.		733
681			734
682			735
683			736
684			737
685			738
686	Yunshui Li, Binyuan Hui, Xiaobo Xia, Jiayi Yang, Min Yang, Lei Zhang, Shuzheng Si, Ling-Hao Chen, Junhao Liu, Tongliang Liu, Fei Huang, and Yongbin Li. 2024b. One-shot learning as instruction data prospector for large language models .		739
687			740
688			741
689			742
690			743
			744
			745
			746
			747
			748

749	Shree Pandya, Siamak Shakeri, Soham De, Ted Klimentko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology .	
750		
751		
752		
753		
754		
755		
756		
757		
758		
759		
760	Teknium. 2023. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants .	
761		
762	Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2024. Function vectors in large language models . In <i>The Twelfth International Conference on Learning Representations</i> .	
763		
764		
765		
766		
767	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	
768		
769		
770		
771		
772	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2024. Jailbroken: How does llm safety training fail? <i>Advances in Neural Information Processing Systems</i> , 36.	
773		
774		
775		
776	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners . In <i>International Conference on Learning Representations</i> .	
777		
778		
779		
780		
781	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022b. Emergent abilities of large language models . <i>Transactions on Machine Learning Research</i> . Survey Certification.	
782		
783		
784		
785		
786		
787		
788		
789	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022c. Chain of thought prompting elicits reasoning in large language models . In <i>Advances in Neural Information Processing Systems</i> .	
790		
791		
792		
793		
794	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. WizardLM: Empowering large pre-trained language models to follow complex instructions . In <i>The Twelfth International Conference on Learning Representations</i> .	
795		
796		
797		
798		
799		
800	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i> .	
801		
802		
803		
	Lei Zhang, Yunshui Li, Ziqiang Liu, Jiayi Yang, Junhao Liu, and Min Yang. 2023a. Marathon: A race through the realm of long context with large language models .	804 805 806 807
	Peitian Zhang, Zheng Liu, Shitao Xiao, Ninglu Shao, Qiwei Ye, and Zhicheng Dou. 2024. Soaring from 4k to 400k: Extending llm’s context with activation beacon. <i>arXiv preprint arXiv:2401.03462</i> .	808 809 810 811
	Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2023b. Automatic chain of thought prompting in large language models . In <i>The Eleventh International Conference on Learning Representations</i> .	812 813 814 815
	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In <i>International conference on machine learning</i> , pages 12697–12706. PMLR.	816 817 818 819 820
	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, LILI YU, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023a. LIMA: Less is more for alignment . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	821 822 823 824 825 826
	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023b. Instruction-following evaluation for large language models. <i>arXiv preprint arXiv:2311.07911</i> .	827 828 829 830 831

A Target Length Generation Task Details

In this section, we present the experimental details of the *Target Length Generation (TLG)*.

A.1 TLG Dataset

Dataset constructed for the *TLG*, totaling 2,000 entries.

TLG Dataset

```
{
  "id": "0"
  "Instruction": "How can I generate an AI model that can classify articles of clothing as shorts, skirts, or pants based on their descriptions?",
  "TargetLength": "50"
}
[...]
```

```
{
  "id": "1999"
  "Instruction": "You will be given several pieces of information about someone, and you will have to answer a question based on the information given. \nJohn is taller than Bill. Mary is shorter than John. Question: Who is the tallest person?",
  "TargetLength": "30"
}
```

A.2 Models & Prompt Templates

In this appendix, we list the models in the *TLG*, including their fullname, params, context length and vocab size. All models are downloaded from Huggingface² and inference is executed using vllm (Kwon et al., 2023).

Model	Model Full Name	Params	Context Length	Vocab Size
Mistral	Mistral-7B-Instruct-v0.3	7B	32,768	32,768
Gemma	gemma-2b-it	2B	8,192	256,000
	gemma-7b-it	7B	8,192	256,000
Llama3	Meta-Llama-3-8B-Instruct	8B	8,192	128,256
	Meta-Llama-3-70B-Instruct	70B	8,192	128,256
InternLM2	InternLM2-Chat-7B	7B	32,768	92,544
	InternLM2-Chat-20B	20B	32,768	92,544
DeepSeek-LLM	deepseek-llm-7b-chat	7B	4,096	102,400
	deepseek-llm-67b-chat	67B	4,096	102,400
Yi-1.5	Yi-1.5-6B-Chat	6B	4,096	64,000
	Yi-1.5-9B-Chat	9B	4,096	64,000
	Yi-1.5-34B-Chat	34B	4,096	64,000
Qwen1.5	Qwen1.5-7B-Chat	7B	32,768	151,936
	Qwen1.5-14B-Chat	14B	32,768	151,936
	Qwen1.5-32B-Chat	32B	32,768	151,936
	Qwen1.5-72B-Chat	72B	32,768	151,936

Table 7: All models used in *TLG*

²<https://huggingface.co/>

Model	Prompt Template	Eos Tokens
Mistral	<s>[INST] {Instruction} [/INST]	</s>
Gemma	<bos><start_of_turn>user\n{Instruction} <end_of_turn>\n<start_of_turn>model\n	<eos>
Llama3	< begin_of_text >< start_header_id >user < end_header_id >\n\n{Instruction}< eot_id > < start_header_id >assistant< end_header_id >\n\n	< end_of_text >,< eot_id >
InternLM2	<s>< im_start >user\n{Instruction} < im_end >\n< im_start >assistant\n	</s>,< im_end >
DeepSeek-LLM	< begin_of_sentence >User: {Instruction} \n\nAssistant:	< end_of_sentence >
Yi-1.5	< im_start >user\n{Instruction}< im_end > \n< im_start >assistant\n	< im_end >,< endoftext >
Qwen1.5	< im_start >system\nYou are a helpful assistant. < im_end >\n< im_start >user\n{Instruction} < im_end >\n< im_start >assistant\n	< im_end >,< endoftext >

Table 8: Prompt templates and Eos tokens for all models used in *TLG*.

A.3 Results on Different Target Length

Here, we present the FM and PM scores of the models at all target lengths.

A.3.1 *Level:0*

The PM and FM scores for each model at *Level:0* are shown in Table 9.

Model	Params	<i>Level:0</i>							
		10		30		50		80	
		PM	FM	PM	FM	PM	FM	PM	FM
Mistral	7B	30.73	30.73	18.60	18.60	16.87	16.87	15.45	28.64
Gemma	2B	21.56	21.56	30.23	30.23	20.88	20.88	11.36	20.45
	7B	12.39	12.39	18.14	18.14	18.88	18.88	12.27	25.91
Llama3	8B	45.41	45.41	35.35	35.35	33.73	33.73	24.09	<u>46.36</u>
	70B	60.55	60.55	66.05	66.05	61.45	61.45	46.82	70.45
InternLM2	7B	17.89	17.89	6.98	6.98	1.20	1.20	1.36	3.64
	20B	20.64	20.64	8.84	8.84	2.81	2.81	4.55	8.18
DeepSeek-LLM	7B	<u>58.26</u>	<u>58.26</u>	25.12	25.12	17.67	17.67	13.18	26.36
	67B	46.79	46.79	20.47	20.47	22.09	22.09	19.09	32.73
Yi-1.5	6B	39.91	39.91	23.72	23.72	20.08	20.08	10.91	20.45
	9B	47.71	47.71	23.72	23.72	17.27	17.27	13.64	29.55
	34B	45.41	45.41	27.44	27.44	20.48	20.48	23.18	42.73
Qwen1.5	7B	31.19	31.19	25.58	25.58	22.89	22.89	17.73	30.45
	14B	45.87	45.87	28.84	28.84	26.51	26.51	12.27	25.45
	32B	46.79	46.79	33.95	33.95	29.32	29.32	20.91	35.91
	72B	39.45	39.45	<u>41.86</u>	<u>41.86</u>	<u>32.53</u>	<u>32.53</u>	<u>29.09</u>	45.91

Table 9: Results of different LLMs of *TLG* at *Level:0*. The best-performing model in each target length is **in-bold**, and the second best is underlined.

A.3.2 *Level:1*

The PM and FM scores for each model at *Level:1* are shown in Table 10.

Model	Params	Level:1					
		150		300		500	
		PM	FM	PM	FM	PM	FM
Mistral	7B	17.86	41.84	14.77	70.04	17.94	30.94
Gemma	2B	17.35	32.65	7.17	33.33	2.69	7.17
	7B	18.88	42.35	12.24	51.90	4.93	13.00
Llama3	8B	<u>38.27</u>	<u>70.92</u>	27.00	<u>78.90</u>	25.11	47.09
	70B	55.10	85.71	22.36	88.61	<u>35.43</u>	59.64
InternLM2	7B	9.18	20.92	5.91	37.55	11.21	22.42
	20B	9.69	22.96	9.28	45.99	13.90	32.29
DeepSeek-LLM	7B	15.31	37.24	18.14	60.76	19.28	33.18
	67B	9.18	34.69	19.83	71.73	21.08	39.01
Yi-1.5	6B	18.88	46.94	12.66	62.45	18.39	35.87
	9B	12.76	33.16	12.66	53.59	26.46	44.39
	34B	25.51	58.67	<u>24.05</u>	78.48	28.70	<u>57.40</u>
Qwen1.5	7B	9.69	29.59	7.17	61.60	26.01	44.39
	14B	5.61	16.84	10.97	56.12	37.67	54.71
	32B	20.92	43.37	14.77	53.59	31.39	50.22
	72B	13.27	35.20	12.66	64.98	28.70	46.19

Table 10: Results of different LLMs of *TLG* at *Level:1*. The best-performing model in each target length is **in-bold**, and the second best is underlined.

A.3.3 Level:2

The PM and FM scores for each model at *Level:2* are shown in Table 11.

Model	Params	Level:1			
		150		300	
		PM	FM	PM	FM
Mistral	7B	3.04	6.96	4.25	4.25
Gemma	2B	0.00	0.00	0.47	0.47
	7B	0.87	0.87	0.00	0.00
Llama3	8B	16.09	21.74	20.28	20.28
	70B	<u>24.35</u>	33.48	49.53	49.53
InternLM2	7B	18.70	23.91	20.75	20.75
	20B	17.39	22.61	17.45	17.45
DeepSeek-LLM	7B	9.13	13.48	12.74	12.74
	67B	9.13	13.91	9.91	9.91
Yi-1.5	6B	12.61	16.96	24.06	24.06
	9B	22.17	<u>31.74</u>	<u>26.89</u>	<u>26.89</u>
	34B	22.17	30.87	20.28	20.28
Qwen1.5	7B	12.17	17.83	5.66	5.66
	14B	15.22	21.30	6.60	6.60
	32B	23.91	31.30	18.87	18.87
	72B	6.09	10.43	1.42	1.42

Table 11: Results of different LLMs of *TLG* at *Level:2*. The best-performing model in each target length is **in-bold**, and the second best is underlined.

Algorithm 1 \mathcal{D}_{MLT} Data Creation

Require: Word count function $L(\cdot)$, meta length tokens $MLTs = \{MLT_0, MLT_1, \dots\}$
Input: Initial dataset \mathcal{D}
Output: \mathcal{D}_{MLT}

```

1:  $\mathcal{D}_{MLT} \leftarrow \{\}$ 
2: for each tuple  $(x, y)$  in  $\mathcal{D}$  do
3:    $mlt \leftarrow \text{None}$ 
4:   for each  $MLT$  in  $MLTs$  do
5:     if  $L(y) > lb_{MLT}$  and  $L(y) \leq ub_{MLT}$  then
6:        $mlt \leftarrow MLT$ 
7:       break
8:     end if
9:   end for
10:  if  $mlt$  is not None then
11:     $\mathcal{D}_{MLT} \leftarrow \mathcal{D}_{MLT} \cup \{(x, mlt, y)\}$ 
12:  end if
13: end for
14: return  $\mathcal{D}_{MLT}$ 

```

C Experiments Details

C.1 MLT in Datasets

To obtain data with varying response lengths for composing \mathcal{D}_{MLT} , particularly those responses exceeding 500, we integrate data from OpenHermes2.5 (Teknum, 2023), LongForm (Köksal et al., 2023) and ELI5 (Fan et al., 2019). We calculate the word count for each response in every dataset, allowing us to statistically analyze the MLT distribution, shown in Table 12.

MLT	OpenHermes2.5 (Teknum, 2023)	LongForm (Köksal et al., 2023)	ELI5 (Fan et al., 2019)
[MLT:10]	28,552	586	3,280
[MLT:30]	16,860	1,428	14,143
[MLT:50]	18,867	1,236	17,597
[MLT:80]	18,014	852	15,926
[MLT:150]	37,515	1,037	19,103
[MLT:300]	7,526	252	2,555
[MLT:500]	1,495	140	682
[MLT:700]	193	101	203
[MLT:800]	1,809	2,465	3,808

Table 12: MLT distribution in each dataset. The OpenHermes2.5 excludes the data utilized in TLG . The LongForm and ELI5 employs its training, validation, and test sets simultaneously. When multiple answers are available in the dataset, the longest answer is selected as the final response.

C.2 More Details of Training

More details of training. We use 4*A100 with 80GB Nvidia GPUs to train the models. The training utilizes both bf16 and tensor tf32 precision formats. The per-device training batch size is set to 4, with gradient accumulation is 8 steps. A cosine learning rate scheduler is applied, starting with an initial

learning rate of $2e-5$ and a warmup ratio of 0.05. All models are trained for 3 epochs. Additionally, log is set to print every 5 steps.

Loss. We document the changes in training loss for all models, as shown in Figure 5.

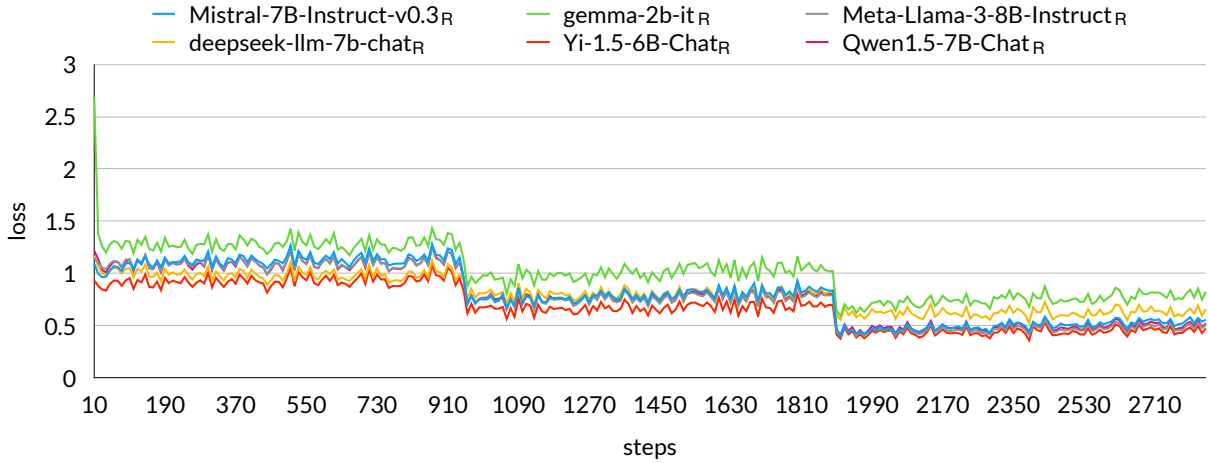


Figure 5: Training loss for models.

C.3 Multi *MLT* generation experiment

Here is the results in multi *MLT* generation experiment.

Model	FM of Different Target Length										Avg FM
	10	30	50	80	150	300	500	700	>800		
Mistral-7B-Instruct-v0.3 +RULER	0.5	0.0	0.5	2.0	18.5	50.5	20.5	3.0	2.5	10.89	
gemma-7b-it +RULER	13.0	17.0	15.5	26.0	54.5	76.5	17.5	0.0	0.0	24.44	
Llama-3-8B-Instruct +RULER	23.5	18.0	12.5	28.0	50.5	76.5	57.0	25.5	30.5	35.78	
deepseek-llm-7b-chat +RULER	36.5	16.0	12.5	17.5	23.5	60.5	36.5	16.0	22.5	26.83	
Yi-1.5-6B-Chat +RULER	26.5	16.5	14.5	14.5	18.5	42.5	35.0	33.5	28.5	25.56	
Qwen1.5-7B-Chat +RULER	13.5	17.0	9.5	16.0	6.5	51.0	57.5	22.5	4.5	22.00	

Table 13: Results in multi *MLT* generation experiment. Generally, the FM scores obtained via RULER surpass those of the baseline models.

C.4 More Details of Other Tasks

We tested the RULER on four benchmarks (HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), TruthfulQA (Lin et al., 2022), and Winogrande (Sakaguchi et al., 2019)) to examine whether the performance of the fine-tuned models varies on different tasks. We employ a 10-shot setting in HellaSwag, 5-shot setting in MMLU, 0-shot setting in TruthfulQA and 5-shot setting in Winogrande.