From Sub-Ability Diagnosis to Human-Aligned Generation: Bridging the Gap for Text Length Control via MARKERGEN

Anonymous ACL submission

Abstract

Despite the rapid progress of large language models (LLMs), their length-controllable text generation (LCTG) ability remains below expectations, posing a major limitation for practical applications. Existing methods mainly focus on end-to-end training to reinforce adherence to length constraints. However, the lack of decomposition and targeted enhancement of LCTG sub-abilities restricts further progress. To bridge this gap, we conduct a bottom-up decomposition of LCTG sub-abilities with human 011 patterns as reference and perform a detailed er-012 ror analysis. On this basis, we propose MARK-ERGEN, a simple-yet-effective plug-and-play approach that: (1) mitigates LLM fundamental deficiencies via external tool integration; (2) 017 conducts explicit length modeling with dynamically inserted markers; (3) employs a threestage generation scheme to better align length 019 constraints while maintaining content quality. Comprehensive experiments demonstrate that MARKERGEN significantly improves LCTG across various settings, exhibiting outstanding effectiveness and generalizability.

1 Introduction

037

041

As a fundamental attribute of text generation, ensuring controllability over text length is of great importance (Liang et al., 2024). Different text types (e.g., summary, story), user needs (e.g., preference for detailed or concise writing), and external requirements (e.g., social media character limits) shape varied length constraints, which are widely present in real-world scenarios (Zhang et al., 2023a). With the rapid development of LLMs, their expanding range of applications has made length-controllable text generation (LCTG) even more crucial in current era (Foster et al., 2024; Gu et al., 2024b).

However, the ongoing enhancements in LLM capabilities have yet to deliver the expected performance in LCTG while ensuring semantic integrity (Foster et al., 2024; Wang et al., 2024; Song et al., 2024). Yuan et al. (2024) reports that even advanced LLMs (e.g., GPT-4 Turbo (OpenAI, 2023)) violate the given length constraints almost 50% of the time. To address this, training-based methods (Park et al., 2024; Yuan et al., 2024; Jie et al., 2023; Li et al., 2024b) have been studied to reinforce LLMs' adherence to length constraints, yet they face two key challenges: (1) Limited generalization: Since text types are diverse and length constraints vary widely (e.g., ranging from an exact 500 words to coarse intervals like 500-600 words or below 500 words), training-based methods often fail to generalize effectively across different settings, as demonstrated in Appendix E.1. (2) Inferior controllability: These methods strengthen LCTG by enforcing implicit length modeling during generation in a top-down manner via training, lacking the decomposition and targeted enhancement of LCTG sub-capabilities, thereby limiting their progress (Retkowski and Waibel, 2024).

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

082

To fill this gap, we take humans as a reference and conduct a bottom-up decomposition of subcapabilities for LCTG. When writing a 500-word story, humans typically begin by planning the content and word allocation for each section. During writing, they continuously track the word count and compose the text in alignment with the plan. This process progressively tests four key abilities: (1) **Identifying** and splitting the words correctly. (2) **Counting** the words precisely. (3) **Planning** the word counts of each part to meet the length constraints. (4) **Aligning** the generated text with length constraints while ensuring semantic integrity.

On this basis, we conduct a decoupled error analysis of LCTG. The experimental results indicate that counting error > perception error > aligning error \gg planning error. This suggests that deficiencies in fundamental capabilities are the primary cause of LCTG's inferior performance. Meanwhile, it further explains why training-based approaches struggle to enhance LCTG effectively,



Figure 1: Sub-ability decomposition of LCTG and corresponding error analysis in LLMs.

as they are unable to provide fine-grained supervision signals for these fundamental capabilities.

084

100

101

103

105

106

109

110

111 112

113

114

115

116

Building upon this, we propose MARKERGEN, a simple-yet-effective, plug-and-play method for achieving high-quality LCTG. Specifically, to address LLMs' weaknesses in identifying and counting, we integrate external tokenizer and counter to track exact length information. To effectively convey these information to LLMs, we design an decaying interval insertion strategy that dynamically injects length markers during the generation process, enabling explicit length modeling while minimizing disruptions to semantic modeling. Furthermore, to mitigate alignment issues, we propose a three-stage decoupled generation paradigm that decouples semantic constraints from length constraints, ensuring that length constraints are better met without compromising content quality.

We conduct experiments with five LLMs on five benchmarks to validate the generalizability of MARKERGEN, covering cross-task (summarization, story generation, QA, heuristic generation), cross-scale (from 10+ to 1000+ words), crosslingual (English and Chinese) and cross-granularity (precise and rough constraints) settings. Experimental results demonstrate that under precise length constraints, MARKERGEN reduces length errors by 12.57% compared to baselines (with an average absolute error of 5.57%), while achieving higher quality scores and incurring only 67.6% of the cost. In range-based length constraints, MARK-ERGEN achieves a 99% acceptance rate, further validating its effectiveness. Finally, we probe into the working mechanism of MARKERGEN through

attention analysis: shallow layers primarily handle length modeling through markers, whereas deeper layers concentrate more on semantic modeling. 117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

2 Preliminaries

We model the LCTG process of LLMs by drawing an analogy to human patterns in this task. Specifically, the model first performs content and length planning based on task requirements and length constraints. Under this plan, the semantic space expands progressively at the word level during generation, accompanied by an implicit counting process. Meanwhile, length estimation acts as a real-time constraint, dynamically regulating further extension. Ultimately, the model strives to align the length constraints while preserving semantic integrity. From this perspective, the overall LCTG ability of LLMs can be systematically decomposed in a bottom-up manner into Identifying, Counting, Planning, and Aligning sub-capabilities (Figure 1). Below we explore LLMs' mastery of these abilities through detailed error analysis on TruthfulQA dataset (Lin et al., 2021).

2.1 Identifying Error

Identifying error refers to the misidentification of fundamental length units (e.g., words), leading to discrepancies between the model's estimated and actual text length. To systematically analyze this error, we instruct the model to recognize the length units of given text one by one. If we define a word as the length unit, the model should output like: "The [1 word] quick [2 words] fox [3 words] ...". On this basis, we calculate the identifying



Figure 2: Error analyses of fundamental abilities in LCTG across LLMs.

error rate e_{I} as follows:

149

150

151

152

153

154

155

157

159

160 161

162 163

164

165

166

167

168

169

171

172

174

175

179

180

181

$$e_{\rm I} = \frac{|N_{\rm pred}^1 - N_{\rm true}|}{N_{\rm true}} \tag{1}$$

where N_{pred}^1 is the model's predicted final count with 1 as count interval, and N_{true} is the actual count. We subtract the error rate obtained when replacing each word with the letter "A" (which barely assess the identifying ability) from e_{I} to further eliminate the influence of other potential factors. We explore the word and token as length unit respectively, as shown in Figure 2a.

Finding 1. *LLMs exhibit notable* e_I with both word and token as unit, showcasing their deficiencies in fundamental identifying ability.

Finding 2. Word yields lower e_1 than token, indicating that LLMs conduct length modeling primarily based on semantic perception rather than decoding mechanics.

2.2 Counting Error

Counting error refers to the inaccurate enumeration of length units in a given sequence, leading to deviations from the intended length. We analyze this error by prompting LLMs to count sequences with varied interval n. The case of n = 1 corresponds to identifying error (see §2.1). A larger n poses a greater challenge for counting accuracy. To decompose counting error from identifying error, we calculate $e_{\rm C}^n$ as follows:

176
$$e_{\rm IC}^n = \frac{|N_{\rm pred}^n - N_{\rm true}|}{N_{\rm true}}$$

$$e_{\rm C}^n = e_{\rm IC}^n - e_{\rm I} \tag{3}$$

Since LLMs exhibit negligible identifying error at the letter level (Figure 2a), error of counting letter serves as a direct measure of pure counting

Models	e_{P}	$s_{ m P}$	$\Delta E\left(\downarrow\right)$	$\Delta S\left(\uparrow\right)$
GPT-40	0.06	4.28	-5.31	0.05
GPT-40 mini	0.33	3.90	+2.11	0.03
Llama-3.1-70B-Instruct	0.00	3.90	-0.63	0.04
Qwen2.5-32B-Instruct	0.04	4.22	-8.93	0.02

Table 1: e_P and s_P denote planning error and planning quality score of LLMs. ΔE and ΔS quantify the LCTG error reduction and text quality gain from two-stage generation over one-stage generation.

ability. We also include a commonly used baseline where the LLMs conduct implicit counting (directly output the length of the entire given text). The results are shown in Figure 2b. 182

183

184

185

186

187

188

189

190

191

193

194

195

196

197

198

199

200

201

202

203

Finding 3. Naive implicit counting can lead to significant errors.

Finding 4. *Explicit counting combined with finegrained intervals leads to better length modeling. At smaller n, the error of explicit counting is significantly lower than that of implicit counting.*

2.3 Planning Error

Planning error refers to the misallocation of word counts across different sections, leading to a discrepancy from target length. For given query and precise length constraint N_{target} , we prompt LLMs to explicitly plan both content and length for each part of the response. We assess the quality ¹ of the plan s_{P} , and calculate the planning error rate e_{P} as:

$$e_{\rm P} = \frac{|N_{\rm plan} - N_{\rm target}|}{N_{\rm target}} \tag{4}$$

where N_{plan} denotes the total word count allocated by the model. Meanwhile, we calculate the reduction in final length error (ΔE) and the improve-

(2)

¹We use Qwen-Plus (Yang et al., 2024) as the judge with a scoring range of [1, 5]. See corresponding prompts in Appendix B.3



Figure 3: Absolute contribution of LCTG sub-capability deficiency on overall LCTG error across LLMs.

ment in content quality (ΔS) achieved by **planning followed by generation** compared to **direct generation**. The results are shown in Table 1.

Finding 5. *LLMs exhibit strong planning ability.* The generated plan effectively meets the length constraints while achieving a quality score of around 4,
 demonstrating well-structured content allocation.

Finding 6. *Planning before generation brings better results. Compared to direct generation, executing planning and generating sequentially for decomposition reduces length deviations while enhancing semantic quality.*

2.4 Aligning Error

211

212

213

214

216

218

219

221

224

225

226

227

231

234

235

Aligning error refers to the discrepancy between the model's perceived length and the target length, arising from the challenge of maintaining semantic integrity while adhering to length constraints. We calculate aligning error as follows:

$$e_{\rm A}^n = \frac{|N_{\rm pred}^n - N_{\rm target}|}{N_{\rm target}}$$
(5)

where N_{pred}^n represents the model's perceived length with counting interval n, i.e., the length the model assumes it has generated. We calculate and show the e_A^n in Figure 4.

Finding 7. Smaller counting intervals introduce greater aligning error. By closely analyzing cases, we find that frequent explicit counting interferes with semantic modeling, causing early termination of generation and poor alignment. In contrast, larger length intervals approximate implicit counting, preserving a more natural generation process.

2.5 LCTG Error

LCTG error refers to the discrepancy between the actual length of generated text and the target length:

$$E = \frac{|N_{\text{true}} - N_{\text{target}}|}{N_{\text{target}}}$$
(6)



Figure 4: Aligning Error across varied length intervals.

238

239

240

241

242

243

245

247

248

250

251

252

254

255

256

257

258

259

260

261

263

As established above, this error is systematically composed of four components: **Identifying Error** (§2.1), **Counting Error** (§2.2), **Planning Error** (§2.3), and **Aligning Error** (§2.4). To investigate the key factors influencing LLMs' LCTG error, we calculate their absolute contributions \dot{e}_i^n under different length interval n as follows:

$$\dot{e}_i^n = \frac{e_i^n}{e_{\mathrm{I}} + e_{\mathrm{C}}^n + e_{\mathrm{P}} + e_{\mathrm{A}}^n} \times E^n, \ i \in [\mathrm{I,C,P,A}] \ (7)$$

The results are shown in Figure 3. Further details can be found in B.

Finding 8. LCTG error is primarily attributed to fundamental deficiencies in length modeling, following the order of Counting Error > Perception Error > Aligning Error > Planning Error. Thus, as counting interval increases, the accumulation of counting errors leads to a corresponding rise in LCTG error.

3 Methodology

Based on the analyses and findings above, we propose MARKERGEN, a simple-yet-effective plugand-play method to help LLMs attain better LCTG performance, as shown in Figure 5. This method consists of two key modules: (1) **Auxiliary Marker Insertion Decoding** mechanism, which explicitly enhances length modeling during generation; (2) **Three-Stage Decoupled Generation**

4



Figure 5: Overview of MARKERGEN.

scheme, which decouples length constraints from semantic content generation to further improve LCTG performance.

264

265

268

269

273

274

276

277

278

281

3.1 Auxiliary Marker Insertion Decoding

External Tool Invocation. Our analysis in §2 reveals that LLMs exhibit significant identifying and counting errors, which directly contribute to inaccuracies in length modeling. To mitigate these fundamental deficiencies, we introduce external tokenizer and counter for unit recognition and counting, respectively. As Finding 1 indicates that LLMs perceive words better than tokens, we select words as the length unit.

Length Information Injection. With precise length information, we consider feeding it into the model for length modeling. Since Finding 3 indicates that LLMs' inherent implicit length modeling leads to significant errors and is inconvenient for incorporating external length information, we actively insert precise length markers during generation to enable explicit length modeling:

Len(x) = Counter(Tokenizer(x))

$$x_{t+1} = \begin{cases} \text{Marker}(\text{Len}(x_{\leq t})), & \text{if } \mathcal{S}(\text{Len}(x_{\leq t}), N) \\ \text{Sampling}(P(x_{t+1}|x_{\leq t})), & \text{else} \end{cases}$$
(8)

where $P(x_{t+1}|x_{\leq t})$ is the LLM's probability distribution for next token, Marker defines the marker format (e.g., [20 words], we discuss the effects of varied marker formats in Appendix C.1), S is the strategy that determines whether to insert a marker based on current length Len(x) and target length N. By treating the inserted markers as anchors, LLMs can continuously adjust the expected length of content to be generated during the generation process, thereby reducing the final LCTG error.

Decaying Interval Marker Insertion Strategy. The most naive insertion strategy involves placing markers at uniform intervals, which we denote as S_{uni} . However, according to Findings 4 and 7, a smaller insertion interval n improves length modeling but compromises semantic modeling, whereas a larger n exhibits the opposite effect. Considering this, we propose a strategy S_{dec} , where n decays exponentially during the generation process:

$$\mathcal{S}_{dec}(x,N) = \begin{cases} \text{True}, & \text{if } x \in \{N - \text{int}(2^{-i} \times N)\}_{i \in \mathbb{N}} \\ \text{False}, & \text{else} \end{cases}$$

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

Taking N = 200 as an example, the maker will be inserted behind the 100th, 150th, 175th, ... words. At the early stage of generation, the model primarily focuses on semantic modeling. As the generation progresses, it increasingly emphasizes length control, ultimately leading to a smaller LCTG error. Consequently, S_{dec} effectively balances semantic modeling and length modeling.

3.2 Three-Stage Decoupled Generation

Finding 7 validates that aligning error primarily arises from the inferior semantic modeling, which

317causes premature termination of the generation pro-
cess. While the planning before generation scheme
alleviates interference in semantic modeling by de-
coupling the planning process (Finding 6), it still
entangles length modeling with semantic model-
ing. To mitigate this, we introduce a three-stage
decoupled generation scheme to further reduce the
alignment error and improve the text quality, as
illustrated in Figure 5.

Stage One: Planning. The model generates a reasonable plan based on the input query and length constraints, including the content of each section and the word allocation.

Stage Two: Ensuring Semantic Integrity. The model focuses on semantic modeling to generate a high-quality response per the plan without being strictly required to adhere to length constraints.

Stage Three: Aligning Length Constraints. 334 Responses generated in stage two are usually of high quality but may not meet length restrictions. To refine them, we use these non-compliant responses 337 as input and apply the Auxiliary Marker Insertion Decoding mechanism for rewriting. The rewrit-339 340 ing requirements include: (1) Retaining the highquality semantic modeling of the input content. (2) Strictly adhering to the specified length constraints. 342 In terms of **workflow**, the model is required to: (1) Firstly analyze the previous stage's response for 345 potential improvements; (2) If its output does not meet the length constraints, it will be regenerated up to T times or until the constraints are met. 347

See Appendix 2 for prompts of each stage.

4 Experiments

We conduct comprehensive experiments to examine MARKERGEN. Specifically, we validate its effectiveness in §4.2, analyze its generalizability in §4.3, explore the impact of its key components in §4.4, and provide further insights into its mechanism in §4.5. Hyperparameter choices and additional analyses are provided in Appendix E.

4.1 Experimental Settings

Benchmarks We choose five benchmarks for experiments, where HelloBench includes two subsets, as shown in Table 2. See details in Appendix D.

361 Baselines

353

354

327

330

331

332

Benchmarks	Ability Tested	Length (words)
CNN/DailyMail (Nallapati et al., 2016)	Summarization	18-165
HANNA (Chhun et al., 2022)	Story Generation	139-995
TruthfulQA (Lin et al., 2021)	Question Answering	101-294
HelloBench (Que et al., 2024)	Heuristic LCTG& Open-ended QA	489-1450
GAOKAO (Zhang et al., 2023b)	History Open-ended QA	71-901

Table 2: Benchmarks Introduction.

• **Ruler** (Li et al., 2024b): A training-based² method that defines length control templates to regulate generation at the range level.

362

363

364

365

366

367

369

370

371

372

373

374

375

376

378

379

380

381

383

384

385

388

389

390

391

392

393

394

395

396

397

• **Implicit** (Bai et al., 2024): Conduct a planand-generate process without explicit counting. To ensure a fair comparison, the model generates multiple responses until token count outperforms MARKERGEN and the candidate with the smallest LCTG error is selected.

We conduct extensive experiments using Details Qwen2.5 series (Qwen2.5-7B/14B/32B-Instruct) (Yang et al., 2024) and the Llama3.1 series (Llama-3.1-8B/70B-Instruct) (Dubey et al., 2024), with sampling temperature as 0.5. We experiment under coarse-grained length constraints on the Openended QA subset of HelloBench and assess the LCTG error rate under precise length constraints on other benchmarks, following Eq. (6). To evaluate the text quality, we use GPT-40 mini (Hurst et al., 2024) as the judge, with a calibration algorithm to mitigate the length bias (Zheng et al., 2023) (See details in Appendix E). For precise constraints, we set the length of ground truth response as desired target length. We run each setting for three times and report the average results.

4.2 Main Results

As shown in Table 3, the commonly used two-stage implicit counting baseline results in a substantial LCTG error rate E of 18.32% on average, even if the best response is chosen across multiple attempts. This intuitively demonstrates the impact of the inherent limitations of LLM's LCTG subcapability. The training-based baseline Ruler, as observed in our preliminary experiments (Appendix E.1), benefits from training on test sets that matches the training domain, while performs poorly on our

²Ruler is the only training-based baseline for which we can find that releases the code and training set.

	Qwen2.5 Series						Llama3.1 Series					
Benchmarks	Methods	7	В	14	В	32	B	8	В	70)B	Costs
		$\overline{E}(\downarrow)$	$S\left(\uparrow ight)$	$E(\downarrow)$	$S\left(\uparrow ight)$	$\overline{E}(\downarrow)$	$S\left(\uparrow ight)$	$\overline{E}(\downarrow)$	$S\left(\uparrow ight)$	$E(\downarrow)$	$S\left(\uparrow ight)$	
CNN/DoilyMoil	Implicit	30.31	3.04	12.54	3.15	11.05	3.21	15.12	3.04	11.07	3.09	$1.30\times\delta$
CININ/Dallylviall	MarkerGen	9.92	3.07	6.06	3.16	4.82	3.25	3.36	3.18	3.18	3.36	δ
	Implicit	28.55	3.47	14.86	3.55	12.03	3.67	16.68	3.54	10.44	3.61	$2.37 \times \delta$
HANNA	MarkerGen	8.49	3.50	5.22	3.55	3.57	3.72	2.98	3.60	2.58	3.63	δ
	Implicit	16.7	4.29	17.9	4.44	8.7	4.45	7.21	4.22	7.64	4.46	$1.75 \times \delta$
ITuullulQA	MarkerGen	9.08	4.33	7.59	4.43	4.48	4.54	3.82	4.25	2.80	4.48	δ
	Implicit	35.69	3.42	21.34	3.80	12.02	3.80	21.91	3.72	27.89	3.74	$1.06 \times \delta$
	MARKERGEN	8.51	4.13	6.35	4.00	5.34	4.14	6.03	4.03	5.03	3.98	δ

Table 3: Overall Performance of MARKERGEN on Various Benchmarks. E denotes LCTG error rate (%) and S denotes the text quality ([1, 5]) given by LLM judge. δ denotes the token cost of MARKERGEN under each setting.

		Target Length Scales								
Model	Methods	100		100 200		300		400		Costs
		$E(\downarrow)$	$S\left(\uparrow ight)$	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$	$E(\downarrow)$	$S\left(\uparrow ight)$	
Qwen2.5-7B-Instruct	Implicit MarkerGen	30.97 8.26	3.45 3.92	22.91 9.06	3.53 4.00	26.12 7.67	3.28 3.75	29.63 5.10	3.08 3.55	$\begin{array}{c} 1.26\times\delta\\ \delta\end{array}$
				Leng	gth Cons	straint Ty	pes			
Models	Methods	<100		<100 100-150		160-200		>500		Costs
		$\overline{E_r}\left(\downarrow\right)$	$S\left(\uparrow ight)$	$\overline{E_r}\left(\downarrow\right)$	$S\left(\uparrow ight)$	$\overline{E_r}\left(\downarrow\right)$	$S\left(\uparrow ight)$	$\overline{E_r} (\downarrow)$	$S\left(\uparrow ight)$	
Qwen2.5-7B-Instruct	Implicit MarkerGen	7.50 0.00	3.47 3.94	63.00 0.50	4.03 4.50	66.00 3.00	4.06 4.53	29.50 0.00	2.65 3.13	$\begin{array}{c} 1.07\times\delta\\\delta\end{array}$

Table 4: Experiments with varied length scales and constraint types on Open-ended QA subset of HelloBench.

evaluated benchmarks, highlighting its limited generalizability. In comparison, under strict length constraints, MARKERGEN achieves an absolute reduction of 12.57% in *E* relative to the implicit baseline, bringing the final error down to just 5.57%. In terms of text quality, by decoupling length modeling and semantic modeling during the generation process and employing the decaying insertion strategy to minimize the damage caused by length constraints to semantic integrity, MARKERGEN achieves a higher *S* in average. Meanwhile, this performance is achieved with only 64% of the tokens used by the baseline.

4.3 Generalizability

Across LLMs and Tasks. Table 3 demonstrates the strong generalizability of MARKERGEN to LLMs and generation tasks.

415Across Length Scale.Table 3 also shows MARK-416ERGEN's strong performance across benchmarks417with varying length scale (18–1450).To further418investigate, we analyze progressively increasing419the target length from 100 to 400.The results in

Table 4 show a declining trend in MARKERGEN's error rate, which can be attributed to the auxiliary marker insertion decoding mechanism that prevents error accumulation from implicit modeling.

Across Constraint Types. In addition to exact length constraints, users may impose range-based limits. We evaluate E_r , the proportion of responses violating these constraints. Table 4 shows that MARKERGEN maintains an E_r below 3% in all cases, significantly lower than the baseline.

Across Lingual. We further validate the effectiveness of MARKERGEN in Chinese setting on GAOKAO benchmark, as shown in Table 8.

4.4 Ablation Studies

In this section, we validate the effectiveness of each module in MARKERGEN with Qwen2.5-32B-Instruct on TruthfulQA, as shown in Table 5.

Tool Invocation.When the model is required to437insert markers independently without relying on
an external tokenizer and counter, its fundamental438limitations lead to a significant increase in the error440



Figure 6: Attention matrices of the first (left) and last (right) layers.

T 7					Mark	er Insert	ion Inter	val n				
Variants	1		4		16		32		64		Decaying	
	$\overline{E}(\downarrow)$	$S\left(\uparrow ight)$	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$	E	$S\left(\uparrow ight)$	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$	$E\left(\downarrow ight)$	$S(\uparrow)$
w/o Tool Two Stage Three Stage	15.53 3.10 4.84	4.28 4.03 4.28	32.50 1.49 4.20	4.29 4.03 4.29	34.64 4.04 4.89	4.46 4.20 4.45	32.50 3.26 5.45	4.48 4.23 4.48	20.44 3.93 5.18	4.58 4.32 4.57	 2.66 4.48	4.28 4.54

Table 5: Ablation studies on key components.

rate, exceeding 15%.

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

Decaying Interval Marker Insertion. When using a fixed marker insertion interval n, since length control is inversely proportional to n, while semantic modeling is directly proportional to n (which induces alignment errors), we observe unstable LCTG error rate. In contrast, by adopting a sparse-to-dense insertion approach, the Decaying Interval Marker Insertion strategy ensures explicit length modeling while maximizing semantic integrity, leading to lower E and superior S.

Three-Stage Decoupled Generation. The introduction of explicit length markers in the two-stage scheme leads to a substantial reduction in LCTG error relative to the implicit baseline $(8.7 \rightarrow 2.66)$. However, this scheme places greater emphasis on length modeling, which consequently diminishes text quality $(4.45 \rightarrow 4.28)$. In comparison, the three-stage scheme achieves a better balance by decoupling semantic and length modeling, thereby improving both length control and text quality.

4.5 Working Mechanism of MARKERGEN

To better understand how LLMs leverage the inserted length markers in MARKERGEN, we visualize the attention matrices of the first and last layers of Llama-3.1-8b-Instruct (Figure 6). In the shallow layers, the attention distribution reveals a clear focus on the length information represented by the length markers (in the red box). As the model progresses to the deeper layers, attention shifts from the length information to the adjacent semantic content (in the orange box). This pattern demonstrates that at shallow layers, the model uses markers to establish length modeling and encode precise length information. At deeper layers, it relies on this length information for semantic modeling, producing tokens that align with the length constraints while maintaining semantic integrity. 468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

Conclusions

To improve the performance of LLMs in lengthcontrollable text generation, we conduct a bottomup error analysis of relevant sub-abilities. The results reveal that deficiencies in identifying, counting, and aligning are key limitations. To fill this gap, we propose MARKERGEN, which leverages external tools to compensate for fundamental deficiencies. Additionally, it introduces Decaying Interval Marker Insertion Strategy to facilitate explicit length modeling and employs Three-Stage Decoupled Generation mechanism to balance semantic coherence and length control. Comprehensive experiments demonstrate the strong generalizability and effectiveness of MARKERGEN in enhancing length control and preserving semantic integrity.

495 Limitations

In this work, we conduct a bottom-up sub-496 capability analysis in the LCTG ability and propose 497 the MARKERGEN method, achieving strong LCTG 498 performance. One major limitation of MARKER-GEN is that it is currently only applicable to open-500 source models and cannot yet be used with closed-501 source models. To address this, we will release our code, allowing closed-source model providers in-503 terested in adapting MARKERGEN to benefit from our method in enhancing LCTG performance.

06 Ethics Statement

507All of the datasets used in this study were publicly508available, and no annotators were employed for our509data collection. We confirm that the datasets we510used did not contain any harmful content and was511consistent with their intended use (research). We512have cited the datasets and relevant works used in513this study.

References

514

515

516

517

518

519

520

521

523

524

525

526

527

528

530

531

532

533

534

535

536

537

539

541

542

- Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. Longwriter: Unleashing 10,000+ word generation from long context llms. arXiv preprint arXiv:2408.07055.
- Bradley Butcher, Michael O'Keefe, and James Titchener. 2024. Precise length control in large language models. *arXiv preprint arXiv:2412.11937*.
- Yingshan Chang and Yonatan Bisk. 2024. Language models need inductive biases to count inductively. *arXiv preprint arXiv:2405.20131*.
- Cyril Chhun, Pierre Colombo, Chloé Clavel, and Fabian M Suchanek. 2022. Of human criteria and automatic metrics: A benchmark of the evaluation of story generation. *arXiv preprint arXiv:2208.11646*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators. arXiv preprint arXiv:2404.04475.
- Kiannah Foster, Andrew Johansson, Elizabeth Williams, Daniel Petrovic, and Nicholas Kovalenko. 2024. A token-agnostic approach to controlling generated text length in large language models.

Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. 2024a. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*. 543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

- Yuxuan Gu, Wenjie Wang, Xiaocheng Feng, Weihong Zhong, Kun Zhu, Lei Huang, Tat-Seng Chua, and Bing Qin. 2024b. Length controlled generation for black-box llms. *arXiv preprint arXiv:2412.14656*.
- Chenyang Huang, Hao Zhou, Cameron Jen, Kangjie Zheng, Osmar R ZaÃane, and Lili Mou. 2025. A decoding algorithm for length-control summarization based on directed acyclic transformers. *arXiv* preprint arXiv:2502.04535.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-40 system card. *arXiv preprint arXiv:2410.21276*.
- Renlong Jie, Xiaojun Meng, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Prompt-based length controlled generation with reinforcement learning. *CoRR*, abs/2308.12030.
- Juseon-Do Juseon-Do, Hidetaka Kamigaito, Manabu Okumura, and Jingun Kwon. 2024. Instructcmp: Length control in sentence compression through instruction-based large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 8980–8996.
- Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy. 2024. The impact of positional encoding on length generalization in transformers. *Advances in Neural Information Processing Systems*, 36.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, et al. 2024a. From generation to judgment: Opportunities and challenges of llm-as-a-judge. *arXiv preprint arXiv:2411.16594*.
- Jiaming Li, Lei Zhang, Yunshui Li, Ziqiang Liu, Yuelin Bai, Run Luo, Longze Chen, and Min Yang. 2024b. Ruler: A model-agnostic method to control generated length for large language models. In *Findings of the Association for Computational Linguistics: EMNLP* 2024, Miami, Florida, USA, November 12-16, 2024, pages 3042–3059. Association for Computational Linguistics.
- Xun Liang, Hanyu Wang, Yezhaohui Wang, Shichao Song, Jiawei Yang, Simin Niu, Jie Hu, Dan Liu, Shunyu Yao, Feiyu Xiong, and Zhiyu Li. 2024. Controllable text generation for large language models: A survey. *CoRR*, abs/2408.12599.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *CoRR*, abs/2109.07958.

- 598 599 601 606 607 610 611 612 613 615 616 617 618 619
- 625 626 627
- 631
- 634

- 647 648

- Sangjun Moon, Jingun Kwon, Hidetaka Kamigaito, Manabu Okumura, et al. Length representations in large language models.
- Ramesh Nallapati, Bing Xiang, and Bowen Zhou. 2016. Sequence-to-sequence rnns for text summarization. CoRR, abs/1602.06023.
- OpenAI. 2023. GPT-4 technical report. CoRR. abs/2303.08774.
- Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. 2024. Disentangling length from quality in direct preference optimization. In Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, pages 4998-5017. Association for Computational Linguistics.
 - Haoran Que, Feiyu Duan, Liqun He, Yutao Mou, Wangchunshu Zhou, Jiaheng Liu, Wenge Rong, Zekun Moore Wang, Jian Yang, Ge Zhang, et al. 2024. Hellobench: Evaluating long text generation capabilities of large language models. arXiv preprint arXiv:2409.16191.
 - Fabian Retkowski and Alexander Waibel. 2024. Zeroshot strategies for length-controllable summarization. arXiv preprint arXiv:2501.00233.
 - Seoha Song, Junhyun Lee, and Hyeonmok Ko. 2024. Hansel: Output length controlling framework for large language models. arXiv preprint arXiv:2412.14033.
 - Zekun Wang, Feiyu Duan, Yibo Zhang, Wangchunshu Zhou, Ke Xu, Wenhao Huang, and Jie Fu. 2024. Positionid: Llms can control lengths, copy and paste with explicit positional awareness. arXiv preprint arXiv:2410.07035.
 - An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. arXiv preprint arXiv:2412.15115.
 - Peiwen Yuan, Shaoxiong Feng, Yiwei Li, Xinglin Wang, Yueqi Zhang, Jiayi Shi, Chuyi Tan, Boyuan Pan, Yao Hu, and Kan Li. 2025. Llm-powered benchmark factory: Reliable, generic, and efficient. arXiv preprint arXiv:2502.01683.
 - Weizhe Yuan, Ilia Kulikov, Ping Yu, Kyunghyun Cho, Sainbayar Sukhbaatar, Jason Weston, and Jing Xu. 2024. Following length constraints in instructions. CoRR, abs/2406.17744.
 - Hanging Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023a. A survey of controllable text generation using transformer-based pre-trained language models. ACM Computing Surveys, 56(3):1-37.
 - Xiaotian Zhang, Chunyang Li, Yi Zong, Zhengyu Ying, Liang He, and Xipeng Qiu. 2023b. Evaluating the performance of large language models on gaokao benchmark. arXiv preprint arXiv:2305.12474.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan 653 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, 654 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, 655 Joseph E. Gonzalez, and Ion Stoica. 2023. Judging 656 llm-as-a-judge with mt-bench and chatbot arena. In 657 Advances in Neural Information Processing Systems 658 36: Annual Conference on Neural Information Pro-659 cessing Systems 2023, NeurIPS 2023, New Orleans, 660 LA, USA, December 10 - 16, 2023. 661 670

671

672

673

674

675

676

677

686

692

A Related Work

LCTG Methods Text length is a fundamental aspect of natural language that carries semantic information, making LCTG a task of balancing length and semantic constraints. Achieving precise length control remains a challenge for LLMs due to limitations in their architecture, such as position encoding (Butcher et al., 2024; Kazemnejad et al., 2024; Chang and Bisk, 2024) and decoding mechanisms(Huang et al., 2025). Consequently, existing methods focus on injecting length information to help LLMs model length accurately, which can be categorized into training-based and inferencebased approaches.

Training-based methods inject varying levels of length signals during fine-tuning or reinforcement learning. For instance, Jie et al. (2023); Li et al. (2024b) use prompt templates to teach LLMs the mapping between length and textual content, while Song et al. (2024); Wang et al. (2024) design fine-grained datasets to guide correct length modeling. Other methods, like Yuan et al. (2024); Jie et al. (2023), utilize reward functions to align length preferences during training. While effective in certain scenarios, these methods suffer from limited generalization across diverse LCTG tasks, including varying length constraints and instructions.Inference-based methods adjust inputs multiple times during generation to inject, such as through prompt-based Automated Revisions and Sample Filtering (Retkowski and Waibel, 2024; Juseon-Do et al., 2024), or length-controlled importance sampling during decoding (Gu et al., 2024b). Although these approaches can better generalize length alignment, they still struggle with achieving precise control.

698While both approaches enhance LCTG, they of-699ten apply a top-down strategy that lacks deep un-700derstanding and targeted enhancement of LCTG701sub-capabilities. This limits progress in meeting702length constraints accurately. Furthermore, many703methods neglect semantic constraints, and injecting704length information may degrade text quality. There-705fore, we propose MARKERGEN to bridge this gap706for precise length control and preserving semantic707integrity.

B Detailed Sub-ability Error analyses in LCTG

B.1 Identifying Error

Identifying error refers to the misidentification of fundamental length units. To systematically analyze this error, we design a counting experiment in which the model is prompted to sequentially recognize and accumulate length units, then compare its predicted count with the ground truth. Experimental results confirm that in the one-by-one accumulation setting, counting errors do not occur, meaning that the final length error entirely arises from identifying error (as shown in Figure 7).

B.2 Counting Error

Counting error refers to the inaccurate enumeration of units in a given sequence, leading to deviations from the intended length. Therefore, in the setting where n > 1 in the counting experiment, the final counting result error is caused by both identifying error and counting error. In this case, counting error can be decoupled by resolving identifying errors in the accumulation process, where errors result from the accumulation step. We also conducted the same counting experiment as in Section 2.2 on the CNN/DailyMail summarization dataset, as shown in Figures 8.

From the figure, we can further validate the same conclusions as in Findings 1, 2, 3, and 4 in Section 2, revealing that the length errors in the generated results of the LCTG task stem from significant errors in the LLM's perception and modeling of length.

B.3 Planning Error

Planning error refers to the misallocation of word counts across different sections, leading to a discrepancy from target length. The planning ability of LLMs encompasses not only length planning but also semantic planning. To effectively assess the quality of LLM's semantic planning, we use Qwen-Plus (Yang et al., 2024) as the judge, with a scoring range of [1, 5]. The specific evaluation prompt is as follows:

You are tasked with evaluating the quality of a generated answer plan for a TruthfulQA question. The evaluation should focus on the truthfulness, logical coherence, and adherence to the given prompt and instructions. Rate the answer plan on a 5-point scale as follows: 709 710

711

712

713

714

715

716

717

708

718 719 720

721 722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

738

739

740

741

742

743

744

745

746

747

750

751

752

753

754

755

756



Figure 7: Schematic diagram of counting experiment under the condition of n = 1

- **5: Outstanding** The plan is highly truthful, logically coherent, and perfectly adheres to the prompt and instructions.
- 4: Very Good The plan is mostly truthful and coherent, with only minor issues in details or adherence to instructions.
- **3: Good** The plan is acceptable but has noticeable shortcomings in truthfulness or coherence.
- 2: Poor The plan has significant issues in truthfulness or logical coherence and does not adequately follow the instructions.
- 1: Unacceptable The plan is largely untruthful, incoherent, or fails to follow the prompt instructions entirely.
- Please provide the overall score in the following format: ###score X
- 778 Question:
 779 + prompt
 780 Generated Answer Plan:
 781 + generated_plan
 782 Evaluate the answer plan based on the above criteria.

758

759

771

772

774

784

Since the LCTG task requires meeting both length and semantic constraints, utilizing the LLM's superior planning ability for explicit planning before generation, as opposed to direct generation, helps to clearly define the modeling space for length and the semantic extension range. This not only contributes to improved text generation quality but also reduces length errors. 789

790

791

792

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

B.4 Aligning Error

Aligning error refers to the discrepancy between the model's perceived length and the target length, arising from the challenge of maintaining semantic integrity while adhering to length constraints.As shown in Figure 4, aside from Finding 7, we observe significant differences in aligning error across models. Qwen2.5-32B-Instruct and GPT-40 mini exhibit larger alignment errors under finegrained length modeling. As discussed in Section 2, "length estimation acts as a real-time constraint, dynamically regulating further extension. Ultimately, the model strives to align the length constraints while preserving semantic integrity." Highfrequency length perception updates pose greater challenges for the natural expansion of the semantic space, which explains why some models with weaker robustness in semantic expansion show significant alignment errors. These errors become a primary source of LCTG inaccuracies (as shown in Figure 9). This further emphasizes that LCTG is a task of balancing length and semantic constraints.

B.5 LCTG Error

Based on the above decomposition of sub-abilities in LCTG and the corresponding error analysis, we can clearly quantify the contribution of each decoupled error to the final LCTG error. As shown in Figure 1, the quantification results in the right figure represent the average values of four models under various n conditions. The conclusion we can draw is that the primary cause of significant length



Figure 8: Error analyses of fundamental abilities in LCTG on CNN/DailyMail.



Figure 9: Absolute contribution of LCTG sub-capability deficiency of GPT 40-mini.

errors in current mainstream LLMs on LCTG tasks is the lack of bottom-up identification and counting capabilities required for accurate length modeling.

C Exploration of Interval Marker Insertion Strategy Variants

C.1 Length Marker Forms

823

824

825

827

828

830

832

835

837

842

843

844

We explored the impact of different forms of length marker insertion on performance, such as the number of words generated "[k]", the semantic marker "[k words]", and the remaining words to be generated "[$N_{target} - k$]" (remaining words). As shown in Table 6, we used Llama-3.1-8B-Instruct on the CNN/DailyMail dataset to investigate the effects of various marker forms under multiple *n* conditions on generation error and text quality. The results show that using a semantic length marker representing the number of words generated achieved the best performance in both length error and text quality.

D Detailed Benchmarks Introduction

The benchmarks used in our experiments are as follows:

Marker Form	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$
[k]	18.28	3.10
[k words]	15.74	3.14
$[N_{\mathrm{target}} - k]$	27.92	3.09

 Table 6: Comparison of Length Marker Forms and Their

 Performance

• **CNN/DailyMail**(Nallapati et al., 2016): A summarization dataset of news articles, with 500 randomly sampled items. (*18–165 words*)

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

- HANNA(Chhun et al., 2022): A long-form story generation dataset with 200 selected items. (*139–995 words*)
- **TruthfulQA**(Lin et al., 2021): A benchmark for factual accuracy in open-domain QA. (*101–294 words*)
- **HelloBench**(Que et al., 2024): A long-text generation benchmark. We selected subsets from *heuristic text generation* (e.g., argumentative and roleplaying writing, covering five types) and *open-ended QA* (spanning ten domains). (489–1450 words)
- GAOKAO-Bench(Zhang et al., 2023b): A benchmark collected from the Chinese college entrance examination (GAOKAO). We selected the 2010-2022 History Open-ended Questions subset. (71–901 words)

E Detailed Experimental Results

E.1 Performance and Generalization Study of Training-based Methods

To investigate the performance and generalization of training-based methods in diverse, real-world LCTG task scenarios, we selected Ruler, a trainingbased method that defines length control templates

TLG dataset									
benchmark	Method	Model	$PM\left(\uparrow ight)$	$FM\left(\uparrow ight)$					
TLG dataset	before training	Llama-3.1-8B-Instruct	5.55	10.20					
	RULER	Llama-3.1-8B-Instruct-ruler	41.75	55.10					
Precise Length Constraint Benchmark									
Benchmarks	Methods	Models	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$					
CNN/DailyMail	Ruler	Llama-3.1-8B-Instruct-ruler	78.21	3.10					
	MarkerGen	Llama-3.1-8B-Instruct	3.36	3.18					
HANNA	Ruler	Llama-3.1-8B-Instruct-ruler	68.21	2.87					
	MarkerGen	Llama-3.1-8B-Instruct	2.98	3.60					
TruthfulQA	Ruler	Llama-3.1-8B-Instruct-ruler	44.93	3.27					
	MarkerGen	Llama-3.1-8B-Instruct	3.82	4.25					
Heuristic Generation	Ruler	Llama-3.1-8B-Instruct-ruler	66.17	2.94					
	MarkerGen	Llama-3.1-8B-Instruct	6.03	4.03					

Table 7: Combined Benchmark Evaluation Table

Methods	Models	$E\left(\downarrow\right)$	$S\left(\uparrow ight)$
Implicit	Qwen2.5-14B-Instruct	27.41	3.47
MARKERGEN	Qwen2.5-14B-Instruct	7.71	3.55

Table 8: GAOKAO-History Chinese Dataset Results

to regulate generation at the range level. This choice is based on the fact that Ruler is the only training-based baseline for which the code and training set are publicly available. We followed the exact setup provided in the repository and verified the correctness of our replication by achieving significant performance improvements on the given test set, as shown in Table 7.

872

874

876

878

882

887

891

892

893

894

895

Next, we tested the trained model, referred to as Llama-3.1-8B-Instruct-ruler, across four selected benchmarks with varying tasks, length scales, and instructions, under cost-alignment conditions. The experimental results revealed substantial errors and a decline in text quality, even when compared to the implicit method's results without training (as shown in Table 3). This finding demonstrates the limited generalization capability of the method, highlighting its struggle to cope with the complexity and diversity of real-world LCTG scenarios.

E.2 Length Bias Correction in LLMs-as-a-Judge

It has been demonstrated that LLMs-as-a-judge exhibit a noticeable length bias (Li et al., 2024a; Gu et al., 2024a). To evaluate the quality of generated text objectively and accurately for LCTG tasks, it is essential to correct for this length bias. We adopt the length-controlled evaluation method outlined in AlpacaEval (Dubois et al., 2024) and Yuan et al. (2025).

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

To derive unbiased judge scores, we use a Multiple Regression model. Specifically, we set the judge score as the dependent variable, with the generator categories as dummy variables, and the length of the generated text as a covariate. The model is formulated as follows:

$f(i) = \beta_0 + \beta_M \cdot C(\text{Method}) + \beta_m \cdot C(\text{Model}) + \beta_l \cdot \text{Length} + \beta_m \cdot C(\text{Model}) + \beta_l \cdot \text{Length} + \beta_m \cdot C(\text{Model}) + \beta_l \cdot \text{Length} + \beta_l \cdot \beta_m \cdot C(\text{Model}) $	$+\epsilon$
(1	0)

where f(i) denotes the judge score for the generated text G_i , C(Method) and C(model) are categorical variables representing the method and the model used, respectively, and Length is the actual length of the generated text. The coefficients β_M , β_m , and β_l represent the adjustments made to the raw judge score f(i), helping to eliminate length bias. These adjusted scores serve as the metrics for faithfulness and alignment.

E.3 Residual Length Error Analysis in MARKERGEN

This subsection focuses on analyzing the residual length errors in the MARKERGEN framework. Building upon the sub-decomposition of LCTG errors presented in Section 2, we eliminate identifying and counting errors through Auxiliary Length Marker Insertion Decoding 3.1. Moreover, by employing the Three-Stage Decoupled Generation strategy 3.2, we effectively reduce aligning errors, thus improving the robustness of all models in semantic expansion under precise length modeling with explicit length markers. This approach ensures



Figure 10: Attention Entropy across layers.

semantic integrity while enhancing text generation quality through a clearer, more in-depth analysis of LLM's LCTG sub-capabilities. Ultimately, resid-932 ual LCTG errors are primarily driven by minimal aligning errors.

930

933

935

937

939

942

943

946

947

951

952

956

958

962

Cross-layer Attention Analysis from the E.4 MARKERGEN Perspective

In this section, we perform a cross-layer attention analysis from the MARKERGEN perspective. By examining attention patterns across different layers of the model, we aim to gain a better understanding of how length and semantic information are processed at various stages of generation, providing insights into improving the accuracy of LCTG tasks.

Combining the analyses from Figures 6 and 10, we infer that in the shallow layers, attention is primarily focused on the length information represented by the length markers. This suggests that the model's early stages prioritize processing and understanding the input length. The higher entropy in these layers indicates that the model needs to integrate various details and information to effectively comprehend the input. As the model progresses to deeper layers, attention shifts from the length information to the adjacent semantic content. The lower entropy in these layers indicates that the model refines its focus, extracting key features and generating more relevant output.

This pattern of attention distribution aligns with the findings from (Moon et al.), which emphasize that length modeling in the early layers serves as a foundation for semantic processing in the later layers. Our analysis further supports the notion 963 that LCTG tasks depend on a dynamic interaction 964 between length control and semantic generation, 965 where early layers focus on length constraints and 966 deeper layers prioritize semantic coherence. 967