Anonymous Author(s)*

ABSTRACT

Recent research on query generation has focused on using Large Language Models (LLMs), which despite bringing state-of-the-art performance, also introduce issues with hallucinations in the generated queries. In this work, we introduce relevance hallucination and factuality hallucination as a new typology for hallucination problems brought by query generation based on LLMs. We propose an effective way to separate content from form in LLMgenerated queries, which preserves the factual knowledge extracted and integrated from the inputs and compiles the syntactic structure, including function words, using the powerful linguistic capabilities of the LLM. Specifically, we introduce a model-agnostic and training-free method that turns the Large Language Model into a Pointer-Generator (LargePiG), where the pointer attention distribution leverages the LLM's inherent attention weights, and the copy probability is derived from the difference between the vocabulary distribution of the model's high layers and the last layer. To validate the effectiveness of LargePiG, we constructed two datasets for assessing the hallucination problems in query generation, covering both document and video scenarios. Empirical studies on various LLMs demonstrated the superiority of LargePiG on both datasets. Additional experiments also verified that LargePiG could reduce hallucination in large vision language models and improve the accuracy of document-based question-answering and factuality evaluation tasks. The source code and dataset are available at https://anonymous.4open.science/r/LargePiG-7674.

KEYWORDS

Query Generation, Hallucination, Pointer Generator

ACM Reference Format:

1 INTRODUCTION

Query generation is an automatic process of generating queries according to the content presented in documents or videos, which not only facilitates information retrieval from documents [12, 35, 48] but also serves applications like short video platforms by creating queries that attract user engagements. There has been notable advancement in query generation using LLMs [5, 12, 36, 39]. However, employing LLMs for query generation often introduces hallucination issues. **Factuality hallucination** refers to inaccuracies in the facts presented in the generated queries, often occurring when the inputs include knowledge not covered by the LLM's pre-training data. For example, being misled by the latest facts in the news documents can make LLMs generate queries that conflict with actual events. **Relevance hallucination** occurs when the generated queries, although factually correct, are irrelevant to the inputs [15]. Both types of hallucinations are not mutually exclusive, with some generated queries exhibiting both issues (see appendix A.1 for the experimental validation of hallucination classification). 59

60

61 62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

Previous research has primarily focused on reducing relevance hallucinations through post-processing methods [5, 12, 15], without addressing hallucinations at the source of generation. With the expanding range of applications for query generation on shortvideo platforms, generating "related search" based on video content to attract user clicks and enhance user engagement has become crucial for these platforms. Figure 1 presents some examples of "related search" on short-video platforms, each of which has hundreds of millions of users ¹. If a generated query exhibits relevance hallucinations, users may not click the query as clicking on "related search" will not find content related to the video, diminishing user interest. Conversely, if a query demonstrates factuality hallucinations (without relevance hallucinations), it might initially attract users' interest through clickbait but fail to deliver content related to the hallucinatory facts, thereby degrading the user experience. Therefore, the queries we generate need to be relevant to the video content, factually accurate, and sufficiently novel to attract user clicks and improve user engagement.

Unlike other generation tasks, query generation primarily relies on the inputs. Thus, decoupling the content and form at the output end of LLMs, ensuring that the factual content of the generated queries mainly comes from the inputs and that the syntax and other forms are organized by LLMs, is key to keeping the generated query truthful and reducing hallucination issues. To this end, we propose to use the Pointer Generator (PG) technology, a sequenceto-sequence model that integrates extraction (pointing to words in the input) and generation (creating new words) strategies to enhance the accuracy and relevance of the generated text [42, 46]. The PG model, combines pointer attention distribution (determining the model's focus on different parts of the inputs), vocabulary distribution (the probability distribution for choosing the next word from a fixed vocabulary), and copy probability (deciding whether to generate a word from the vocabulary distribution or copy directly from the input), not only increases the probability of mentioning facts presented in the inputs and decreases the likelihood of generating unrelated facts but also ensures the correctness of syntax and other forms generated by LLMs. Although PG technology has

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

⁵ Conference'17, July 2017, Washington, DC, USA

^{56 © 2024} Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-x-xxxx-xXYY/MM

⁵⁷ https://doi.org/10.1145/nnnnnnnnnn
58

¹TikTok: www.tiktok.com; Kwai: www.kuaishou.com; Xiaohongshu: www. xiaohongshu.com.

Conference'17, July 2017, Washington, DC, USA

Anon



(a) Relevance Hallucination: The video is from TikTok, where the "related search" at the top presents a certain relevance hallucination, as the person in the video is playing an electric piano rather than an electric guitar.



(b) Factuality Hallucination: The video is from Kwai, where the "related search" presents a certain factuality hallucination. Singapore itself is a country, so it is illogical to ask which country's nationality it belongs to.



(c) Truthful Query: The video is from Xiaohongshu, where the "related search" at the top present is relevant and factual.

Figure 1: Examples of query generation in real applications across different short video platform.

been applied in query generation tasks with traditional language models [19, 47], considering the enormous parameter size and training resource consumption of LLMs, adopting the traditional PG scheme, which requires learning pointer attention distribution and copy probability, may not only disrupt the original representations of LLMs but also diminish their generalization capability.

Facing the above challenge, we propose a novel PG implementation that can achieve PG functionality within LLMs without requiring additional training. Our method is based on two core observations: (1) Attention modules are more 'truthful' than other modules in LLMs (e.g., FFN modules), allowing the intrinsic attention weights towards the input sequence within LLMs to serve as the PG's pointer attention distribution; (2) LLMs generate different types of words (function words and factual knowledge words) with distinct patterns [10, 41]. When generating function words, the vocabulary distribution obtained from the high layers of LLMs is relatively consistent, whereas, for factual knowledge words, the vocabulary distribution from the high layers of LLMs shows significant differences. Further analyzing the internal mechanism behind the occurrence of different patterns in LLMs, we find that this pattern is rooted in the difference in the amount of information between function words and factual knowledge words in human linguistics. We relaxed the requirement for LLMs to generate the correct words, only needing them to identify the type of word to be generated and calculate the copy probability through the difference between the vocabulary distribution of the model's high layers and the last layer.

Based on this concept, we propose that Large Language Models can essentially act as an implicit Pointer-Generator (LargePiG ...), better addressing the hallucination issues in query generation. Our method has several notable advantages: Firstly, it preserves LLMs' powerful capabilities and generalizability, as it does not require significant modifications to the model architecture or additional training. Secondly, by simplifying the implementation process of PG, our method reduces additional computational and resource requirements, making it more efficient and easy to implement. Lastly, this approach retains the advantages of PG, achieving decoupling of content and form at the output end of LLMs, making the generated content faithful to the inputs.

To better assess the capability of LargePiG in solving hallucination issues within query generation, we introduce TruthfulVQG and TruthfulDQG, two challenging Truthful Query Generation benchmarks gathered from video and document scenarios, respectively. Experiments on these datasets demonstrate that LargePiG is capable of increasing the factuality and relevance of various LLM-based query generation methods across different LLMs. More experiments on the LLaVA [24] family validate the effectiveness of LargePiG in addressing hallucination issues in query generation within multimodal scenarios. Further experiments on relevance testing and factuality evaluation demonstrate that LargePiG can individually address relevance hallucination and factuality hallucination. Efficiency analysis shows that LargePiG causes negligible latency in

the query generation process, proving the practical applicability of LargePiG.

We summarize the major contributions of this paper as follows: (1) We identify the relevance and factuality hallucination issues in query generation, which are crucial for ensuring effective "related search" in short-video platforms.

(2) We propose **LargePiG**, a training-free, and model-agnostic decoding method that mitigates query generation hallucinations without modifying LLM architectures, ensuring ease of deployment.

(3) We introduce two truthful query generation benchmarks, TruthfulVQG and TruthfulDQG, and demonstrate through extensive experiments the effectiveness of LargePiG in reducing hallucinations while maintaining efficiency.

2 RELATED WORK

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

290

Large language models based query generation. Query gener-249 ation is vital for improving information retrieval systems and user 250 experience on short video platforms. Doc2Query [35] implements 251 this concept using a sequence-to-sequence model for generating 252 queries based on document contents. Advancing this, UDP [40] 253 254 utilizes LLMs in a zero-shot setting to predict query likelihood 255 from text passages. Building on this, PQGR [12] and InPars [5] introduce few-shot and contrastive example approaches, enhanc-256 ing the contextual awareness of query generation. AOG [25] fur-257 ther develops LLM adaptability to query generation by employing 258 LoRA [17] for fine-tuning with real user queries and context, along-259 side other parameter-efficient methods like soft-prompt tuning and 260 adapters [36, 37]. Additionally, UDAPDR [39] explores efficiency by 261 combining large and small models to generate and refine queries. 262 Our work addresses hallucination in query generation, introduc-263 ing LargePiG, a novel decoding method applicable to LLM-based 264 query generation approaches to reduce relevance and factuality 265 hallucination. 266

267 Hallucination mitigation in large language models. Large 268 Language Models exhibit a critical tendency to produce hallucinations, resulting in content that is inconsistent with real-world 269 facts or user inputs. Hallucination mitigation strategies can be data-270 271 driven, involving more refined filtering of pretraining data [28] or high-quality instruction-tuning datasets [52] to reduce the likelihood of LLMs learning hallucinatory knowledge. Alternatively, 273 approaches from the input side, such as Retrieval Augmented Gen-274 eration, utilize data to reduce LLM-generated hallucinations by 275 grounding the model with an external knowledge base [14]. How-276 277 ever, Retrieval Augmented Generation is not well-suited for tasks 278 like query generation, as there is no explicit need for external retrieval content. Our LargePiG method focuses on reducing halluci-279 280 nation for the query generation task from the generation side, transforming the LLM into a pointer generator by leveraging intrinsic 281 features of the LLM to separate content and form in LLM-generated 282 queries. Unlike DoLa [10], which contrasts between transformer 283 layers to correct the next word's probability, LargePiG derives the 284 copy probability from the difference between the vocabulary distri-285 bution of the model's high layers and the last layer. Moreover, these 286 hallucination mitigation methods are orthogonal to the LargePiG 287 288 approach taken in this paper and could potentially be used in conjunction to mitigate hallucinations further. 289

3 METHOD

Current Large Language Models are fundamentally based on the Transformer decoder-only architecture. Initially, the input text is tokenized and transformed into numerical vectors by the embedding layer. Given a sequence of input tokens as $X = \{x_1, x_2, \ldots, x_{t-1}\}$, where the input tokens may include the instruction $I = \{x_1, \ldots, x_{m-1}\}$, the source document $D = \{x_m, \ldots, x_n\}$, and part of generated query $\tilde{Q} = \{x_{n+1}, \ldots, x_{t-1}\}$, the embedding layer first converts these tokens into a series of vectors $H_0 = \{h_1^{(0)}, \ldots, h_{t-1}^{(0)}\}$. After passing through multiple Transformer Decoder Layers, H_N is processed by a Classification Layer, usually composed of a layer of linear layers and softmax, mapping to the vocabulary distribution.

To address the hallucination issues present in LLM-based query generation, we propose to incorporate the mechanism of the Pointer-Generator to enhance the model's faithfulness to the factual knowledge contained within the source document *D*. The Pointer-Generator combines the original decoding vocabulary distribution P_{vocab} of the LLM with the newly introduced pointer attention distribution P_{source} , the latter representing the probability distribution over the source document *D*. Furthermore, the Pointer-Generator includes a copy probability p_{copy} , which determines whether the model selects the next word from a predefined vocabulary or directly copies a word from the source document. We propose to use this mechanism to ensure that the factual content in the generated query mainly comes from *D* and that the syntax and other forms are organized by LLMs, significantly reducing the occurrence of hallucinations.

Unlike previous approaches that required retraining the pointergenerator model to learn the pointer attention distribution and copy probability, we propose **LargePiG**, a plug-in and training-free method, to implement pointer-generator decoding within LLMs (see Figure 2). The pointer attention distribution can utilize the LLM's intrinsic attention weights towards the source document (§ 3.1); the vocabulary distribution comes from the output of the original LLM, ensuring the generative capability of the model (§ 3.2); and the copy probability is derived from the difference between the vocabulary distribution of the model's high layers and the last layer (§ 3.3). Finally, we delve into the rationality of why LargePiG can implicitly transform LLM into a pointer generator (§ 3.4).

3.1 LargePiG: Pointer Attention Distribution

The core module of Large Language Models consists of N stacked Transformer layers. Each Transformer layer contains a self-attention module and feedforward neural networks (FFN) to process the embedded vectors, allowing the model to focus on the most relevant parts of the input dynamically. As the vectors in H_0 pass through each Transformer layer, they are successively transformed, with the output of the layer j represented as H_j . In this process, taking the layer j as an example, H_{j-1} , the output of the layer (j - 1), first passes through the j-th layer's self-attention module. Here, we take Multi-Head Attention (MHA) as an example, which can be easily generalized to Multi-Query Attention [43] and Grouped-Query Attention [2]:

$$MHA = Concat(head_1, \dots, head_M)W^O,$$
(1)

head_i =
$$A_i(H_{j-1}W_i^Q, H_{j-1}W_i^K, H_{j-1}W_i^V)$$
, (2)

348

Figure 2: The architecture of the proposed plug-in and training-free method LargePiG. Pointer Attention Distribution (§ 3.1) from the LLM's self-attention weights, Vocabulary Distribution (§ 3.2) from the output of the original LLM, Copy Probability (§ 3.3) from the difference between the vocabulary distribution of the model's high layers and the last layer.

$$A_i(Q, K, V) = A_i^w V, \ A_i^w = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d}}\right), \tag{3}$$

where A_i^w denotes the attention weights of MHA, with M as the number of heads, $W^{Q/K/V/O}$ are learnable parameters and \sqrt{d} are scaling factor. Since each head captures a unique attention pattern, we aggregate these by averaging: $A^w = \frac{1}{M} \sum_{i=1}^M A_i^w$, enabling a unified representation of attention mechanisms across heads.

In the context of LargePiG, computing the pointer attention distribution P_{source} primarily focuses on the attention weights from the last token in H_{j-1} (i.e., A_{t-1}^w) to the tokens of the source document D. As the source document D corresponds to tokens from m to n in the input sequence, we use $A_{t-1,m:n}^w$ to compute P_{source} . First, for the values in $A_{t-1,m:n}^w$, we normalize them to ensure their sum equals one, forming a probability distribution. Since we are only concerned with the tokens corresponding to the source document in A^w and we already know this is extracted from a larger softmax function, direct normalization suffices. Let this normalized vector be $P_{m:n}$:

$$\mathbf{P}_{m:n} = \frac{A_{t-1,m:n}^{w}}{\sum_{i=m}^{n} A_{t-1,i}^{w}} \tag{4}$$

Next, we construct the probability distribution to match the vocabulary distribution. We depart from traditional PG by not considering new word emergence, focusing on maintaining LLM generation fidelity to input while acknowledging the prevalent use of sentence-piece tokenization [23]. Let \mathcal{V} be the vocabulary of the LLM. The probability distribution for each token x_i in P_{source} within \mathcal{V} comes from the corresponding attention weight in $\mathbf{P}_{m:n}$. Therefore, for each token x_i in the vocabulary \mathcal{V} , its pointer attention distribution $P_{\text{source}}(x_i)$ is defined as:

$$P_{\text{source}}(x_i) + = \begin{cases} \mathbf{P}_{m:n}[j] & \text{for all } j \text{ where } x_j = x_i \text{ and } x_j \in D \\ 0 & \text{otherwise} \end{cases}$$

(5)

Thus, the probability $P_{\text{source}}(x_i)$ for each $x_i \in D$ directly corresponds to the normalized attention weight $\mathbf{P}_{m:n}$, while the probability for vocabulary token not in D is 0.

3.2 LargePiG: Vocabulary Distribution

The generation of the vocabulary distribution in the LargePiG model is seamlessly integrated with the output of the original LLM. This integration is achieved through the model's final component, an affine transformation layer commonly called the classification layer. This layer maps the output of the last Transformer layer H_N , to the vocabulary distribution P_{vocab} over the vocabulary set \mathcal{V} . The probability distribution for the next token x_t given the preceding sequence $x_{< t}$, is computed by applying a softmax function to the affine-transformed output:

$$P_{\text{vocab}}(x_t) = q_N(x_t \mid x_{< t}) = \operatorname{softmax}\left(\phi\left(h_{t-1}^{(N)}\right)\right)_{x_t}, \quad x_t \in \mathcal{V}$$
(6)



524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

where $h_{t-1}^{(N)}$ is the output vector from the last Transformer layer for the position (t-1) in H_N , and $\phi(\cdot)$ performs the affine transformation to project this vector into the vocabulary space. The subscript x_t indicates that we extract the probability corresponding to the token x_t from the softmax output. This approach ensures that the generative capabilities of the underlying LLM are preserved within our LargePiG framework. Through this methodology, LargePiG leverages the extensive linguistic and syntactic knowledge of the LLM, thereby significantly retaining the richness and fluency of the generated query.

3.3 LargePiG: Copy Probability

465

466

467

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

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

The copy probability in our LargePiG model leverages the difference between the vocabulary distribution of the LLM's high layers and the last layer. For the layer *j*, we also compute the vocabulary distribution using $\phi(\cdot)$ as follows, where \mathcal{J} is a set of candidate layers and this operation is called early exiting [41, 44]:

$$q_{j}\left(x_{t} \mid x_{\leq t}\right) = \operatorname{softmax}\left(\phi\left(h_{t-1}^{(j)}\right)\right)_{x_{t}}, \quad j \in \mathcal{J}.$$
(7)

Based on the findings of Chuang et al. [10] and early exit decoding research [13, 41], when LLMs generate function words (e.g., auxiliary verbs, prepositions, conjunctions), the vocabulary distribution $q_i(x_t \mid x_{\le t})$ stabilizes at high layers. In contrast, when generating factual knowledge words (e.g., names, places, dates), the vocabulary distribution continues to evolve at high layers. In the query generation task, we expect the factual content in the generated query primarily comes from the source document, while syntax and other forms are organized by the LLM. This implies we can use the vocabulary distribution $q_N(x_t \mid x_{\le t})$ from the last transformer layer as an anchor layer, and by calculating the distributional differences with the vocabulary distributions from other high layers, determining whether LLM is generating factual knowledge words or function words. A larger distributional difference suggests a higher likelihood of generating factual knowledge words. Since our goal is to ensure that the factual content of the generated query mainly comes from the input document, the copy probability should be higher in such cases, and vice versa. Therefore, the copy probability can be calculated as follows:

$$p_{\rm cp} = O_{j \in \mathcal{J}} d\left(q_N \left(x_t \mid x_{< t} \right), q_j \left(x_t \mid x_{< t} \right) \right), \tag{8}$$

where *O* can be an average $\frac{1}{|\mathcal{J}|} \sum$, a max, or a min operation, $d(\cdot, \cdot)$ is a distributional distance measure such as Jensen-Shannon Divergence [10, 32], and \mathcal{J} is the set of high-layers around the anchor layer. We can control the intensity of copying by adjusting *O* and \mathcal{J} . A larger range of \mathcal{J} and *O* being max increases the likelihood of copying, while a smaller range of \mathcal{J} and *O* being min decreases it.

The final distribution generated by LargePiG is given by:

$$P_{\text{LargePiG}}(x_t) = p_{\text{cp}} P_{\text{source}}(x_t) + (1 - p_{\text{cp}}) P_{\text{vocab}}(x_t).$$
(9)

3.4 Exploring the Internal Mechanisms of LargePiG

The key to LargePiG's functionality lies in LLM's ability to correctly reflect the current generated token's attention weights towards

the source document and generate factual knowledge words and function words in the pattern we mentioned in § 3.3.

Regarding the pointer attention distribution, we analyzed the causes of hallucinations in query generation in § 1, concluding that the attention modules in LLMs are more 'truthful' than the FFN modules and classification laver. The factuality hallucination mainly arises from the LLM's insufficient knowledge about the source document. Some studies have shown that knowledge is mainly stored in the FFN module of the transformer layer in pretrained language model [11]. Even if the self-attention module correctly focuses on the relevant token, the FFN module may still produce factuality hallucinations due to insufficient pre-training [31]. Moreover, Jiang et al. [20] found that MLP modules have a more significant impact on incorrect outputs than attention modules, indicating that in the transformer layers of LLMs, attention modules are more 'truthful' than FFN modules. The relevance hallucination can be attributed to the softmax bottleneck issue inherent in LLMs [7, 51], where the model predicts the probability of each word across the entire vocabulary, struggling to differentiate between words that are almost equally likely in a given pre-training context but have different meanings in the current situation. The softmax bottleneck primarily stems from the final classification layer, which is structurally unrelated to the attention module in the transformer layer we use. In Appendix A.2, we further experimentally verify that the attention modules in LLMs are more 'truthful' than the FFN modules and the classification layer.

Regarding the copy probability, we delve deeper into the findings of [10, 41], questioning why LLM predictions for function words stabilize at high layers' vocabulary distributions, while predictions for factual knowledge words do not. Research on early exit decoding [13, 41, 44] has demonstrated that different data samples (tasks) possess varying complexities. For multi-layer stacked deep models, such as ResNet [16] and LLaMA [45], simple tasks may only require shallow layers for completion, whereas complex tasks demand the involvement of all layers. The scaling law [22] and the emergence ability [49] also testify to this, with the model's ability to solve more complex tasks increasing alongside its size and layer number. Returning to our task, predicting function words can exit at shallower layers, while predicting factual knowledge words requires deeper layers, indicating that predicting function words is simpler, whereas predicting factual knowledge is more complex.

Why is predicting function words simpler, and predicting factual knowledge more complex? Achille et al. [1] demonstrated that tasks with greater information content are more complex. Since LLMs learn from human language, if we can verify that factual knowledge words in human language convey more information than function words, then the pattern mentioned above is determined by the nature of human language itself. Our experimental analysis within our TruthfulVQG and TruthfulDQG benchmarks investigated the semantic impact of removing factual knowledge words versus function words, with experimental details provided in Appendix A.3. The results show that on both datasets, removing factual knowledge words causes a greater decrease in semantic similarity scores with the original sentence compared to function words. These findings confirm that factual knowledge words contribute more significantly to the sentence's informational content

			Qwen1.5	7B Cha	t				LLaMA2	7B Chat	t	
Model	TruthfulVQG			Tr	TruthfulDQG		TruthfulVQG		TruthfulDQG			
Model	MC1	MC2	MC3	MC1	MC2	MC3	MC1	MC2	MC3	MC1	MC2	MC3
Base	40.35	66.97	37.70	27.34	85.77	39.83	52.94	75.12	46.01	33.72	71.61	34.29
+ CD	35.79	63.43	36.49	24.25	85.60	37.99	-	-	-	-	-	-
+ DoLa	37.97	64.73	35.68	23.52	85.05	37.09	52.79	75.25	46.10	35.09	69.97	33.19
+ LargePiG	41.49^{\dagger}	68.12^{\dagger}	38.92^\dagger	29.91 [†]	89.33 †	42.18^{\dagger}	54.56^{\dagger}	76.15^{\dagger}	47.20^{\dagger}	37.23^{\dagger}	70.95	36.93 [†]
PQGR	43.61	70.08	41.26	25.86	77.23	36.86	52.22	74.21	45.60	32.28	65.74	31.41
+ CD	41.71	66.10	40.69	23.84	77.90	35.58	-	-	-	-	-	-
+ DoLa	40.13	66.50	38.24	23.79	76.51	35.67	51.83	73.69	44.54	31.92	64.41	31.52
+ LargePiG	45.52^\dagger	70.79 [†]	42.54^\dagger	27.12^{\dagger}	79.20^{\dagger}	38.35^\dagger	52.87^{\dagger}	74.87^{\dagger}	46.27^{\dagger}	34.66 [†]	68.34^{\dagger}	34.21^{\dagger}
InPars	44.35	70.77	41.56	26.09	78.82	37.37	52.53	74.53	45.85	30.66	64.43	30.32
+ CD	43.91	68.90	39.82	24.06	77.20	35.69	-	-	-	-	-	-
+ DoLa	40.35	66.90	38.48	24.48	77.57	36.96	51.59	74.33	44.86	29.87	63.97	29.52
+ LargePiG	46.26^\dagger	71.51^\dagger	42.82^{\dagger}	27.34^{\dagger}	81.17^\dagger	38.53^\dagger	53.03 [†]	74.74	46.20^{\dagger}	33.70 [†]	67.30 [†]	33.36 †
AQG	40.50	67.26	37.85	27.41	85.86	39.93	54.00	75.92	46.87	34.82	71.62	34.42
+ CD	36.79	63.36	33.44	24.23	83.56	37.96	-	-	-	-	-	-
+ DoLa	37.99	64.65	35.62	25.59	85.28	39.21	52.79	75.25	46.10	33.02	70.96	33.17
+ LargePiG	41.56^{\dagger}	68.13^{\dagger}	39.06 [†]	29.99 [†]	89.58^{\dagger}	42.35^{\dagger}	54.84^{\dagger}	76.73^{\dagger}	47.76^{\dagger}	37.09 [†]	71.04	36.82 [†]

Table 1: Performance comparisons between LargePiG and the baselines. The boldface represents the best performance. '†' means improvements are significant (paired t-test at *p*-value < 0.05).

than function words, highlighting the complexity of predicting factual knowledge words. Verifying that the pattern found in [10, 41], rooted in the linguistic properties of human language, is a principle that holds true across multiple languages, even though initial studies focused on English scenarios. Our subsequent experiments expanded this understanding to multiple languages, validating the feasibility of employing this pattern for calculating copy probability in LargePiG. For further analysis of the effectiveness of copy probability in LargePiG, see Appendix A.3.

4 EXPERIMENT

4.1 Experimental Settings

Datasets. To quantitatively assess the truthful query generation capabilities of LargePiG in both video (e.g., TikTok) and document (e.g., Bing Search) scenarios, considering the absence of relevant datasets, we constructed two challenging benchmarks named Truth-fulVQG and TruthfulDQG. These benchmarks correspond to for-mats similar to TruthfulQA [29], crafted from video (Chinese corpus) and document (English corpus) respectively, to validate the model's query generation truthfulness. The construction of the benchmarks utilized a combination of LLM and manual methods. The completed data format is shown in Table 8 of Appendix A.5, where "Bad queries" are those containing either relevance hallucina-tions or factuality hallucinations or both, "Good queries" are those without any hallucinations, and "Best query" represents the optimal query. The construction process is detailed in Appendix A.4 and

 35'
 54.84'
 76.73'
 47.76'
 37.09'
 71.04
 36.82'

 Appendix A.5, and the statistical results of the datasets are shown in Table 9.
 Metrics. To evaluate LLMs in truthful query generation, we inde

Metrics. To evaluate LLMs in truthful query generation, we independently compute each reference query's log-probability. Drawing inspiration from the evaluation metrics of TruthfulQA-MC [10, 29], the metrics used to assess the truthfulness of the model-generated queries include MC1 (the percentage of all data where the best query log-probability is greater than all bad queries log-probability), MC2 (normalized total probability assigned to the set of good queries), and MC3 (the percentage of all good queries where each good query log-probability is greater than all bad queries log-probability).

Models and Baselines. We employed two types of backbone LLMs, Qwen1.5 7B chat [3] and LLaMA2 7B chat [45], and utilized four LLM-based query generation approaches, including (1) Base: using the backbone LLMs to directly generate queries in a zero-shot manner; (2) PQGR [12]: prompting the LLM with 8 incontext examples to generate queries, which achieves more suitable queries compared to the Base approach; (3) Inpars [5]: includes not only good queries in the in-context examples but also bad queries to enable the model to generate better queries through comparison; (4) AQG [25]: employ LoRA [17] to fine-tuning the LLM using real-world user-input queries and context data to enhance the model's query generation capability. The implementation details of these LLM-based query generation approaches are in Appendix A.6. Our approach, LargePiG, is model-agnostic and can be applied to different LLM-based query-generation methods, reducing the relevance and factuality hallucinations associated with model-generated queries. The implementation details of LargePiG

Conference'17, July 2017, Washington, DC, USA

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

Table 2: Experimental results on multimodal data.

Model	MC1	MC2	MC3
LLaVA-7B	58.40	80.45	51.54
+ LargePiG	59.80	81.74	52.9 4
LLaVA-13B	57.20	79.18	50.74
+ LargePiG	58.10	79.93	51.26

are provided in Appendix A.7. For baseline models, we compared LargePiG with recent closely related work aimed at reducing hallucinations in LLMs: DoLa [10], which enhances factuality in LLMs by decoding through contrasting layers, and Contrastive Decoding (CD) [27], which improves factuality in LLMs' generations by leveraging the contrasts between LLMs of different sizes, selecting tokens that maximize their log-likelihood difference. For Qwen1.5 7B chat, we chose Qwen1.5 1.8B chat [3] as the contrast model for CD. Since there is no smaller-sized model for LLaMA2 7B chat, we could not perform CD experiments on this model. DoLa, CD, and LargePiG are all training-free decoding methods for reducing hallucinations in LLM generation, making them fair for comparison.

4.2 Results

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722 Main result. As shown in Table 1, LargePiG has demonstrated improvements across two datasets, various backbone methods, and 723 different metrics, validating LargePiG's ability to enhance the truth-724 fulness of LLM-based query generation methods. The effectiveness 725 observed across datasets in different languages further corrobo-726 rates the analysis presented in Section 3.4. Moreover, our method 727 has surpassed CD and DoLa, which even exhibited negative gains 728 729 on some datasets. The primary reason is that query generation primarily relies on the factual knowledge in the inputs, requiring 730 731 less generated factual knowledge from the model, whereas DoLa 732 and CD stimulate the model's knowledge by contrasting shallow layers' logits with deep layers' logits or contrasting large LLM's 733 logits with small LLM's logits, which may lead to the generation of 734 735 facts that do not align with the context. In the following analysis experiments, we will further discuss the respective advantages of 736 CD, DoLa and LargePiG, and analyze in detail from the perspec-737 tives of relevance hallucinations and factuality hallucinations how 738 739 LargePiG can improve the truthfulness of LLM generation. Further verification in Appendix A.13 confirmed that queries generated by 740 LargePiG exhibit higher similarity to the real queries. Additionally, 741 742 human evaluation and case studies in Appendix A.13 showed that 743 LargePiG not only reduced relevance and factual hallucinations in 744 the generated queries but also made them more appealing to users.

745 Multimodal result. LargePiG is effective not only on large language models but can also be applied to Large Vision-Language 746 Models (LVLM), further enhancing the truthfulness of query gener-747 748 ation that integrates both vision and language modalities. We selected the recently popular large vision-language model LLaVA [24] 749 as the backbone model. Detailed method descriptions about the im-750 plementation can be found in Appendix A.8. To validate LargePiG's 751 752 ability to address hallucination issues in multimodal query generation tasks, we compiled a multimodal version of the TruthfulVQG

Table 3: Experimer	t results	on FACTOR
--------------------	-----------	-----------

	LLaMA-7B		LLaMA-13B		
Model	News	Wiki	News	Wiki	
Base	58.3	58.6	61.1	62.6	
+ CD [26]	-	-	62.3	64.4	
+ DoLa [10]	62.0	62.2	62.5	66.2	
+ LargePiG	71.0	60.4	72.1	63.1	
+ DoLa + LargePiG	63.4	64.7	65.3	68.8	

dataset, named TruthfulVQG-M. Experimental results on LLaVA-7B/13B, shown in Table 2, indicate that the truthfulness of queries generated by LargePiG surpasses those produced by the original decoding method, confirming the effectiveness of LargePiG in multimodal tasks. We also observed that LLaVA-13B performs less effectively than LLaVA-7B, a potential reason being that in the video query generation task, due to the high noise level in video content, the more complex LLaVA-13B model might be more sensitive to noise. Furthermore, short videos contain some new content not present in the pre-training data, which could lead to easier overfitting to the training data in a zero-shot scenario, thus resulting in suboptimal performance compared to LLaVA-7B.

4.3 Analysis

LargePiG's ability to reduce factuality hallucinations. To specifically validate LargePiG's capability to address factual hallucinations, we selected the News and Wiki categories of FACTOR dataset [33], which assesses LLMs' factuality in long-paragraph settings by completion task. The News' ground-truth answers are based on facts from news content, which LLMs may not have sufficiently learned during training; the Wiki contains general facts well-learned during pre-training, allowing LLMs to respond based on pre-trained knowledge and also to learn from the context. To ensure a fair comparison with DoLa, we chose LLaMA-7B and LLaMA-13B as the backbone LLMs following DoLa's setting.

The experimental results shown in Table 3 demonstrate that on the News dataset, LargePiG successfully enhanced the copy ability of Base models to address hallucinations, thereby significantly outperforming other methods that solely rely on the model's intrinsic pre-trained knowledge and original context understanding capabilities. Given the feature of the Wiki dataset, although the results for LargePiG on Wiki do not surpass other methods that stimulate the model's own pre-trained knowledge, they still exceed the base model, validating the contribution of LargePiG's copy ability to resolving hallucinations. Moreover, LargePiG can be combined with state-of-the-art methods that are based on the model's pretrained knowledge, achieving advancements beyond the current state of the art (i.e., +DoLa + LargePiG > +DoLa). This suggests that LargePiG's copy ability can be synergistically integrated with the model's inherent pre-trained knowledge.

LargePiG's ability to reduce relevance hallucinations. To independently verify LargePiG's capability to resolve relevance hallucinations, we generated queries using different models and

Anon.



Figure 3: Semantic similarity win rate of Qwen1.5-7B-Chat with LargePiG vs without LargePiG on TruthfulVQG.

 Table 4: Relevance win rate comparison between Qwen1.5-7B

 Chat with LargePiG and without LargePiG on TruthfulVQG.

Model	LargePiG Win	Original Model Win	Tie
Base	827	70	103
PQGR	749	181	70
InPars	805	141	54
AQG	831	73	96

then encoded them and the corresponding context using the cur-rent state-of-the-art text representation model BGE [50] to calcu-late their cosine semantic similarity. The pairwise comparisons of cosine similarity are presented on Figure 3, demonstrating that LargePiG notably outperforms the baseline models. The results on TruthfulDQG are detailed in Appendix A.9, which presents similar conclusions to those found in the experiments on TruthfulVQG. This indicates that LargePiG effectively reduces the relevance hallucinations of query generation. In addition, we used GPT-40 (from OpenAI) to assess LargePiG's ability to reduce relevance halluci-nations (see prompt in Appendix A.11). Considering time and API cost factors, we sampled 1000 data points from TruthfulVQG for evaluation. The experiments in Table 4, judged by GPT-40, further confirm that LargePiG can mitigate relevance hallucinations.

LargePiG's ability to copy. To validate whether LargePiG has a stronger copy ability compared to the original LLM decoder, we tested the performance of LLM with and without the addition of LargePiG on tasks that require copying from the inputs. Following the setting of Jelassi et al. [18] for validating LLMs' copy capability, we selected the SQuAD question-answering dataset [38], which pro-vides text paragraphs along with several questions pertaining to the text and features various inputs lengths. We conducted experiments on Qwen1.5-7B-Chat, reported the F1 score, and classified ques-tions into short and long categories based on whether their length exceeded 200 words. The results on Figure 4 show that LargePiG significantly improved the F1 score on Qwen1.5-7B-Chat, with more pronounced improvements for scenarios with long inputs, indicat-ing that LargePiG indeed enhances the copy ability of LLMs. Similar results on LLaMA2-7B-Chat are shown in Appendix A.10.



Figure 4: Comparison of the Copying Ability between Qwen1.5-7B-Chat and Qwen1.5-7B-Chat with LargePiG on the SQuAD dataset.

Table 5: Decoding latency (ms/token).

	Baseline	DoLa	LargePiG
Base / AQG	95.9 (×1.00)	99.9 (×1.04)	101.8 (×1.06)
InPars	135.1 (×1.00)	142.4 (×1.05)	139.8 (×1.03)
PQGR	142.0 (×1.00)	148.3 (×1.04)	149.1 (×1.05)

Efficiency analysis. We use NVIDIA V100-32G GPUs and 52core Intel(R) Xeon(R) Gold 6230R CPUs at 2.10GHz machine to analyze the efficiency of original decoding (baseline), DoLa, and LargePiG when applied across different query generation models. The decoding time of LargePiG in LLaMA2-7B models increases by a maximum of 6% compared to the baseline and is on par with the decoding time of DoLa, as shown in Table 5 (experiments on Qwen1.5-7B are detailed in Appendix A.12). The results demonstrate that LargePiG can enhance the truthfulness of query generation with negligible additional time consumption, proving the practical applicability of LargePiG.

5 CONCLUSIONS

LLM-based query generation significantly improves query quality and user experience in information retrieval systems, but it also introduces hallucination challenges, hindering its application in emerging use cases such as "related search". To address these, we propose LargePiG, a training-free method transforming an LLM into a Pointer-Generator. LargePiG separates content and form in LLMgenerated queries, using input knowledge for fact generation and LLM capabilities for syntactic structure. It combines self-attention weights for pointer attention distribution, LLM original output as vocabulary distribution, and high-layer vocabulary distribution for copy probability. Our empirical evaluations on the proposed TruthfulVQG and TruthfulDQG datasets confirm LargePiG's effectiveness in reducing hallucination on query generation tasks.

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

929 **REFERENCES**

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

- Alessandro Achille, Giovanni Paolini, Glen Mbeng, and Stefano Soatto. 2021. The information complexity of learning tasks, their structure and their distance. Information and Inference: A Journal of the IMA 10, 1 (2021), 51–72.
- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebron, and Sumit Sanghai. 2023. GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints. In <u>Proceedings of the 2023</u> Conference on Empirical Methods in Natural Language Processing. 4895–4901.
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen Technical Report. arXiv preprint arXiv:2309.16609 (2023).
 - [4] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. 2016. Ms marco: A human generated machine reading comprehension dataset. arXiv preprint arXiv:1611.09268 (2016).
 - [5] Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and Rodrigo Nogueira. 2022. Inpars: Unsupervised dataset generation for information retrieval. In Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval. 2387–2392.
 - [6] Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. 2024. Internlm2 technical report. arXiv preprint arXiv:2403.17297 (2024).
 - [7] Haw-Shiuan Chang, Zonghai Yao, Alolika Gon, Hong Yu, and Andrew Mccallum. 2023. Revisiting the Architectures like Pointer Networks to Efficiently Improve the Next Word Distribution, Summarization Factuality, and Beyond. In <u>Findings</u> of the Association for Computational Linguistics: ACL 2023. 12707–12730.
 - [8] Xiang Chen, Chenxi Wang, Yida Xue, Ningyu Zhang, Xiaoyan Yang, Qiang Li, Yue Shen, Jinjie Gu, and Huajun Chen. 2024. Unified Hallucination Detection for Multimodal Large Language Models. arXiv preprint arXiv:2402.03190 (2024).
 [9] I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting
 - [9] I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. 2023. FacTool: Factuality Detection in Generative AI-A Tool Augmented Framework for Multi-Task and Multi-Domain Scenarios. arXiv preprint arXiv:2307.13528 (2023).
 - [10] Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R Glass, and Pengcheng He. 2023. DoLa: Decoding by Contrasting Layers Improves Factuality in Large Language Models. In <u>The Twelfth International Conference on Learning</u> Representations.
 - [11] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge Neurons in Pretrained Transformers. In <u>Proceedings of the 60th</u> <u>Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 8493–8502.</u>
 - [12] Zhuyun Dai, Vincent Y Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith Hall, and Ming-Wei Chang. 2022. Promptagator: Few-shot Dense Retrieval From 8 Examples. In <u>The Eleventh International Conference on</u> Learning Representations.
 - [13] Siqi Fan, Xin Jiang, Xiang Li, Xuying Meng, Peng Han, Shuo Shang, Aixin Sun, Yequan Wang, and Zhongyuan Wang. 2024. Not all Layers of LLMs are Necessary during Inference. arXiv preprint arXiv:2403.02181 (2024).
 - [14] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997 (2023).
 - [15] Mitko Gospodinov, Sean MacAvaney, and Craig Macdonald. 2023. Doc2query--: When less is more. In <u>European Conference on Information Retrieval</u>. Springer, 414–422.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In <u>Proceedings of the IEEE conference on</u> computer vision and pattern recognition. 770–778.
- [17] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In International Conference on Learning Representations. https://openreview.net/forum?id=nZeVKeeFYf9
- [18] Samy Jelassi, David Brandfonbrener, Sham M Kakade, and Eran Malach. 2024. Repeat after me: Transformers are better than state space models at copying. arXiv preprint arXiv:2402.01032 (2024).
- [19] Xin Jia, Wenjie Zhou, Xu Sun, and Yunfang Wu. 2021. Eqg-race: Examinationtype question generation. In Proceedings of the AAAI conference on artificial intelligence, Vol. 35. 13143–13151.
- [20] Che Jiang, Biqing Qi, Xiangyu Hong, Dayuan Fu, Yang Cheng, Fandong Meng, Mo Yu, Bowen Zhou, and Jie Zhou. 2024. On Large Language Models' Hallucination with Regard to Known Facts. arXiv preprint arXiv:2403.20009 (2024).

- [21] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. <u>IEEE Transactions on Big Data</u> 7, 3 (2019), 535–547.
- [22] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <u>arXiv preprint arXiv:2001.08361</u> (2020).
- [23] Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations. 66–71.
- [24] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2024. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. <u>Advances</u> in Neural Information Processing Systems <u>36</u> (2024).
- [25] Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing. 2023. Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting over heterogeneous sources. In <u>The</u> <u>Twelfth International Conference on Learning Representations.</u>
- [26] Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2022. Contrastive decoding: Open-ended text generation as optimization. <u>arXiv preprint arXiv:2210.15097</u> (2022).
- [27] Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. Contrastive Decoding: Open-ended Text Generation as Optimization. In <u>Proceedings of the 61st</u> <u>Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 12286–12312. https: //doi.org/10.18653/v1/2023.acl-long.687</u>
- [28] Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Textbooks are all you need ii: phi-1.5 technical report. arXiv preprint arXiv:2309.05463 (2023).
- [29] Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 3214– 3252.
- [30] Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2023. Hallusionbench: You see what you think? or you think what you see? an image-context reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models. <u>arXiv preprint</u> <u>arXiv:2310.14566</u> (2023).
- [31] Ang Lv, Kaiyi Zhang, Yuhan Chen, Yulong Wang, Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan. 2024. Interpreting Key Mechanisms of Factual Recall in Transformer-Based Language Models. arXiv preprint arXiv:2403.19521 (2024).
- [32] ML Menéndez, JA Pardo, L Pardo, and MC Pardo. 1997. The jensen-shannon divergence. Journal of the Franklin Institute 334, 2 (1997), 307–318.
- [33] Dor Muhlgay, Ori Ram, Inbal Magar, Yoav Levine, Nir Ratner, Yonatan Belinkov, Omri Abend, Kevin Leyton-Brown, Amnon Shashua, and Yoav Shoham. 2023. Generating benchmarks for factuality evaluation of language models. <u>arXiv</u> preprint arXiv:2307.06908 (2023).
- [34] Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, KaShun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. RAGTruth: A Hallucination Corpus for Developing Trustworthy Retrieval-Augmented Language Models. In <u>Proceedings</u> of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Bangkok, Thailand.
- [35] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document expansion by query prediction. <u>arXiv preprint arXiv:1904.08375</u> (2019).
- [36] Zhiyuan Peng, Xuyang Wu, and Yi Fang. 2023. Soft prompt tuning for augmenting dense retrieval with large language models. <u>arXiv preprint arXiv:2307.08303</u> (2023).
- [37] Zhiyuan Peng, Xuyang Wu, Qifan Wang, Sravanthi Rajanala, and Yi Fang. 2024. Q-PEFT: Query-dependent Parameter Efficient Fine-tuning for Text Reranking with Large Language Models. <u>arXiv preprint arXiv:2404.04522</u> (2024).
- [38] Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know What You Don't Know: Unanswerable Questions for SQuAD. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Iryna Gurevych and Yusuke Miyao (Eds.). Melbourne, Australia, 784–789. https: //doi.org/10.18653/v1/P18-2124
- [39] Jon Saad-Falcon, Omar Khattab, Keshav Santhanam, Radu Florian, Martin Franz, Salim Roukos, Avirup Sil, Md Sultan, and Christopher Potts. 2023. UDAPDR: Unsupervised Domain Adaptation via LLM Prompting and Distillation of Rerankers. In <u>Proceedings of the 2023 Conference on Empirical Methods in</u> <u>Natural Language Processing</u>. 11265–11279.
- [40] Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving Passage Retrieval with Zero-Shot Question Generation. In <u>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</u>. 3781–3797.
- 1042 1043 1044

1046

1047

1048

1049

1050

1051

1052

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

- [41] Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Tran, Yi Tay, and Donald Metzler. 2022. Confident adaptive language modeling. <u>Advances</u> in Neural Information Processing Systems 35 (2022), 17456–17472.
- [42] Abigail See, Peter Liu, and Christopher Manning. 2017. Get To The Point: Summarization with Pointer-Generator Networks. In <u>Association for Computational</u> Linguistics. https://arxiv.org/abs/1704.04368
- [43] Noam Shazeer. 2019. Fast transformer decoding: One write-head is all you need. arXiv preprint arXiv:1911.02150 (2019).
- [44] Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. 2016.
 Branchynet: Fast inference via early exiting from deep neural networks. In 2016 23rd international conference on pattern recognition (ICPR). IEEE, 2464–2469.
- [45] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <u>arXiv</u> preprint arXiv:2307.09288 (2023).
 - [46] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. Advances in neural information processing systems 28 (2015).
 - [47] Siyuan Wang, Zhongyu Wei, Zhihao Fan, Yang Liu, and Xuanjing Huang. 2019. A multi-agent communication framework for question-worthy phrase extraction and question generation. In <u>Proceedings of the AAAI conference on artificial</u> intelligence, Vol. 33. 7168–7175.
 - [48] Yujing Wang, Yingyan Hou, Haonan Wang, Ziming Miao, Shibin Wu, Qi Chen, Yuqing Xia, Chengmin Chi, Guoshuai Zhao, Zheng Liu, et al. 2022. A neural corpus indexer for document retrieval. <u>Advances in Neural Information Processing</u> Systems 35 (2022), 25600–25614.
 - [49] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent Abilities of Large Language Models. <u>Transactions on Machine Learning Research</u> (2022).
 - [50] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighof. 2023. C-pack: Packaged resources to advance general chinese embedding. <u>arXiv preprint</u> arXiv:2309.07597 (2023).
 - [51] Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W Cohen. 2018. Breaking the Softmax Bottleneck: A High-Rank RNN Language Model. In International Conference on Learning Representations.
 - [52] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. <u>Advances in Neural Information Processing Systems</u> 36 (2024).

A APPENDIX / SUPPLEMENTAL MATERIAL

A.1 Experimental Verification of Hallucination Classification

Our hallucination classification for generated query is grounded in real-world observations, aiming to help readers better understand the distinct types of hallucinations present in query generation. This categorization is intended to offer valuable insights for future research in this domain. To further validate our classification, we conducted an analysis experiment using the TruthfulVQG dataset (Detailed in Section 4.1). In this experiment, we encoded both the generated queries and corresponding video content using BGE [50] and computed the cosine similarity to obtain a **Semantic Similarity score**. The results demonstrate that factual hallucinations can occur independently, even in the absence of relevance hallucinations, highlighting the need to decouple these two types of hallucinations for more precise handling.

The experimental results in Table 6 confirm that factual hallucinations can persist even with high semantic similarity scores. This finding underscores the importance of treating relevance and factual hallucinations separately to improve query generation and retrieval quality.

A.2 Evaluation of attention modules are more 'truthful' than FFN modules

In section 3.4, we propose that Attention modules are more 'truthful' than other modules in LLMs (e.g., FFN modules). To validate this core observation, we conducted experiments on the RAGTruth dataset, which is a word-level hallucination corpus in various tasks within the Retrieval-augmented generation (RAG). The RAGTruth dataset contains responses from various LLMs that exhibit hallucinations in RAG scenarios, with manually annotated hallucination spans, hallucination types, and hallucination reasons. The RAG scenario is similar to the query generation scenario, as both involve generation based on input content, approximating the query generation context.

Specifically, in our validation experiments, we selected data from RAGTruth [34] where LLaMA2-7B-Chat exhibited hallucinations. Using LangChain ², a widely used open-source toolkit, we applied the RecursiveCharacterTextSplitter to segment the input retrieved document into different spans. We then calculated whether the attention module attended span (mean pooling the attention scores then selecting the input span with the highest score) of LLaMA2-7B-Chat during the generation of hallucination spans could identify the hallucination spans in the response. This evaluation was based on GPT4-0 (from OpenAI) using the following prompt:

Prompt: { external context + query}

Respond: {response}

- Conflict Span: { Conflict Span} Conflict Type: { Conflict Type} Reason: {Reason}
- Given the following context information: "{Attend Span → }", can this support the existence of a → conflict in the response? Please answer with " → Yes" or "No" and give the reason on the

 \hookrightarrow newline.

The results are shown in Table 7. We found that, in most cases, the attention modules of LLMs can attend to the correct input spans to identify hallucinations in the response. This indicates that hallucinations in LLMs are caused by other modules in LLMs (e.g., FFN modules), thereby proving the core observation that Attention modules are more 'truthful' than other modules in LLMs.

A.3 Implementation Details of Words Information

In the experiments concerning word information, we conducted tests using the TruthfulVQG and TruthfulDQG benchmarks constructed in this paper. For English in the TruthfulDQG benchmark, we used Spacy ³ for tokenization and part-of-speech tagging, while for TruthfulVQG (Chinese corpus), we employed Jieba ⁴ and Hanlp ⁵. Factual knowledge words include organizations, personal names, locations, and dates. Function words include auxiliary verbs, prepositions, determiners, conjunctions, and coordinating conjunctions. Subsequently, on both datasets, we removed an equal number of factual knowledge words and function words and then utilized BGE embeddings [50] to align and compare the cosine similarity

https://www.langchain.com/.	

³https://spacy.io ⁴https://github.com/fxsjy/jieba ⁵https://www.hanlp.com

1103

1104

1105

1106

1135

1136

1137

1138

1139

1154

1155

1156

1157

1158

1159

Table 6: Semantic similarity scores across query types, highlighting that factual hallucinations can occur despite high similarity with relevant content.

Туре	Max Semantic Similarity	Min Semantic Similarity	Average Semantic Similarity
Facticity hallucination queries	0.8492	0.3792	0.6479
Facticity truth queries	0.8607	0.2647	0.6482
Random similarity	N/A	N/A	0.2709

Table 7: Proportion of data where Llama2 7B Chat attentionheads attend to the correct information.

1161

1162

1163

1164

1165 1166

1167

1168

1169

1170

1171

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1218

Attention heads attend	Attention heads mis-attend
77.5%	22.5%

between the modified sentences and the original sentences. The results are shown below:

- In the TruthfulDQG benchmark, removing factual knowledge words resulted in a similarity score of 0.7741, while removing function words led to a higher similarity score of 0.9296.
 - In the TruthfulVQG benchmark, the removal of factual knowledge words produced a similarity score of 0.7415, compared to 0.9477 when function words were eliminated.

The results show that on both datasets, removing factual knowledge words causes a greater decrease in semantic similarity scores with the original sentence compared to function words. These findings confirm that factual knowledge words contribute more significantly to the sentence's informational content than function words, highlighting the complexity of predicting factual knowledge words. Verifying that the pattern found in [10, 41], rooted in the linguistic properties of human language, is a principle that holds true across multiple languages, even though initial studies focused on English scenarios.

Why Can LLM Identify Factual Knowledge Words and Function Words?

Considering that LLMs can only directly learn to predict the next word in the natural language training corpus, they may not have an intuitive concept of what constitutes factual knowledge words and function words. Therefore, we conducted an intrinsic frequency analysis of factual knowledge and function words on the TruthfulDQG benchmark. The statistical results are shown below:

- Number of different words in function words: 228
- Number of different words in factual knowledge words: 3263
- Total number of words in function words: 33849
- Total number of words in factual knowledge words: 6026
- Average occurrence of function words: 148.46
- Average occurrence of factual knowledge words: 1.85

These results show that function words appear much more frequently than factual knowledge words, particularly evident from their average occurrences. It is evident that due to the substantially larger training data of function words compared to factual knowledge words, LLMs can predict function words at shallower layers while predicting factual knowledge words need deeper layers.

Another Perspective on the Effectiveness of Copy Probability in LargePiG.

Besides the pattern we mentioned above, Jiang et al. [20] observes that in hallucinated cases, the output token's information rarely shows abrupt increases and maintains consistent superiority in high layers of the LLMs. This corresponds to cases in LargePiG where there is a higher copy probability, thus enabling the reduction of hallucinations by copying factual knowledge words from the source document. This further demonstrates the capability of the copy probability in LargePiG to address the issue of hallucinations.

A.4 Details about Dataset Collection

The TruthfulVQG dataset is collected from a real short video platform used by over one billion users. The TruthfulDQG dataset is adapted from the MS-MARCO dataset [4]. The data processing for TruthfulVQG is more complex than TruthfulDQG's. Thus, we will use TruthfulVQG as an example to illustrate the process.

Data Collection:

The raw data was collected from Search Click Data and Post-Watch Search Data, and the final processed public data does not include any user search information, only video content, and LLMgenerated queries.

• Collected Data Source:

- Search Click Data (30,000 samples): We collect 30,000 samples of users' clicked videos after searching the corresponding queries with data flowing from query to video.
- Post-Watch Search Data (10,000 samples): We collect 10,000 samples of users' searched queries after watching the corresponding videos, which is a smaller subset compared to click data, with data flowing from video to query.
- Criteria for Inclusion:
 - Search Click Data: Include only data with positions greater than one and less than twenty to mitigate position bias of the top results and low relevance of farther results.
 - Post-Watch Search Data: Include only data with total count numbers greater than five to ensure relevance to previously viewed videos.

Components of Video Content:

- Title: Accurate representation of video content.
- Video Dialogue Text (ASR): Prone to noise but contain detailed information about the video.
- Video Text Information (OCR): More reliable than ASR and contains more information than Title.

1274

1275

1276

1219

1220

1221

Data Preprocessing: Remove examples lacking textual features,
 containing sensitive words, or background music that affect ASR
 results.

Next, we will use LLMs to generate multiple queries for data annotation of all videos. To enable the LLMs have the ability to generate high-quality queries, we first fine-tuned these LLMs. Then, we combined them with the original LLMs to generate queries.

Model Fine-Tuning: • Models Used:

 Qwen1.5 7B Chat [3] and InternLM 7B Chat [6] ⁶: Among the strongest for Chinese language capabilities.

Purpose:

 Employing multiple LLMs ensures diversity in generated queries, reducing the risk of repetitive queries that single model sampling might produce.

Data Utilization and Query Generation

- Sort data by video quality scores and select the top 10,000 samples for query generation (Generation is time-intensive, approximately 40 hours per week. Hence, only the top entries are used).
 - Approximately 20+ queries are generated per video using the following prompt.

Query Generation Prompt:

instruction : Based on the video's title , dialog text , and					
\hookrightarrow text information within the video, generate a					
\hookrightarrow relevant and engaging search query. This query					
\hookrightarrow should accurately reflect the video content,					
\hookrightarrow adhere to factual information, and stimulate user					
\hookrightarrow interest to drive clicks . Ensure the query is					
\hookrightarrow concise and contains key information points.					
input: Title: { Title content }					
Dialog text : {Dialog text content}					
Text information: {Text information content}					
Query:					
output: {Query content}					

This prompt is also used in our experiments to generate queries ⁷.

A.5 Details about Dataset Annotation.

During the data annotation section, we first performed further cleaning and filtering of the data. We utilized a combination of LLM and manual annotation to label TruthfulDQG and Truthful-VQG. This hybrid approach of LLM and manual annotations has been employed in numerous works on hallucination benchmark annotation [8, 30].

A.5.1 Phase One: Filter Dataset. Remove sensitive words and per form heuristic query quality filtering based on repetitiveness and
 length scores.

A.5.2 Phase Two: Relevance Assessment. This phase focuses on
 detecting relevance hallucination by measuring the relevance of
 generated queries to the video content.

Similarity Calculations

- (1) **Embedding-Based Similarity:** Utilizes BAAI BGE Embedding [50] and cosine similarity to compute similarity scores between text embeddings.
- (2) Word-Based Similarity: Employs Jieba for text segmentation and calculates similarity using the Jaccard similarity ⁸.

Weighting Method Adjusts relevance scoring based on the ASR noise level:

ASR Score = $0.6 \times \cos(ASR, OCR) + 0.4 \times \cos(ASR, Title)$

	$(0.34 \times (Query, Title) +$	
	$0.33 \times (Query, ASR) +$	
	$0.33 \times (Query, OCR),$	if ASR Score > 0.5
	$0.4 \times (Query, Title)+$	
Query Scoring =	$0.2 \times (Query, ASR) +$	
	$0.4 \times (Query, OCR),$	if ASR Score > 0.3
	$0.5 \times (Query, Title)+$	
	$0.1 \times (Query, ASR) +$	
	$0.4 \times (Query, OCR),$	otherwise

A.5.3 *Phase Three: Factuality Assessment.* Detecting the factuality hallucination of the generated queries by using LLM-based fact-checking methods–Self-Check (4-shot CoT) and FacTool [9].

Self-Check (4-shot CoT). We implement Self-Check (4-shot CoT) using the larger and more powerful LLM Qwen1.5-72B-Chat [3] to detect queries' factuality hallucination. The prompt is shown below ⁹:

You will receive a query generated by another model. Your					
\hookrightarrow task is to check whether this query contains any					
\hookrightarrow factual errors. Please refer to the examples and					
\hookrightarrow guidelines below when evaluating the query:					
- If the query accurately reflects verifiable facts, it					
\hookrightarrow should be considered factually correct.					
- If the query contains misleading or inaccurate					
\hookrightarrow information, it should be considered factually					
\hookrightarrow incorrect.					
- If you cannot determine the accuracy of the query, or					
\hookrightarrow if the query requires more context for					
\hookrightarrow evaluation, it should be considered					
\hookrightarrow indeterminate.					
- Your response must follow the specified format,					
\hookrightarrow containing two keys: "reasoning" (the process					
\hookrightarrow of reasoning) and " factuality " (the judgment					
\hookrightarrow of factuality , where True if the query is					
\hookrightarrow factually correct or does not involve factual					
\hookrightarrow information; False if the query contains					
\hookrightarrow factual errors; No if indeterminate).					
You must respond only in the format described below.					
\hookrightarrow Do not reply in any other form. Adding any					
\hookrightarrow content that violates the response format is					
\hookrightarrow prohibited. Start your response with '{{'.					
https://goilit loam arg/gtable/modules/generated/glearm metrics is seend as an					
HUUS.//SUINT=TEATTLUTP/STAUTE/HUUUTES/PETETATEU/SKTEATTLITTEUTUS.TAUGATU_SCOFE.					

⁸https://scikit-learn.org/stable/modules/generated/sklearn.metrics.jaccard_score. html

Anon

⁶We replace InternLM 7B Chat with LLaMA2 7B Chat on TruthDQG.

 ¹³³² ⁷As the TruthfulVQG is a Chinese Dataset, we translate the prompt from Chinese
 using ChatGPT-4.
 1334

⁹As the TruthfulVQG is a Chinese Dataset, we translate the prompt from Chinese using ChatGPT-4.

```
[Response Template]:
{{
  "reasoning ": "Reason whether the query is factual.
        \hookrightarrow Think through step by step .",
  " factuality ": "True if the query is factually correct
        \hookrightarrow or does not involve factual information:
        \hookrightarrow False if the query contains factual errors;
        \hookrightarrow No if indeterminate."
}}
Examples:
1. [Query]: "Collapse of a tunnel in Antarctica"
{{
  "reasoning ": "This query contains a factual error.
        \hookrightarrow Given the extremely low temperatures in
        \hookrightarrow Antarctica, constructing tunnels is
        \hookrightarrow extremely difficult, and based on current
        \hookrightarrow knowledge, there are no tunnels in
        \hookrightarrow Antarctica, thus a collapse cannot occur.",
  " factuality ": False
}}
2. [Query]: "The Asian Games in Hangzhou will open on
      \hookrightarrow September 23, 2023"
{{
   "reasoning": "The factuality of this query cannot be
        \hookrightarrow determined with the information at hand; it
        \hookrightarrow requires consultation of the latest official
        \hookrightarrow announcements or news sources to verify the
        \hookrightarrow specific opening date .",
  " factuality ": No
}}
3. [Query]: "How to make scrambled eggs with tomatoes"
{{
  "reasoning": "This query is not about the
        \hookrightarrow truthfulness of a statement but requests a
        \hookrightarrow recipe, therefore it does not involve
        \hookrightarrow factual errors .",
  " factuality ": True
}}
4. [Query]: "Messi is Argentine"
{{
  "reasoning": "This query is factually correct. Lionel
        \hookrightarrow Messi is a well-known football player born
        \hookrightarrow in Argentina, a fact that is widely known
        \hookrightarrow and can be verified through reliable sources
        → .",
    factuality ": True
}}
Below is the given query -
[Query]: {}
```

Advanced Fact-Checking. For indeterminate cases after Self-Check, we use advanced fact-checking tools FacTool [9] with Qwen1.5 72B Chat [3] and Serper ¹⁰ to further check queries' factuality based on external data sources from Google Search. The prompt is shown below ¹¹:

You are an excellent assistant.
You will receive a piece of text. Your task is to
\rightarrow identify any factual errors within this text
When judging the factuality of the given text you may
\hookrightarrow refer to provided evidences if necessary
These evidences could be helpful. Some evidences might
\hookrightarrow contradict each other. You must be
careful when using evidences to assess the factuality
\hookrightarrow of the given text
The response should be a dictionary containing three
here = "reasoning" " factuality "
"error" and "correction " corresponding to the
() reasoning whether the given text is
\rightarrow reasoning, whether the given text is
the (boolean value - flue of Faise), the factual effor
\rightarrow present in the text, and the
corrected text.
Below is the given text
[text]: {query}
Below is the provided evidence
[evidences]: {evidence}
You should respond only in the format described below.
\hookrightarrow Do not return any other content.
Start your response with '{{'.
[response format]:
{{
"reasoning ": "Why is the given text factual or not?
\hookrightarrow Be careful when you claim
something is not factual . When you claim something is
↔ not factual , you must provide
multiple pieces of evidence to support your decision
\hookrightarrow .",
"error ": " If the text is factual, then None;
\hookrightarrow otherwise, describe the error .",
" correction ": " If there is an error , then the
\hookrightarrow corrected text .",
" factuality ": " If the given text is factual , then
\hookrightarrow True; otherwise, False ."
}}

Finally, the completed data format is shown in Table 8, and the statistics of TruthfulVQG and TruthfulDQG are shown in Table 9.

Human Assessment. To further ensure the relevance and factual accuracy of the query, we request three annotators with graduatelevel qualifications to manually evaluate the "good queries" to confirm factuality and relevance to the context, ensuring they are both engaging and appropriate.

¹⁰The website of Serper is https://serper.dev/.

¹¹As the TruthfulVQG is a Chinese Dataset, we translate the prompt from Chinese using ChatGPT-4.

Table 8: Description of data fields in TruthfulVQG and TruthfulDQG.

	Video / Document Content	Best query	Good Queries	Bad Queries
Data Type	string	string	[string]	[string]
Description	Description of the video / document content	Best query (factual and most relevent)	Array of good queries	Array of bad queries

Table 9: Statistics of TruthfulVQG and TruthfulDQG. # denotes the average number.

Dataset	Data Count	# Good Queries	# Bad Queries	# Total Queries	Language
TruthfulVQG	4,148	3.82	4.75	8.56	Chinese
TruthfulDQG	2,718	4.04	4.00	8.05	English

A.6 Implementation Details of LLM-based Query Generation Approaches

The prompts used on TruthfulDQG for different LLM-based query generation approaches are shown below (The prompts used on TruthfulVQG are just different in the instruction, which has been demonstrated on Appendix A.4):

Base / AQG:

Given the following document, generate a concise, factual → and relevant query that a user might type into a → search engine to find this information. Document: {Document contents}. Related Query:

PQGR:

Given the following document, generate a concise, factual → and relevant query that a user might type into a → search engine to find this information. Example 1: Document: {Document contents}. Related Query: {The query relevant and factual to document → contents}. ... Example 9: Document: {Document contents}. Related Query:

InPars:

Given the following document, generate a concise, factual
\hookrightarrow and relevant query that a user might type into a
\hookrightarrow search engine to find this information.
Example 1:
Document: {Document contents}.
Related Query: {The query relevant and factual to document
\hookrightarrow contents}.
Hallucination Query: {The query irrelevant and unfactual to
\hookrightarrow document contents}.
Example 4:
Document: {Document contents}.
Related Query:

The size of the dataset for LoRA fine-tuning AQG is 10,000 pairs. The fine-tuning targets the q_proj and v_proj within the transformer layers. The learning rate is set to 5e-5, the per-device train batch size is 4, and the gradient accumulation steps are 4.

A.7 Implementation Details of LargePiG.

We run all the experiments on machines equipped with NVIDIA V100 GPUs and 52-core Intel(R) Xeon(R) Gold 6230R CPUs at 2.10GHz. We utilize the Huggingface Transformers package to conduct experiments. During the decoding of responses from the language models, we employ random sampling with a temperature of 0.8 and a maximum of 256 new tokens to generate responses. The rest of the parameters use the models' default settings. As for selecting the layer to calculate the pointer attention distribution, we used the last layer's attention weights by comparing them with other layers. As for selecting the words to calculate the pointer attention distribution, we recommend filtering the function words in the input using tools detailed in Appendix A.3. Considering that the Jensen-Shannon divergence is usually small in the high-dimensional space of vocabulary distribution, we scale the copy probability p_{cp} in LargePiG by a factor of α . To ensure that the scaled p_{cp} remains within a reasonable range, we clip its value to be less than 0.5, thus maintaining a balance between copy and generation. The value of α is selected from the set [100, 500, 1000]. The $O_{j \in \mathcal{J}}$ in Equation 8 is selected as max, and \mathcal{T} comprises the last 8 or 16 layers of the backbone LLMs, excluding the anchor layer which is the last layer (for increased efficiency, either even or odd numbered layers may be selected). We use two-fold validation to select the hyper-parameters. The LLaMA2-7B-Chat can be downloaded from https://huggingface.co/metallama/Llama-2-7b-chat-hf. The Qwen1.5-7B-Chat can be down-

llama/Llama-2-7b-chat-hf. The Qwen1.5-7B-Chat can be downloaded from https://huggingface.co/Qwen/Qwen1.5-7B-Chat. Due to the limited Chinese training corpus of LLaMA2-7B-Chat, we used Llama2-Chinese-7b-Chat on TruthfulVQG, which can be downloaded from https://huggingface.co/LinkSoul/Chinese-Llama-2-7b.

A.8 Details about LargePiG Applied to LLaVA

The architecture of LLaVA [24] is straightforward, comprising only a Vision Encoder, Projection, and Language Model, with training



(a) Video cover one. Map tokens: Cat, gray, black, elve, eyes, sitting ...



(b) Video cover two. Map tokens: pandas, white, fang, gry, Chinese ...

Figure 5: An example of two video covers mapped to tokens, where we have ignored other irrelevant words and the "_" character before some tokens.

conducted in two stages: Stage 1: Pre-training for Feature Alignment, and Stage 2: Fine-tuning End-to-End. A key issue when applying LargePiG to LLaVA concerns how to map image tokens to text tokens, thus establishing an attention distribution based on the source content. Considering during the Feature Alignment stage, the primary task is aligning the image features H_v with the pre-trained LLM word embeddings, we propose mapping each image token to the closest text token in the embedding space when computing the Pointer Attention Distribution. In the implementation, we utilize the faiss vector database [21] to store text token embeddings and retrieve the corresponding tokens using image token embeddings, allowing for rapid retrieval of relevant tokens. Case studies shown in Figure 5 reveal that this retrieval method can accurately reveal the main information in the images, although many noise tokens are also retrieved. Therefore, we apply rule-based filtering to remove tokens with low similarity to the text part and construct the attention distribution using the remaining tokens together with the text tokens.

A.9 More results on LargePiG's Ability to Reduce Relevance Hallucinations

More results on LargePiG's ability to reduce relevance hallucinations are shown in Figure 6 and the left of Figure 7, both LLaMA27B-Chat and Qwen1.5-7B-Chat with LargePiG can generate more
semantic relevance queries with the document / video contents,
indicating that LargePiG is effective in reducing the relevance hallucinations of query generation.

A.10 More Results on LargePiG's Copy Ability

The results of LargePiG with LLaMA2-7B-Chat on the SQuAD question-answering dataset are shown on the right of Figure 7, which also show that LargePiG significantly improved the F_1 score on LLaMA2-7B-Chat, with more pronounced improvements for scenarios with long inputs, indicating that LargePiG indeed enhances the copy ability of LLMs.

A.11 Prompt for relevance evaluation by GPT4-0

The evaluation prompt are shown below:

You will be given a description of a video and queries						
\hookrightarrow generated by the baseline model and the LargePiG						
\hookrightarrow model based on the video description . Your task is						
\hookrightarrow to determine which model generates higher-quality						
\hookrightarrow queries. When evaluating the queries, please						
\hookrightarrow refer to the following guidelines :						
1. Are the queries relevant to the video description?						
2. You must reply only in the format described below.						
\hookrightarrow Adding any extra content that violates the reply						
\hookrightarrow format is prohibited.						
Begin your reply with '{{'.						
[Reply Template]:						



Figure 6: Results of LLaMA2-7B-Chat without LargePiG vs with LargePiG on TruthfulDQG. Left: Overall semantic similarity scores. Right: Win rate with LargePiG compared against without LargePiG.



Figure 7: Left: Overall semantic similarity scores of Qwen1.5-7B-Chat with LargePiG vs without LargePiG on TruthfulVQG. Right: Performance of LLaMA2-7B-Chat vs LLaMA2-7B-Chat + LargePiG on SQuAD.

'win_model': 'LargePiG (LargePiG generated query is more \hookrightarrow relevant) or Baseline (Baseline generated query is \hookrightarrow more relevant) or Tie (both models generated \hookrightarrow similar queries)', 'reason ': 'The reason for determining the winning model in → the previous statement' }} Video description : `{}` Baseline generated query: `{}` LargePiG generated query: `{}`

More Results on Efficiency Analysis A.12

Table 10 shows the decoding latency for different models on Qwen1.5-7B. It is evident that compared to LLaMA2-7B, the inference latency of Qwen1.5-7B is significantly reduced. Therefore, the addition of DoLa or LargePiG, although increasing the time cost compared

to LLaMA2-7B, still shows a relatively small overall increase. The maximum increase in time cost is about 10%, which is within an acceptable range.

Generated Query Quality Evaluation A.13

To verify the quality of queries generated by LargePiG compared with the baseline models, we first encode the generated queries and the corresponding user-input queries using BGE Embedding. Subsequently, we compute the cosine similarity to compare the semantic similarity between queries generated by different models and those input by users. As can be observed from Table 11 and 12, queries generated by LargePiG exhibit higher similarity to actual user input queries, thereby confirming the high quality of LargePiGgenerated queries from a semantic relevance perspective.

To further validate the performance of LargePiG in real-world scenarios, we compared it with the online query generation model (A 7 billion parameter transformer decoder-only model) of a shortvideo platform with billions of users, observing the performance of the online model after integrating LargePiG. In the offline evaluation, we retrieved the results generated by the online model along with the corresponding input data (totaling 1108 samples) and

Anon.

	Baseline	DoLa	LargePiG
Base / AQG	64.21 (×1.00)	70.67 (×1.10)	68.56 (×1.07)
InPars	69.54 (×1.00)	75.82 (×1.09)	76.50 (×1.10)
PQGR	82.26 (×1.00)	92.45 (×1.12)	89.49 (×1.09)

Table 10: Decoding latency (ms/token) on Qwen1.5-7B.

 Table 11: Semantic Evaluation of Generated Query Quality on the Qwen1.5 7B Chat.

	Qwen1.5 7B Chat			Qwen	1.5 7B C	hat + Lar	gePiG	
	Base	PQGR	InPars	AGQ	Base	PQGR	InPars	AGQ
TruthfulDQG	0.6751	0.6629	0.6586	0.6858	0.7086	0.6893	0.6725	0.7077
TruthfulVQG	0.6143	0.5983	0.5957	0.6102	0.6260	0.6074	0.6056	0.6247

Table 12: Semantic Evaluation of Generated Query Quality on the LLaMA2 7B Chat.

		LLaMA2-7B-Chat			LLaM	A2-7B-C	hat + Lar	gePiG
	Base	PQGR	InPars	AGQ	Base	PQGR	InPars	AGQ
TruthfulDQG	0.6574	0.6561	0.6561	0.6568	0.6850	0.6897	0.6979	0.6870
TruthfulVQG	0.5898	0.5687	0.5674	0.5898	0.6003	0.5792	0.5875	0.5953

Table 13: Human and LLM evaluations of the queries generated by LargePiG and the original online model.

	Baseline Win	LargePiG Win	Tie
Count Number	54	1031	23

conducted generation after adding LargePiG to the online model. Subsequently, we employed a collaborative approach of LLM and human evaluation. Initially, the Qwen1.5-72B-Chat was used to determine whether LargePiG wins, the base model wins, or if it is a tie, providing reasons for each. Then, two human evaluators with graduate-level qualifications reviewed the LLM's outputs, correct-ing any erroneous assessments made by the LLM, thereby enhanc-ing the overall efficiency and accuracy of the evaluations. Combined with experimental results on Table 13 and case studies (translated from Chinese) below, it was demonstrated that adding LargePiG not only reduced the relevance and factual hallucinations in the generated queries but also made them more attractive to users, fur-ther validating the effectiveness of LargePiG. From the analysis of case studies, we found that the reason why LargerPiG can generate queries that are more attractive to users may be the interpolation of vocabulary distribution, which can reduce the probability of generating an end token. Moreover, during the query generation process, there is a consistent high alignment with the video content. Consequently, the generated queries are more detailed and specific, thereby attracting more user clicks. **Evaluation Prompt:**

You will receive a video's description, along with queries \hookrightarrow generated by the baseline model and the LargePiG \hookrightarrow model based on that video description . Your task \hookrightarrow is to determine which model produced the higher \hookrightarrow quality query. When evaluating the queries, please refer to the following guidelines : - Whether the query is relevant to the video \hookrightarrow description - Whether the query is factually accurate - Whether the query can attract user interest - You must reply only in the format described \hookrightarrow below. Do not respond in any other \hookrightarrow form. Adding any extra content that \hookrightarrow violates the reply format is \hookrightarrow prohibited. Start your reply with \hookrightarrow '{{'. [Reply Template]: łł "win_model": "LargePiG (if LargePiG generated a \hookrightarrow better query) or Baseline (if \hookrightarrow Baseline generated a better query) or \hookrightarrow Tie (if both models generated similar \hookrightarrow queries) ", "reason ": "The reason for the previous winning ← model decision' }} Video description : {}

1973 Query generated by baseline : {}	4. Farewell Song IV series online viewing,	2031
1974 Query generated by Largeric: {}	5. An Duo Xiao Zhe	2032
1975 Case studies:	7. Where to watch Farewell Song TV series	2033
1977 Example One:	8 Farewell Song finale	2034
1978	9 An Duo Xiao Zhe arguing	2035
video_content: Video Category: Film and Television. Video	LargePiG	2037
\hookrightarrow Tags: Mermaid. Video Title : The most terrifying	1 Farewell Song TV series	2038
\hookrightarrow hunters often appear in the guise of their prey,	2 Farewell Song An Duo and Xiao Zhe	2039
\hookrightarrow horror thriller mermaid series #FastTalkingFilm #	3 An Duo and Xiao Zhe arguing	2040
\hookrightarrow KuaishouCinemaSuperstarPlan #FilmCommentary.	4 An Duo and Xiao Zhe diyorce.	2041
\hookrightarrow Video Cover: Fish Cost. Video Author: Firebird Films	5. An Duo Farewell Song	2042
1985 \hookrightarrow ,	6. Xiao Zhe An Duo.	2043
Baseline :	7. Xiao Zhe Farewell Song An Duo finale.	2044
1. Mermaid Movie,	8. Xiao Zhe An Duo divorce truth.	2045
2. Mermaid Movie Commentary,	Judge Result :	2046
3. Mermaid Movie Episode One,	win model: LargePiG.	2047
4. Commentary on the First Episode of Mermaid,	reason: The LargePiG model generated queries are more	2048
5. Complete Mermaid Movie Online Viewing,	\hookrightarrow specific, incorporating key events from the	2049
6. First Part of Mermaid,	\hookrightarrow video description (arguing, divorce) which	2050
7. Second Part of Mermaid Movie,	\hookrightarrow better reflect the video content and engage	2051
8. Complete Mermaid Movie Online Viewing Uncut,	\hookrightarrow user interest, such as 'An Duo and Xiao Zhe	2052
LargePiG:	\hookrightarrow arguing', 'An Duo and Xiao Zhe divorce', 'Xiao	2053
1. Mermaid Movie,	\hookrightarrow Zhe An Duo divorce truth '. In contrast, the	2054
2. Mermaid Movie Online Viewing,	\hookrightarrow Baseline queries, although relevant, do not	2055
3. Complete Mermaid Video,	\hookrightarrow sufficiently highlight the conflicts and plot	2056
4. Horror Mermaid Movie,	\hookrightarrow developments in the video. such as ' finale '	2057
5. High Definition Online Viewing of Mermaid Movie,	\hookrightarrow and 'online viewing', which are more general	2058
6. Free High Definition Online Viewing of Horror	\hookrightarrow and factually incorrect.	2059
\hookrightarrow Mermaid Movie,	,	2060
7. Uncut Super High Definition Online Viewing of	Example Three	2061
\hookrightarrow Mermaid Movie,		2062
8. Mermaid Movie Commentary,	video content: Video Category: Anime Video Tags: Ultraman	2063
Judge Result :	\hookrightarrow Video Title : Episode 4 Go Ahead, Ultraman Zero!	2064
win_model: LargePiG,	\hookrightarrow Performed by: Shota Minami / voyager Lyrics by:	2065
reason: The LargePiG's queries contain more specific	→ Hideki Tamiya/Koki Konishi Composed by: Tomohiro	2066
\hookrightarrow information, such as 'horror movie' and 'free	→ Yamaguchi Edited by: Ultraman Information Master	2067
\hookrightarrow high definition online viewing, which more	→ Ultraman Golden Songs Zero Edition #Ultraman #	2068
\hookrightarrow accurately reflect the video's theme (horror	→ ZeroUltraman 15th Anniversary Video Cover: Go	2069
\hookrightarrow thriller) and are likely to generate more user	→ Ahead Ultraman Zero Performance by: voyager / Shota	2070
\hookrightarrow interest . Additionally, it offers a variety of	→ Minami Composed by: Tomohiro Yamaguchi Lyrics by:	2071
\hookrightarrow viewing options, such as the uncut super high	→ Hideki Tamiya/Koki Konishi Edited by: Ultraman,	2072
\hookrightarrow definition version, which may be more	→ Information Master Video Author: Ultraman	2073
\hookrightarrow appealing to users.	\hookrightarrow Information Master,	2074
2017	Baseline :	2075
2018 Example Two:	Go Ahead, Ultraman Zero,	2076
2019 video_content: Video category: Film and TV show; Video tags	Ultraman Zero,	2077
2020 \hookrightarrow : Farewell Song; Video title : Xiao Zhe argues with	Go Ahead, Ultraman Zero Song,	2078
2021 \hookrightarrow An Duo, An Duo proposes divorce, and they are	Complete Lyrics of Ultraman Zero Theme Song,	2079
2022 \hookrightarrow destined to break up! #CatchTheNewDrama #	Go Ahead, Zero,	2080
2023 → WebDramaFarewellSong; Video cover: Attending an	Original Singer of Go Ahead, Ultraman Zero,	2081
2024 → international music festival ; Video creator : Old	How to Sing Go Ahead, Ultraman Zero Song,	2082
2025 \hookrightarrow Friend Qi (recruiting apprentices),	Original Sound of Go Ahead, Zero,	2083
2026 Baseline :	Ultraman Zero Go Ahead,	2084
2027 1. Farewell Song TV series,	LargePiG:	2085
2028 2. Farewell Song episode 36 finale,	Go Ahead, Ultraman Zero,	2086
3. Farewell Song An Duo and Xiao Zhe,	Original Singer of Go Ahead, Ultraman Zero Song,	2087

- 3. Farewell Song An Duo and Xiao Zhe,

		·····
2089 2090 2091 2092 2093 2094 2095 2096 2097 2098	Go Ahead, Ultraman Zero Theme Song Lyrics, Go Ahead, Ultraman Zero Anime Episode One, Go Ahead, Ultraman Zero Theme Song, Go Ahead, Ultraman Zero Ultraman Zero Song, Ultraman Zero, How to Sing Go Ahead, Ultraman Zero Song., Ultraman Zero Go Ahead, Original Singer of Go Ahead, Ultraman Zero Ultraman ← Zero Song., Judge Result :	win_model: LargePiG, 2147 reason: LargePiG's queries are more specific , 2148 → containing more information related to the 2149 → video content such as 'Go Ahead, Ultraman Zero 2150 → Anime Episode One', which can stimulate user 2151 → interest and provide a richer background 2152 → related to the video. In contrast , Baseline 's 2153 → queries , while related to the video theme, are 2154 → more generic and do not specify details such 2155 → as the original singer or anime episodes . 2156
2099		2157
2100		2158
2101		2159
2102		2160
2103		2161
2104		2162
2105		2163
2106		2104
2107		2103
2109		2167
2110		2168
2111		2169
2112		2170
2113		2171
2114		2172
2115		2173
2116		2174
2117		2175
2118		21/6
2119		21// 2178
2120		2170
2122		2180
2123		2181
2124		2182
2125		2183
2126		2184
2127		2185
2128		2186
2129		2187
2130		2188
2131		2109
2132		2170
2134		2192
2135		2193
2136		2194
2137		2195
2138		2196
2139		2197
2140		2198
2141		2199
2142		2200
2143		2201
2145		2202 2203
2146		19 2204