

# Increasing the Difficulty of Automatically Generated Questions via Reinforcement Learning with Synthetic Preference

Anonymous ACL submission

## Abstract

The demand for high-quality question-answering (QA) datasets has surged with the proliferation of language models and conversational agents in various emerging domains. As these models become ever more capable, the possibility of applying them to more challenging tasks is growing. Manual dataset annotation is costly and time-consuming, necessitating a more efficient approach. Automatically generated questions often suffer from a lack of quality or difficulty; hence, we propose a methodology to increase the difficulty of automatically generated questions using synthetic preference data, derived from SQuAD, to fine tune a question generation model using reinforcement learning. We empirically show an improvement in question difficulty over a supervised-finetuned model with minimal impact on question validity and perform an extensive error analysis. We believe our methodology provides a feasible approach to creating high quality synthetic datasets in emerging domains.

## 1 Introduction

Question-answering (QA) datasets serve diverse purposes, from providing educational materials for students (Das et al., 2021) to serving as crucial resources for model training and evaluation (Rajpurkar et al., 2016). As new domains begin to incorporate language models into workflows and customer service tasks based around information extraction and content reasoning, the need for challenging, in-domain datasets has become increasingly evident. Difficult datasets are crucial for advancing the capabilities of language models, pushing them to handle complex tasks and enhancing their performance in these real-world, challenging scenarios. This growing need is underscored by the rapid proliferation of QA datasets, with over 80 new datasets emerging within the last two years alone (Rogers et al., 2023). Despite this, many QA

Shark Tank (based on the Dragon's Den reality format) also became a midseason sleeper hit on Sundays in the spring of 2010; the following season, it became the tentpole of the network's Friday night schedule, gradually helping make ABC a strong competitor...

### SFT

What reality format did Shark Tank come from?  
(answer: the Dragon's Den)

✗ Short Range Dependencies   ✗ No Entity Disambiguation   ✗ Passage Duplication

### PPO

What show became the tentpole of ABC's Friday night schedule? (answer: Shark Tank)

✓ Long Range Dependencies   ✓ Entity Disambiguation   ✓ Sequence Modification

$\frac{42}{\text{word span}} > \frac{9}{\text{word span}}$    "The network" = ABC   Shark Tank ↓ show

Figure 1: Example generated questions from supervised-fine-tuned question generation model and one fine-tuned with PPO from synthetic difficulty samples.

datasets suffer from a lack of quality or difficulty while economically scaling in size.

One major challenge faced in developing QA datasets is cost. Annotation cost for QA datasets is especially high because of the time and cognition required to write questions and validate them. To exemplify this, the popular question-answering dataset SQuAD (Rajpurkar et al., 2016) recommended workers to take 4 minutes for every 5 questions at a rate of \$9/ hour. This amounts to roughly \$12,000 just to write the dataset's 100,000 questions; moreover, the cost is likely much higher when considering answer validation, and discarded samples due to duplication or poor quality.

Automatic Question Generation (AQG) systems present a remedy to these challenges given their efficiency and scalability compared to human annotators. Even in a zero-shot setting, language models are able to generate coherent questions (Sachan et al., 2022; Wang et al., 2023b); as such, we argue that writing coherent questions is no longer the main goal of AQG systems. Controlling more

abstract attributes such as question difficulty, desirable for improving model performance, remains challenging as the concept is somewhat subjective and hard to manipulate. However, recent innovations in reinforcement learning for language models now enable these human-like ideals to be injected into the model learning process (Ouyang et al., 2022).

Pinning down a definitive description of question difficulty is near impossible as it depends on many factors. Common syntactic measurements of question difficulty include: question length; the average frequency of question terms in the English language (AlKhuzayy et al., 2023; Beinborn et al., 2014); and the syntactic difference between the dependency parse trees of a question and answer sentence (Rajpurkar et al., 2016). Semantic measurements may consider the relatedness between an answer span and the surrounding context (Beinborn et al., 2015), or the cosine similarity between distractors and the correct answer (Hsu et al., 2018). We argue that difficult questions also require: reasoning over long spans of text; disambiguation of entities; and the use of synonyms to distance the question from the source text. A combination of all of these features is incredibly challenging to directly incorporate into the model training process.

We initially attempted to define such a task to encourage Large Language Models (LLMs) to rank samples with respect to difficulty. We extensively explored defining a set of criteria for difficulty for zero-shot models, tasking the model with selecting the more difficult sample between two question-answer pairs. To validate the proficiency of the model, we aimed to maximise the kappa agreement between the LLM and human annotators; however, the results were very poor, achieving a Cohen’s  $\kappa$  of only 0.14. These results led to the understanding that textually specifying the full scope of difficulty would become an intractable problem. Therefore, we pivoted to leveraging the feature extraction capabilities of transformer models to infer the components of difficulty.

In this paper we present a methodology for increasing the difficulty of automatically generated questions using synthetic preference data. We derive this preference data from the ability of question-answering models to correctly identify answer spans in a subset of SQuAD, assigning to each question a score based on the number of models that incorrectly answered the question. We assume that more challenging questions are answered cor-

rectly less frequently, and use this as the basis for our comparisons.

We summarise this paper’s contributions as follows:

1. A methodology for increasing the difficulty of automatically generated questions using PPO and synthetic data;
2. Empirical evidence of the methodology’s efficacy including human evaluation;
3. An in-depth error analysis and study of interesting phenomena that emerge as part of this approach.
4. An open-source code base and set of models to recreate and adapt our work<sup>1</sup>

## 2 Related Work

A similar question generation approach to ours is employed by Zhang et al. (2022) who adopt a pipeline structure. However, their primary objective is to generate suitable questions rather than specifically focusing on difficulty. An important distinction lies in their extensive pre-processing applied to identify candidate answers before feeding them to the question generation model. We argue that pre-identifying answers may limit diversity and prevent the inclusion of potentially complex answer types.

**Analyzing and Controlling Question Difficulty** Understanding and managing question difficulty holds significant importance, especially in tasks involving the creation of exams and assessments (AlKhuzayy et al., 2023). One approach, as presented by Loginova et al. (2021), involves modelling the difficulty of multiple-choice questions through the use of softmax scores obtained from a pre-trained QA model. The variance in these scores is then calculated, with higher variance indicating greater difficulty.

Lin et al. (2015) controls the difficulty of quiz questions through the selection of distractor answers based on semantic similarity between linked data items. This involves collecting both structured RDF data and unstructured text, computing similarity scores through K-means clustering, and generating questions and answers via template-based methods. Importantly, the semantic similarity plays a role in determining the difficulty level, with more

<sup>1</sup>We release all code and models on [GitHub](#).

challenging questions featuring distractors exhibiting higher semantic similarity.

**Reinforcement Learning with Human Feedback** RLHF is a machine learning paradigm that combines reinforcement learning with human-provided guidance to steer language models to meet the needs of users, finding frequent use in chatbot and AI assistant settings (Ouyang et al., 2022). The basis for most modern methods is the Proximal Policy Optimisation (PPO) algorithm (Schulman et al., 2017), which iteratively enhances the language model’s policy to maximize cumulative rewards through interactions with a dataset or language simulation. It collects experiences, evaluates advantages, and updates the policy with a clipped surrogate objective to ensure stability, gradually improving the model’s performance.

**Automatic Question Generation** Chen et al. (2019) introduce a cross-entropy loss with a reinforcement learning-based loss function when training a gated bi-directional neural network for question generation. In this context, the reward model is optimising the semantic and syntactic quality of the question. BLEU-4, as a reward function, optimises the model for the evaluation metrics and the negative Word Movers Distance component is used to ensure semantic quality by maximising the similarity between a generated sequence and a ground truth sequence. Although question quality is maintained, other factors such as question difficulty are not considered.

Self-critic sequence training (SCST) (Rennie et al., 2017) uses a classical policy gradient method, REINFORCE, which is a Monte Carlo method. SCST computes rewards with n-gram token overlap as sub-sentence level rewards. Since training sets often have limited questions, these training rewards are arguably sparse, hindering the question generation model from extrapolating beyond the training distribution.

Liu et al. (2019) adopt a two-component reward for refining ill-formed questions. Question wording is used as a measure of short-term reward, and alignment between the question and answer represents a long-term component.

### 3 Method

To challenge the high cost of manual annotation while maintaining quality and increasing difficulty, we design and implement a robust system capable of generating contextually relevant, coherent, and

challenging question-answer pairs from textual input. The process follows the core methodology of RLHF, deviating only in the use of synthetic preference data to train a reward model. Rather than explicitly defining the characteristics of difficulty and risking failure to capture certain aspects, we exploit the ability of leading question-answer models to derive which questions are challenging, and allow a reward model to extract the component features of the task.

We task three models with answering all questions in our validation split of SQuAD. These questions are assigned a score based on the number of times they were answered incorrectly, which are in turn used to generate pairwise preference data. These pairwise samples enable the training of a reward model for use in fine-tuning a supervised model for the task of question generation.

We embed this synthetic RLHF process into a greater pipeline for generating samples, shown in Figure 2. This ensures the quality of the final dataset. The pipeline also contains a set of rule-based critics which are used to exclude samples that are malformed and those with non-unique answers in the source text. Samples are then deduplicated using exact string matching.

The remainder of this section discusses each of the relevant components of the pipeline and the RLHF process.

#### 3.1 Supervised Fine-Tuning

In our training process for question generation and response formatting, we begin by employing a reformatted version of the SQuAD v1 training split (see Table 1). The reformatting converts SQuAD to the task of question-answer pair generation, as shown in Figure 3. We select the "correct" answer as the one that appears most frequently in the list of answers for each question in the dataset, selecting randomly among the most common if there is no victor. To ensure model robustness without overfitting, the model undergoes a single epoch of training, enabling it to effectively capture the nuances of the task.

#### 3.2 Reward Modelling

To control the difficulty of our model, we leverage the intrinsic properties present in challenging questions from SQuAD. To extract these attributes, we employ three question answering models that almost match or exceed human performance on SQuAD v2 to evaluate our development split: a

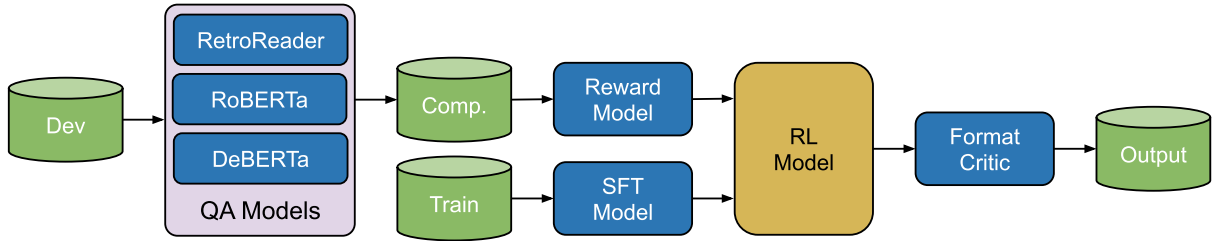


Figure 2: Depiction of our dataset generation pipeline. Question-Answering models are first used to create pairwise comparison data to train a reward model. An SFT model is trained on the train split of SQuAD and then fine-tuned using the reward model, producing the RL model. When generating question-answer pairs for the final dataset, generations are passed through the format critics to ensure data quality.

<p><b>Instruction</b> Write 1 answerable span extraction question and provide the correct answer based on the text.</p> <p><b>Input</b> ... Upon its arrival in Canberra, the Olympic flame was presented by Chinese officials to local Aboriginal elder Agnes Shea, of the Ngunnawal people. She, in turn, offered them a message stick ...</p> <p><b>Response</b> Who received the flame from Chinese officials in Canberra? (answer: <u>Agnes Shea</u>)</p>	<p>Samples that pass these critics are then deduplicated using exact matching.</p>
--	--

Figure 3: Example training sample from the reformatted SQuAD dataset for use in supervised fine-tuning.

262 RoBERTa-large model<sup>2</sup>, a DeBERTa-large model<sup>3</sup>  
 263 and RetroReader (Zhang et al., 2020). Each ques-  
 264 tion is assigned a score based on the number of  
 265 models that failed to correctly answer the ques-  
 266 tion. These scores are used to place questions into  
 267 a pairwise ranking setup against other questions  
 268 for the same input context. Where a question’s  
 269 scores are equal, they are considered ties, and no  
 270 pairwise sample is created. We also record the mar-  
 271 ginal, defined as the difference in score between the  
 272 chosen and rejected samples, to experiment with  
 273 the marginal ranking loss, as defined in Touvron  
 274 et al. (2023).

### 275 3.2.1 Format Critics

276 To ensure the quality of the final dataset, we utilise  
 277 a collection of rule-based critics which we call *For-*  
 278 *mat Critics*. These critics have two main functions:  
 279 they remove questions that don’t adhere to the de-  
 280 sired format of  $Q?$  (answer:  $A$ ); they ensure the  
 281 provided answer is unique in the text, minimising  
 282 the number of ambiguous or impossible questions.

<sup>2</sup>deepset/roberta-large-squad2

<sup>3</sup>deepset/deberta-v3-large-squad2

283 Samples that pass these critics are then deduplica-  
 284 ted using exact matching.

### 285 3.3 Reinforcement Training

286 We use Proximal Policy Optimisation (Schulman  
 287 et al., 2017) with multiple sets of adapters to reduce  
 288 the memory overhead during training, implemented  
 289 using the Transformers Reinforcement Learning  
 290 library (von Werra et al., 2020). A single base  
 291 model is used with separate LoRA adapters for the  
 292 policy, reference, and reward model components;  
 293 each is switched to perform the relevant aspect of  
 294 the reinforcement training process.

295 During early experiments, we found that training  
 296 was often very unstable or resulted in low pass  
 297 rates at the format critic. To combat this, we added  
 298 a rule-based reward component to penalise gener-  
 299 ations that did not pass the format critic. This  
 300 simple function converts the reward to be the nega-  
 301 tive absolute reward in the case that samples are  
 302 malformed. Using a rule-based reward that manipu-  
 303 lates the original reward prevents the instability  
 304 caused by hard coding a fixed penalty and saves  
 305 the computational complexity and imperfection of  
 306 a second adapter-based reward model:

$$R_i = \begin{cases} -|R_i| & \text{if malformed} \\ R_i & \text{otherwise} \end{cases}$$

## 307 4 Experimental Setup

### 308 4.1 Models

309 We conduct our experiments with LLaMA2-7B-  
 310 chat and apply LoRA adapters to all linear layers  
 311 on both SFT and RM models to enable training on  
 312 a single A100 80GB GPU using Flash Attention 2  
 313 (Dao, 2023). All LoRA adapters share the same hy-  
 314 perparameters: LoRA rank of 16,  $\alpha$  of 32, dropout  
 315 of 0.05, no bias and BrainFloat (BF16) datatype.



Split	# Contexts	# Questions
Train	18,891	87,599
Dev	1,567	8,038
Test	500	2,532
Human Test	50	50
Train comp.	1,107	8,394
Dev comp.	123	950

Table 1: Split of contexts and questions from SQuAD. The *comp.* splits are derived from the dev split, used to evaluate the performance of the reward model during training.

We experiment with marginal ranking loss to help distinguish between slight and significant differences in question difficulty while training the reward model. Under the hypothesis that the difficulty of a question is not independent of the associated passage of text, we also experiment with training a reward model with and without the input text attached. Results of these experiments can be found in Appendix A.

## 4.2 Generation Settings

During generation, the model is tasked with producing a single output for each question in the training set using nucleus sampling (Holtzman et al., 2020). We maintain the original configuration for LLaMa-2 with a repetition penalty of 1.1, top P of 0.7, and top K of 0 but increase the temperature from 0.6 to 0.9 to increase the diversity of generations.

## 4.3 Data Splits

We base our splits off the original SQuAD to minimise the risk of data leakage. We maintain the full train split unchanged as any model previously trained on SQuAD will have seen the full train split. We extract a test split of 500 contexts from the dev split, ensuring no contexts appear in both the dev and test splits. We extract 50 unique contexts from the test split for a human evaluation of question quality and answerability. In all cases, context-question pairs were only kept if they fit into the context length of LLaMa-2 when formatted in the correct prompt format. All samples were formatted into the three instruction components: *instruction*, *input*, *response* as shown in Figure 3.

Only the dev set of our SQuAD dataset was used to derive difficulty comparison data, to ensure the reward model never sees the samples used for evaluation. To evaluate the reward model, we extract 10% of the comparison contexts. Full dataset statistics can be found in Table 1.

## 4.4 Evaluation Metrics

As our goal is to evaluate the difficulty of answerable questions, we provide the input passage, question and answer to GPT-4o<sup>4</sup> and Gemini-1.5-pro<sup>5</sup> and ask whether the sample meets our specification of validity. We take samples to be answerable if they were unanimously labelled as such, and reject all other samples. GPT-based evaluations have demonstrated a robust alignment with human preferences across various complex tasks in reference-free settings (Fu et al., 2023; Liu et al., 2023). The results of this analysis can be found in Appendix C.

To assess the quality of generated questions relative to our SQuAD test split, we *intentionally avoid*  $n$ -gram based metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and more modern alternatives such as Q-Metrics (Nema and Khapra, 2018), as we believe they restrict diversity of generation, constraining the model to reference questions and answers. We instead adopt the following reference-free metrics:

**Syntactic Divergence** provides a distance measure between two dependency paths which acts as a measure of difficulty. Word-lemma anchors, common to both the question and answer sentence, are first detected. A dependency path from the anchor to the interrogative word (who, what, etc.) in the question is compared to the dependency path between the anchor and the answer span in the answer sentence using Levenshtein distance (Levenshtein et al., 1966).

**RQUGE** calculates an *acceptability-score* by generating an answer for the candidate question and predicting the semantic similarity between the predicted answer and the gold answer provided by the user. In our setup, this metric acts as an assessment of both the question and answer quality (Mohammadshahi et al., 2023).

**QAScore** attempts to align AQG evaluation to human judgements. Question-answer pairs are evaluated by summing log-probabilities of RoBERTa correct token predictions for all words in the answer when masked individually. QAScore claims to show strong correlation with human judgement (Spearman  $r = 0.864$ ) (Ji et al., 2022).

**Self-BLEU** assesses how similar questions are to other questions generated for a given context. Each question is taken as a hypothesis and the others as a reference for the BLEU calculation. The

<sup>4</sup>gpt-4o as of 1st June 2024

<sup>5</sup>gemini-1.5-pro as of 1st June 2024

Model	Total Valid ( $\uparrow$ )	DeBERTa ( $\downarrow$ )	RoBERTa ( $\downarrow$ )	RetroReader ( $\downarrow$ )
SQuAD	2,532 (-)	0.68	0.68	0.65
ZeroShot	357 $\pm$ 14 (0.14)	0.644 $\pm$ 0.007	0.650 $\pm$ 0.007	0.629 $\pm$ 0.009
SFT	1252 $\pm$ 2 (0.49)	0.654 $\pm$ 0.012	0.653 $\pm$ 0.005	0.616 $\pm$ 0.015
PPO-input	<b>1375 <math>\pm</math> 18 (0.54)</b>	<b>0.601 <math>\pm</math> 0.004</b>	<b>0.606 <math>\pm</math> 0.003</b>	<b>0.582 <math>\pm</math> 0.007</b>
PPO-input-margin	1373 $\pm$ 4 (0.54)	0.612 $\pm$ 0.001	0.608 $\pm$ 0.005	0.587 $\pm$ 0.002

Table 2: Question-Answering model performance on each set of samples. Models were only supplied samples which passed the format critics and were unanimously deemed answerable by GPT-4o and Gemini-1.5-pro. The *Total Valid* column indicates this number of valid samples used during question answering. Accuracy is based on exact match and results are mean and standard deviation across three sets of generated samples. Lower accuracy indicates harder questions.

self-BLEU is taken as the average BLEU for the question collection (Zhu et al., 2018).

## 5 Results & Discussion

**Model Accuracy** To measure performance, we observe the difference in prediction accuracy for QA models on each dataset. Table 2 shows that in all cases of PPO training, we observe a decrease in average model prediction accuracy and an increase in the total number of valid generations. The consistent decrease in absolute prediction accuracy for all models when using the PPO trained models over both zero-shot and SFT signifies an increase in average question difficulty. The SFT process vastly improves the model’s ability to generate valid questions. The PPO process further bolsters this capability which illustrates that the model is learning the intrinsic properties of high-quality questions. The performance of the reward models, shown in Appendix A, is reflected here, showing lesser degrees of improvement for those models fine-tuned without access to the input passage.

**External Metrics** Figure 4 shows results for the reference-free metrics. RQUGE is clearly effective at discriminating between human-written SQuAD samples, those generated by the fine-tuned models and the zero-shot examples, but it is unable to separate out the SFT and PPO results. The particularly high score for SQuAD could in part be due to data leakage as the answer generation model for the metric was trained on SQuAD (Khashabi et al., 2022). This would indicate why our newly generated questions might score lower as it cannot have memorised the answer. Syntactic divergence results for the SQuAD test split and all trained model generations follow a consistent distribution but the zero-shot results appear much better, despite having a higher average prediction accuracy than the SFT and PPO models. Zero-shot obtaining

higher syntactic divergence could stem from the general purpose language generation objective of LLaMa-2-chat. This could cause the model to generate boilerplate text which distances the structure of the question from that of the answer sentence but doesn’t necessarily result in a more difficult question. QAScore proves uninformative, only being able to subtly identify SQuAD samples from model generated samples. Self-BLEU indicates that SQuAD samples are the most diverse, which is to be expected, but that zero-shot samples exhibit a distinct lack of diversity when compared with fine-tuned models. This result is, in part, misleading as Self-BLEU was only calculable for input passages with at least two valid questions. As the number of valid generations was so low for the zero-shot model, the cases where multiple valid questions were generated for a context was disproportionately in favour of identical generations.

In general we find the reference-free metrics to show limited correlation with model prediction accuracy and an ability differentiate human written samples from model generations. We believe this is evidence for the continued need for more reliable, reference-free evaluation tools for question generation.

**Human Evaluation** To evaluate question quality, we conduct a human evaluation on a subset of 50 passages from the test split. Each input passage and question is filtered through the format critic then provided to two annotators who select either the correct answer span or indicate that the question cannot be answered. In the case of annotator disagreement or the annotated answers differing from the model generated answer, the annotator responses and the model answer are provided to two new annotators who both select which responses are appropriate. We allow annotators to select multiple responses as correct but only include those

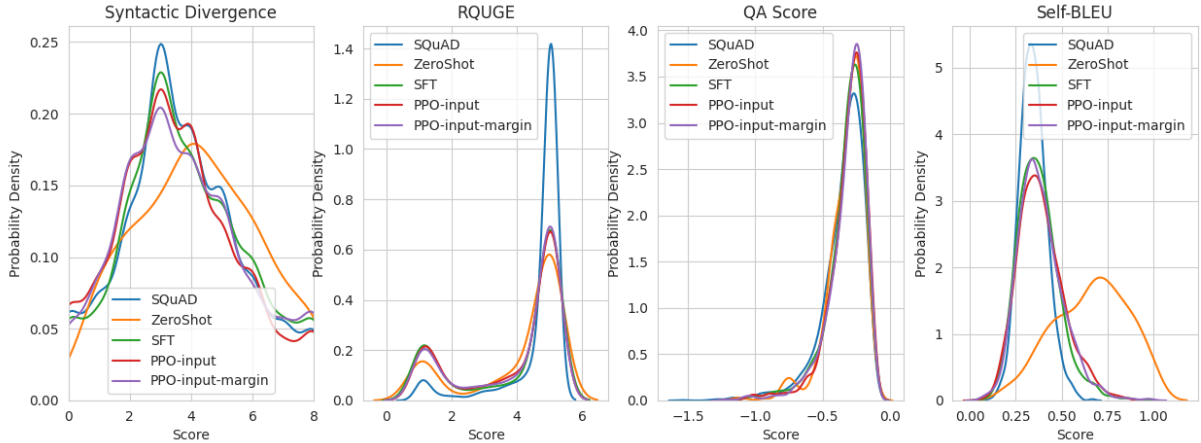


Figure 4: Distribution of reference free metrics results for each model’s generations based on our SQuAD test set.

Model	Full	Partial
<b>ZeroShot</b>	0.10	0.14
<b>SFT</b>	0.52	0.60
<b>PPO-input</b>	0.52	0.64
<b>PPO-input-margin</b>	<b>0.56</b>	<b>0.64</b>

Table 3: Results of human evaluation for question quality. *Full* indicates that the model generated answer was a valid answer according to the format critics and identified by human annotators and *Partial* indicates that the sample passed format critics and a valid answer was identified for the question but the model generated answer did not match.

480 that were selected unanimously by both annotators  
 481 as valid. We observe an agreement of  $\kappa = 0.7975$   
 482 between annotators. The results of this evaluation,  
 483 shown in Table 3, displays an equivalent or  
 484 improved rate of answerability when fine-tuning  
 485 with PPO; the answerability proportions for each  
 486 dataset are roughly equivalent to those presented in  
 487 Table 2. This further corroborates the efficacy of  
 488 our approach.

### 489 5.1 Error Analysis

490 **Failure Modes** At a high level, we can observe the  
 491 reasons for sample rejection for each model. As  
 492 shown in Figure 5, the zero-shot model is generally  
 493 unable to generate samples that have a single  
 494 answer span in the text, despite exactly specifying  
 495 this in the prompt. The high number of incorrectly  
 496 formatted samples was a result of only a question  
 497 being generated or neither a question nor answer  
 498 being generated. For all the trained model variants,  
 499 the dominant failure mode was unanswerable  
 500 questions. As shown in Appendix C, each of the  
 501 fine-tuned models show a similar proportion of

502 otherwise valid samples being unanswerable. The  
 503 answerability rate could potentially be improved by  
 504 generating candidate answers, as in (Zhang et al.,  
 505 2022), and passing an input passage and answer to  
 506 the question generation model.

507 **Positional Bias** One interesting phenomenon  
 508 is the positional bias in where the model chooses  
 509 to generate answers. To calculate positional bias,  
 510 we treat the full answer span as a single "word"  
 511 and calculate the proportion through the input para-  
 512 graph in which the answer word appears. As seen  
 513 in Figure 6, the zero-shot positional bias is less  
 514 severe than in the other datasets. The positional  
 515 bias of SQuAD is clearly seen as, after training on  
 516 the dataset, all models exhibit this same preference  
 517 for the beginning of input passages. The clear bias  
 518 observed in the zero-shot model, despite not being  
 519 fine-tuned, is documented in other tasks such as  
 520 LLM ranking (Wang et al., 2023a; Li et al., 2023)  
 521 and in summarisation where introductory content is  
 522 favoured (Ravaut et al., 2023). A potential remedy  
 523 is to supply the model with a sliding window of  
 524 sentences across the context paragraph to force the  
 525 model to generate questions throughout the text.  
 526 While this would improve the diversity of a final  
 527 dataset, it may have the adverse effect of limiting  
 528 the range of dependencies, restricting potentially  
 529 challenging questions across the whole text.

530 **Hallucinated External Knowledge** Where am-  
 531 biguous references to specific entities exist in the  
 532 input passage such as *the museum collection*, the  
 533 models frequently attempt to fill in which entity  
 534 is being referred to. From a context containing  
 535 ambiguous references to an unnamed museum, the  
 536 questions *What year did the Tate acquire the statue*  
 537 *of St John the Baptist?*, *How many works does*

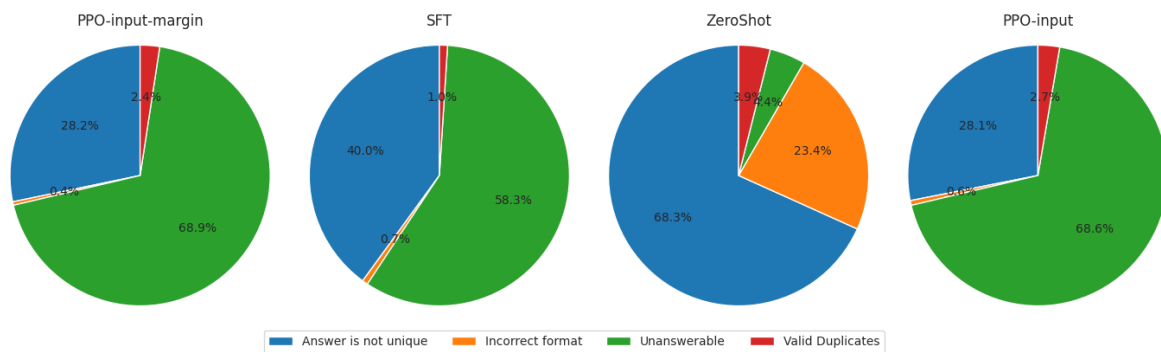


Figure 5: Error distribution of questions for SFT, ZeroShot, and the two best performing PPO variants.

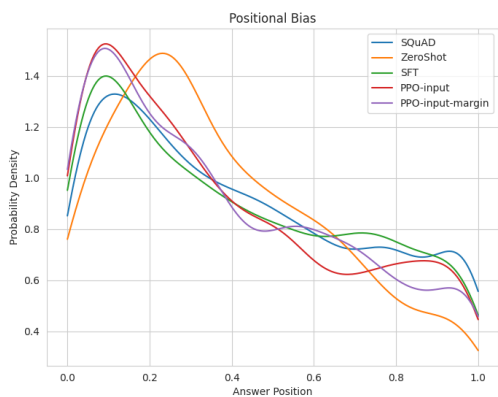


Figure 6: Position of answer span, merged to be a single word, as a proportion of the way through the input passage when split into words. SQuAD positions are selected from our test split and answers are chosen to be the most common from the list of suitable answers. Neither invalid nor exact duplicate questions are considered.

538 *Rodin have in the British Museum’s collection?*  
 539 were generated across both the SFT and PPO mod-  
 540 els; the examples consistently passed LLM evalua-  
 541 tions of answerability. This suggests the solution  
 542 to this problem is more holistic and requires im-  
 543 provements at a foundational model level to resolve.  
 544 We could resolve this at a critic level through more  
 545 careful prompting, however, this returns to our origi-  
 546 nal and intractable task of textually describing a  
 547 complex task. A more holistic solution could be to  
 548 adapt PPO with functional grounding (Carta et al.,  
 549 2023) to be a pure text task. However, this may  
 550 lower the quality of questions as it could discourage  
 551 the use of implicit or complementary knowledge.

552 **Unidirectional Relationships** A strategy to in-  
 553 crease the difficulty of questions is to invert re-  
 554 lationships found in the text. The models some-

555 times misappropriate this tool, resulting in invalid  
 556 questions such as the question *What did the Ming*  
 557 *dynasty represent?* from a passage containing *...ex-*  
 558 *plorer Zheng He representing the Ming Dynasty...*  
 559 Knowledge graph assisted generation could help  
 560 to resolve these logical inconsistencies (Lin et al.,  
 561 2015). However, expecting our target demograph-  
 562 ics, emerging domains, to possess high-quality  
 563 knowledge graphs is an unreasonable assumption.

## 6 Conclusion

564  
 565 In this paper we have introduced a robust ap-  
 566 proach for automatically generating question-  
 567 answer pairs from textual input. Using existing,  
 568 high-performing question answer-models, we are  
 569 able to determine which questions are most chal-  
 570 lenging and develop synthetic pairwise data for  
 571 training a reward model. Rather than explicitly  
 572 defining the characteristics of question difficulty,  
 573 we allow the reward model to extract these fea-  
 574 tures, leading to a significant increase in question  
 575 difficulty when used to fine-tune the SFT model.

576 Furthermore, we have conducted an extensive  
 577 analysis of the current issues with this approach and  
 578 provide potential remedies which may be explored  
 579 in future work.

580 We believe this technique may be extended to ad-  
 581 dress further abstract properties of question genera-  
 582 tion such as ambiguity, completeness and relevance.  
 583 This method may also be adapted to tackle multi-  
 584 ple aspects at once through the use of multi-reward  
 585 model setups as in Wu et al. (2023).

586 All code and models from this project is made  
 587 available for adaptation and reuse.



## 588 Limitations

589 This project only shows the suitability of the  
590 method on a single model. In future work, we  
591 seek to address this by performing a more compre-  
592 hensive review of the approach across a range of  
593 model sizes and architectures. We also acknowl-  
594 edge that this method currently only addresses an-  
595 swerable questions while most contemporary QA  
596 datasets utilise both answerable and unanswerable  
597 questions. Finally, despite using LoRA and multi-  
598 adapter training, we still required approximately 15  
599 GPU hours on an A100 80GB which restricts the  
600 potential audience for this approach. Evaluating  
601 smaller models or quantisation will enable greater  
602 access to this project’s benefits.

## 603 Ethics Statement

604 This project has been approved by the relevant in-  
605 stitution’s ethics committee. We use LLaMa2 in  
606 accordance with Meta’s license<sup>6</sup>. All annotators  
607 were located through word of mouth are paid £12  
608 per hour - above the UK National Living Wage of  
609 £11.44

## 610 References

- 611 Samah AlKhuzaey, Floriana Grasso, Terry R. Payne,  
612 and Valentina Tamma. 2023. [Text-based Question  
613 Difficulty Prediction: A Systematic Review of Auto-  
614 matic Approaches](#). *International Journal of Artificial  
615 Intelligence in Education*.
- 616 Lisa Beinborn, Torsten Zesch, and Iryna Gurevych.  
617 2014. [Predicting the Difficulty of Language Pro-  
618 ficiency Tests](#). *Transactions of the Association for  
619 Computational Linguistics*, 2:517–530.
- 620 Lisa Beinborn, Torsten Zesch, and Iryna Gurevych.  
621 2015. [Candidate evaluation strategies for improved  
622 difficulty prediction of language tests](#). In *Proceed-  
623 ings of the Tenth Workshop on Innovative Use of NLP  
624 for Building Educational Applications*, pages 1–11,  
625 Denver, Colorado. Association for Computational  
626 Linguistics.
- 627 Thomas Carta, Clément Romac, Thomas Wolf, Sylvain  
628 Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer.  
629 2023. [Grounding large language models in interac-  
630 tive environments with online reinforcement learn-  
631 ing](#).
- 632 Yu Chen, Lingfei Wu, and Mohammed J. Zaki. 2019.  
633 [Reinforcement learning based graph-to-sequence  
634 model for natural question generation](#). *CoRR*,  
635 abs/1908.04942.

<sup>6</sup><https://ai.meta.com/llama/license/>

- Tri Dao. 2023. [FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning](#). ArXiv:2307.08691 [cs]. 636  
637  
638
- Bidyut Das, Mukta Majumder, Santanu Phadikar, and Arif Ahmed Sekh. 2021. [Automatic question generation and answer assessment: a survey](#). *Research and Practice in Technology Enhanced Learning*, 16(1):5. 639  
640  
641  
642
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. [GPTScore: Evaluate as You Desire](#). ArXiv:2302.04166 [cs]. 643  
644  
645
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The Curious Case of Neural Text Degeneration](#). ArXiv:1904.09751 [cs]. 646  
647  
648
- Fu-Yuan Hsu, Hahn-Ming Lee, Tao-Hsing Chang, and Yao-Ting Sung. 2018. [Automated estimation of item difficulty for multiple-choice tests: An application of word embedding techniques](#). *Information Processing & Management*, 54(6):969–984. 649  
650  
651  
652  
653
- Tianbo Ji, Chenyang Lyu, Gareth Jones, Liting Zhou, and Yvette Graham. 2022. [QAScore—An Unsupervised Unreferenced Metric for the Question Generation Evaluation](#). *Entropy*, 24(11):1514. 654  
655  
656  
657
- Daniel Khashabi, Yeganeh Kordi, and Hannaneh Hajishirzi. 2022. [UnifiedQA-v2: Stronger Generalization via Broader Cross-Format Training](#). ArXiv:2202.12359 [cs]. 658  
659  
660  
661
- Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union. 662  
663  
664  
665
- Zongjie Li, Chaozheng Wang, Pingchuan Ma, Daoyuan Wu, Shuai Wang, Cuiyun Gao, and Yang Liu. 2023. [Split and Merge: Aligning Position Biases in Large Language Model based Evaluators](#). ArXiv:2310.01432 [cs]. 666  
667  
668  
669  
670
- Chenghua Lin, Dong Liu, Wei Pang, and Edward Apeh. 2015. [Automatically Predicting Quiz Difficulty Level Using Similarity Measures](#). In *Proceedings of the 8th International Conference on Knowledge Capture, K-CAP 2015*, pages 1–8, New York, NY, USA. Association for Computing Machinery. 671  
672  
673  
674  
675  
676
- Chin-Yew Lin. 2004. [ROUGE: A Package for Automatic Evaluation of Summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics. 677  
678  
679  
680
- Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. [G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment](#). ArXiv:2303.16634 [cs]. 681  
682  
683  
684
- Ye Liu, Chenwei Zhang, Xiaohui Yan, Yi Chang, and Philip S. Yu. 2019. [Generative question refinement with deep reinforcement learning in retrieval-based QA system](#). *CoRR*, abs/1908.05604. 685  
686  
687  
688

689	Ekaterina Logina, Luca Benedetto, Dries Benoit, and Paolo Cremonesi. 2021. <a href="#">Towards the Application of Calibrated Transformers to the Unsupervised Estimation of Question Difficulty from Text</a> . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 846–855, Held Online. INCOMA Ltd.	
690		
691		
692		
693		
694		
695		
696	Alireza Mohammadshahi, Thomas Scialom, Majid Yazdani, Pouya Yanki, Angela Fan, James Henderson, and Marzieh Saeidi. 2023. <a href="#">RQUGE: Reference-Free Metric for Evaluating Question Generation by Answering the Question</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 6845–6867, Toronto, Canada. Association for Computational Linguistics.	
697		
698		
699		
700		
701		
702		
703		
704	Preksha Nema and Mitesh M. Khapra. 2018. <a href="#">Towards a Better Metric for Evaluating Question Generation Systems</a> .	
705		
706		
707	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	
708		
709		
710		
711		
712		
713	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a Method for Automatic Evaluation of Machine Translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	
714		
715		
716		
717		
718		
719		
720	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">SQuAD: 100,000+ Questions for Machine Comprehension of Text</a> . ArXiv:1606.05250 [cs].	
721		
722		
723		
724	Mathieu Ravaut, Shafiq Joty, Aixin Sun, and Nancy F. Chen. 2023. <a href="#">On Context Utilization in Summarization with Large Language Models</a> . ArXiv:2310.10570 [cs].	
725		
726		
727		
728	Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 7008–7024.	
729		
730		
731		
732		
733	Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2023. <a href="#">QA Dataset Explosion: A Taxonomy of NLP Resources for Question Answering and Reading Comprehension</a> . <i>ACM Computing Surveys</i> , 55(10):197:1–197:45.	
734		
735		
736		
737		
738	Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. <a href="#">Improving Passage Retrieval with Zero-Shot Question Generation</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3781–3797, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
739		
740		
741		
742		
743		
744		
745		
	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. <a href="#">Proximal Policy Optimization Algorithms</a> . ArXiv:1707.06347 [cs].	746 747 748
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <a href="#">Llama 2: Open Foundation and Fine-Tuned Chat Models</a> . ArXiv:2307.09288 [cs].	749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771
	Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. <a href="#">Trl: Transformer reinforcement learning</a> . <a href="https://github.com/huggingface/trl">https://github.com/huggingface/trl</a> .	772 773 774 775 776
	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023a. <a href="#">Large Language Models are not Fair Evaluators</a> . ArXiv:2305.17926 [cs].	777 778 779 780
	Zhenduo Wang, Yuancheng Tu, Corby Rosset, Nick Craswell, Ming Wu, and Qingyao Ai. 2023b. <a href="#">Zero-shot Clarifying Question Generation for Conversational Search</a> . In <i>Proceedings of the ACM Web Conference 2023, WWW '23</i> , pages 3288–3298, New York, NY, USA. Association for Computing Machinery.	781 782 783 784 785 786 787
	Zequ Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A. Smith, Mari Ostendorf, and Hannaneh Hajishirzi. 2023. <a href="#">Fine-Grained Human Feedback Gives Better Rewards for Language Model Training</a> . ArXiv:2306.01693 [cs].	788 789 790 791 792
	Cheng Zhang, Hao Zhang, Yicheng Sun, and Jie Wang. 2022. <a href="#">Downstream transformer generation of question-answer pairs with preprocessing and post-processing pipelines</a> . In <i>Proceedings of the 22nd ACM Symposium on Document Engineering, DocEng '22</i> , pages 1–8, New York, NY, USA. Association for Computing Machinery.	793 794 795 796 797 798 799
	Zhuosheng Zhang, Junjie Yang, and Hai Zhao. 2020. <a href="#">Retrospective reader for machine reading comprehension</a> . <i>CoRR</i> , abs/2001.09694.	800 801 802

Model	Accuracy (%)
RM	63.66
RM-input	<b>70.69</b>
RM-margin	62.39
RM-input-margin	70.38

Table 4: Accuracy of reward model variants based on the test split of the comparisons dataset. *input* indicates that the model was trained with the question and associated text passage as input and *margin* indicates that marginal ranking loss was used.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. [Taxygen: A Benchmarking Platform for Text Generation Models](#). ArXiv:1802.01886 [cs].

## A Reward Model Performance

To understand the relative contributions of marginal ranking loss and the use of the input when training reward models to discriminate based on difficulty, we trained all four permutations of settings on the whole training split of the comparisons dataset and evaluated on the test split. As shown in Table 4, the inclusion of the input text had a very significant impact on performance. This was expected as the difficulty of a question is not independent of the related passage. Surprisingly, marginal ranking loss had a very slight negative impact on reward model performance. We believe this could be due to the fact that features of difficulty are very subtle and the marginal component may have caused too significant adjustments due to higher loss values.

## B Obtaining Zero-Shot Model Generations

To obtain zero-shot generations, we adopted a slightly different approach. To not constrain the output of the model too much, thus harming generation performance, we adopted a two-stage process. LLaMa-2-7b-chat was first tasked with generating a question-answer pair based on the text, unconstrained. We then passed this output back into the model with the task of extracting the question and answer components and placing them into a JSON file with the keys *question* and *answer*. We used the same, high temperature of 0.9 for generating the samples and a much lower temperature of 0.2 for extracting into a JSON to reduce the chance of models altering the generated sequences while structuring them.

## C API-Based LLM Answerability Annotation

To ensure that we evaluate performance on as high-quality questions as possible, we extract only those questions deemed *answerable*, by our definition, by both GPT-4o and Gemini-1.5-pro. Table 5 shows that the zero-shot samples had the highest rate of predicted answerability; each other variant shows very consistent rates of answerability. This outcome should be tempered by the results in Figure 5 which indicates that the zero-shot model had an extremely high failure rate in many other regards.

Following is a text, a question and an answer. You must determine whether the provided answer is a correct span-extraction response to the question. If there are multiple plausible answers in the text, the answer should be the most relevant or accurate one. If there are multiple equally plausible answers in the text, respond "NO". If the provided answer is incomplete or contains excess information, respond "NO". If the answer does not correctly answer the question, respond "NO". Only if the answer is correct and does not breach the aforementioned requirements, respond with "YES".

**Text:** ... Upon its arrival in Canberra, the Olympic flame was presented by Chinese officials to local Aboriginal elder Agnes Shea, of the Ngunnawal people. She, in turn, offered them a message stick ...

**Question:** Who received the flame from Chinese officials in Canberra?

**Answer:** Agnes Shea

Respond with only "YES" or "NO" in response to this task. Do NOT provide any other text or reasoning.

Figure 7: Example prompt and response to GPT-4o (gpt-4o as of 1st June 2024) and Gemini-1.5-pro (gemini-1.5-pro as of 1st June 2024).

Model	Answerable ( $\uparrow$ )	Unanswerable ( $\downarrow$ )	Undetermined ( $\downarrow$ )	Cohen's $\kappa$ ( $\uparrow$ )
<b>ZeroShot</b>	<b>0.73</b>	<b>0.14</b>	<b>0.13</b>	0.61
<b>SFT</b>	0.64	0.20	0.16	<b>0.62</b>
<b>PPO</b>	0.64	0.20	0.16	<b>0.62</b>
<b>PPO-input</b>	0.62	0.20	0.18	0.58
<b>PPO-margin</b>	0.62	0.19	0.19	0.56
<b>PPO-input-margin</b>	0.63	0.21	0.16	<b>0.62</b>

Table 5: Results of answerability task posed to GPT-4o and Gemini-1.5-pro. Results represent the proportion of samples that are answerable, unanswerable and undecided, taken from those samples which passed the format critic.