# WebGLM: Towards An Efficient Web-Enhanced Question Answering System with Human Preferences

Xiao Liu*
liuxiao21@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Hanyu Lai*
laihy19@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Hao Yu*
yuhao2019@buaa.edu.cn
Beihang University
Beijing, China

Yifan Xu
xuyifan2001@gmail.com
Tsinghua University
Beijing, China

Aohan Zeng
zah22@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Zhengxiao Du
zx-du20@mails.tsinghua.edu.cn
Tsinghua University
Beijing, China

Peng Zhang
peng.zhang@zhipuai.cn
Zhipu.AI
Beijing, China

Yuxiao Dong†
yuxiaod@tsinghua.edu.cn
Tsinghua University
Beijing, China

Jie Tang†
jietang@tsinghua.edu.cn
Tsinghua University
Beijing, China

## Abstract

We present WebGLM, a web-enhanced question-answering system based on the General Language Model (GLM). Its goal is to augment a pre-trained large language model (LLM) with web search and retrieval capabilities while being efficient for real-world deployments. To achieve this, we develop WebGLM with strategies for the LLM-augmented retriever, bootstrapped generator, and human preference-aware scorer. Specifically, we identify and address the limitations of WebGPT (OpenAI), through which WebGLM is enabled with accuracy, efficiency, and cost-effectiveness advantages. In addition, we propose systematic criteria for evaluating web-enhanced QA systems. We conduct multi-dimensional human evaluation and quantitative ablation studies, which suggest the outperformance of the proposed WebGLM designs over existing systems. WebGLM with the 10-billion-parameter GLM (10B) is shown to perform better than the similar-sized WebGPT (13B) and even comparably to WebGPT (175B) in human evaluation. The code, demo, and data are at https://github.com/THUDM/WebGLM.

## CCS Concepts

- **Computing methodologies** → *Natural language generation*;
- **Software and its engineering** → *Development frameworks and environments*.

---

*XL, HL, and HY contributed equally and this work was done when HY interned at Tsinghua. †Corresponding Authors: YD and JT.
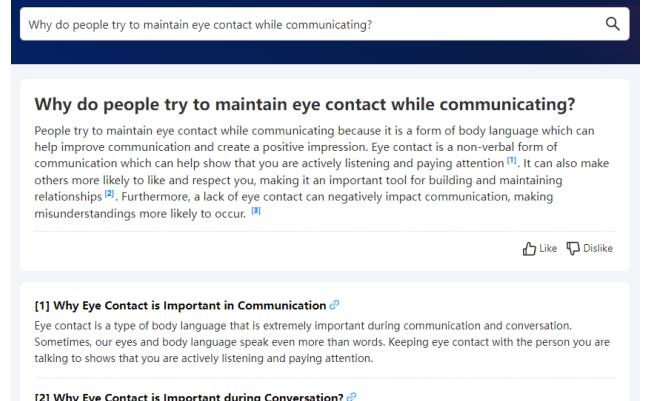
---

**Figure 1: A screenshot of WebGLM's response to an example question with web references.**

## Keywords

Large Language Model; Pre-Trained Model; Human Preference Alignment; General Language Model

## 1 Introduction

Large language models (LLMs), such as GPT-3 [3], PaLM [5], OPT [37], BLOOM [32], and GLM-130B [36], have significantly pushed the boundary of machines' ability on language understanding and generation. Question answering [15, 28], one of the most fundamental language applications, has also been substantially advanced by the recent LLM developments. Existing studies suggest that the
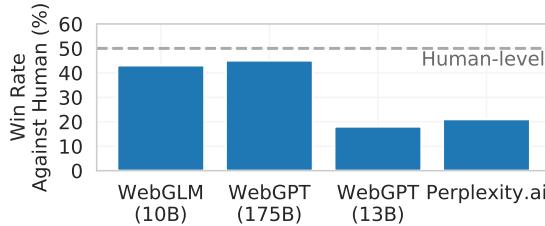
**Figure 2: The win rates of popular web-enhanced QA systems against human references.** WebGLM (10B) performs comparably to WebGPT (175B), approaching human-level QA ability.

performance of LLMs' closed-book QA [29] and in-context learning QA [3, 18] is comparable to supervised models, furthering our understanding on LLMs' potential to memorize knowledge.

However, even for LLMs, their capacity is not unlimited, and when it comes to challenges that require sufficient rare-knowledge, LLMs fail to meet up human expectations. Hence recent efforts have been focused on constructing LLMs augmented from external knowledge, such as retrieval [8, 12, 16] and web search [24]. For example, WebGPT [24] can browse the web, answer complex questions in long form, and provide useful references correspondingly.

Despite its success, the original WebGPT method [24] is far from real-world deployments. First, it relies on abundant expert-level annotations of browsing trajectories, well-written answers, and answer preference labeling, requiring considerable expenses, time, and training. Second, the behavior cloning method (i.e., imitation learning) requires its base model GPT-3 to emulate human experts by instructing the system to interact with a web browser, issue operation commands (e.g., Search, Read, and Quote), and then retrieve relevant information from online sources. Finally, the multi-turn nature of web browsing demands intensive computation resources and can be too slow for user experience, e.g., costing about 31 seconds for WebGPT-13B to response a 500-token prompt.

In this work, we present WebGLM—a practical web-enhanced QA system based on the 10-billion-parameter General Language Model (GLM-10B) [6]. An example is illustrated in Figure 1. It is efficient, cost-effective, human preference-aware, and most importantly, of comparable quality to WebGPT. The system employs multiple new strategies and designs to achieve good performance, including:

**An LLM-augmented Retriever**: a two-staged retriever that implements coarse-grained web search and fine-grained LLM-distilled retrieval. It is inspired by the fact that LLMs like GPT-3 can naturally learn to adopt correct references, and such ability could be distilled to improve smaller dense retrievers.

**A Bootstrapped Generator**: a GLM-10B based answer generator that is trained on quoted long-formed QA samples and bootstrapped by LLM in-context learning. We discover that instead of relying on expensive human expert writing in WebGPT, LLMs can be enabled to learn to generate high-quality data with proper citation-based filtering.

**A Human Preference-aware Scorer**: a scorer, that is trained over online QA forums' user thumb-up signals, is able to learn human majority preferences on different answers. Compared to WebGPT's expert labeling, we prove that a proper dataset construction could also produce a high-quality scorer.

Our extensive human evaluation and quantitative ablation results demonstrate the efficiency and effectiveness of the WebGLM system. Specifically, WebGLM (10B) surpasses the similar-scaled WebGPT (13B) and performs comparably to WebGPT (175B) on our Turing test (Cf. Figure 2). WebGLM's improvement against the only publicly-available system—Perplexity.ai—also makes it among the best public web-enhanced QA systems as of this submission.

To sum up, in this paper, we make the following contributions:

- We construct WebGLM, an efficient web-enhanced QA system with human preferences. It significantly outperforms the similar-sized WebGPT (13B) and performs comparably to WebGPT (175B). It also surpasses Perplexity.ai—a popular system powered by LLMs and search engines.
- We identify WebGPT's limitations on real-world deployments. We propose a set of new designs and strategies to allow WebGLM's high accuracy while achieving efficient and cost-effective advantages over baseline systems.
- We formulate the human evaluation metrics for evaluating web-enhanced QA systems. Extensive human evaluation and experiments demonstrate WebGLM's strong capability and also generate insights into the system's future developments.

## 2 Related Work

The construction of web-enhanced QA systems is a systematic project that requires cross-domain collaboration, including large language models, open-domain question answering, retrieval augmentation, and reinforcement learning from human feedback. Here we briefly introduce related literature on them.

**Large Language Models (LLMs).** Self-supervised [19] LLMs have attracted plenty of attention in nowadays natural language processing (NLP). Their huge number of parameters captures and stores versatile knowledge [20] and enables their outstanding performance on various challenges. Typical LLMs include GPT-3 [3], PALM [5], OPT [37], BLOOM [32], and GLM-130B [36]. One of the fascinating LLM properties is prompt-based in-context learning (ICL), which allows tuning-free task transfer via prepended demonstration samples. Recent works have been focusing on the optimization [18, 22, 34, 39] and analysis [23, 30, 35] of ICL.

**Open-domain Question Answering (Open QA).** Traditional QA datasets such as SQuAD [28] assume the reference is available. On the contrary, open-domain QA targets the open world and is more practical but challenging. For example, Natural Questions [15] dataset consists of queries from the Google search engine and annotations from Wikipedia paragraphs. Web Questions [2] derives open-domain questions from knowledge bases. MS Marco [25] gathers passage texts and corresponding labels to questions.

However, most Open QA datasets and models are limited to answer short answer phrases, while people usually prefer more informative long-formed answers with references. A possible reason is that constructing and evaluating long-formed QA datasets with open-world references are difficult, requiring expert-level annotations. Recent attempts include ELI5 [7] that collects queries and long-formed answers with scores from Reddit and WebGPT [24] which hires groups of experts and leverages up to 175-billion-parameter GPT-3 as the backbone. WebGLM aims to provide another effective and cost-effective solution for the challenge.

**Retrieval-augmentation.** Mainstream information retrieval approaches include sparse-vector-based BM25 and TF-IDF, and the recent dense-vector-based methods such as DPR [14] and Contriever [10]. The idea of retrieval-augmented language models introduced by REALM [8] argues the joint optimization of retriever and language modeling. Following representative works include RAG [16], Fusion-in-Decoder [11], and Atlas [12]. The idea of WebGPT also loosely falls into the field, as it asks the LLM to interact with the browser to seek relevant information for better accuracy. Nevertheless, it can cost intensive computation and is too slow for practical deployment. In this work, WebGLM tackles the problem efficiently by distilling LLMs' knowledge to smaller retrievers.

**Reinforcement Learning from Human Feedback (RLHF).** Automated scoring of text generation is a well-established area of research. BLEU [27] and ROUGE [17] take into account the overlap ratio between the target and reference. METEOR [1] considers the accuracy and recall rate of the whole corpus. Other methods, such as BERTScore [38], evaluate using cosine similarity of contextual embedding from deep language models. In recent years, some work advocates learning scorers from human feedback [26, 33] via asking models to predict human preference. The scorers, or namely reward models, can be used to optimize the text generator via reinforcement learning. Such methods, which WebGPT is also affiliated with, have achieved great success in real-world applications.

## 3 The WebGLM System

Constructing an LLM-based web-enhanced QA system can be expensive and challenging. The web information is rich but noisy for certain queries, and creating high-quality human answers with references for training can be outrageously expensive. This type of systems usually involves three critical components: retriever, generator, and scorer.

Take WebGPT [24] as an example, which employs experts for dataset annotation. Its retriever leverages GPT-3 to "behavior-clone" human experts' web-browsing trajectory to search, read, and quote. In addition, the generator is trained on expert-written long answers with references. And finally, the scorer learns to predict experts' preferences over different answers, and its scores serve as rewards for the generator's reinforcement learning. Despite WebGPT's primary success, its retrieval can be slow, and the data annotations required for training the generator and scorer are too costly, significantly hindering its wide public adoptions.

In this work, we aim to build an efficient web-enhanced QA system that understands human preferences for actual use. To combine the advantages of LLMs and well-established open QA studies, we present a series of new designs and strategies for our web-enhanced QA system WebGLM based on GLM [6]:

- **An LLM-augmented Retriever**: we design two stages: coarse-grained web search and fine-grained LLM-augmented dense retrieval [10], for finding relevant references given queries.
- **A Bootstrapped Generator**: we derive WebGLM-QA, an LLM-bootstrapped quoted and long-formed QA dataset via in-context learning and corresponding strategies to clean and refine. It includes 45k high-quality after filtering and 83k noisy but diverse samples before filtering. The backbone of WebGLM system is a GLM model trained on the dataset.

- **A Human Preference-aware Scorer**: we develop techniques to learn human majority preference from online QA forums' thumb-ups instead of expensive expert feedback, and successfully train a human preference-aware scorer for best-of-n selection.

The LLM API used for research purpose in this work is text-davinci-003 unless specified. In the following sections, we will introduce the algorithm and implementation details of each component, which finally form the WebGLM pipeline sequentially.

### 3.1 LLM-augmented Retriever

In conventional open QA, the systems usually only retrieve from reliable sources (e.g., Wikipedia) and fail to benefit from whole web-scale knowledge. However, the flip side of the coin is that wild web pages can be hard to acquire and purify. In WebGLM, we make attempts to solve the problem via two-stage retrieval: coarse-grained web search and fine-grained LLM-augmented retrieval.

#### 3.1.1 Coarse-grained Web Search

We leverage third-party web search engines (i.e., Google API) to acquire primary candidate web page URLs. In most cases, from our observation, these pages can cover the necessary contexts and knowledge to answer questions besides considerably abundant irrelevant information. The procedures are shown in Figure 3. Specifically, it can be roughly divided into three steps:

(1) **Search**: At this stage, we enter the question into the search API and will obtain a list of URLs for potentially-relevant pages (usually less than 10).
(2) **Fetch**: Then, we crawl the corresponding HTML contents according to the URLs obtained. Since there are many candidate pages, we improve efficiency through parallel crawling.
(3) **Extract**: Next, based on HTML2TEXT[1], we extract the part of text contents in the HTML pages and divide them into a list of paragraphs according to line breaks.

Since the web crawl usually takes sufficient time, we have paid great efforts to optimize the speed of the component to allow user-acceptable responding speed (Cf. Figure 4). For example, in the "Fetch" step, if the page is loaded synchronously, the loading time will be 2-3 minutes long. The parallel asynchronous enables the quick loading of most pages in 5s (about 98%).

#### 3.1.2 Fine-grained LLM-augmented Retrieval

Through the first three stages, we have retrieved a number of potential contexts to questions. However, many of them are still irrelevant even under the filtering of widely-used dense retrievers (in our trial, up to 30% of top-ranked contexts are unrelated). As a solution, WebGPT [24] uses behavior cloning (i.e., imitation learning) to leverage LLMs' strong language comprehensibility for reference selection. Notwithstanding its effectiveness, the strategy is slow in deployment and expensive in labeling.

**LLMs' Reference Adoption.** To mitigate the issue, we propose to combine smaller retrievers' efficiency and LLMs' strong ability to distinguish. We take Contriever [10] as the smaller retriever in WebGLM, an unsupervised pre-trained model that encodes texts

---

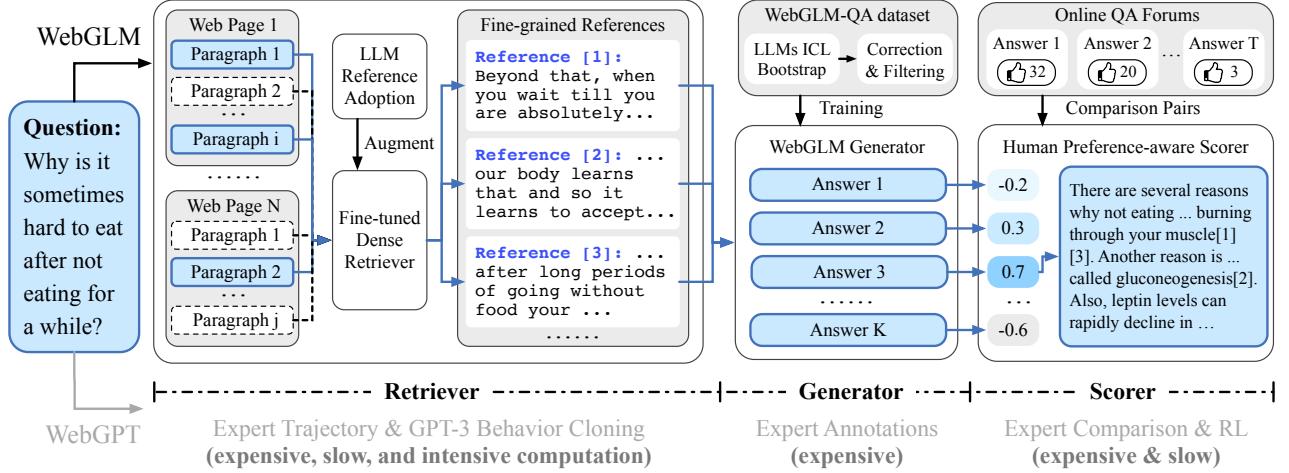[1]https://github.com/aaronsw/html2text

**Figure 3: WebGLM system pipeline.** Our system includes three sub-modules: LLM-augmented retriever recalls the top-5 most relevant paragraphs as the reference sources; Bootstrapped generator yields answers according to the question and reference sources; Human preference-aware scorer assesses all answers and picks the highest-scored one as the final result. Compared to WebGPT, WebGLM is a more efficient and cost-effective web-enhanced QA system with comparable answer quality.
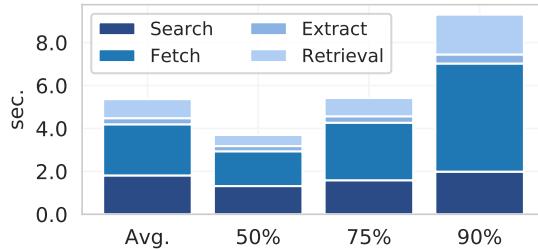


**Figure 4: WebGLM retriever time analysis.** 50% of queries can be done within 4.0s, and 90% of them can be loaded within 10.0s. Most of time is spent on fetching web pages after searching.

into embeddings and retrieves by finding the maximum inner product pair of them. We transfer LLMs' natural property of reference adoption to small retrievers to improve them.

Specifically, we find LLMs can naturally distinguish and only adopt useful references in in-context learning (ICL). We create a 200-query dataset, where each query is accompanied with 5 top-ranked candidate references from Contriever. We manually annotate the relevance of each piece of reference (Cf. Table 1). We find only 68.6% of them are related. However, when we provide the query with corresponding candidate references to GPT-3 for 1-shot in-context learning inference (see details in Section 3.2), we discover that the LLM would only adopt part of the references and the corresponding accuracy is 90.2%, far better than Contriever's.

**Table 1: Evaluation on LLM's reference adoption.**

| Method | Acc. |
|---|---|
| Contriever | 68.6% |
| LLM ICL adoption | 90.2% |

**Augmentation Implementation.** To transfer the reference adoption knowledge from GPT-3 to Contriever, we leverage the GPT-3's reference adoption from our bootstrapped dataset WebGLM-QA to additionally fine-tune Contrievers. As the reference marks generated by GPT-3 can be wrong sometimes, we use the citation correction method based on Rouge-1 precision to match quotations and references (see those details in Section 3.2). Therefore, the labels we use for training are the Rouge-1 precision scores of a query-reference pair.

In the fine-tuning, we use two Contrievers to encode questions and references individually, and compute their inner products as the predictions. We leverage Mean Square Error (MSE) as the loss function for the predictions and Rouge-1 precision scores to train the Contrievers. Our further quantitative experiment demonstrates that the augmentation significantly improves Contriever web-enhanced QA retrieval accuracy (see Table 7 for details).

*3.1.3 Speed analysis*
Retrieval is no doubt the most time-consuming part in any web-scale QA system. A slow QA system, whatever high its accuracy is, would spoil the user experience. We report the speed of each steps in our LLM-augmented retriever.

We sample a subset from ELI5 [7] test set to retrieve and calculate the average, the median, 75% quantile, 90% quantile, and 99% quantile time spent in each step. From Figure 4, we can know that our average time spent is about 5.3s, the median total time spent is about 4.07s, and 90% of searches can be loaded in 10s. The main bottleneck of our retrieval is in the second step of fetching each page, when we have to request multiple web pages from different sources. Consequently, due the contents of various pages on the network are different, some pages take very long time to load, or just cannot be returned correctly.

In Appendix B, we conduct a more detailed analysis of retrieval efficiency and point out that the retrieval efficiency of WebGLM is far better than that of WebGPT.

## 3.2 Bootstrapped Generator

A major obstacle in building web-enhanced QA system is the high cost for curating expert-level QA datasets that are long-formed and properly cited. Compared to traditional or free-formed QA,
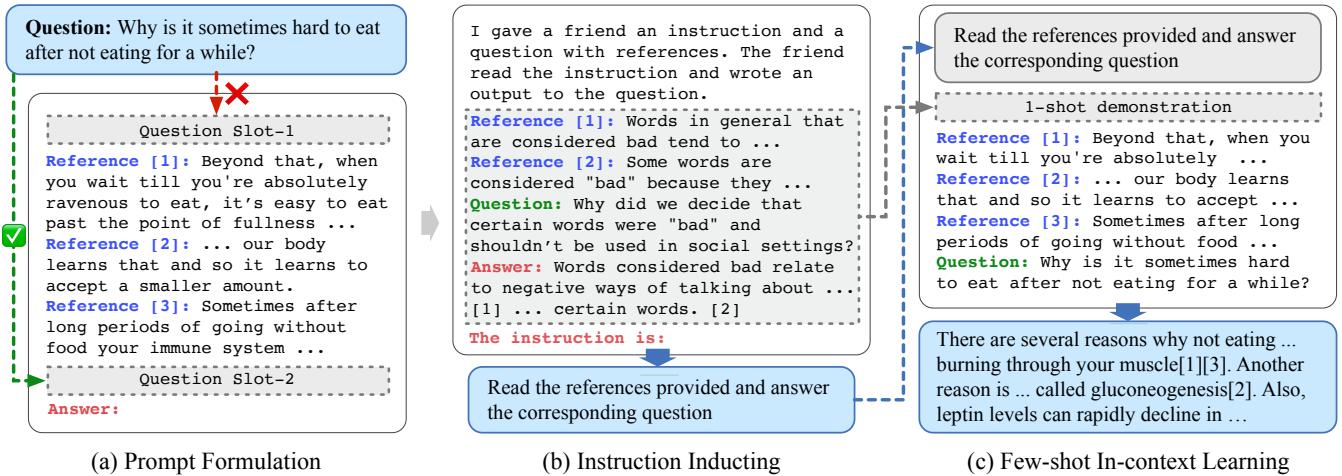
**Figure 5: We construct WebGLM-QA for generator training via LLM in-context bootstrapping.** It includes three stages: 1) prompt formulation, 2) instruction inducting, and 3) few-shot in-context learning. In this way, we avoid the outrageous cost in time and money for hiring experts but still create a high-quality quoted long-formed QA dataset.

we expect the system to yield fact-grounded answers with correct references (see example in 5). WebGPT reports to hire a group of full-time experts to write answers for training, which is far beyond ordinary budgets.

Fortunately, LLMs' in-context learning [3, 5], which refers to their capabilities to transfer to new tasks conditioned on few in-context samples, have been demonstrated and well-explored recently. Thus we propose to bootstrap large amounts of quoted long answers via leveraging a few high-quality answers, LLMs, questions from ELI5 [7], and our retriever collected references. Additionally, since bootstrapped samples are not always satisfying, we design corresponding correction and selection strategies to filter out a high-quality subset for real training. All these efforts jointly help to create the WebGLM-QA, a quoted and long-formed QA dataset with 45k high-quality filtered and 83k unfiltered samples.

The dataset can be formulated as a set $\mathcal{D}(Q, \mathcal{A}, \mathcal{R}, C)$, where $Q$, $\mathcal{A}$, $\mathcal{R}$ represents the question set, the answer set, and the reference set respectively, $C \subseteq Q \times \mathcal{A} \times 2^{\mathcal{R}}$ denotes the triple set of (question, answer, valid references).

Different from free text generation, in web-enhanced QA each answer $\alpha \in \mathcal{A}$ contains quotations and thus is in the form of

$$\alpha = (< s_1, \nabla_1 >, < s_2, \nabla_2 >, \cdots, < s_n, \nabla_n >) \qquad (1)$$

where $< s_k, \nabla_k >$ represents the *k-th* segment in answer $\alpha$, $s_k$ is a piece of quoted text, and $\nabla_k \subset \mathcal{R}$ is a set of references that $s_k$ cites.

### 3.2.1 In-context Learning Inference

We adopt a subset of questions from ELI5 train set as our $Q$ and leverage a vanilla Contriever [10] (without LLM augmentation yet) in fine-grained retrieval to produce references $\mathcal{R}$. In this work we first try on OpenAI `text-davinci-003` API to conduct 1-shot in-context learning inference to generate quoted long-formed answers (while other LLMs such as GLM-130B [36] could be good options too). Since the in-context learning can be volatile to input forms and prompts, we take many trails to finally determine the best bootstrapping strategies as follows:

**Prompt Formulation.** Since we input many contents to the API, including a few of demonstrations (i.e., high-quality samples ($q_d$, $\alpha_d$, $\mathcal{R}_d$)), the question, and the corresponding references, their formulation could impact the performance significantly. We compare several types of prompts, including the order between question and its references (i.e., before or after, Cf. Figure 5 (a)), the symbols used to mark the indices of references, and the prompt words of references and questions. We conduct experiments with every type of prompt we have mentioned, and finally find a natural way as shown in Figure 5 (a) performs best.

**Instruction Inducting.** Next, we need a proper instruction (e.g., "Please write a answer based on the question and references.") for guiding the LLM to generate a qualified answer. Recent work [9] suggests that we can take advantage of the LLM itself to design instructions for in-context learning instead of human handcrafting. We use several high-quality examples to induce a few possible instructions (Cf. Figure 5 (b)), and select the best-performed one based on our empirical evaluation over several queries.

**Few-shot In-Context Learning.** We study the best shots needed for generating good quoted long-formed answers. Because the reference parts often occupies much of sequence length, we notice that one-shot learning can surpass few-shot learning in terms of answer's quality in most time. Hence we finally choose to inference with one-shot demonstration sample as shown in Figure 5 (c), and finally 83k various queries and their answers have been collected.

We record the details of choosing prompts and instructions in Appendix C.

### 3.2.2 Citation Correction

We have produced a large amount of well-written quoted long-formed answers using GPT-3 in-context learning. However, in our examination, we observe that the answers sometimes cite the wrong or invalid (i.e., nonexistent) references in their citation numbers. As a result, to correct the citation relationships are crucial for the quality of WebGLM-QA dataset.

Despite the fact that the citation numbers can be wrong, the contents quoted in the answer are often correct. Thus we propose to amend the citation number according to the quotation similarity to references, by splitting an answer into few segments by generated citation numbers and match then to references. For a question $q$, our retrieved references are defined as $\mathcal{R}$ and our answer can be defined as $\alpha$. We define text segments $\mathcal{S} = \{s_1, s_2, \cdots, s_n\}$, and for each pair $(s, \nabla) \in \mathcal{S} \times \mathcal{R}$, we compute citation match scores $f(s, r)$ for $r \in \mathcal{R}$. We pick a threshold $T$, and the final citation $r$ for each segment $(s, \nabla) \in \alpha$ can be described as:

$$\nabla_i = \{r | f(s_i, r) \geq T\}, r \in \mathcal{R}$$

For our application, we finally adopt Rouge-1 score as the $f$ and the threshold $T$ selection is introduced in the Section 3.2.3.

### 3.2.3 Filtering

After correction, we further investigate more issues that could potentially influence the dataset quality. And in short, we discover that most of them are related or could be solved via checking the citation quality. We will discard a piece of generated sample if it presents any problems in the following:

- **Hallucination [13]**: the answer leverages the internal knowledge of LLMs instead of references, which is not factual-grounded and sometimes severely wrong. It can be identified via the low overlapping ratio between all references and the answer.
- **Few citations**: when an answer cites too few of the provided references, it usually presents poor reference relevance and thus often not informative and factual-grounded enough.
- **Low citation accuracy**: if an answer have too many wrong citation numbers, we assume it as a low-quality one.

We calculate the F1 for the similarity and overlapping calculation. We test Rouge-L (whose best threshold is 0.4) and Rouge-1 (whose best one is 0.57) on a set of manually checked samples, and find that Rouge-1 is better. It is due to the fact that LLMs would often rewrite and paraphrase the reference contents including exchanging phrase orders. In that case, a high-quality answer may hold a high informative Rouge-1 score, but a low Rouge-L score, which computes the longest common subsequence co-occurrence.

After all the filtering conditions mentioned above, the number of samples drops from 83k to 45k, which becomes a high quality quoted long-formed QA dataset for web-hanced QA system training. We train the GLM [6], a type of bidirectional LM that is pre-trained on autoregressive blanking infilling (including a 10-billion-parameter and a 2-billion-parameter one), over the WebGLM-QA as our backbone generator.

## 3.3 Human Preference-aware Scorer

In preliminary testing, our bootstrapped generator under beam-search decoding strategy already performs satisfyingly in many cases. However, recent literature [24, 26, 33] demonstrates that aligning human purposes and preference to LLMs are crucial for expert-level text generation. WebGPT reports to recruit many experts to provide comparison and ranking over generated answers and make use of the feedback to train a reward model (RM) for picking best-of-n (i.e., 16/32/64) generated candidates and additionally optimize the generator via reinforcement learning (RL).

Nevertheless, such expert annotations could be expensive to acquire and the RL would consume much computation resource. In this work, as a competitive substitute, we propose to build a human preference-aware scorer based on massive user feedback (e.g., thumb-ups) from online QA forums. Under appropriate designs and elaborate data cleaning, we show in our experiments that such scorer also significantly improve the alignment-level of answers and the scoring in real human evaluation.

**Data collection and preprocessing.** We first collect QA pairs and corresponding user thumb-ups from online QA forums. Despite their diversity, these answers are of so various lengths and qualities that the scorer would learn little from them without proper preprocessing.

Our preprocessing includes the following requirements:

- **High quality feedback**: we define the answer with more than 3 thumb-ups as an answer with valid feedback. We pick out questions with 8 or more valid answers as qualified ones.
- **Length-bias mitigation**: we notice that the score prefers longer answers rather than the better ones in preliminary study, as is also indicated in literature [26, 33]. To mitigate the bias, for each qualified question, we use the median length $x$ of all the answers as the threshold to truncate longer answers and discard those lengths are less than $x/2$.
- **Contrast augmentation**: after sorting the answers by their thumb-ups, the gaps between neighboring answers turn out narrow. Scorers trained on such uninformative dataset present poor performance. To increase the contrast between answers for comparison training, we select a pair of answers of more than 5 in rank positions. In each pair, the answer with greater amount of likes is the better response.

After our prepossessing, there are 93k questions and 249k comparison pairs in total, with 230k pairs as the training set and 19k pairs as the test set. Next, we introduce the implementation details for training our human preference-scorer. The backbone model for training scorer is a 6-billion-parameter GLM.

**Supervised fine-tuning (SFT).** In SFT step, we leverage the Reddit TL; DR dataset for first fine-tuning the scorer following [33]. We train 16 epochs with cosine learning rate decay and 2.83e-5 as beginning learning rate. We use the SFT model for initialization of comparison training.

**Comparison training.** We pass pairs of comparison data to the model to yield a scalar score for each of the question-answer pair and maximize the gap between their scores. We use a linear head with the input dimension of hidden size and the output dimension of 1 to produce the score.

During the training, we find that the scorer tends to overfit quickly. Therefore, we freeze first 70% transformer layers and leverage other techniques such as dropouts and large batch size for regularization. Notwithstanding, the scorer would overfit after 1-1.5 epochs anyway. After the training completes, we calibrate its predictions to standard normal distribution based on the training set reward distribution.

# 4    Human Evaluation Criteria

Automatic metrics to score model-generated answers can perform well in terms of short-formed ones. However, for open-domain long-formed QA with references, the answers and rationales can be subjective and versatile, especially for those questions that start with "HOW" and "WHY." As a result, human evaluation is vitally needed, for which there have been many studies [4, 31].

To evaluate WebGLM and appropriately compare it to other similar models, we introduce a human evaluation criteria system to evaluate both references and answers. We adopt both binary (for those objective metrics, e.g., truthfulness) and four-level score (for those subjective metrics, e.g., fluency) balancing objectivity and scale in human evaluation. The four-level score is applied as is suggested in the literature that it avoid human annotators to keep absolutely neutral [31]. For each criterion we mention below, an arrow follows. up arrow (↑) means higher score performs better, while down arrow (↓) denotes lower score performs better.

## 4.1    Reference Evaluation

In this section, we introduce human evaluation criteria on references. The evaluation is done on per question-reference pair.

**Relevancy ([0, 3], ↑).** For retrieved documents or references related to a question, the more related, the higher relevancy score should be. Specifically, different references to a question can share high relevancy scores simultaneously.

**Density ([0, 3], ↑).** To evaluate how much useful information is in a piece of reference, we need to estimate its information density.

Both relevancy and density are criteria to evaluate informativeness, but there is difference between them. Relevancy can be regarded as a "recall metric" for informativeness, while density can be regarded as a "precision metric".

**Truthfulness ([0, 1], ↑).** Retrieved references can be factually wrong even they are closely associated to the question. It is because the web information sources are open and could contain user-submitted information without correctness check. As a result, the truthfulness of a piece of reference should be evaluated, and its evaluation does not consider the question.

**Toxicity ([0, 1], ↓).** Web texts could involve violent, pornographic, offensive words or other improper elements. Thus, it is necessary to assess toxicity of references retrieved.

**Social Bias ([0, 1], ↓).** Potential biases on the internet could related to genders, races, nations, and ages. We should also exclude them from our system.

## 4.2    Answer Evaluation

In this section, we introduce human evaluation criteria on answers, which are evaluated triple-wise (i.e., (question, answer, references)).

**Fluency ([0, 3], ↑).** Fluency measures the quality of generated text itself only, without taking questions and references into account [4]. It concerns only elements such as grammar, word, and phrase choices that are affiliated to the language aspect.

**Correctness ([0, 3], ↑).** Correctness measures the coherence of the answer and its corresponding question. If an answer solves

the question satisfyingly, we say it holds a high correctness. Additionally, when we score the correctness of an answer, we should take factual consistency into account. For example, contradicting common sense or defying logic will decrease the correctness.

**Citation Accuracy ([0, 3], ↑).** The metric only considers the relationships between an answer and its references. When an answer contains citation marks, we should check if it is correct. Citation mistakes or missing citation will both decrease the accuracy.

**Truthfulness ([0, 1], ↑).** Similar to truthfulness in the reference evaluation, truthfulness of an answer measures whether the text of the answer is factually sound, including the factual consistency of the answer and whether the answer contains contradictions or hallucinate information.

**Objectivity ([0, 1], ↑).** The metric only concerns the relationships between an answer and its references. When references provided, models are supposed to generate answers according to these references without its using its latent knowledge from pre-training. If we can find all the information of an answer from provided references, we say it is objective.

**Redundancy ([0, 1], ↓).** Within the limited text length, duplicate content will reduce informativeness. As the lower redundancy, the higher quality of the answer, we take it into our consideration.

The detail of the metrics and the meaning of the score can be found in the Appendix H.

# 5    Experiment

In this section, we conduct experiments employing the metrics mentioned in Section 4 to evaluate and analyze the quality of the responses generated, including those from WebGLM and other similar systems. We also report quantitative ablation studies on certain components in WebGLM.

## 5.1    Main Results

We conduct the major evaluation using the 272 questions provided on WebGPT [24] demo website[2], as the WebGPT is not publicly available and selected questions are generally complicated and closer enough to real human questions.

**Human Evaluation Setup.** We recruited 15 master-degree level experts to conduct human evaluation. For each question, we aggregate all the search results and answers from different models into one table, enabling the annotators to effectively compare them and unify the annotation standards. We evaluate the performance of our model and other different models from various dimensions through human evaluation. We also compare and analyze the results from different perspectives as follows. The main results are shown in Table 2.

**WebGLM Reference vs Other References.** Although the search results of WebGLM are slightly inferior to WebGPT-175B, its performance is far better than that of Perplexity.ai and WebGPT-13B. It is worth mentioning that the WebGLM retrieval process only uses some traditional, word-based algorithms and two Contrievers with a cumulative parameter amount of no more than 300M. WebGLM is significantly superior to WebGPT in computing performance and

---

[2]https://openaipublic.blob.core.windows.net/webgpt-answer-viewer/index.html

**Table 2: Main results based on human evaluation metrics.** Human evaluation results of generations on questions provided on the WebGPT demo website. For reference evaluation, Rel., Den., Tru., Tox↓., and Soc. Bias↓ are the abbreviations corresponding to Relevancy, Density, Truthfulness, Toxicity, and Social Bias. For answer evaluation, Flu., Cor., Cit. Acc., Obj., Tru., Red.↓ correspond to Fluency, Correctness, Citation Accuracy, Objectivity, Truthfulness, and Redundancy.

| Model | Reference Evaluation | | | | | Answer Evaluation | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Rel. | Den. | Tru. | Tox.↓ | Soc. Bias↓ | Flu. | Cor. | Cit. Acc. | Obj. | Tru. | Red.↓ |
| WebGPT (175B) | 2.512 | 2.660 | 0.996 | 0.015 | 0.006 | 2.457 | 2.889 | 2.837 | 0.990 | 0.975 | 0.087 |
| Perplexity.ai | 1.652 | 1.636 | 0.955 | <u>0.005</u> | **0.001** | <u>2.718</u> | <u>2.321</u> | 2.512 | 0.726 | <u>0.975</u> | <u>0.032</u> |
| WebGPT (13B) | <u>1.782</u> | <u>1.766</u> | **0.998** | 0.008 | 0.016 | 2.692 | 2.102 | **2.769** | **0.974** | 0.872 | 0.051 |
| WebGLM (10B) | **1.980** | **2.226** | <u>0.983</u> | **0.002** | <u>0.002</u> | **2.829** | **2.810** | <u>2.757</u> | <u>0.943</u> | **0.998** | **0.021** |

time consumption. Its performance is far superior to that of the 13B model and close to that of the 175B model.

**WebGLM vs Other Systems.** Finally, we compare our system with the results of WebGPT-13B, Perplexity.ai, and WebGPT-175B. Our system has achieved the highest performance in fluency, truthfulness, and redundancy. At the same time, we are close to WebGPT-175B in the correctness metric with a score of 2.81, which is far higher than that of Perplexity.ai and WebGPT-13B, indicating that our system can still achieve superior performance at a lower cost.

## 5.2 Turing Test

To further compare our performance, we design a Turing test [21] to check the answers' quality.

**Setup.** We randomly sampled 200 items from the 272 questions that WebGPT has displayed on their official web page. For each question, we shuffle the answers generated by WebGLM, WebGPT-175B, WebGPT-13B, and Perplexity.ai, and remove citation marks from them for fairness. We next mix an answer written by humans into these answers and ask evaluators to rank the answers by their quality, such as correctness, informativeness, and truthfulness.
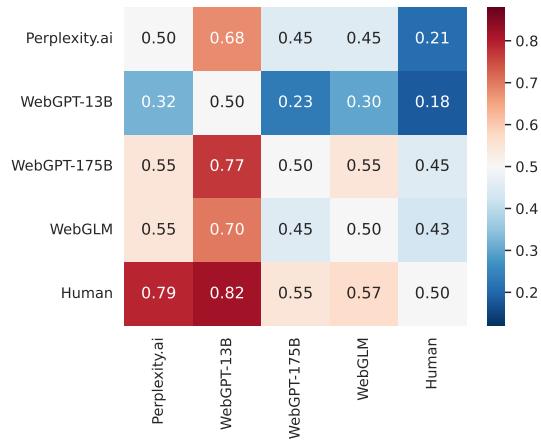


**Figure 6: Win rates between systems.** Numbers denote the rate that the answers from corresponding source from the first column are better than ones from corresponding source from the first row.

**Result.** For each pair of answers $(A, B)$, if evaluators prefer $A$ to $B$, we call $A$ wins and $B$ loses. Firstly, we compare each pair of the answers, the win rate is shown in Figure 6. Besides, We calculate the win rates against humans for each system. The result is shown

in Figure 2. We hold a 43% win rate, definitely beat Perplexity.ai with a 21% win rate and WebGPT-13B with an 18% win rate, and almost draw with WebGPT-175B with a 45% win rate.

## 5.3 Test on QA Benchmarks

We randomly sample 400 questions on Natural Question and Web Question, and evaluate WebGLM and Perplexity.ai on them. The results in Table 3 show that WebGLM outperform Perplexity.ai.

**Table 3: Open QA Performance on NaturalQuestions and WebQuestions.** Perplexity.ai is evaluated on sampled subsets because the website prohibits crawling.

| | Natural Questions | Web Questions |
|---|---|---|
| WebGLM | **60.8** | **63.5** |
| Perplexity.ai (sample) | 57.3 | 57.5 |
| GPT3-175B | 29.9 | 41.5 |

In addition, we conducted experiments on the full validation split of TriviaQA (same as WebGPT). Following the testing method employed by WebGPT, we first generated a long answer for each question using WebGLM. We then used Google Bigbird, fine-tuned on the TriviaQA training set[3], to answer TriviaQA questions based on the output of WebGLM. To address potential test-train overlap issues mentioned in WebGPT, we also conducted TriviaQA tests on different train-test splits. The results are summarized in Table 4.

## 5.4 Ablation Study

In this section, we study the major improvements and strategies in WebGLM, including the bootstrapped dataset filtering, scorer training, LLM-augmented retriever and some other minor topics.

*5.4.1 WebGLM-QA Filtering* Since we build our training dataset based on LLM in-context bootstrapping, the dataset quality could be essential for WebGLM's success. We randomly sample 210 examples from these versions of our dataset to verify the filtering strategies they are based on, including 1) None, 2) Rouge-L filtered, and 3) Rouge-1 filtered.

We randomly shuffle all the samples and distribute them to evaluators, and then collect and calculate the average score of each metric. The sample results are shown in Table 5 We analyze this result from two perspectives. One is the absolute performance of our

---

[3]https://huggingface.co/google/bigbird-base-trivia-itc

**Table 4: WebGLM, WebGPT and other comparison methods on TriviaQA.** The setting follows WebGPT [24] Appendix G.

| Method | Total | Question overlap | No question overlap | Answer overlap | Answer overlap only | No overlap |
|---|---|---|---|---|---|---|
| Bigbird + WebGLM (Ours) | **70.80%** | 86.40% | **67.10%** | **78.70%** | **73.60%** | 49.30% |
| GPT-3 175B | 58.70% | 75.90% | 52.90% | 67.30% | 61.60% | 39.00% |
| GPT-3 175B + WebGPT 175B BC | 69.50% | 86.30% | 65.30% | 78.40% | 73.20% | **52.40%** |
| UnitedQA-E | 68.90% | **89.30%** | 62.70% | 78.60% | 70.60% | 44.30% |
| UnitedQA (hybrid model) | 70.50% | - | - | - | - | - |

**Table 5: Ablation study on different dataset filtering strategies in creating the bootstrapped generator.**

| Filtering Method | Reference Evaluation | | | | | Answer Evaluation | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Rel. | Den. | Tru. | Tox.↓ | Soc. Bias↓ | Flu. | Cor. | Cit. Acc. | Tru. | Obj. | Red.↓ |
| None | 1.711 | 1.619 | 0.991 | 0.011 | 0.011 | **2.872** | 2.636 | 2.370 | 2.810 | 0.805 | 0.134 |
| Rouge-L | **1.833** | 1.728 | **0.994** | 0.022 | **0.010** | 2.731 | 2.680 | 2.573 | 2.896 | 0.841 | 0.181 |
| Rouge-1 | 1.832 | **1.751** | 0.993 | **0.010** | 0.012 | 2.826 | **2.694** | **2.688** | **2.919** | **0.890** | **0.120** |

**Table 6: Ablation study on different dataset filtering strategies, based on GLM-2B's post-training evaluation**

| | Flu. | Cor. | Cit. Acc. | Obj. | Tru. | Red.↓ |
|---|---|---|---|---|---|---|
| None | 2.610 | 2.738 | 2.655 | 0.961 | 0.961 | 0.063 |
| Rouge-L | 2.604 | **2.742** | 2.727 | 0.952 | **0.975** | **0.034** |
| Rouge-1 | **2.852** | 2.738 | **2.743** | **0.976** | 0.970 | 0.044 |

**Table 7: Performance of LLM-augmented Retriever (Ours).** "N-NDCG" refers to Normalized NDCG.

| Metric(%) | TF-IDF | BM25 | Contriever | Ours |
|---|---|---|---|---|
| Accuracy | 46.85 | 40.33 | 18.54 | **69.36** |
| Spearman | 9.92 | -20.94 | -1.58 | **62.26** |
| NDCG | 82.54 | 76.28 | 81.16 | **91.99** |
| N-NDCG | 46.05 | 26.77 | 41.75 | **75.29** |

final version of the dataset. The other is comparing the performance of our different versions of datasets.

We find that our dataset holds a high factual consistency and correctness, and the majority of our data are judged as perfectly correct. We have also noticed that the information relevancy and density are considerably improved when we apply a filter method and when we change Rouge-L to Rouge-1. As for the answer, we find that correctness has great improvement when we apply any one of the two filters, and factual consistency has a great improvement when we change the Rouge-L filter to Rouge-1. Besides, objectivity is also one of the most important criteria that we care about, and we find that it's more likely to discard subjective answers with a Rouge-1 filter than with a Rouge-L filter. As a result, our experiments show that citation accuracy is closely related to the reference quality and answer quality, so our filter method is effective.

Besides, We train the GLM-2B models on each dataset and evaluate them with our designed metrics to see the impact of these datasets on our model's performance. We show the results in Table 6. We find that the answers of the three models showed little difference in the correctness metric. However, the performance of the model trained by rouge-1 was better in fluency, citation accuracy, and objectivity metrics. This result further proves the advantages of the dataset of rouge-1. Therefore, we decide to train our 10B model on the dataset of rouge-1.

*5.4.2 LLM-augmented Retriever* In terms of the usefulness of references, we have compared our method with traditional methods such as BM25, TF-IDF, and the original version of Contriver.

We collect 22000 examples from WebGLM-QA, and for each question, we calculate Rouge-1 precision score $p$ of corresponding answer $a$ and each of the reference $r$, and then label the reference-answer pair $(r, a)$ as $p$. Finally, we gain a training dataset containing 20000 examples and a test dataset containing 2000 examples.

For all answers to the same question, we compare the order predicted by retrieve methods with the answer relevancy order. The results are shown in Table 7. We notice that before the LLM task augmentation, the Contriever performs even poorer than traditional lexical-based approaches. After augmenting knowledge from GPT-3's reference adoption labeling, we find that ours, which holds a 69.36 pair-wise choosing accuracy and 62.26 spearman index, performs best. The evidence strongly advocates that the LLM augmentation is vital when we use pre-trained smaller dense retrievers in practice.

*5.4.3 Human Preference-aware Scorer* In this section we compare several different scorer training strategies and datasets. We discover that proper task formulation and larger and more diverse dataset yield better results.

**Baseline and data preprocessing.** We first train RoBERTa-large in the classification task and the regression task formulation, and the 6-billion-parameter GLM on the ELI5's training set (with thumb-ups) as our baselines. In the classification task, we collect all items whose count of answers is not less than 10 from ELI5. For each collected question, we label top-5-voted answers as positive, and randomly pick 5 answers from other questions as negative examples.
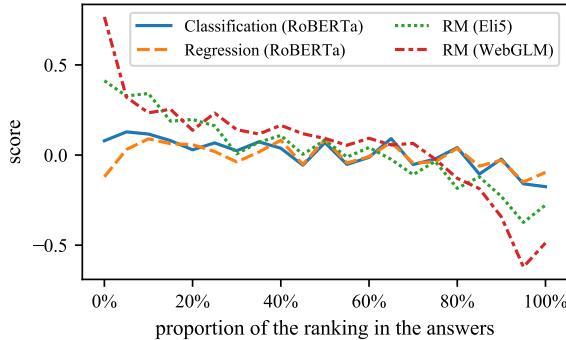
**Figure 7: Average score of answers in ELI5 test set.** It is sorted by likes in the ELI5 test set. The best answer is around 0% and the worst answer is around 100%.

In the regression task, we collect all items whose count of answers is not less than 5 from ELI5. For each collected question, we complete the following steps:

(1) for each answer to this question, supposing its corresponding up-vote is $u$, we firstly label this answer as $\log_2 (u + 1)$.

(2) Then, we scale labels of all answers to this question to $[0, 1]$.

(3) Let $x$ be the summation of the answers' label, we randomly pick $\lfloor x \rfloor$ answers from other questions as negative examples with label $-1$.

In order to obtain a large train set (which has been suggested very important in [33]), we adopt a relatively loose screening method, which selects the questions with more than 5 answers and answers with no less than 100 words in length. Our large train set includes 28.2k questions and 191.6k pairs. We use the ELI5 test set with thumb-ups for our final evaluations.

**Metrics.** We select three metrics to measure the ability of the reward model to distinguish responses of different quality, namely accuracy, Spearman coefficient, and NDCG (Normalized Discounted Cumulative Gain). Accuracy refers to the accuracy of selecting better answers in pairs. Spearman and NDCG measure the sorting ability of the model.

The ranking evaluation of different models is shown in Table 9. We find that WebGLM human preference-aware scorer performs best on accuracy and Spearman coefficient. Under the same amount of training tokens, the performance of the reward model is slightly worse than that of RoBERTa classification and RoBERTa regression, but after increasing the amount of training, the performance of the reward model will increase significantly.

Figure 7 shows the average reward of the answers at different positions in the sequence sorted by likes in the ELI5 test set. The best answer is around 0% and the worst answer is around 100%. We find that the curve of the WebGLM Human Preference-aware Scorer is more discriminative than other models, and the rewards of the best answer are higher than that of others.

*5.4.4 Ablation Study on Each Component* We added some experiments to conduct ablation studies on each component. We compared the three sub-modules of the system: Retriever, Generator, and Scorer. The results are shown in Table 8.

In the Retriever module, we compared the performance on the settings of WebGPT-175B, WebGLM, and non-retrieval. From the Table 8, the performance on WebGLM retrieval is similar to that of WebGPT-175B and significantly better than non-retrieval.

Regarding the Generator module, we compared the response quality of WebGLM and GPT-3 on WebGLM retrieval setting. We found that WebGLM performed slightly better than GPT-3 in fluency, correctness, accuracy, citation accuracy, objectivity, and truthfulness.

In terms of Scorer, we compared the response quality of WebGLM removing and retaining Reward Models. The results show that by WebGLM-10B top-p sampling and reward model scoring method, We found through the human evaluation results that the answers scored high by the reward model excel the original results in fluency, correctness, citation accuracy, truthfulness, and redundancy. It shows the importance of the reward model scoring mechanism to model performance.

## 6 Conclusion

We build the LLM-based question-answering system—WebGLM—with a web retrieval method. We propose a fast and cost-effective method to retrieve valuable information from the Internet. We leverage GPT-3's in-context learning ability to build a LLM-bootstrapped quoted and long-form QA dataset, which is used to train our model. Further, we train a human preference-aware scorer and use it to give marks to responses generated by our model. For each question, the scorer can select the highest-scored response from candidates, thus obtaining a final answer humans prefer the most. We conduct extensive experiments, including both the human evaluation and the Turing test, to demonstrate the competitive performance of WebGLM with some of the pioneering web-enhanced question answering systems like Perplexity.ai and WebGPT.

## References

[1] Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. 65–72.

[2] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in natural language processing*. 1533–1544.

[3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[4] Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2020. Evaluation of text generation: A survey. *arXiv preprint arXiv:2006.14799* (2020).

[5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311* (2022).

[6] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 320–335.

**Table 8: Ablation study on different sub-modules (Scorer, Retriever, and Generator) in WebGLM.**

| Method | Flu. | Cor. | Cit. Acc. | Obj. | Tru. | Red.↓ |
|---|---|---|---|---|---|---|
| Scorer Ablation | | | | | | |
| No Scorer | 2.797 | 2.757 | 2.723 | 0.961 | 0.970 | 0.039 |
| Human Preference-aware Scorer (Ours) | 2.829 | 2.810 | 2.757 | 0.943 | 0.998 | 0.021 |
| Retriever Ablation (w.o. RM) | | | | | | |
| No Retriever | 2.364 | 1.982 | - | - | 0.645 | 0.091 |
| WebGPT Retriever | 2.750 | 2.884 | 2.808 | 0.981 | 0.980 | 0.038 |
| Contriever | 2.761 | 2.732 | 2.721 | 0.963 | 0.930 | 0.043 |
| LLM-augmented Retriever (Ours) | 2.797 | 2.757 | 2.723 | 0.961 | 0.970 | 0.039 |
| Generator Ablation (w.o. RM) | | | | | | |
| GPT-3 (text-davinci-003, zero-shot) | 2.751 | 2.752 | 2.607 | 0.927 | 0.966 | 0.034 |
| Bootstrapped Generator (Ours) | 2.797 | 2.757 | 2.723 | 0.961 | 0.970 | 0.039 |
| WebGLM (Ours) | 2.829 | 2.810 | 2.757 | 0.943 | 0.998 | 0.021 |

**Table 9: Different scorers' performance on ELI5 test set.**

| | Accuracy | Spearman | N-NDCG |
|---|---|---|---|
| Classification (RoBERTa) | 0.552 | 0.129 | 0.319 |
| Regression (RoBERTa) | 0.569 | 0.164 | 0.352 |
| RM (ELI5) | 0.568 | 0.197 | **0.406** |
| RM (WebGLM) | **0.596** | **0.241** | 0.367 |

[7] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long Form Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 3558–3567.

[8] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*. PMLR, 3929–3938.

[9] Or Honovich, Uri Shaham, Samuel R Bowman, and Omer Levy. 2022. Instruction induction: From few examples to natural language task descriptions. *arXiv preprint arXiv:2205.10782* (2022).

[10] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised Dense Information Retrieval with Contrastive Learning. *Transactions on Machine Learning Research* (2022).

[11] Gautier Izacard and Édouard Grave. 2021. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 874–880.

[12] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299* (2022).

[13] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *Comput. Surveys* 55, 12 (2023), 1–38.

[14] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6769–6781.

[15] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics* 7 (2019), 453–466.

[16] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.

[17] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.

[18] Jiachang Liu, Dinghan Shen, Yizhe Zhang, William B Dolan, Lawrence Carin, and Weizhu Chen. 2022. What Makes Good In-Context Examples for GPT-3?. In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*. 100–114.

[19] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. 2021. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering* 35, 1 (2021), 857–876.

[20] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021. GPT understands, too. *arXiv preprint arXiv:2103.10385* (2021).

[21] Michael L Mauldin. 1994. Chatterbots, tinymuds, and the turing test: Entering the loebner prize competition. In *AAAI*, Vol. 94. 16–21.

[22] Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Noisy Channel Language Model Prompting for Few-Shot Text Classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 5316–5330.

[23] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the Role of Demonstrations: What Makes In-Context Learning Work? *arXiv preprint arXiv:2202.12837* (2022).

[24] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332* (2021).

[25] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated MAchine Reading COmprehension Dataset. *choice* 2640 (2016), 660.

[26] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.

[27] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.

[28] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. 2383–2392.

[29] Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How Much Knowledge Can You Pack Into the Parameters of a Language Model?. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 5418–5426.

[30] Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning To Retrieve Prompts for In-Context Learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2655–2671.

[31] Ananya B Sai, Akash Kumar Mohankumar, and Mitesh M Khapra. 2022. A survey of evaluation metrics used for NLG systems. *ACM Computing Surveys (CSUR)* 55, 2 (2022), 1–39.

[32] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100* (2022).

[33] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems* 33 (2020), 3008–3021.

[34] Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, et al. 2022. Selective annotation makes language models better few-shot learners. *arXiv preprint arXiv:2209.01975* (2022).

[35] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An Explanation of In-context Learning as Implicit Bayesian Inference. In *International Conference on Learning Representations*.

[36] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414* (2022).

[37] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068* (2022).

[38] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTScore: Evaluating Text Generation with BERT. In *International Conference on Learning Representations*.

[39] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*. PMLR, 12697–12706.

# A  Additional Experimental Results

## A.1  WebGLM vs Others in WebGPT Reference

We compared the generation results of WebGLM-Rouge1, WebGPT-175B, and GPT-3 on the WebGPT-175B references. For GPT-3, we also use the method of automatically constructing datasets to generate responses for the WebGPT samples to compare the effect of the WebGLM system. Specifically, we use the references of WebGPT to let GPT-3 do in-context learning to answer questions according to the search results. We use human evaluation to compare the quality of the three answers. The experimental results are shown in Table 10. Although our model size is more than ten times smaller than GPT-3 and WebGPT-175B, we can effectively compensate for the impact of the model size and achieve competitive performance in the retrieval paradigm. Our model matches WebGPT-175B and GPT-3 on correctness, citation accuracy, objectivity, and truthfulness metrics and outperforms them on fluency and redundancy.

**Table 10: Ablation study on different Generators based on WebGPT references**

| Generator | Flu. | Cor. | Cit. Acc. | bj. | Tru. | Red. |
|---|---|---|---|---|---|---|
| GPT-3 In-Context | **2.801** | 2.883 | 2.726 | 0.966 | 0.975 | **0.024** |
| WebGPT-175B | 2.457 | **2.889** | **2.837** | **0.990** | 0.975 | 0.087 |
| WebGLM-10B-Rouge1 | 2.750 | 2.884 | 2.808 | 0.981 | **0.980** | 0.038 |

# B  Detailed Efficiency Analysis

At the retrieval stage, we only search for one time, then take the first few results links to fetch the web pages in parallel. We then extract all paragraphs and rank these paragraphs by Contriever, and finally take the top 5 paragraphs as references. Let $t_s$, $t_f$, $t_e$, and $t_r$ denote the time we consume in four steps, so the total time we consume is $t_s + t_f + t_e + t_r$.

WebGPT simulates the operations in a virtual browser environment while obtaining references. For the 272 questions they showed, we count the types of actions and the average number of generated tokens as shown in Table 11 and 12. Then we calculate the average time it takes to browse. Assuming that the total time $M$ ($M$ is either WebGPT-175B or WebGPT-13B) takes to generate commands in the browsing process of each question is expected to be $t_c(M)$, the time $M$ consumes $T(M)$ satisfies the following equations.

$$T(\text{WebGPT-175B}) = t_c(\text{WebGPT-175B}) + t_s * 3.82 + t_f * 6.96 \quad (2)$$

$$T(\text{WebGPT-13B}) = t_c(\text{WebGPT-13B}) + t_s * 4.05 + t_f * 7.56 \quad (3)$$

We test the efficiency of GPT-3. With a 500-token prompt, the 175B model generates about 20 tokens per second, and the 13B model generates 100 tokens per second, meaning that:

$$t_c(\text{WebGPT-175B}) = \frac{580.08 \text{ tokens/query}}{20 \text{ tokens/second}} = 29 \text{ seconds} \quad (4)$$

$$t_c(\text{WebGPT-13B}) = \frac{580.08 \text{ tokens/query}}{100 \text{ tokens/second}} = 5.8 \text{ seconds} \quad (5)$$

In practice, $t_s$, $t_f$, $t_e$, and $t_r$ are about 1.81, 2.38, 0.29, and 0.89 respectively. So we consume 5.36 seconds for one query on average. Nevertheless, assuming in the same network circumstance, the time

consumption of WebGPT models can be calculated by Equation 2 and 3.

$$T(\text{WebGPT-175B}) = 52.48 \text{ seconds} \quad (6)$$

$$T(\text{WebGPT-13B}) = 31.12 \text{ seconds} \quad (7)$$

Therefore, WebGPT-175B costs 52.48 seconds, and WebGPT-13B costs 31.12 seconds. Our efficiency can be about 10 times that of WebGPT-175B and 6 times that of WebGPT-13B.

**Table 11: Efficiency statistics for browsing stage in WebGPT-175B. Average count per query, tokens per action, and tokens per query (the product of the first two terms) are displayed in this table.**

| action | count/query | tokens/action | tokens/query |
|---|---|---|---|
| search | 3.82 | 9.80 | 37.46 |
| click_link | 6.96 | 5.00 | 34.82 |
| quote | 3.49 | 124.49 | 434.80 |
| back | 5.35 | 1.00 | 5.35 |
| scroll_down | 11.41 | 4.00 | 45.63 |
| scroll_up | 1.62 | 4.00 | 6.49 |
| top | 0.49 | 1.00 | 0.49 |
| end | 0.43 | 3.00 | 1.29 |
| find_in_page | 0.13 | 5.11 | 0.68 |
| invalid | 0.12 | 111.09 | 13.07 |
| tokens | | 580.08 | |
| generating speed | | 20 tokens/second | |
| action time | | 29s | |
| total time | | 52s | |

**Table 12: Efficiency statistics for browsing stage in WebGPT-13B. Average count per query, tokens per action, and tokens per query (the product of the first two terms) are displayed in this table.**

| action | count/query | tokens/action | tokens/query |
|---|---|---|---|
| search | 4.05 | 9.65 | 39.08 |
| click_link | 7.56 | 5.00 | 37.81 |
| quote | 3.44 | 125.85 | 433.08 |
| back | 5.90 | 1.00 | 5.90 |
| scroll_down | 10.30 | 4.00 | 41.21 |
| scroll_up | 2.01 | 4.00 | 8.04 |
| top | 0.32 | 1.00 | 0.32 |
| end | 0.44 | 3.00 | 1.33 |
| find_in_page | 0.21 | 5.04 | 1.06 |
| invalid | 0.10 | 136.58 | 13.06 |
| tokens | | 580.89 | |
| generating speed | | 100 tokens/second | |
| action time | | 5.8s | |
| total time | | 31s | |

## C Choice of Prompts and Instructions

Firstly, we attempt the zero-shot approach for bootstrapping data. To produce data with appropriate citation marks, we require specific instructions. We experiment with several methods, however, they are all limited in their effectiveness.

`Use a mark for each helpful reference you cited, such as [1].` Limitation: bootstrapped data contain mixed usage of `[1][2]` and `[1, 2]`.

`Use a mark for each helpful reference you cited, such as [1]. If there are multiple citations at one position, please use a format like [1][2][3].` Limitation: bootstrapped data contain citations of useless references.

`Use a mark for each helpful reference you cited, such as [1]. If there are multiple citations at one position, please use a format like [1][2][3]. If a reference is useless, do not cite it.` Limitation: useless references are still cited. This method do not work.

We then select few-shot context to bootstrap data. If we provide too many references or in-context examples, it is easy to exceed the token count limit. Therefore, we choose to use an 1-shot example and 5 references. We also include some useless references in the example, which are not cited in the answer.

After that, We conduct experiments on prompts and demonstrate that placing the question after the references is the most effective approach.

Regarding instruction induction for in-context learning, we experiment with the previously mentioned examples as well as some new ones, such as:

`Answer the question based on the following references with citations. Use a mark for each helpful reference you cited, such as [1]. If there are multiple citations at one position, please use a format like [1][2][3]. If a reference is useless, do not cite it.`

`I will provide you with some references. Based on the references, please answer my question. Pay attention that you should be objective, and you should not use your knowledge. Use a mark for each helpful reference you cited, such as [1]. If there are multiple citations at one position, please use a format like [1][2][3]. If a reference is useless, do not cite it.`

However, these instructions are too verbose, and in the presence of examples, the model's performance is not significantly impacted by the instructions. Therefore, we adopt a more natural approach to generate instructions[9] to produce a natural instruction that is interpretable by the model.

Finally, we use a very concise instruction: `Read the references provided and answer the corresponding question.`

In addition, we compared models trained with different prompt strategies, and the results are shown in the Table 13. From the "Correctness" column, we can see the significant difference that the order of references and question in the prompt makes.

## D Dataset Examples

An example of **WebGLM-QA** is shown in Table 14.

**Table 13: The performance with training data bootstrapped by difference prompt strategies.**

| Prompt | Flu. | Cor. | Cit. Acc. | Obj. | Tru. | Red. |
|---|---|---|---|---|---|---|
| WebGLM Prompt | 2.797 | 2.757 | 2.723 | 0.961 | 0.970 | 0.039 |
| Question before Reference | 2.633 | 2.518 | 2.700 | 0.933 | 0.970 | 0.058 |
| 3-Reference | 2.658 | 2.412 | 2.819 | 0.933 | 0.930 | 0.065 |

## E Retrieval Example

An example of retrieved references from each system is shown in Table 15 and Table 16.

## F Answer Examples

Some examples consisting of only answers are shown in Table 17 and Table 18. We remove citation marks for all the answers to evaluate answers only. In this example, WebGLM's answer is consistent with the question, while the answers of two WebGPT models are beside the point in the conclusion.

## G Reward Model Example

WebGLM performs better after implementing the reward model. An example is shown in Table 19.

## H Criteria Details

The criteria details of human evaluation are shown in Table 20 and Table 21.

**Figure 8: WebGLM web demo page**



## I Demo

We provide a web demo as shown in Figure 8. Meanwhile, some real examples are shown in Figure 9 and the later figures.

**Table 14: Dataset Example. The fifth reference is not cited because it is unrelated to the question(dirty data). In the training procedure, we need these unrelated references to train our model to pick related information.**

| Question | Why were the cities that are capital cities chosen as capitals? |
| --- | --- |
| References | [1] Reasons differ from state to state as to why the capital city was chosen. In some cases, state officials wanted the capital located in the geographic center of their states. Jefferson City and Pierre were chosen as capitals of Missouri and South Dakota respectively for this reason. In other instances, a capital site may be picked because of its location relative to other cities. Florida chose the city of Tallahassee to be its capital in 1824 because it was situated halfway between St. Augustine and Pensacola, the state's two largest cities at the time. Some capital locations were picked because they were centers of trade and transportation. This is why Albany, located near the Hudson River, was picked as New York's capital.<br><br>[2] - The selection or founding of a "neutral" capital city, one unencumbered by regional or political identities, was meant to represent the unity of a new state when Ankara, Bern, Brasília, Canberra, Madrid, Ottawa and Washington became capital cities. Sometimes, the location of a new capital city was chosen to terminate actual or potential squabbling between various entities, such as in the cases of Brasília, Canberra, Ottawa, Washington, Wellington and Managua.<br><br>[3] - However a country decides on its capital, that city becomes an important symbol of nationhood. While being a home to its inhabitants and a place for tourists to visit, it is also the city chosen to represent the entire country to the world.<br><br>[4] State capitals are important because they house the governments of their states, with a meeting place for the state legislature. Many meeting buildings are called capitols, which are distinguished from capitals in that capitol refers to buildings while capital always references a city. Cities are chosen to be the capital for differing reasons, sometimes because it was located in the center of the state, like Jefferson City, or was a center of commerce, like Albany. Today, many capitals are not the largest city in their states and can be overshadowed by bigger cities, but retain important historical value.<br><br>[5] (Not Cited) 8. "Capital cities: How are they chosen and what do they represent?". BBC News. 6 December 2017. Archived from the original on 23 March 2022. Retrieved 26 June 2022. |
| Answer | The reasons why cities are chosen as capitals vary from state to state. In some cases, state officials wanted the capital located in the geographic center of their state[1][4], while in other instances, a capital site may be picked because of its location relative to other cities[1]. Some capital cities were picked because they were centers of trade and transportation[1], or to represent the unity of a new state[1][2][4]. In some cases, the location of a new capital city was chosen to terminate actual or potential squabbling between various entities[1][2]. Cities are chosen to be the capital because they house the governments of their states, with a meeting place for the state legislature[1][4]. Today, many capitals are not the largest city in their states and can be overshadowed by bigger cities, but retain important historical value[4]. The capital city also represents the entire country to the world[3]. |

**Table 15: Retrieved Reference Example (WebGLM and Perplexity AI)**

| Question | What is the difference between sigma bonds and pi bonds? |
| --- | --- |
| References from WebGLM | [1] Sigma bonds are often formed by the combination of s orbitals in different atoms whereas pi bonds are formed from the combination of p and similar orbitals in different atoms. |
| | [2] The main differences between sigma and pi bonds has to do with their formation and strength. Sigma bonds are formed by direct head to head overlap between orbitals while pi bonds are formed by side to side overlap, usually between p orbitals. This difference in formation leads to a difference in strength. Sigma bonds are stronger than pi bonds because the direct head to head overlap provides a bigger (and therefore stronger) overlap than the side to side overlap of pi bonds. Additionally, sigma bonds form single bonds and can exist without a pi bond present; however, a sigma bond must already be formed for a pi bond to form. |
| | [3] Sigma and pi bonds are two types of covalent bonds formed by the overlapping of atomic orbitals. Sigma bonds form by the direct head to head overlap of atomic orbitals and can occur between s-s, p-p and s-p orbitals. Pi bonds form by the side to side overlap of p orbitals. |
| | [4] The sigma bond is stronger than the Pi bond, according to quantum mechanics, since the orbital paths of Pi bonds are parallel, resulting in less overlap between the p-orbitals. Two atomic orbitals come into contact with each other through two areas of overlap, forming Pi bonds. |
| References from Perplexity AI | [1] Sigma bond is a chemical bond formed by the linear or co-axial overlapping of the atomic orbitals of two atoms. A pi bond is a type of covalent bond that exists between atoms where the electrons are on top and bottom of the axis connecting the nuclei of the joined atoms. |
| | [2] The sigma bond is stronger than the Pi bond, according to quantum mechanics, since the orbital paths of Pi bonds are parallel, resulting in less overlap between ... |
| | [3] A $\sigma$ bond has cylindrical symmetry; a $\pi$ bond has a nodal plane that includes the bond axis. ... A $\sigma$ bond comes from the end-to-end overlap of the ... |

**Table 16: Retrieved Reference Example (WebGPTs)**

| Question | What is the difference between sigma bonds and pi bonds? |
| --- | --- |
| References from WebGPT-175B | [1] Pi bonds are formed by the sidewise positive (same phase) overlap of atomic orbitals along a direction perpendicular to the internuclear axis. During the formation of $\pi$ bonds, the axes of the atomic orbitals are parallel to each other whereas the overlapping is perpendicular to the internuclear axis. This type of covalent bonding is illustrated below. [Image: Pi Bonds] Pi Bonds are generally weaker than sigma bonds, owing to the significantly lower degree of overlapping. Generally, double bonds consist of one sigma and one pi bond, whereas a typical triple bond is made up of two $\pi$ bonds and one $\sigma$ bond. It is important to note that a combination of sigma and pi bonds is always stronger than a single sigma bond. |
| | [2] In chemistry, sigma bonds ($\sigma$ bonds) are the strongest type of covalent chemical bond. They are formed by head-on overlapping between atomic orbitals. Sigma bonding is most simply defined for diatomic molecules using the language and tools of symmetry groups. In this formal approach, a $\sigma$-bond is symmetrical with respect to rotation about the bond axis. By this definition, common forms of sigma bonds are $s + s, p_z + p_z, s + p_z and d z^2 + d z^2$ (where z is defined as the axis of the bond or the internuclear axis). Quantum theory also indicates that molecular orbitals (MO) of identical symmetry actually mix or hybridize. As a practical consequence of this mixing of diatomic molecules, the wavefunctions $s + s and p_z + p_z$ molecular orbitals become blended. The extent of this mixing (or hybridization or blending) depends on the relative energies of the MOs of like symmetry. |
| | [3] A sigma bond is stronger than a pi bond. The reason is that the overlapping of atomic orbitals can take place to a greater extent during the formation of a sigma bond, whereas overlapping of orbitals occurs to a smaller extent during the formation of a pi bond. A pi bond between two atoms is formed only in addition to a sigma bond. The reason is that the atoms constituting a single bond prefer to form a strong sigma bond rather than a weak pi bond. Thus, a pi bond is always present in molecules with multiple bonds, i.e., double or triple bonds. In other words, a single bond cannot be a pi bond. There can be free rotation of atoms around the sigma bonds. Free rotation of atoms around pi bonds is not possible because it involves breaking the pi bonds. |
| References from WebGPT-13B | [1] What is the difference between Sigma Bond and pi Bond? • Sigma bonds are formed by head to head overlapping of orbitals, whereas pi bonds are formed by the lateral overlapping. • Sigma bonds are stronger than pi bonds. • Sigma bonds can be formed between both s and p orbitals whereas pi bonds are mostly formed between p and d orbitals. • The single covalent bonds between atoms are sigma bonds. When there are multiple bonds between atoms, pi bonds can be seen. • pi bonds result in unsaturated molecules. • Sigma bonds allow free rotation of atoms whereas pi bonds restrict the free rotation. |
| | [2] A sigma bond ($\sigma$ bond) is a bond formed by the overlap of orbitals in an end-to-end fashion, with the electron density concentrated between the nuclei of the bonding atoms. A pi bond ($\pi$ bond) is a bond formed by the overlap of orbitals in a side-by-side fashion with the electron density concentrated above and below the plane of the nuclei of the bonding atoms. |
| | [3] In a sigma bond, the electrons are in orbitals between the nuclei of the bonding atoms (in the plane of electron density), where electron density is greatest. The end-to-end axial overlap of electrons attracts the nuclei to form a strong bond. In a pi bond, the p-orbitals overlap side-to-side above and below the nuclei (equatorial plane), so this attraction is more of a "sideways" pull as opposed to a direct line of pulling. Additionally, the election density is spread out over a greater area. Thus, sigma bonds are stronger than pi bonds. |

**Table 17: Answer Example (1)**

| Question | When electronics show black, the pixels are not off. Wouldn't that use less energy, and show truer black? |
|---|---|
| WebGLM | displaying black on OLED screens does use less power than displaying white, because the LEDs that make up each pixel are off when displaying black. However, this is not the case for LCD monitors, where white is actually slightly more energy efficient than black. This is because LCD monitors rely on an array of thin-tube fluorescent bulbs that provide a constant source of light to create a white screen, and to make it black, LCDs rely on a diffuser to block this light. Therefore, displaying black does not necessarily use less energy, and may not even show truer black than white. |
| WebGPT-175B | Most computer and mobile device displays do not use black pixels to save energy. Instead, they use an array of thin-tube fluorescent bulbs to provide a constant source of light to create a white screen. To make it black, LCDs rely on a diffuser to block this light. As a result, LCDs use more energy than CRTs to display a black screen. Measuring a 17-inch (43-centimeter) LCD monitor, Schindler found that white required 22.6 watts, while black came in a tad higher at 23.2 watts. With a 20-inch (50.8-centimeter) LCD, black required 6 percent more energy than white. In contrast, AMOLED displays don't have a backlight at all. Instead, each pixel is like a tiny red, blue, or green colored light. If the screen is instructed to display black, it doesn't need to block any light, it simply doesn't light up any of the little colored sub-pixels. As a result, black pixels save you a lot of power because those pixels can be turned completely off. |
| WebGPT-13B | It depends on the resting state of the LCD as to whether they require energy to stop light or to allow light to pass through," Heap explains. "This is why screen tests show that some CCFL (cold cathode fluorescent lamp) LCD screens save energy and some use a bit more. All of the scientific test data we have come across shows a slight saving on black LCD screens, which suggests that the rest state on many LCD screens does not allow light through." Heap also points out that a large number of Blackle users come from India and South America, where CRTs are still commonly sold. Most LCD devices use a backlight which is always on. When an LCD is black, it is doing full light output, and then also blocking all of that light to create the appearence of black. If a monitor is displaying all-black, most will reduce power to the backlight accordingly, primarily to produce a deeper black, and of course also to save power. |
| Perplexity AI | While it is true that black pixels use less energy than white pixels on some displays, such as OLEDs, this is not always the case. On LCDs, for example, black pixels still require power to be lit up. Additionally, the amount of energy used by a pixel depends on its color intensity, so even if a pixel is black it may still use some energy. |
| Human Writen | It's not true for all types of screens. For LCD screens, energy cost is not definitely related to brightness. The principle of the LCD screen is that there is a white LED panel at the back, and the front uses the LCD panel to cover the pixels one by one. Turn down the backlight brightness to save power. The screen is adjusted to internal black, but the LCD panel is shielded to black, which will not save power. The power saving has nothing to do with the background color. |

**Table 18: Answer Example (2)**

| | |
|---|---|
| Question | Why does it seem like when I watch something the second time around, it goes by faster than the first time I watched it? |
| WebGLM | It seems like when you watch something the second time around, it goes by faster than the first time you watched it because you are focusing on the little details instead, which might make it feel like time is going by faster. Another reason could be that when you know what's coming, your mind breaks it up and you have intervals/milestones, and know when it's going to end, which helps you digest it easier. However, if you go back a third time, it may be ruined because all the things you could discover have been ruined by the second rerun. |
| WebGPT-175B | There are several different ideas about why this happens. One theory is that when you are cognitively engaged, meaning you are thinking about the task at hand and are busy doing it, you don't have the opportunity to notice the passage of time. When you're not engaged, like when you are bored, you have the time to think about the passage of time and so it seems like it goes by slower. Another theory is that the chemicals in your brain release dopamine when you are having fun, and it impacts how your brain perceives time. When you are having fun and your brain is releasing dopamine, time seems to pass faster. |
| WebGPT-13B | There are two main reasons time seems to pass faster the second time around. First, information our brains absorbs and processes increases with each viewing, so the more information there is, the slower time goes. Second, as we age, the rate at which we process visual information slows down, and this is what makes time 'speed up' as we grow older. The complexity of the information we examine and the networks of neurons in our brain increase as we age, and electrical signals must traverse greater distances, causing signal processing to take more time. This is what leads to time passing more rapidly. |
| Perplexity AI | The phenomenon of time seeming to pass faster when watching something for a second time is likely due to the recency effect, which amplifies the importance of recent stimuli or observations. Additionally, paying more attention to time passing can make it seem like it's going by faster, as can familiarity with the content. Finally, life-threatening situations, eye movements, tiredness, hypnosis, age and emotions can all warp our perception of time. |
| Human Writen | When humans repeatedly look at the same thing, the brain will quickly recall the previous memory when it sees the same content. After people recall quickly in the brain, they can quickly grasp its contents, and feel the time has become faster. |

**Table 19: Comparison Example of GLM-10B-RM and GLM-10B-Rouge1**

| Question | Why is it that soap can wash almost anything off of your hands but doesn't wash its fragrance off of you? |
|---|---|
| References | [1] Soap is able to clean hands and dishes because of some pretty nifty chemistry. Soap molecules have on one end what's known as a polar salt, which is hydrophilic, or attracted to water. The other end of the molecule is a nonpolar chain of fatty acids or hydrocarbons, which is hydrophobic—meaning that it's repelled by water but attracted to grease and other oily substances. When you wash your hands, the soap forms something like a molecular bridge between the water and the dirty, germ-laden oils on your hands, attaching to both the oils and the water and lifting the grime off and away. Soaps can also link up with the fatty membranes on the outside of bacteria and certain viruses, lifting the infectious agents off and even breaking them apart. Once the oily dirt and germs are off your hands, the soap molecules thoroughly surround them and form tiny clusters, known as micelles, that keep them from attaching to anything else while they wash down the drain. <br><br>[2] First, the hydrophilic ends of the soap molecule attach to the water, then the hydrophobic ends attach to the oils, pathogens, or other debris on the hand. After the soap has bonded with the germs on your hand the water can then wash it down the drain. This works because the soap molecule can attach to the germs more strongly than the germs can hold onto your skin. <br><br>[3] Therefore, when you wash with soap and water the soap molecules grab on to the fat molecules and pulls them off the skin into the rinse water. When your rinse the skin the fat and germs are removed from the skin. Therefore, the main reason why soap works to cleanse and sanitise is it literally washes everything off the skin. |
| WebGLM-10B-RM | Soap is able to cleanse and sanitise hands and bodies because of its unique chemistry. Soap molecules have on one end a polar salt, which is attracted to water, and on the other end a nonpolar chain of fatty acids or hydrocarbons, which is hydrophobic and attracted to grease and other oily substances. When you wash your hands, the soap forms a molecular bridge between the water and the dirty, germ-laden oils on your hands, attaching to both the oils and the water and lifting the grime off and away.[1] Soap molecules can also link up with the fatty membranes on the outside of bacteria and certain viruses, lifting the infectious agents off and even breaking them apart. Once the oily dirt and germs are off your hands, the soap molecules thoroughly surround them and form tiny clusters, known as micelles, that keep them from attaching to anything else while they wash down the drain.[2] Therefore, the main reason why soap works to cleanse and sanitise is it literally washes everything off the skin.[3] However, while soap can wash almost anything off of your hands, it cannot wash its fragrance off of you because fragrances are not dirt and are not removed by rinsing. |
| WebGLM-10B-Rouge1 | The main reason why soap works to cleanse and sanitise is because of its molecules, which have on one end a polar salt, which is hydrophilic, or attracted to water, and on the other end a nonpolar chain of fatty acids or hydrocarbons, which are hydrophobic, meaning they are repelled by water but attracted to grease and other oily substances[1]. Soap molecules form a molecular bridge between the water and the dirty, germ-laden oils on your hands, attaching to both the oils and the Water and lifting the grime off and away[1][2]. Then, the soap molecules surround the oily dirt and germs and form tiny clusters, known as micelles, that keep them from attaching to anything else while they wash down the drain[1]. However, the fragrance of soap is not removed this way because it is not bound to the fat molecules that are removed from the skin when you wash with soap and water[3]. |

**Table 20: Criteria Details For References**

| | |
|---|---|
| Relevancy | Evaluating how a reference text is related to the question and how much information of the question can be found in the reference. |
| 0 | The reference is almost not related to the question. |
| 1 | A few information or answers related to question can be found in the reference. |
| 2 | Most of the answers to the question can be found in the reference, but a small amount of reasoning or external information is still needed. |
| 3 | The complete answer to the question can be found in the reference |
| Density | Evaluating the proportion of information related to question in a reference text. |
| 0 | Reference content is almost irrelevant to questions and answers. |
| 1 | Reference has most content irrelevant to the problem, but there is still a small amount of content related to the problem. |
| 2 | Reference has a few irrelevant contents. |
| 3 | Reference is closely related to the problem, and almost the whole reference focuses on the problem. |
| Truthfulness | Evaluating whether the text itself violates common sense, logic or contradictions. There is no need to pay attention to whether it is relevant to the question. |
| 0 | The text contains statements that are contrary to the facts, and the authenticity is low. |
| 1 | The text has no statements contrary to the facts and is highly authentic. |
| Keep Blank | If the evaluator cannot determine the authenticity through his or her own knowledge, leave it blank. |
| Toxicity ↓ | Evaluating whether It involves violence, pornography, indecent words and other similar elements. |
| 0 | The above element does not exist. |
| 1 | One or more of the above elements exist. |
| Social Bias ↓ | Evaluating whether there are prejudices against specific social groups, such as racial discrimination and gender discrimination. |
| 0 | The above element does not exist. |
| 1 | One or more of the above elements exist. |

**Table 21: Criteria Details For Answers**

| | |
|---|---|
| Fluency | Evaluating whether grammar, spelling, word usage, etc. conform to people's grammar habits (please do not consider any semantic factors). |
| 0 | There are major errors in grammar and spelling, which make the text difficult to read. |
| 1 | There are small errors in grammar and spelling, which will slightly affect understanding. |
| 2 | There are a few grammatical, spelling or case errors that do not affect understanding. |
| 3 | Fluent language, correct grammar, no mistakes, easy to read. |
| Correctness | Evaluating whether the question is correctly answered. |
| 0 | No answer, or the answer is irrelevant or wrong. |
| 1 | A few answers are given, but they are particularly incomplete or fragmented. The question is basically not answered. |
| 2 | Basically answer the questions, but there are a few mistakes or omissions. |
| 3 | Answer the question perfectly. |
| Citation Accuracy | Evaluating whether the reference marks in the answer are accurate. |
| 0 | The reference marks are basically wrong or there is no reference label. |
| 1 | There are a large number of missing and wrong marks. |
| 2 | There are a few missing and wrong marks. |
| 3 | The reference marks are completely accurate. |
| Objectivity | Evaluating whether all the answers come from references. |
| 0 | There is external knowledge in the answer which does not come from references. |
| 1 | All answers can be based on the reference. |
| Truthfulness | Evaluating whether the text itself violates common sense, logic or contradictions. There is no need to pay attention to whether it is relevant to the question. |
| 0 | The text contains statements that are contrary to the facts, and the authenticity is low. |
| 1 | The text has no statements contrary to the facts and is highly authentic. |
| Keep Blank | If the evaluator cannot determine the authenticity through his or her own knowledge, leave it blank. |
| Redundancy ↓ | Evaluating whether there is redundancy in the answer, such as repeating the same sentence or the same fact repeatedly. |
| 0 | There is no redundancy. |
| 1 | There is redundancy. |

**Figure 9: Real Example: How to balance career and hobbies?**

How to balance career and hobbies?                                                                🔍

## How to balance career and hobbies?

Balancing career and hobbies is important to achieve a full, enriched life [3].It is important to make a routine and prioritize healthy behaviors [1] in order to find balance. Additionally, it can help to change how much time you spend on hobbies by making time for them and prioritizing them [2].It is also important to draw a sharp line between work and everything else [4], and to remember to enjoy time with family and hobbies as much as work. [4]

### [1] How to Balance Work and Hobbies - Real Recovery 🔗

Do you have things you enjoy that bring happiness and fill blank spaces in your day? Participating in pleasurable activities decreases stress and improves mental health. However, balance is essential. You might start to notice your hobbies cutting into time meant for work or other tasks. Learn to balance your work responsibilities and hobbies by making a routine and prioritizing healthy behaviors.

### [2] How to Balance Work and Hobbies - Real Recovery 🔗

Have you noticed that your hobbies take up more of your space, time, and energy? The good news is that finding balance can be quick and straightforward. Change how much time you spend on hobbies by doing the following:

### [3] How to Find and Make Time for Your Passion Even When You ... 🔗

Working a day job while discovering and pursuing your passions demands a lot of time and focus, but they are equally fundamental human needs—one for economic stability and the other for self-actualization. For many of us—especially those with a lot of creative drive and entrepreneurial spirit—balancing our career and pursuing our dreams are critical to leading a full, enriched life.

### [4] The Balancing Act: Work, family, hobbies—how does it all fit ... 🔗

Maybe balance is the ability to draw a sharp line between work and everything else. If I weren't able to do this, I wouldn't enjoy spending time with my family as much. I wouldn't enjoy my hobbies as much, either. Recently I was able to achieve a lifelong dream. I recorded a hip-hop track with my kids. It's amateur, but I had so much fun doing it—and I have so much fun listening to it—that I'll be making music with them until I can't pick up a mic any more. You can find our songs on SoundCloud—search for "Mad Daddy"—or on my LinkedIn profile.

**Figure 10: Real Example: FL Studio and Cubase, which is better?**

FL Studio and Cubase, which is better? 🔍

## FL Studio and Cubase, which is better?

It is difficult to say which digital audio workstation (DAW) is better, FL Studio or Cubase. [5] Both are excellent tools and have made music production, editing, and creation easier than ever before, with professional-sounding audio results [2].FL Studio is thought to be less expensive than Cubase [4] and is more approachable and simpler to use [3], while Cubase is better for vocals and creating full songs easily [1].Additionally, FL Studio has iOS app compatibility, which is convenient [3].Ultimately, both DAWs are excellent and have their own pros and cons, so it is up to the user to decide which one is right for them.

**[1] Cubase vs FL Studio Compared: Which Software is Your Best ...** 🔗

- FL Studio is great for music production and MIDI composition, whereas Cubase is best for vocals and creating full songs easily.

**[2] Cubase vs FL Studio Compared: Which Software is Your Best ...** 🔗

FL Studio and Cubase are two of the best DAWs that you can buy right now. Fans of both software attest that these two have made music production, editing, and creation a whole lot easier. Plus, the audio quality of the results is professional sounding and can rival the quality of those created in studios.

**[3] Cubase vs FL Studio: How to Choose the Best DAW** 🔗

Both are excellent tools, however, for me, it is more approachable and simpler to use FL Studio. The Cubase package is still a very good option but feels a little more restrictive and doesn't have the same iOS app compatibility that you have with FL Studio, which I think is convenient that can't be ignored.

**[4] FL Studio Vs Cubase: Which Music Production Software Is ...** 🔗

FL Studio is often thought to be less expensive than Cubase. FL Studio includes several pricing options, including a free demo version and three paid versions: Fruity Edition, Producer Edition, and Signature Bundle.

**[5] Cubase vs FL Studio Compared: Which Software is Your Best ...** 🔗

While both Cubase and FL Studio can stand their ground as the best digital audio workstation package out there, there are some differences.

**Figure 11: Real Example: Is attention better than CNN?**

Is attention better than CNN?                                                                                    🔍

## Is attention better than CNN?

It depends on the application. Self-attention layers have the advantage of being able to consider the impact of every other pixel of the image, as opposed to convolution layers whose receptive field is the $K\times K$ neighborhood grid [1].This can be beneficial in certain applications, such as image recognition. Attention based CNNs can also be used to model a pair of sentences, used in answer selection, paraphrase identification and textual entailment tasks [2].Furthermore, attention models can outperform the baseline across all depths while having 12% fewer FLOPS and 29% fewer parameters [3].However, it is important to note that while local attention with a window of size 5 has the same receptive field as a CNN with kernel size 5, the learned weights are applied to the inputs differently. [4]

### [1] How a self-attention layer can learn convolutional filters? 🔗

The main difference between CNN and self-attention layers is that the new value of a pixel depends on every other pixel of the image. As opposed to convolution layers whose receptive field is the $K\times K$ neighborhood grid, the self-attention's receptive field is always the full image. This brings some scaling challenges when we apply transformers to images that we don't cover here. For now, let's define what is a multi-head self-attention layer.

### [2] Attention in Neural Networks - Towards Data Science 🔗

In this paper Yin et al presented ABCNN — Attention Based CNN to model a pair of sentences, used in answer selection, paraphrase identification and textual entailment tasks. The key highlight of the proposed attention based model was that it considers the impact/relationship/influence that exists between the different parts or words or whole of one input sentence with the other, and provides an interdependent sentence pair representation that can be used in subsequent tasks. Let's take a quick look at the base network first before looking at how attention was introduced into it.

### [3] How To Go Beyond CNNs With Stand-Alone Self-Attention ... 🔗

The results show that the attention models outperform the baseline across all depths while having 12% fewer FLOPS and 29% fewer parameters.

### [4] Are there any advantages of the local attention against ... 🔗

It is true that when using local attention with a window of size 5, the "receptive field" is the same as a CNN with kernel size 5 (or two CNN layers with kernel size 3). However, there is a key difference in how the learned weights are applied to the inputs.

**Figure 12: Real Example: How to survive in the first-tier cities without a high-salary work?**

How to survive in the first-tier cities without a high-salary work?                                         🔍

## How to survive in the first-tier cities without a high-salary work?

Surviving in first-tier cities without a high-salary work can be difficult, as most entry-level jobs do not pay enough to support someone living in a large city [1].One way to survive is to be great at what you do so that you stand out and can find a job in a smaller economy [2].Additionally, it can be important to make sure you have a steady job with a consistent salary, as the cost of college, healthcare, and housing have gone up while real wages have gone down [1].It is also important to make sure you have a good understanding of the cost of living in the city you are moving to, as the cost of rent, food, and other expenses can quickly add up. Lastly, it can be important to make sure you are realistic about your spending, as it can be easy to fall into the trap of wanting to have a luxury lifestyle, but not having the funds to support it.

### [1] Strategies for survival on an entry-level salary 🔗

Most entry-level jobs do not pay enough to support someone living in a large city. This is a problem for recent grads. They imagine life in a big city as lots of entertainment, crowds of young people for fun, and a great dating scene. But it's a hard life to fund: The cost of college, healthcare, and housing have gone up, all while real wages have gone down. This generation is facing a gap between wages and the cost of living that their parents never did.

### [2] Strategies for survival on an entry-level salary 🔗

I think the key to surviving in a smaller economy is to be great at what you do so that you stand out. If there are not a lot of jobs you can take, you need to make sure you're a top candidate.

**Figure 13: Real Example: What do you think of the 3.5 version of Genshin Impact?**

What do you think of version 3.5 of Genshin Impact?                                                                                              🔍

## What do you think of version 3.5 of Genshin Impact?

Genshin Impact version 3.5 is the latest update coming to the live servers on March 1st. [1] [2] It will bring back familiar events like the Windblume Festival, as well as various adjustments and optimizations [1]. The update is available for PlayStation 5, PlayStation 4, PC via client and Epic Games Store, iOS via the App Store, and Android via Google Play [3] [4]. The Genshin Impact 3.5 preload size is 2.03GB on iOS and can be downloaded by tapping the pre-install resource package from the title screen [4]. It is a fun game to play with a controller, and the update seems to keep the game running well. [4]

### [1] Genshin Impact 3.5 Release Date and Details - IGN 🔗

Genshin Impact Version 3.5 is the next Genshin Impact update coming to the live servers on March 1st. Version 3.5 will herald a return to Mondstadt, and bring back familiar events like the Windblume Festival.

### [2] 'Genshin Impact' Version 3.5 Update Pre-Installation Is Now ... 🔗

Genshin Impact (Free) version 3.5 update pre-installation has finally gone live on iOS, Android, and PC platforms ahead of its release date this Wednesday for all platforms. Genshin Impact version 3.5 'Windblume's Breath' arrives on March 1st for iOS, Android, PS5, PS4, and PC platforms worldwide bringing in the Windblue Festival, a new Archon Quest, two new characters, updated rules for Genius Invokation TCG, one extra Intertwined Fate for each completed Archon Quest, and more. If you missed the previously-announced Prime Gaming collaboration, read this. Watch the Genshin Impact version 3.5 update Dehya character trailer below:

### [3] Genshin Impact version 3.5 update now available - Gematsu 🔗

Genshin Impact is available now for PlayStation 5, PlayStation 4, PC via client and Epic Games Store, iOS via the App Store, and Android via Google Play. A Switch version is also planned.

### [4] 'Genshin Impact' Version 3.5 Update Pre-Installation Is Now ... 🔗

The Genshin Impact 3.5 preload size is 2.03GB on iOS. You can download this by tapping the pre-install resource package from the title screen as usual or from the Paimon menu in-game under other. PC pre-installation details are here. If you haven't checked out Genshin Impact yet, you can download it for free on the App Store for iOS here and on Google Play for Android here. The PC version is available on the official website here and the Epic Games Store. If you play on iOS, with iOS 14.5 or iPadOS 14.5 and later, you can use PS5 and Xbox Series X|S controllers to play Genshin Impact. We featured Genshin Impact as our Game of the Week when it released and awarded it our 2020 Game of the Year. I also featured it as one of the best iOS games to play with a controller. What do you think of what we've seen of Genshin Impact 3.5 so far?

**Figure 14: Real Example: transformers are originated from NLP, but why they can be applied in CV?**

transformers are originated from NLP, but why they can be applied in CV?                    🔍

## transformers are originated from NLP, but why they can be applied in CV?

Transformers were originally developed for Natural Language Processing (NLP) and have since been applied to Computer Vision (CV) due to their effectiveness in improving the performance of language processing models [1].They have also been used in CV to improve the performance of image classification on ImageNet and object detection on COCO [1].Additionally, integration of NLP and CV techniques have been used to improve results in both fields [2].Furthermore, Transformers make use of attention mechanisms which, unlike RNNs, do not have a recurrent structure, allowing them to match the performance of RNNs with attention with enough training data [3].Finally, attempts have been made to apply Transformers in CV by substituting word embeddings with flattened image patches. [4]

### [1] Year 2020: the Transformers Expansion in the CV - Medium 🔗

In the precedent years transformers have perform greatly in the field of NLP. They've significantly improved the performance of the language processing models and the effect is comparable to what had been done to image understanding since 2012 by convolution neural networks. Now at the end of 2020 we have transformers entering the top quartile of well-known computer vision benchmarks, such as image classification on ImageNet and object detection on COCO.

### [2] Natural Language Processing techniques in Computer Vision 🔗

Self Attention and Transformer-based architectures have recently boosted results not only in the NLP domain, but also in the CV domain. Integration of NLP and CV techniques have also inspired different creative approaches in both fields. These approaches have improved SotA, and the potential for even greater results is possible. Thus, every DS/ML/DL practitioner should be aware of these recent developments to successfully implement them to applied tasks.

### [3] Transformer (machine learning model) - Wikipedia 🔗

Before transformers, most state-of-the-art NLP systems relied on gated RNNs, such as LSTMs and gated recurrent units (GRUs), with added attention mechanisms. Transformers also make use of attention mechanisms but, unlike RNNs, do not have a recurrent structure. This means that provided with enough training data, attention mechanisms alone can match the performance of RNNs with attention.

### [4] Natural Language Processing techniques in Computer Vision 🔗

This model design represents the first attempts to apply Transformers in CV. As you can see, it is fairly straightforward - word embeddings are substituted with flattened image patches.

**Figure 15: Real Example: Who proposed Music Transformer? How does it work?**

Who proposed Music Transformer? How does it work? 🔍

## Who proposed Music Transformer? How does it work?

The Music Transformer paper was proposed by Huang et al. from Google Magenta [1].It is a state-of-the-art language-model based music generation architecture that uses Transformers, which were successful in the NLP field, to the symbolic music generation domain [1].The paper also includes optimization work on implementing Transformer with relative attention [2].Some related work of using Transformer architecture on generating music include MuseNet (from OpenAI) [3] and Pop Music Transformer [3].The paper shows the first successful use of Transformer in generating music that exhibits long-term structure. [4]

### [1] Transformer (machine learning model) - Wikipedia 🔗

The Music Transformer paper, authored by Huang et al. from Google Magenta, proposed a state-of-the-art language-model based music generation architecture. It is one of the first works that introduce Transformers, which gained tremendous success in the NLP field, to the symbolic music generation domain.

### [2] Transformer (machine learning model) - Wikipedia 🔗

In my opinion, the Music Transformer paper is not only an application work, but its crux also includes optimization work on implementing Transformer with relative attention. We will delve into this part below.

### [3] Transformer (machine learning model) - Wikipedia 🔗

Some related work of using Transformer architecture on generating music include MuseNet (from OpenAI), and also Pop Music Transformer. It is evident that the Transformer architecture would be the backbone of music generation models in future research.

### [4] Music Transformer: Generating Music with Long-Term Structure 🔗

TL;DR: We show the first successful use of Transformer in generating music that exhibits long-term structure.

**Figure 16: Real Example: What is the backbone of Toolformer?**

What is the backbone of Toolformer?                                                                                          🔍

## What is the backbone of Toolformer?

The backbone of Toolformer is a pretrained 6.7b parameter 'GPT-J' model [1].This model is used to decide which APIs to call, when to call them, what arguments to pass, and how to best incorporate the results into future token prediction [1].The model is finetuned to learn to use APIs [2] by using a language model to build Toolformer's dataset [2].Specifically, the dataset is augmented with API calls in the text and then checked to make sure they worked and were useful [2].The resulting dataset is used to finetune the model so it can learn to use APIs. [2]

### [1] Import AI 🔗

What is Toolformer? "A model trained to decide which APIs to call, when to call them, what arguments to pass, and how to best incorporate the results into future token prediction" . The model is based on a pretrained 6.7b parameter 'GPT-J' model and, despite its small size, outperforms many much larger models, including

### [2] Import AI 🔗

How they did it: They use a language model to build Toolformer's dataset. Specifically, they take a dataset of plain text, augment that data with API calls in the text, then check if the calls a) worked and b) were useful and if they were, then weave that back into the dataset. They use the resulting dataset to finetune the model so it can learn to use APIs. "Moreover, as API calls are inserted in exactly those positions and with exactly those inputs that help M predict future tokens, finetuning... enables the language model to decide when and how to use which tool, based purely on its own feedback."

**Figure 17: Real Example: Why CyGames succeed? What games have they launched?**



Why CyGames succeed? What games have they launched?

## Why CyGames succeed? What games have they launched?

CyGames succeed because they have acquired some of the best names and talent in the Japanese games industry, as they are more than willing to pay for them [1].They have also been successful in the mobile and social game marketplace, with an audience of over 20 million users worldwide [2] .Some of their known games include Dragalia Lost, Granblue Fantasy, Princess Connect! [4] Rage of Bahamut, Shadowverse, Uma Musume: Pretty Derby, and World Flipper series [4].They have also been successful due to their high return on investment in Japan [5].However, some of their strategies are outdated, such as their gacha system of high generosity + bad pull mechanics. [3]

**[1] Cygames And Platinum Games Announce 'Granblue Fantasy ...** 🔗

Due to the success of games like Granblue Fantasy in Japan, Cygames is doing very well financially. This has resulted in the company acquiring some of the best names and talent in the Japanese games industry, as they are more than willing to pay for them.

**[2] Rage of Bahamut developer Cygames breaks 20 million users ...** 🔗

Say what you will about the games themselves, but it's pretty fair to conclude at this point that mobile and social games are where nearly all of the economic growth in the Japanese game marketplace is coming from at this point. There are few better examples of this than Tokyo-based Cygames, a company that was founded in May 2011, released its first game four months later, and now enjoys an audience of over 20 million users worldwide.

**[3] CyberAgent, the mother company of Cygames, released their ...** 🔗

Cygames strategies are just outdated. Their gacha system of high generosity + bad pull mechanics simply doesn't cut it anymore. Their character stories are still the best in the business but they have some of the worst overall narrative stories. Their games require a significant time investment - fairly grindy with low reward-to-time ratio. Modern games have great pity/dupe clauses and are generous enough to cultivate a healthy base of small spenders while still attracting whales. CN and KR games entering the market with high QoL like sweeps and slick UIs redefined the mobage UX. I'm sure they know all of this already, but it's hard to change all of this with current IPs.

**[4] Cygames - Gematsu** 🔗

Cygames is a video game publisher and developer headquartered in Shibuya, Tokyo, Japan. It is known for the Dragalia Lost, Granblue Fantasy, Princess Connect!, Rage of Bahamut, Shadowverse, Uma Musume: Pretty Derby, and World Flipper series.

**[5] Cygames on on why the West requires an entirely different UA ...** 🔗

For Cygames our return on investment is obviously significantly higher in Japan, but it also has to do with what kind of game it is.

**Figure 18: Real Example: When will the COVID-19 disappear?**

When will the COVID-19 disappear? 🔍

## When will the COVID-19 disappear?

Experts do not expect COVID-19 to become eradicated anytime soon, and believe it will become endemic in certain areas according to established patterns over time [2].According to Bill Gates, at some point next year, Covid-19 will become an endemic disease in most places [1].The COVID-19 Scenario Modeling Hub made some pandemic projections spanning August 2022 to May 2023, assuming the new tweaked boosters adding protection for the newest omicron relatives would be available and a booster campaign would take place in fall and winter [3].A February 2021 survey of over 100 scientists found that over half said they did not believe the virus could be eradicated, and 89% of them believe that COVID-19 will become endemic over time. [4]

**[1] Bill Gates: How Covid pandemic ends and becomes endemic ...** 🔗

But that doesn't mean the coronavirus will disappear. In a blog post on Tuesday, Bill Gates laid out one seemingly likely scenario: "At some point next year, Covid-19 will become an endemic disease in most places."

**[2] Will COVID disappear? Scientists say no, so we should think ...** 🔗

Experts expect COVID-19 will someday become endemic, meaning it occurs regularly in certain areas according to established patterns. But they don't think that will be very soon.

**[3] Will COVID disappear? Scientists say no, so we should think ...** 🔗

Experts say COVID will keep causing serious illness in some people. The COVID-19 Scenario Modeling Hub made some pandemic projections spanning August 2022 to May 2023, assuming the new tweaked boosters adding protection for the newest omicron relatives would be available and a booster campaign would take place in fall and winter.

**[4] COVID Endemic vs Pandemic - is it Here to Stay? - Novavax** 🔗

COVID-19 is considered to be in its pandemic phase because the number of infections continues to increase and affect people worldwide.1 A February 2021 survey of over 100 scientists found that over half said they did not believe the virus could be eradicated, and 89% of them believe that COVID-19 will become endemic over time.1

**Figure 19: Real Example: Who is the president of United States now?**

**Figure 20: Real Example: Tell me about the movie Black Panther 2**

Tell me about the movie Black Panther 2 🔍

## Tell me about the movie Black Panther 2

Black Panther 2, officially titled Black Panther: Wakanda Forever, is the sequel to the 2018 movie Black Panther. It follows Queen Ramonda (Angela Bassett), Shuri (Letitia Wright), M'Baku (Winston Duke), Okoye (Danai Gurira) and the Dora Milaje (including Florence Kasumba) as they fight to protect their nation from intervening world powers in the wake of King T'Challa's death [1].Director Ryan Coogler returns to direct the movie, which is expected to both pay homage to the late actor Chadwick Boseman and expand the wonderful world he helped bring to life [4].The cast for Black Panther: Wakanda Forever includes Lupita Nyong'o, Danai Gurira, Florence Kasumba, Martin Freeman, Letitia Wright, Winston Duke, Angela Bassett, Tenoch Huerta, Michaela Coel, and Dominique Thorne, with Lake Bell also cast in an undisclosed role [3].The face of the late Chadwick Boseman – and his beloved character T' Challa – appears early in the sequel, and his presence and legacy is profoundly felt throughout the film, as various characters (and the audience) deal with grief and move on with their lives in the face of adversity – personal and superheroic. [2]

### [1] Black Panther 2 Wakanda Forever | Release date, cast, trailer ... 🔗

The official synopsis for Black Panther: Wakanda Forever reads: "In Marvel Studios' Black Panther: Wakanda Forever, Queen Ramonda (Angela Bassett), Shuri (Letitia Wright), M'Baku (Winston Duke), Okoye (Danai Gurira) and the Dora Milaje (including Florence Kasumba), fight to protect their nation from intervening world powers in the wake of King T'Challa's death.

### [2] 'Black Panther 2' review: 'Wakanda Forever' is a ... - USA Today 🔗

The face of the late Chadwick Boseman – and his beloved character T' Challa – appears early in the sequel "Black Panther: Wakanda Forever," on memorial murals and that nifty Marvel Studios intro. But his presence and legacy is profoundly felt throughout the film, as various characters (and the audience) deal with grief and move on with their lives in the face of adversity – personal and superheroic.

### [3] Black Panther 2: Release Date, Trailer, Plot, and Everything ... 🔗

While the King might be absent from the story, we will be seeing quite a few characters from the first movie returning for the sequel. The cast for Black Panther: Wakanda Forever includes Lupita Nyong'o, Danai Gurira, Florence Kasumba, Martin Freeman, Letitia Wright, Winston Duke, Angela Bassett, Tenoch Huerta, Michaela Coel, and Dominique Thorne. Lake Bell, who voiced Black Widow in the animated series What If...?, has also been cast in an undisclosed role.

### [4] Black Panther 2: Release Date, Trailer, Plot, and Everything ... 🔗

The tragic loss of Chadwick Boseman has left quite a hole in the world, first and foremost the very real loss, as well as his sudden absence in the ongoing Marvel Cinematic Universe. But the show must go on, as they say, and Black Panther 2, officially titled Black Panther: Wakanda Forever, is currently in the works. The movie is expected to both pay homage to the late actor and expand the wonderful world he helped bring to life.

**Figure 21: Real Example: What is Hogwarts Legacy?**

**Figure 22: Real Example: What is google bard?**

**Figure 23: Real Example: What is the most popular AI technology in 2023?**

**Figure 24: Real Example: Tell me the two teams of NBA all-star in 2023.**

Tell me the two teams of NBA all-star in 2023.

## Tell me the two teams of NBA all-star in 2023.

The two teams of NBA All-Star in 2023 will be Team LeBron and Team Giannis [2].The game will be hosted by the Utah Jazz at Vivint Arena and televised nationally by TNT [1].The full list of players, injury replacements and captains can be found at the link provided in reference . [3]

### [1] 2023 NBA All-Star Game - Wikipedia 🔗

The 2023 NBA All-Star Game will be an exhibition game played on February 19, 2023, on the 30th anniversary of the first All-Star Game held in Salt Lake City in 1993. It will be the 72nd edition of the event. The game will be hosted by the Utah Jazz at Vivint Arena. The game will be televised nationally by TNT for the 21st consecutive year.

### [2] 2023 NBA All-Star Game rosters: Full list of players, injury ... 🔗

The 2023 NBA All-Star Game is here. All-Star weekend festivities got started Friday night with the Rising Stars Game and the celebrity game. Mac McClung stole the show at the Dunk Contest and Damian Lillard won the 3-point shootout at All-Star Saturday Night in Salt Lake City. The weekend culminates with Team LeBron vs. Team Giannis in the 2023 NBA All-Star Game on Sunday night. Here's how to watch.

### [3] 2023 NBA All-Star Game rosters: Full list of players, injury ... 🔗

2023 NBA All-Star Game rosters: Full list of players, injury replacements as LeBron, Giannis serve as captains

**Figure 25: Real Example: What is copilot?**

What is copilot?　　　　　🔍

## What is copilot?

Copilot is a cloud-based artificial intelligence tool developed by GitHub and OpenAI to assist users of Visual Studio Code, Visual Studio, Neovim, and JetBrains integrated development environments (IDEs) by autocompleting code [3].It is powered by a brand new AI system named Codex, which is based on the GPT-3 model [2].It is somewhat GPT-3 based, which means it may be able to understand natural language in a plain text file and establish a conversation based on this [1].It is an extremely fun-to-use tool and is extremely useful for coding in Python, JavaScript, TypeScript, Ruby, and Go [3] .However, it also presents a lot of bugs [4], so community feedback is appreciated. [4]

### [1] What is GitHub Copilot? An AI Pair Programmer for Everyone 🔗

Copilot is somewhat GPT-3 based, which means it may be able to understand natural language in a plain text file and establish a conversation based on this.

### [2] What is GitHub Copilot? An AI Pair Programmer for Everyone 🔗

Copilot is the result of a collaboration between GitHub and OpenAI, which is heavily backed by Microsoft. It's powered by a brand new AI system named Codex, which is based on the GPT-3 model.

### [3] GitHub Copilot - Wikipedia 🔗

GitHub Copilot is a cloud-based artificial intelligence tool developed by GitHub and OpenAI to assist users of Visual Studio Code, Visual Studio, Neovim, and JetBrains integrated development environments (IDEs) by autocompleting code. Currently available by subscription to individual developers, the tool was first announced by GitHub on 29 June 2021, and works best for users coding in Python, JavaScript, TypeScript, Ruby, and Go.

### [4] What is GitHub Copilot? An AI Pair Programmer for Everyone 🔗

Copilot is a brand new tool, and it presents a lot of bugs. I created some posts on the GitHub Copilot's discussion page (available for technical preview users), and I've already received feedback from the community.

**Figure 26: Real Example: What is the core technique chatgpt use?**

What is the core technique chatgpt use? 🔍

## What is the core technique chatgpt use?

The core technique used by ChatGPT is Reinforcement Learning with Human Feedback [1] [2].This technique is based on previous academic research and is used to guide the learning process with the specific goal of mitigating the model's misalignment issues [1].Reinforcement Learning with Human Feedback allows the AI to learn what humans expected when they asked a question, going beyond simply training the LLM to predict the next word [2].Additionally, ChatGPT is highly scalable and can be trained to handle new topics and tasks without the need for extensive retraining [3], and is able to carry a conversation through multiple queries in addition to generating software code. [4]

### [1] How ChatGPT actually works - AssemblyAI 🔗

Researchers and developers are working on various approaches to address the alignment problem in Large Language Models. ChatGPT is based on the original GPT-3 model, but has been further trained by using human feedback to guide the learning process with the specific goal of mitigating the model's misalignment issues. The specific technique used, called Reinforcement Learning from Human Feedback, is based on previous academic research. ChatGPT represents the first case of use of this technique for a model put into production.

### [2] ChatGPT: What Is It & How Can You Use It? 🔗

ChatGPT was also trained using human feedback (a technique called Reinforcement Learning with Human Feedback) so that the AI learned what humans expected when they asked a question. Training the LLM this way is revolutionary because it goes beyond simply training the LLM to predict the next word.

### [3] ChatGPT is a new AI chatbot that can answer questions and ... 🔗

One of the key advantages of ChatGPT is its ability to learn and adapt quickly to new information. This means that it can be trained to handle new topics and tasks without the need for extensive retraining. Additionally, ChatGPT is highly scalable, which makes it well-suited for use in large-scale applications.

### [4] ChatGPT is a new AI chatbot that can answer questions and ... 🔗

But with early ChatGPT adopters demonstrating the technology's ability to carry a conversation through multiple queries in addition to generating software code, the world of so-called natural language processing appears to be entering a new phase.

**Figure 27: Real Example: Where does the code data used to train copilot come from?**

Where does the code data used to train copilot come from?                                                    🔍

## Where does the code data used to train copilot come from?

Copilot is trained on massive amounts of open source code from GitHub [3] [4] .This process is repeated for millions of lines of code totaling close to 200Gb of code [1] .The model itself has 12 billion parameters which can be altered during the course of training [1] .Although most code output by Copilot can be classified as a transformative work, some of the code it produces may be copied verbatim from the training data [2] .This has led to fears that the output code is insufficiently transformative to be classified as fair use and may infringe on the copyright of the original owner. [2]

### [1] A Review of Github Copilot - Partee.io 🔗

This process is repeated for millions of lines of code totaling close to 200Gb of code. The model itself has 12 billion parameters which can be altered during the course of training. Below is a visual of what this training process looks like for the model Copilot is based on (GPT-3).

### [2] GitHub Copilot - Wikipedia 🔗

Although most code output by Copilot can be classified as a transformative work, GitHub admits that a small proportion is copied verbatim, which has led to fears that the output code is insufficiently transformative to be classified as fair use and may infringe on the copyright of the original owner. This leaves Copilot on untested legal ground, although GitHub states that "training machine learning models on publicly available data is considered fair use across the machine learning community". The company has also stated that as of June 2022 only a few source codes are taken over completely or partially unchanged. Therefore as the software continues to learn, this figure is expected to drop. Also in June 2022, the Software Freedom Conservancy announced it would end all uses of GitHub in its own projects, accusing Copilot of ignoring code licenses used in training data. In November 2022, a class-action lawsuit was filed, challenging the legality of Copilot.

### [3] A Review of Github Copilot - Partee.io 🔗

Copilots ability and shortcomings both arise from the data it was trained on: open source repositories of code. Think about it. Did the OpenAI developers have time to filter through and only select the code they thought was written well or correct?? Absolutely not. There are millions of repositories and that would take forever. Even though Copilot uses a state of the art AI model, and uses a brilliant training process, it's still guaranteed to occasionally write code that is either incomplete, incorrect, or inefficient. This is because it was trained on code with all of those qualities. Anyone can push their code to GitHub, even the most junior developers. That code, written by those junior developers, is part of the corpus of text that Copilot is trained to produce.

### [4] Everything you need to know about Github Copilot - Medium 🔗

Copilot is built on Codex, a novel model based on GPT3 that has been trained on massive amounts of open source code from GitHub. It's directly connected with VSCode to create suggestions based on a combination of the current context (i.e., your code) and the "knowledge" it's gained during the training process.

**Figure 28: Real Example: What is the model behind Perplexity AI?**