

Tasnim Kabir, Yoo Yeon Sung, Saptarashmi Bandyopadhyay, Hao Zou, Abhranil Chandra, and Jordan Lee Boyd-Graber. **You Make me Feel like a Natural Question: Training QA Systems on Transformed Trivia Questions.** *Empirical Methods in Natural Language Processing*, 2024.

```
@inproceedings{Kabir:Sung:Bandyopadhyay:Zou:Chandra:Boyd-Graber-2024,  
Title = {You Make me Feel like a Natural Question: Training QA Systems on Transformed Trivia Questions},  
Author = {Tasnim Kabir and Yoo Yeon Sung and Saptarashmi Bandyopadhyay and Hao Zou and Abhranil Chandra and  
Booktitle = {Empirical Methods in Natural Language Processing},  
Location = {Miami},  
Year = {2024},  
Url = {http://umiacs.umd.edu/~jbg//docs/2024_emnlp_natural.pdf},  
}
```

Accessible Abstract: Many of the questions for training AIs how to answer questions come from the queries users type into search engines (like Google’s Natural Questions). Is there a cheaper—perhaps even better—way? We propose a “naturalization” technique to turn high-quality, rigorously edited trivia questions into examples that resembles Natural Questions. Training on our naturalized questions and testing on natural questions comes close to the results with using Natural Questions, and we can improve results on MMLU (a standard modern evaluation set) by using our data.

Links:

- ArXiv [<https://arxiv.org/abs/2402.11161>]
- Research Talk [<https://youtu.be/gbBraibisY>]

Downloaded from http://umiacs.umd.edu/~jbg/docs/2024_emnlp_natural.pdf

Contact Jordan Boyd-Graber (jbg@boydgraber.org) for questions about this paper.

You Make me Feel like a Natural Question: Training QA Systems on Transformed Trivia Questions

Tasnim Kabir
Computer Science
University of Maryland
tkabir1@umd.edu

Yoo Yeon Sung
College of Information
University of Maryland
yysung53@umd.edu

Saptarashmi Bandyopadhyay
Computer Science
University of Maryland
saptab1@umd.edu

Hao Zou
Computer Science
Columbia University
zou00080@umn.edu

Abhranil Chandra
Computer Science
University of Waterloo
a23chand@uwaterloo.ca

Jordan Lee Boyd-Graber
Computer Science
University of Maryland
jbg@umiacs.umd.edu

Abstract

Training question answering (QA) and information retrieval systems for web queries require large, expensive datasets that are difficult to annotate and time-consuming to gather. Moreover, while *natural* datasets of information-seeking questions are often prone to ambiguity or ill-formed, there are troves of freely available, carefully crafted question datasets for many languages. Thus, we automatically generate shorter, information-seeking questions, resembling web queries in the style of the Natural Questions (NQ) dataset from longer trivia data. Training a QA system on these transformed questions is a viable strategy for alternating to more expensive training setups showing the F1 score difference of less than six points and contrasting the final systems.¹

1 Introduction

Question answering is a central problem in AI research. One way of understanding *why* people ask questions was explained in [Rodriguez and Boyd-Graber \(2021\)](#)²: questions come from either an information-seeking paradigm ([Voorhees, 2019](#), henceforth Cranfield) or a probing, evaluative paradigm ([Turing, 1950](#), Manchester).

While it is easy to get *questions* in the Cranfield paradigm because the asker creates questions that they do not know the *answer* to, additional annotations to find these answers are expensive. For example, Natural Questions ([Kwiatkowski et al., 2019](#)), a benchmark dataset collected by Google

from questions people asked online, critically does not include “found” correct *answers*. Instead, annotating these answers could be more expensive than their Manchester counterparts, mostly written by QA writing experts (e.g., trivia members).

Moreover, while large corporations can collect large-scale *natural* Cranfield questions *at no cost*, these questions sometimes are of poor quality because of ambiguity ([Min et al., 2020](#)) or false presuppositions ([Yu et al., 2023b](#)). Due to these pit, [Boyd-Graber and Börschinger \(2020\)](#) argue that Manchester questions are more useful for building and evaluating QA systems. Thus, we utilize the Quiz Bowl (QB) samples, a Manchester QA dataset, created by trivia experts (Section 2).³

This paper investigates whether and how we can transform Manchester QB samples into questions that resemble natural, Cranfield questions. To this end, we propose syntactic transformations (NATURALIZATION) that convert QB elicitations into QB-TRANS questions that resemble NQ (Section 3).

To validate the quality of QB-TRANS for training QA systems, we consider two experimental settings: zero-shot and supervised. The zero-shot setting examines whether QB-TRANS is an effective training data for a QA system when compared to NQ (Section 4). We train QA systems with QB-TRANS training data and compare the two systems on the NQ test set. Average F1 scores on NQ test set vary by less than 6 points, which implies that QB-TRANS can replace NQ training data.

¹The codebase and data is available at <https://github.com/Pinafore/qb2nq>

²[Rogers et al. \(2023\)](#) call this probing.

³QB writers are particularly known for understanding what makes for a good QA pair; QB dataset avoids the ambiguity and false presuppositions that are often in NQ.

We also combine NQ with QB-TRANS as training data in our supervised setting (Section 5), improving F1 (tested on NQ test set) by 10 points compared to training on only NQ. QB-TRANS lacks issues that plague NQ: presupposition and ambiguity (Section 6). Moreover, NATURALIZATION generalizes to other datasets (Section 6.4). Our contributions are naturalizing Manchester QB questions into Cranfield QB-TRANS while retaining the positive traits of QB samples, thereby improving QA with a more affordable process. The dataset generated from NATURALIZATION can be used to answer non-NQ data (Section 6.5) which proves the generalization of NATURALIZATION. Section 8 shows how this can ensure a cheaper and more up-to-date alternative to NQ data by generating large-scale Cranfield dataset that benefits training question-answering models and generalizes to other non-NQ datasets.

2 Artful but Arcane QB dataset

This section discusses why we use QB data and how different they are from NQ questions. The next section explains NATURALIZATION (Section 3).

Elicitations from QB dataset Consider this QB example:

A radio mast named for this city was the world’s tallest structure until the mast collapsed in 1991. This capital contains a skyscraper formerly known as the Joseph Stalin Palace of Culture and Science. A landmark called Sigismund’s Column commemorates Sigismund III Vasa, who moved his capital from Kraków to this city on the Vistula River. A 1943 Jewish ghetto uprising occurred in—for 10 points—what Polish capital?

Here, clues are introduced pyramidally—harder, more obscure clues about *Warsaw* appear first (Rodriguez et al., 2019)—so that whoever knows the most about Warsaw should be able to answer the question sooner.⁴

However, we do not need this complexity. Instead, we extract the series of clues that an expert author thought was noteworthy about *Warsaw* (e.g., key sites that commemorate its history and rulers who made it the capital).

We define the source text paragraph as an *elicitation*. As they are combined clues in multiple sentences, they are not grammatical or natural. Thus,

⁴For example, deciding if “moved his capital from Kraków to this city on the Vistula” is when the player should answer requires the ability to decide not just what to answer, enough to answer but also *when* to answer in the quiz bowl tournament (He et al., 2016).

we turn each clue extracted from an elicitation into multiple NQ-like questions, which are short and simple. Ultimately, our goal is NATURALIZING these clues into information-seeking, *natural* questions.

Comparison with NQ datasets We extract an average of seven sentences for each QB elicitation. Each of these sentences is twenty-two words on average. On the other hand, in NQ, the average question length is eight words (Kwiatkowski et al., 2019). The NQ questions were harvested from Google queries based on heuristics.⁵ The number of samples from QB and NQ are comparable (QB: 112,927 elicitations and answers and NQ: 307,373 samples); however, there is a substantial difference in cost, quality, and quantity.

For cost comparison, while the QB elicitations have answers unambiguously created by trivia authors, answers to NQ questions must be laboriously annotated by paid workers. While Google has not officially released costs, the convoluted process and the lack of reproduction since 2019 suggests that its price is high. From the QA researcher’s perspective, the elicitation process is free.

For quality comparison, trivia authors who created QB elicitations understand the importance of discouraging ambiguity and false suppositions in their clues (Boyd-Graber and Börschinger, 2020) while they are prevalent in NQ. Thus, if we can faithfully elicit these clues from QB, the resulting questions may be of higher quality than NQ questions (Detail analysis is in Section 6).

Finally, for quantity comparison, because each QB elicitation contains many clues, the size of a transformed dataset is three-fold larger than NQ. Also, while the NQ dataset may only ask a single question about a rare entity, this is not likely the case for QB: a single elicitation would produce several clues about an entity, allowing a model to understand more about each potential answer.

3 NATURALIZATION

This section outlines NATURALIZATION: converting the elicitations into multiple NQ-like questions (Figure 1).

⁵For example, the questions start with “who”, “when” or “where” followed by a finite form of “do” or a modal verb (Kwiatkowski et al., 2019)

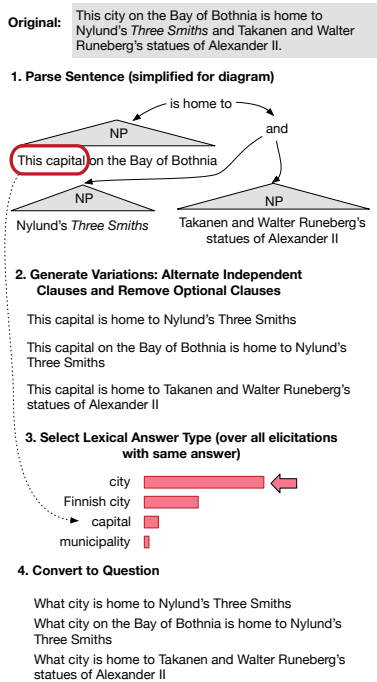


Figure 1: In the process of creating Cranfield style questions from Manchester elicitations, (1) we take each clue sentence from the paragraph-long QB question, and parse it. (2-3) The parsed sentences are transformed into variants, (4) that are finally turned into Cranfield questions.

3.1 Generating Candidates

Many of the transformations depend on an initial dependency parse (Nivre, 2010; Honnibal and Montani, 2017). Some parsed elicitations are statements about a target entity that do not resemble how questions are asked (e.g. statements about the target entity “she was the last Queen of Hawaii” or “this element is mined from bauxite”). To transform these into questions, we find mentions coreferent with the answer.

Conjunction and Removing Clauses Given these candidates, we then extract the minimal facts that could form the basis of a question. For example, if the QB elicitation had “he wrote *Animal Farm* and 1984”, this can become two facts: “he wrote *Animal Farm*” and “he wrote 1984”. Thus, we construct independent clauses by extracting spans that contain the mention (“he”), a verb (“wrote”), and one member of a conjunction (either of the two works). Similarly, we can sometimes remove clauses: “this author who graduated Eton College wrote *Homage to Catalonia*” can be simplified to “this author wrote *Homage to Catalonia*” (Details in Appendix, Algorithm 2).

Canonical Answer Type Next, we identify what kind of answer the question is looking for. This is important because sometimes questions written in QB’s pyramidal style uses oblique references, particularly at the beginning of the question: “substance” for zinc, “creator” for Chinua Achebe, or “polity” for Bangladesh. However, these are rarer than the most straightforward and direct references. For example, zinc is most often asked about using “what element”, Chinua Achebe with “what playwright”, and Bangladesh with “what nation”. Thus, we group all QB elicitations that have the same answer and for each answer find the most frequent string used to refer to about the answer. These canonical answer types then replace the mentions in the original question.

Imperative to Interrogative The most obvious difference between QB elicitations and NQ questions is that QB elicitations are not grammatical questions: rather, they are declarative statements about the answer. For imperative statements such as “name this first prime minister of Canada”, we generate a synthetic mention that makes the object of the imperative verb the question: “who was the first prime minister of Canada” by mapping the canonical answer type to its WORDNET (Fellbaum, 1998) hypernym and applying the appropriate question word (e.g., person.n.01 maps to “who”, time_period.n.01 maps to “when”). The pseudocode is in Algorithms 4 (interrogative) and 5 (answer type).

Additional Heuristics Through observation of the linguistic and grammatical style of NQ we add additional heuristics to further improve the candidates such as **removing punctuation** and **adding subject** (full list in Appendix A).

3.2 LLM Transformation: Llama2 and GPT

As a baseline, we transform QB clues into questions by prompting generative text models— GPT and LLAMA2 (Touvron et al., 2023).⁶ For a fair comparison, we follow a similar process of generating NQ-like questions from QB. We extract clue sentences from the QB elicitations, then prompt GPT and LLAMA2 to produce a natural question. We first identify the lexical answer type (e.g., “this person”) and ask GPT and LLAMA2 to formulate a query that could be used as a Google search. This

⁶<https://huggingface.co/meta-llama/Llama-2-7b-hf>

Original QB	Answer	QB-TRANS	LLAMA2
This "sea Prince of Wales" leaves his home by hiding out in a canoe near a coral reef, and he is mistakenly called "Hedgehog" by a character who offers him a ninetieth lay, a partner of Bildad named Peleg.	Queequeg	Which "sea prince of wales" leaves his home by hiding out in a canoe near a coral reef, and he is mistakenly called "hedgehog" by a character who offers him a ninetieth lay, a partner of bildad named peleg.	What sea did the "sea Prince of Wales" hide out in a canoe near a coral reef?
In one of this man's paintings, one character oddly uses her left hand to grasp the red-cloaked character's chin while her right hand sits at his knee	Jean Auguste Dominique Ingres	In one of which man's paintings, one character oddly uses her left hand to grasp the red-cloaked character's chin while her right hand sits at his knee?	What is the significance of the left-hand grasp and the right-hand placement in Jean Auguste Dominique Ingres' painting featuring a red-cloaked character and another woman?
This substance moves to the liver where it is turned into glucose in the Cori cycle.	Lactic acid	Which substance moves to the liver where it is turned into glucose in the Cori cycle.	What substance is produced when the body breaks down carbohydrates during intense exercise?

Table 1: LLAMA2 converted question can introduce error in question. Conversion of QB clues using rule-based and large language model-based techniques (LLAMA2) is shown here. In the case of LLAMA2, in the first and second examples, the questions changed the entity the QB clue is asked about. In the third example, words were removed from the question and additional random words were added, implying LLAMA2 transformation is worse than that of QB-TRANS.

is an example prompt for GPT and LLAMA2 (text in bold is specific to answer):

This "sea Prince of Wales" leaves his home by hiding out in a canoe near a coral reef, and he is mistakenly called "Hedgehog" by a character who offers him a ninetieth lay, a partner of Bildad named Peleg. This is a trivia question. Turn this into shorter question of fewer than 20 words that start with "what character", and ask about this "character" in the short question. The questions should be natural as a Google query to find out what the answer to the long question is. The shorter questions you write should not include the answer, **Queequeg** and not be confusable with other answers.

LLAMA2 and GPT transformations do not contain all the clues or hallucinate some information (Table 1). In the first example, the question is asking about the character "Queequeg" from the 1851 novel Moby-Dick. However, LLAMA2 did not capture the entity of interest and asked about the "sea" instead of a "character". In the second example, the question is asking about French painter Jean Auguste Dominique Ingres. However, LLAMA2 asked about the significance of the position of the hands in the painting instead of the entity of the painter and also included the answer in the question. In example 3 for LLAMA2 generated question, important clues are removed (e.g. substance moves to the liver where it is turned into glucose in the Cori cycle) and random clues are added (e.g. substance is produced when the body breaks down carbohydrates during intense exercise). More examples can be seen in Table 10. Similarly, GPT also generated questions with hallucinations, including random clues. It also sometimes changes the entity about the original QB-TRANS and includes the "answer" in the generated question (the prompt instructs not to include the answer in the question). For example, for QB clue "*This language uses five cases, though the genitive and dative cases are identical, as are the nominative and accusative.* (Answer:

Daco-Romanian), GPT converts it to "*What are the five cases used in the Romanian language?*". More examples can be seen in Table 11. However, LLAMA2 and GPT have similar generated questions (Examples in Table 12).

4 Zero-shot QA with QB-TRANS training

We ensure we use no NQ data and evaluate on NQ test set which disadvantages our approach as NQ has issues such as presupposition and ambiguity (Section 6.1).

4.1 Challenges in Zero-shot QA System

There are challenges in comparing models for zero-shot QA because some models are based on large language models (LLMs) that do not disclose training data. Thus we do not know whether some zero-shot systems use NQ (Shi et al., 2023). For example, Narayanan (2023); Magar and Schwartz (2022); Sainz et al. (2023a,b) suggest that GPT-3.5 is contaminated with NQ training and development set.

One sign that these models train on NQ is that they give an abnormal probability for tokens in NQ as measured by Min K% probability (Shi et al., 2023). The state-of-the-art LLMs have an average probability of 63% (Detail results in Appendix, Table 7). This indicates that these state-of-the-art LLMs have a high probability of having NQ in the training data.

Another clue that these models have used NQ for training is that they repeat NQ answers to questions even when NQ is wrong (Table 2); this is the clearest signal that the model has seen the NQ data's answers, as annotation errors are less likely to be by coincidence. GPT incorrectly answers those questions, with the answers included in the NQ dataset. Thus, it is likely for GPT's training data to be contaminated (Sainz et al., 2023a; Cotton et al.,

NQ question	NQ answer (wrong)	Gold answer	GPT answer	Comment
Who won the Oscar for best picture in 1976?	Rocky	One Flew Over The Cuckoo’s Nest	Rocky	Rocky won the best picture in 1977. ⁷
Where was held the first session of Muslim league	Dhaka, Bangladesh	Karachi	Dhaka, Bangladesh	The AIME Conference in 1906, held at Dhaka, Bangladesh, laid the foundation of the Muslim League. ⁸
Total number of death row inmates in the us	2,718	2,331	Over 2,400 people	This information is changed over periods.
Who is next in line to be the monarch of England	Charles, Prince of Wales	Prince William	Charles, Prince of Wales	The answer is outdated.

Table 2: To determine whether NQ is in the training data of GPT, we take the answers given by GPT 3.5. If the answer is the same as given in NQ dataset, we can assume it has seen those datasets.

2024) and can no longer be a fair candidate for zero-shot experiments.

4.2 Zero-shot QA systems

Thus, we select two systems with high accuracy on traditional NQ training: Deep Passage Retrieval (Karpukhin et al., 2020b, DPR) and Retrieval-Augmented Language Modeling (Shi et al., 2024, REPLUG). These systems are trained from the ground up. DPR (Karpukhin et al., 2020a) extracts the answer from a context which is extracted using passage retriever models. We train DPR on the questions, answers, and context passages for the NQ-like generated QB-TRANS questions dataset (ours). In training, we generate the positive context by collecting passages that contain answer string, and negative context otherwise (Example in Appendix, Table 13). In REPLUG (Shi et al., 2024), the retrieval model finds the most appropriate passage from a large corpus; then the model produces more accurate answers by augmenting retrieved information to the input context.

4.3 Training Data

We compare all of our generated datasets with the original NQ dataset (NQ). Our goal is to create a QA system with the same accuracy as the original NQ dataset while training on the QB-TRANS dataset, so this is an upper bound. In this zero-shot experiment, we train the model with different percentages derived from QB-generated questions. We compare this traditional training regime with several training sets derived from QB-TRANS (Full results in Appendix, Figure 6). We compare against all transformed sentences from our syntactic-based method (QB-TRANS) to the LLM baselines (QB-GPT and QB-LLAMA2). We also use individual

⁷<https://www.oscars.org/oscars/ceremonies/1976>

⁸https://en.wikipedia.org/wiki/All-India_Muslim_League

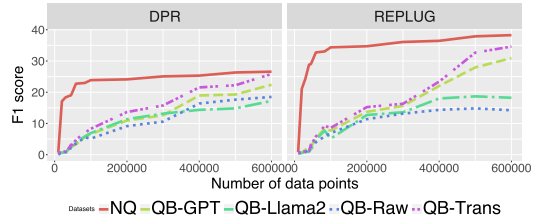


Figure 2: QB-Trans can replace NQ in training QA system and achieve accuracy close to NQ training system. **DPR:** As expected, **QB-TRANS** without any NQ data comes within 5 points of a model trained on NQ. Training on the full QB-TRANS and evaluating it produces the highest F1 score system with DPR. This does better than transformations created by prompting a GPT and LLAMA. **REPLUG:** Again, **QB-TRANS** without any NQ data comes within 7 points of a model trained on NQ.

elicitation sentences from the QB dataset *without* any transformation: QB-RAW.

We used multiple passes when there is a difference in dataset size. For example, because NQ has 307k, we used multiple passes to compare against QB-TRANS dataset of size 800k.

4.4 Results and Analysis

Our transformations lag behind a model trained directly on NQ by only about six points on average, while the LLMs lags by over ten points. QB-TRANS data can be applied to different QA systems and achieve comparable performance (Figure 2). While we expect QB-RAW to do poorly, it shows how much our transformation improves upon the original dataset.

LLM-based transformation (QB-GPT and QB-LLama2) performs worse than syntactic NATURALIZATION. As discussed in Section 3.2, not only does the desired answer change in LLM-based transformation (it is not clear that there is a correct answer), but the answer also sometimes appears in the question (despite prompt instructions).

5 Supervised QA System with QB-NQ training data

We combine all of the naturalized datasets with the original NQ dataset (NQ), with the goal of having the largest NQ-like dataset and highest accuracy.

5.1 Supervised QA systems

As the baseline, we use the top model in the NQ challenge leaderboard **ReflectionNet** (Wang et al., 2020a): an MRC model for answer prediction and Reflection model for answer confidence. We also use the state-of-the-art **GENREAD** (Yu et al., 2023a), which is a *generate-then-retrieve* pipeline QA system that directly generates the contextual documents by using clustering document representations. This method outperforms traditional *retrieve-then-read* methods. We also use the two retrieval-based systems **DPR** (Karpukhin et al., 2020b) and **REPLUG** (Shi et al., 2024) from the previous section, but this time trained with QB-TRANS data along with NQ dataset.

5.2 Training Data

We train the supervised QA systems with our QB-NQ dataset, the combination of original NQ and QB-TRANS questions. Here, QB-NQ-20, represents the filtered and transformed QB-TRANS dataset and 20% percent of the original NQ data. NQ examples are selected uniformly at random. We also use multiple passes when differences in dataset size like zero-shot setting. More detail on the formation of training questions and answers in Appendix C.

5.3 Supervised Classifier

Training our supervised NQ system requires a balance of NQ and NQ-like data. However, the generation process results in many questions that insufficiently resemble the Cranfield questions we want to emulate: some are too short or long, do not make sense, or still look too much like a Manchester QB elicitation. Like how Goodfellow et al. (2014) use a classifier to filter the outputs of an automatic generative process, we identify the best examples from the above process. We use a simple logistic regression classifier (Cox, 1958) trained on the generated NQ-like examples (through the process described in the previous section) as negative examples and with real NQ examples as positive examples. Our features identify question topics and formats that occur frequently in NQ. For example, the bigram “who played”, reflects NQ’s emphasis on popular

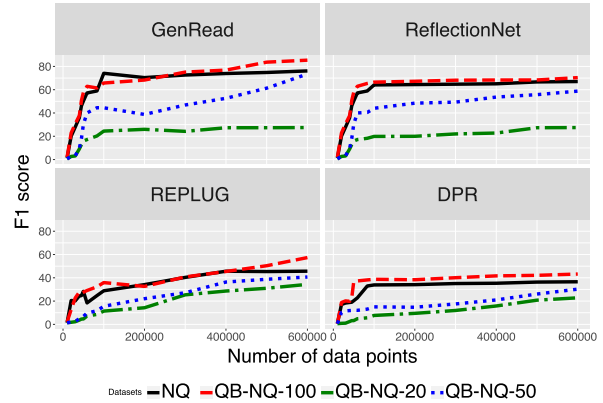


Figure 3: **GENREAD**: Supervised training on **QB-NQ-100** and evaluating on NQ test set produces the highest F1 score system with **GENREAD**. However, the cheaper datasets from our systematic conversion (**QB-NQ-50**), with a noisier but larger dataset, reached within 7 points of the F1 score of NQ training systems. In similar supervised settings, **ReflectionNet**, **REPLUG**, **DPR**: **QB-NQ-100** data crosses the NQ by 12 points on average compared to trained only on NQ, and adding just 50% of NQ data (QB-NQ-50) allows the model to reach within 10 points on average of the F1 score of the model trained on the NQ dataset. QB-Trans adding with NQ in training QA system can achieve F1 much higher (10 points on average on four systems) to NQ training system.

culture; starting questions with “how”, “when”, or “where” recapitulates the process for harvesting NQ; and short questions have the highest feature weight, emphasizing that NQ questions are short. To make use of the answers provided in the dataset, we also include the answers as a feature in the classifier.

The QA system training has early stopping to decide how much NQ-like data to use. At each epoch, we add 50k QB-TRANS data (sorted by classifier score) along with existing QB-TRANS and NQ train set and test it on NQ dev set until the F1 score stops increasing. When the score starts to drop we continue it for five more iterations to avoid local minima. If F1 again starts to increase, we continue. Otherwise, the QB-TRANS data number that has the best F1 score on the dev set is chosen as the optimal train set to be included along with NQ train set. For example, in the first epoch, we take classifier-scored top 50k data. In the second epoch, we use the next best scored 50k data and the previous 50k data along with the NQ and retrain the system.

5.4 Result and Analysis

Section 4 argues that using transformed QB-TRANS data would be cheaper than using NQ data (which

Models	Datasets			
	NQ	QB-NQ-100		
		<i>No classifier</i>	<i>With classifier</i>	
		<i>no early stopping</i>	<i>early stopping</i>	
DPR	39.23	43.54	46.21	49.12
REPLUG	45.75	55.29	49.12	57.56
ReflectionNet	64.01	68.36	73.89	75.87
GenRead	74.31	79.56	85.03	78.01

Table 3: The best F1-score is reported here. The classifier with early stopping helps us to find out the optimal number of data points needed for the model.

is expensive) to gather answers. What if we have access to a *fraction* of the NQ data? Finally, given the best configuration of the previous experiment, we add small amounts of NQ data to see how much is needed to recreate the best NQ result. No data in the training process is changed. Adding half of the NQ brings parity to the result. Therefore, our experiments show the effectiveness of QB-TRANS dataset as an alternative of NQ dataset in the zero-shot setting and an expansion of NQ dataset in supervised QA systems. Similar results can be seen in all the systems (Figure 3). REFLECTIONNET and GENREAD have higher F1 score than DPR and REPLUG because of their usage of large language models and ensemble models in training. The result is summarised in Table 3.

6 Analysis of Transformed Questions

This section discuss the quality of our dataset compared to NQ. We incorporate answer equivalence to the experiment with the goal of improving F1 score. Finally we prove the generalization of NATURALIZATION by showing how well our transformations can apply to non-NQ data (as evaluated on the NQ test set) and how well our transformed data can answer non-NQ data.

6.1 Quality Analysis of QB-TRANS and NQ

To analyze the quality of our dataset, we use CREPE (Yu et al., 2023b) to identify false presuppositions (Table 4). Our dataset has fewer presuppositions than NQ.

NQ has more ambiguous questions, as found using Min et al. (2020)’s AmbigQA binary classifier and GPT-3.5 (Table 4). An example of an ambiguous question from NQ is “*How many nominations does Game of Thrones have?*” This question can ask about the number of nominations “Game of Thrones” has across all its seasons, or it can ask about any particular season or award ceremony. Therefore, no precise answer can be given without additional context. However, QB elicitation gener-

Dataset	Size	% of Presupposition	% of Ambiguity	
			using GPT-3.5	using AmbigQA
NQ	307373	21	63	68
QB-Trans	800000	16	27	25

Table 4: The percentage of harmful presupposition and ambiguous questions in NQ and QBTrans dataset. QB-Trans has fewer presuppositions and significantly fewer ambiguities than NQ.

ally ensures each clue points to a unique answer without any ambiguity (given its rigorous editing).

6.2 Transformation Error Analysis

Not all of the original elicitations are transformed correctly. Consider this original elicitation:

This author created a character who smokes a cigarette before the body of his dead mother, and who vacations with his friend Raymond and shoots an Arab on the beach.

The heuristic “split conjunction” and “no wh-word” are applied and generate questions “This author created a character who smokes a cigarette before the body of his dead mother”, “what author vacations with his friend Raymond,” and “what author shoots an Arab on the beach”. The second and third questions are incorrect. This happens because there is an error in finding relative clauses.

6.3 Answer Equivalence in Zero-shot and Supervised Training

While Section 5 focuses on ensuring that the transformed questions resemble the target NQ data as much as possible, it did not consider the answers. To fully emulate NQ data, the answers need to be comparable. Thus, we expand the answer set provided in the QB dataset (which typically is more formal and verbose than NQ) with the WikiData answer equivalence sets from Si et al. (2021) for both training and evaluation.

For example, NQ has a question “Where do the greasers live in the outsiders?” with the gold answer set comprised of {“Tulsa”, “Oklahoma”}. However, if the QA system answers “Tulsa, Oklahoma”, it will be considered incorrect in the exact match. Thus, we apply an answer equivalence system to change the answer set to {“Tulsa”, “Oklahoma”, “ttown”, “Tulsa”, “Tulsa Oklahoma”, “Wagoner county Tulsa city”}. After adding answer equivalence in the supervised setting, the F1 score for QB-NQ-100 increases by 12 points on average from NQ which is three points more than systems without answer equivalence on the similar experiment on four models (GENREAD, REFLECTIONNET, REPLUG, DPR) from Section 5. Moreover,

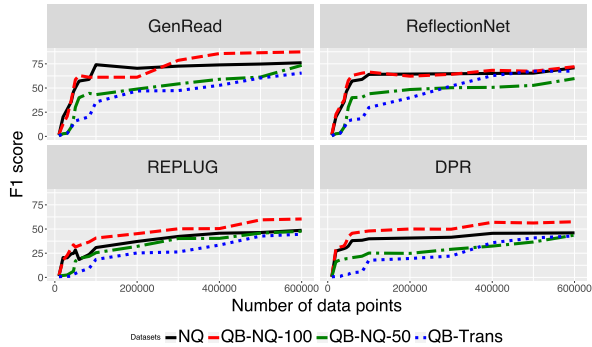


Figure 4: **With answer equivalence: QB-NQ-100** data crosses by 12 points on average of a model trained on NQ, and adding just 50% of the NQ data allows the model to reach within 7 points of the whole NQ with answer equivalence. **QB-TRANS** comes within 4 points of model trained on NQ.

on the same four models, the F1 score for QB-NQ-50 is much closer (two points improvement) to NQ than without answer equivalence. In zero-shot setting, with answer equivalence, the gap between the F1 score for QB-TRANS and NQ closes to four from six (consistent with results in Si et al. (2021)) (Figure 4) on same experiment from Section 4.

6.4 Cost of Heuristics and Generalization

Our NATURALIZATION technique needed multiple iterative cycles to fine-tune and optimize the heuristics. This systematic approach allowed us to acquire accuracy in under one hundred hours ensuring both effectiveness and efficiency.

All these heuristics can be directly applied to other pyramidal and clue-based question-answering datasets and generate NQ-like data at a cheaper cost without going through each clue manually.

To show the generalization of our heuristics, we apply the heuristics to different datasets. For example, *Jeopardy!* has an elicitation:

This small, red summer fruit develops tiny seeds on the outside and often tops shortcake.

After applying the heuristics described in Section 3.1 the question becomes

Which small, red summer fruit develops tiny seeds on the outside?
Which small, red summer fruit often tops shortcake?

We apply these heuristics to similar clue-based datasets *Jeopardy!* (Jeo, 2024), *TriviaQA* (Joshi et al., 2017), *HotpotQA* (Yang et al., 2018) and the Japanese dataset *AI King* (Aik, 2024). Examples of the original questions from these datasets and transformed questions after applying our heuristics

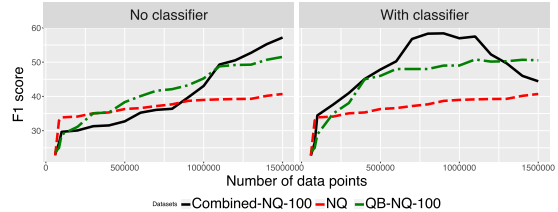


Figure 5: **No classifier:** The combined dataset shows similar performance initially with the model trained on NQ and QB-NQ. However, when we increase the data point, it goes 12 points higher than the model trained on NQ. With the **classifier**, the classifier chose the training data to resemble NQ. Therefore, the data selected earlier produces a better F1 score. However, after 110k data points, the performance starts to deteriorate. That means the data we add does not resemble NQ after that.

Models	Datasets			
	NQ	No classifier		With classifier
		no early stopping		early stopping
DPR	39.23	52.20	53.48	57.54
REPLUG	45.75	58.35	57.10	60.92
ReflectionNet	64.01	75.91	77.96	79.89
GenRead	74.31	80.98	82.90	86.87

Table 5: The best F1-score on NQ test is reported here. The classifier with early stopping based on NQ dev helps us to find out the optimal number of data points.

are in Appendix Table 15 and 16.

This generated combined dataset creates larger training data for models leading to improvement in performance. Figure 5 shows the application of heuristics to other datasets can generate larger datasets and this combined dataset (COMBINED-NQ-100) can improve the F1 score for DPR. We can significantly increase the size of datasets by applying these heuristics automatically to different language and domain datasets which can increase the system’s F1 score compared to the system solely trained on NQ. The results of these datasets are shown in Table 5. Table 14 shows the percentage of error our heuristics have while applying to different domain and language datasets is less than 1%. Our heuristics can also detect errors (e.g. ill-formed sentences, ambiguous clues about the entity, etc.) in the datasets. For example, in the *Jeopardy!* elicitation “Hits hard”, it is impossible to answer that without more context or needs additional category information to answer.⁹ Our heuristics can be applied to identify them.

⁹One needs to know the category “Clothing words” to know the answer is “socks/ belts.”

Models	Before finetuning	After finetuning using generated dataset
GPT-3.5-TURBO (ACHIAM ET AL., 2023)	70.0	72.1
GPT-4o-MINI (ACHIAM ET AL., 2023)	82.0	83.2
LLAMA2-7B (TOUVRON ET AL., 2023)	45.3	49.9
LLAMA2-13B (TOUVRON ET AL., 2023)	54.8	58.3
FLAN-PALM (CHUNG ET AL., 2024)	72.1	75.3
FALCON (ALMAZROUEI ET AL., 2023)	57.1	60.0

Table 6: The best average accuracy on the MMLU dataset is reported here. The LLMs fine-tuned with our generated dataset (QB-NQ-100-Jeopardy-TriviaQA-AI King-HotpotQA) help improve accuracy (an improvement of three points on average).

6.5 Generalization to QA with non-NQ Data

NQ is a part of several alignment datasets (Yang, 2023; Herzig et al., 2021), therefore, we see if this can improve modern LLMs. We experimented with how well our transformed data can answer non-NQ data. We used our transformed dataset to train systems and tested it on the Massive Multi-task Language Understanding (MMLU) benchmark dataset (Hendrycks et al., 2021) which consists of exam questions from 57 tasks ranging from history, mathematics, law, and computer science. We fine-tuned the LLM models such as GPT, LLAMA2, etc with our generated dataset and saw an average accuracy improvement of 3 points (Table 6) on the MMLU set. This points to the generalizability of our dataset in both settings—our NATURALIZATION can be applied to non-NQ data (Section 6.4) and our transformed data can answer non-NQ data.

7 Related Work

This section discusses the question generation and transformation from existing datasets, which is an effective alternative to expensive data collection.

7.1 Generating Questions

Given the expense of gathering these data, an obvious alternative is to generate your data. While we transform one question format into another, Probably Asked Questions (Lewis et al., 2021, PAQ) transforms source documents into questions that *could* be asked. These questions are more formulaic than those carefully crafted by trivia experts in the QB dataset, but an obvious extension would be to see if PAQ questions could help augment the results. Another transformed question class is translated questions that convert datasets like SQUAD into multiple languages (Carrino et al., 2020; d’Hoffschmidt et al., 2020). A frequent research thrust has been to create methods to generalize these datasets, either by merging datasets

together (Artetxe et al., 2019; Khashabi et al., 2020) or by QA-driven slot-filling (Du et al., 2021b) or event extraction via QA (Lyu et al., 2021) by creating algorithms that explicitly generalize (Munteanu et al., 2004; Munteanu and Marcu, 2005) or use existing algorithm for different use cases (Liang et al., 2023; Gou et al., 2023).

7.2 Transforming Questions

Our approach of transforming the form of QB elicitation is inspired by a long line of research. Machine translation models are used to transform questions to resemble the text where the answer would be found (Wang et al., 2007; Chen et al., 2013) or to transform a context-dependent question into more closely resembled NQ question (Demszky et al., 2018). More related work about other QA dataset, large language models, and zero-shot QA system is in Appendix, Section F.

8 Conclusion and Future Work

Transformed NQ-like questions from the QB data is an alternative to expensive datasets like NQ. The transformed data itself is not as good as NQ by itself, but is competitive; this is a reasonable option if the resources are not available to curate a dataset like NQ. However, the dataset is getting old with obsolete questions and out-of-date answers (Zhang and Choi, 2021). If there is a budget to create a dataset comparable to NQ, a small amount of this data augmented with transformed data from a dataset like QB can surpass a model trained on the NQ dataset. This can act as a continuous flow of new natural questions. Moreover, as no new NQ tests sets are published, this can provide an alternative benchmark to the obsolete eval NQ data. For future work, we can apply this conversion technique to other languages’ probing datasets (Han et al., 2023) where transformation heuristics can be learned using human data.

There are methods like reinforcement learning from human feedback (RLHF) that use NQ along with other datasets (Li et al., 2024; Feng et al., 2023) or create new datasets aligning NQ with other datasets for LLMs (Yang, 2023) or create adversarial dataset (Eisenschlos et al., 2021) or rank skill with complex questions (Joshi et al., 2017). Our work shows that there are additional sources of information that are cheaper and more recent that can feed into these datasets instead of NQ.

9 Limitations

Focus on Natural Questions We focus on NQ, a popular and respected dataset. It contains real user questions from Google on a variety of topics and they are natural queries. This diversity helps in training QA models and is suitable as a benchmark for the evaluation of QA systems. Other datasets are different, and we do not know how well our transformations would generalize to other datasets. However, we suspect that similar transformations would also succeed.

Errors hidden by Correct Answers While our transformed data often gets to the right answer, we have not systematically verified that the produced questions are themselves correct. It could be that enough of the necessary contents within the conversions remain that systems can reach the correct answer but that the questions contain errors (either factual or grammatical). From our inspection of the questions, we do not believe this to be the case, but a systematic evaluation would be needed to confirm this. However, this would dramatically raise the cost of the dataset, obviating one of the motivations for this approach.

Distribution Shift QB and NQ have very different distributions: QB is more academic, while NQ has more questions about sports and pop culture. Thus, solely evaluating on NQ potentially says little about how well our conversion process works for the topics that are over-represented in QB compared to NQ. While NQ does have some questions about literature and science, they are under-represented; it could be that our transformations are particularly brittle on questions about equations or works of fiction but NQ evaluation does not expose that weakness.

Ethical Considerations

The most important ethical consideration of this paper is that we are using the data from the trivia community to train a model. In contrast to datasets like SearchQA (Dunn et al., 2017) or TriviaQA (Joshi et al., 2017) where it is unclear how the original trivia authors feel about the use of the data, the QB community explicitly welcomes the sharing and dissemination of the data to train QB players: datasets are covered by a creative commons license (and the norm of sharing indeed predates the formal creation of creative commons). While computer QA systems are a different kind of trivia player

(machine rather than human), we believe that this would be in the spirit of the community.

Acknowledgement

This work is supported by NSF Grant IIS-2403436. We thank the University of Maryland Institute for Advanced Computer Studies (UMI-ACS) for their continuous help and support with the computational resources for the project. We are thankful to the reviewers for the valuable comments and helpful suggestions which helped us improve our paper’s clarity and quality. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

References

2024. AI King Japan Quiz AI Championship. <https://sites.google.com/view/project-ai0/home?authuser=0>.
2024. Jeopardy! dataset. <https://huggingface.co/datasets/jeopardy-datasets/jeopardy>.
- OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, and Shyamal Anadkat. 2023. [Gpt-4 technical report](#).
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, M rouane Debbah,  tienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al. 2023. The falcon series of open language models. *CoRR*.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2019. [On the cross-lingual transferability of monolingual representations](#). *CoRR*, abs/1910.11856.
- Pratyay Banerjee and Chitta Baral. 2020. [Self-supervised knowledge triplet learning for zero-shot question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 151–162, Online. Association for Computational Linguistics.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv:2004.05150*.
- Jordan Boyd-Graber and Benjamin B rschinger. 2020. [What question answering can learn from trivia nerds](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7422–7435, Online. Association for Computational Linguistics.
- Jordan Boyd-Graber and Benjamin B rschinger. 2020. [What question answering can learn from trivia nerds](#). In *Association for Computational Linguistics*.
- 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*, 33:1877–1901.
- Casimiro Pio Carrino, Marta R. Costa-juss , and Jos  A. R. Fonollosa. 2020. [Automatic Spanish translation of SQuAD dataset for multi-lingual question answering](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5515–5523, Marseille, France. European Language Resources Association.
- Boxing Chen, George Foster, and Roland Kuhn. 2013. Adaptation of reordering models for statistical machine translation. In *Proceedings of the Association for Computational Linguistics*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Debby RE Cotton, Peter A Cotton, and J Reuben Shipway. 2024. Chatting and cheating: Ensuring academic integrity in the era of chatgpt. *Innovations in Education and Teaching International*, 61(2):228–239.
- David R Cox. 1958. The regression analysis of binary sequences. *Journal of the Royal Statistical Society: Series B (Methodological)*, 20(2):215–232.
- Dorottya Demszky, Kelvin Guu, and Percy Liang. 2018. Transforming question answering datasets into natural language inference datasets. *arXiv preprint arXiv:1809.02922*.
- Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathy Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V. Le, Yonghui Wu, Zhifeng Chen, and Claire Cui. 2021a. [Glam: Efficient scaling of language models with mixture-of-experts](#). *CoRR*, abs/2112.06905.
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. [Glam: Efficient scaling of language models with mixture-of-experts](#). In *International Conference on Machine Learning*, pages 5547–5569. PMLR.
- Xinya Du, Luheng He, Qi Li, Dian Yu, Panupong Papat, and Yuan Zhang. 2021b. Qa-driven zero-shot slot filling with weak supervision pretraining. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 654–664.

- Matthew Dunn, Levent Sagun, Mike Higgins, V Ugur Guney, Volkan Cirik, and Kyunghyun Cho. 2017. Searchqa: A new qa dataset augmented with context from a search engine. *arXiv preprint arXiv:1704.05179*.
- Martin d’Hoffschmidt, Wacim Belblidia, Quentin Heinrich, Tom Brendlé, and Maxime Vidal. 2020. Fquad: French question answering dataset. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1193–1208.
- Julian Eisenschlos, Bhuwan Dhingra, Jannis Bulian, Benjamin Börschinger, and Jordan Boyd-Graber. 2021. Fool me twice: Entailment from Wikipedia gamification. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 352–365, Online. Association for Computational Linguistics.
- Falcon. 2024. Falcon 7B. <https://falconllm.tii.ae/falcon-models.html>. [Online; accessed 8-May-2024].
- C. Fellbaum. 1998. *WordNet : An Electronic Lexical Database*, chapter A semantic network of English verbs. MIT Press, Cambridge, MA.
- Tao Feng, Zifeng Wang, and Jimeng Sun. 2023. Citing: Large language models create curriculum for instruction tuning. *arXiv preprint arXiv:2310.02527*.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. **Generative adversarial nets**. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.
- Qi Gou, Zehua Xia, Bowen Yu, Haiyang Yu, Fei Huang, Yongbin Li, and Nguyen Cam-Tu. 2023. Diversify question generation with retrieval-augmented style transfer. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1677–1690.
- HyoJung Han, Jordan Boyd-Graber, and Marine Carpuat. 2023. **Bridging background knowledge gaps in translation with automatic explicitation**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9718–9735, Singapore. Association for Computational Linguistics.
- He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. 2016. Opponent modeling in deep reinforcement learning. In *International conference on machine learning*, pages 1804–1813. PMLR.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. **Measuring massive multitask language understanding**. In *ICLR*. OpenReview.net.
- Jonathan Herzig, Thomas Mueller, Syrine Krichene, and Julian Eisenschlos. 2021. Open domain question answering over tables via dense retrieval. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 512–519.
- Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. **TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension**. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020a. **Dense passage retrieval for open-domain question answering**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020b. **Dense passage retrieval for open-domain question answering**.
- Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Han-naneh Hajishirzi. 2020. **UNIFIEDQA: Crossing format boundaries with a single QA system**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1896–1907, Online. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. **Natural questions: A benchmark for question answering research**. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenetorp, and Sebastian Riedel. 2021. Paq: 65 million probably-asked questions and what you can do with them. *Transactions of the Association for Computational Linguistics*, 9:1098–1115.
- Lei Li, Yekun Chai, Shuohuan Wang, Yu Sun, Hao Tian, Ningyu Zhang, and Hua Wu. 2024. **Tool-augmented reward modeling**. In *The Twelfth International Conference on Learning Representations*.

- Yuanyuan Liang, Jianing Wang, Hanlun Zhu, Lei Wang, Weining Qian, and Yunshi Lan. 2023. Prompting large language models with chain-of-thought for few-shot knowledge base question generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4329–4343.
- Qing Lyu, Hongming Zhang, Elicor Sulem, and Dan Roth. 2021. Zero-shot event extraction via transfer learning: Challenges and insights. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 322–332.
- Inbal Magar and Roy Schwartz. 2022. Data contamination: From memorization to exploitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 157–165.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. **AmbigQA: Answering ambiguous open-domain questions**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5783–5797, Online. Association for Computational Linguistics.
- Dragos Stefan Munteanu, Alexander Fraser, and Daniel Marcu. 2004. Improved machine translation performance via parallel sentence extraction from comparable corpora. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 265–272.
- Dragos Stefan Munteanu and Daniel Marcu. 2005. Improving machine translation performance by exploiting non-parallel corpora. *Computational Linguistics*, 31(4):477–504.
- Arvind Narayanan. 2023. Gpt-4 and professional benchmarks: the wrong answer to the wrong question. <https://www.aisnakeoil.com/p/gpt-4-and-professional-benchmarks>.
- Joakim Nivre. 2010. Dependency parsing. *Language and Linguistics Compass*, 4(3):138–152.
- OpenOrca. 2024. OpenOrca - Mistral - 7B - 8k. <https://huggingface.co/Open-Orca/Mistral-7B-OpenOrca>. [Online; accessed 8-May-2024].
- 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*, 35:27730–27744.
- Jiezhong Qiu, Hao Ma, Omer Levy, Wen-tau Yih, Sinong Wang, and Jie Tang. 2020. **Blockwise self-attention for long document understanding**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2555–2565, Online. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. **Exploring the limits of transfer learning with a unified text-to-text transformer**. *Journal of Machine Learning Research*, 21(140):1–67.
- Pedro Rodriguez and Jordan Boyd-Graber. 2021. **Evaluation paradigms in question answering**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9630–9642, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Pedro Rodriguez, Shi Feng, Mohit Iyyer, He He, and Jordan Boyd-Graber. 2019. **Quizbowl: The case for incremental question answering**. *CoRR*, abs/1904.04792.
- Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2023. Qa dataset explosion: A taxonomy of nlp resources for question answering and reading comprehension. *ACM Computing Surveys*, 55(10):1–45.
- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023a. **NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark**. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787, Singapore. Association for Computational Linguistics.
- Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, and Eneko Agirre. 2023b. **Did chatgpt cheat on your test?**
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. **Detecting pretraining data from large language models**.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. Replug: Retrieval-augmented black-box language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8364–8377.
- Chenglei Si, Chen Zhao, and Jordan Boyd-Graber. 2021. **What’s in a name? answer equivalence for open-domain question answering**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9623–9629, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hao Sun, Xiao Liu, Yeyun Gong, Anlei Dong, Jingwen Lu, Yan Zhang, Daxin Jiang, Linjun Yang, Rangan Majumder, and Nan Duan. 2023. **Beamsearchqa: Large language models are strong zero-shot qa solver**.

- Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2021. CommonsenseQA 2.0: Exposing the limits of ai through gamification. In *Proceedings of Advances in Neural Information Processing Systems*.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, and Shruti Bhosale. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *ArXiv*, abs/2307.09288.
- A. M. Turing. 1950. [Computing machinery and intelligence](#). *Mind*, LIX(236):433–460.
- Ellen M Voorhees. 2019. *The evolution of cranfield*, pages 45–69. Springer.
- Mengqiu Wang, Noah A Smith, and Teruko Mitamura. 2007. What is the jeopardy model? a quasi-synchronous grammar for qa. In *Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL)*, pages 22–32.
- Xuguang Wang, Linjun Shou, Ming Gong, Nan Duan, and Daxin Jiang. 2020a. [No answer is better than wrong answer: A reflection model for document level machine reading comprehension](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4141–4150, Online. Association for Computational Linguistics.
- Xuguang Wang, Linjun Shou, Ming Gong, Nan Duan, and Daxin Jiang. 2020b. [No answer is better than wrong answer: A reflection model for document level machine reading comprehension](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4141–4150.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. [Finetuned language models are zero-shot learners](#). In *International Conference on Learning Representations*.
- BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Lucicioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, and Pawan Sasanka Ammanamanchi. 2023. [Bloom: A 176b-parameter open-access multilingual language model](#).
- Hui Yang, Lekha Chaisorn, Yunlong Zhao, Shi-Yong Neo, and Tat-Seng Chua. 2003. Videoqa: question answering on news video. In *Proceedings of the eleventh ACM international conference on Multimedia*, pages 632–641.
- Jianxin Yang. 2023. Longqlora: Efficient and effective method to extend context length of large language models. *arXiv preprint arXiv:2311.04879*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, S Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023a. Generate rather than retrieve: Large language models are strong context generators. In *International Conference on Learning Representations*.
- Xinyan Yu, Sewon Min, Luke Zettlemoyer, and Hananeh Hajishirzi. 2023b. Crepe: Open-domain question answering with false presuppositions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10457–10480.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297.
- Hang Zhang, Yeyun Gong, Yelong Shen, Weisheng Li, Jiancheng Lv, Nan Duan, and Weizhu Chen. 2021. Poolingformer: Long document modeling with pooling attention. In *International Conference on Machine Learning*, pages 12437–12446. PMLR.
- Michael Zhang and Eunsol Choi. 2021. Situatedqa: Incorporating extra-linguistic contexts into qa. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7371–7387.

Appendix A provides a complete list of all the heuristics used in the NATURALIZATION method. These heuristics are applied based on preconditions that ensure only relevant modifications are made to the elicitation (Appendix B). We transform answers of the QB dataset to resemble the answer structure of the NQ (Appendix C). For our zero-shot experiment, we explain the zero-shot system and how we have computed whether LLMs have seen NQ in training data which supports the removal of the state-of-the-art LLMs from zero-shot QA (Appendix D). We also provide the details of DPR training and their results in Appendix D. Appendix E explains the comparison in results of two baseline LLMs in our experiment—LLAMA2 and GPT. Finally, we give some more related work in this field in Appendix F.

A Heuristics List

Through observation of the linguistic and grammatical style of NQ we add additional heuristics to further improve the candidates such as **removing punctuation** and **adding subject**:

- **punctuation**: Natural questions typically do not include punctuation, so we remove punctuation at the boundary of a generated question.
- **Adding subject**: If a question is missing a subject (e.g., “wrote *Burmese Days*”, we add “which” answer_type (in this example, author) to the beginning of the question.

Full list of heuristics in Table 8 and 9.

B Process of Application of heuristics

We have applied all the heuristics to all the questions with some precondition to determine the applicability of those heuristics. For example, when we apply “remove conjunctions” heuristics, we determine whether that particular question has a conjunction (via a dependency parse). If it has a conjunction, only then that heuristics will be applied. Otherwise, the question goes to the next heuristics unchanged. Similarly, for “Imperative to Interrogative” heuristic checks whether the subject of that question is imperative and if it is, converts it to interrogative. The algorithm is given in Algorithm 1.

C Answer Formation in QB

We also transform answers from the QB dataset to look like the NQ data. For example, one of the QB questions after transformation “Which ethnic

group’s language and customs were adopted by a majority of the uru people?” with the answer “Aymara people (the Quechua were the larger group targeted by the genocide)”. However, if we observe the NQ answer list, there is no description given using the parenthesis. Therefore, we convert the answer set to also include “Aymara people” to make the answer set look like NQ formatted.

D Zero-shot QA with QB-TRANS Data

D.1 What is a zero-shot system?

Zero-shot systems enables the models to answer the questions without explicitly trained on them. Under zero-shot setting for the NQ dataset, there can be no training on NQ data— not with questions and their answers and not with their contextual documents. Therefore, when given any NQ test data, the zero-shot systems directly encode the given question and predict the answer. A question q is given to the model as the input. Based on that input, the model generates the answer a denoted by $p(a|p, \theta)$ where θ is the model parameters (Yu et al., 2023a).

The state-of-the-art zero-shot QA system AL-LIES (Sun et al., 2023) framework generates additional questions through an iterative process. In this process, an LLM is used to generate queries based on existing query-evidence pairs and score the answer. This iteration process continues until the score reaches a predefined threshold. Therefore, this system decomposes the original question into multiple sub-questions and achieves state-of-the-art performance on the zero-shot setting for the NQ dataset. Another state-of-the-art zero-shot model GENREAD Yu et al. (2023a) uses the large language model InstructGPT (Ouyang et al., 2022) to directly generate contextual documents from a given question.

D.2 Min K% probability

To design a fair zero-shot system to compare NQ with QB, we first detect whether NQ data exists in the training data of an LLM by using Shi et al. (2023)’s Min K% probability technique. This technique utilizes minimum token probabilities of a text for detecting data in pertaining. The hypothesis is that a member example in training data does not have words with a high negative log-likelihood. The average log-likelihood of K-% tokens is computed using

LLM name	Min K% probability
GLAM (Du et al., 2021a)	71.1%
FLAN (Wei et al., 2022)	62.9%
PALM (Chowdhery et al., 2023)	68.3%
LLAMA (Chowdhery et al., 2023)	57.0%
T-5 (RAFFEL ET AL., 2020)	77.9%
BLOOM (WORKSHOP ET AL., 2023)	64.4%
MISTRALORCA (OPENORCA, 2024)	47.1%
FALCON (FALCON, 2024)	55.2%

Table 7: We validate if NQ is present in their pretraining data by MIN-K(K=60)% PROB (Shi et al., 2023). A high average probability suggests that the NQ is likely part of the pertaining data. We can see for all the state-of-the-art LLMs, the probability is 63% on average. Thus, we can say, these models likely have NQ in their training data.

$$Min-K(\%)Prob(x) = \frac{1}{E} \sum_{x_i \in Min-K\%(x)} \log P(x_i | x_1, \dots, x_{i-1}) \quad (1)$$

After feeding in an NQ sample into the model, we use the technique to yield Min K% probability by taking k% tokens with minimum probabilities with K=60 and calculating their average log-likelihood. Based on the hypothesis in Shi et al. (2023), if the log-likelihood is high, then NQ is likely to exist in the model’s training data.

D.3 DPR Training

The passages that contain any of the answer strings are positive examples, while the passages that do not are negative examples. One example is shown in Table 13.

D.4 Zero-shot Training and Results

We use individual elicitation sentences from the QB dataset *without* any transformation: **QB-Raw**. While we expect this to do poorly, it shows how much our transformation improves upon the original dataset.

E Comparison of LLMs and Error in Transformation

We use LLAMA2 baseline because of the cost efficiency. Both GPT and LLAMA2 showed similar conversion (Table 12). However, LLAMA2 baseline results are comparable to the GPT models. For example, training with the first 10000 examples ends

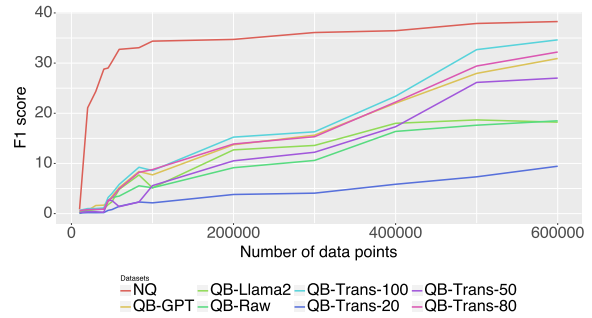


Figure 6: QB-Trans can replace NQ in training QA system and achieve accuracy close to NQ training system. As expected, **QB-Trans-100** without any NQ data comes within 5 points of a model trained on NQ. Training on the full QB-Trans and evaluating it produces the highest accuracy system with DPR. However, the percentage of that dataset from our systematic conversion (**QB-Trans-80**) reaches a substantial fraction of the accuracy. This does better than conversions created by prompting a LLM.

with an accuracy of 0.58 for GPT and 0.45 accuracy for LLAMA2. Similarly, when we have 50000 samples for both models, the accuracy is 3.13 for GPT and 2.64 for LLAMA2. We can see both the language models perform worse than the rule-based conversion in the QA systems. That is why we can say, the rule-based system (QB-TRANS) performs better irrespective of language model choice as the baseline (Figure 6).

F Related Work

F.1 An Explosion of Datasets

The last few years have seen a flurry of datasets. Some of these datasets are created at great expense through crowdsourcing to capture common sense, numerical reasoning, visual QA (Antol et al., 2015), video QA (Yang et al., 2003), common sense questions (Talmor et al., 2021) or multicultural questions (Clark et al., 2020); Rogers et al. (2023) gives a thorough summary. Less common are datasets focusing on found data, although there is nonetheless a panoply of questions harvested from educational resources, civil service exams, users, and trivia games.

F.2 Large Language Models and Transformer-based Models

Due to the increasing sequence length, the transformer uses sparse attention to handle the complexity of long document modeling (Zhang et al., 2021). In this method, each token is made to attend to a

more important context or local context (Qiu et al., 2020). Another approach uses a sliding window pattern to capture local information that includes Longformer (Beltagy et al., 2020), BigBird (Zaheer et al., 2020). Lastly, PoolingFormer (Zhang et al., 2021) uses full self-attention in a two-level attention schema—the first one works as a sliding window attention pattern and the second level increases the receptive field. Wang et al. (2020b) uses a machine reading comprehension (MRC) model for answer prediction and a Reflection model for answer confidence. This achieves state-of-the-art performance on the NQ dataset in the leader board of NQ challenge.

F.3 Zero-shot QA

In a zero-shot setting, the large language model generates new questions. In BeamSearchQA (Sun et al., 2023), new questions are generated using LLM by iterative refining and expanding the scope of the question to achieve a state-of-the-art EM score of 38.0, there are some approaches without the retriever. The in-context learning approach is applied using GPT-3 (Brown et al., 2020), cost-efficient Generalist Language Model (GLaM) GPT-3 (Du et al., 2022), instruction-tuned model (Wei et al., 2022) in zero-shot setting. Self-supervised knowledge learning is applied in zero-shot QA, for example, heuristic-based graph (Banerjee and Baral, 2020). However, we are creating nq-like questions from qb questions in our work. The main difference between our work from the previous work is that we are using a different dataset to train the model in a zero-shot to make it compatible with the NQ dataset. With a proper classifier and carefully chosen heuristics, we introduce a conversion of different domain datasets as a replacement of the NQ dataset.

In rewriting elicitations into questions, we need to replace uncommon, odd answer mentions (e.g., “this polity”) with more traditional ones (e.g., “this country”). Thus, we count all mentions used to refer to an answer a , then store the most frequent in M . This becomes the canonical mention in Algorithm 5 which we will always use for rewriting questions.

Algorithm 1 Transform QB Questions to NQ-like Question. We split clues from QB questions into elicitation questions (QB_E) and applied various heuristics to transform them and maintain proper syntax.

```

1: Split each clue in QB questions into QB elicitation ( $QB_E$ ) by splitting them through period(.)
2: procedure APPLY HEURISTICS FOR TRANSFORMER( $QB_E$ )
3:   Heuristics list ( $H$ ) $\leftarrow$ {Split Conjunction, Imperative to Integrative, No Wh-words, . . . }
4:   for each  $QB_e \in QB_E$  do
5:     for each  $h \in H$  do
6:       Flag $\leftarrow$ PreCondition( $QB_e$ )                                      $\triangleright$  Check if heuristic can be applied to  $QB_e$ 
7:       if Flag is True then
8:          $QB_e \leftarrow h(QB_e)$                                         $\triangleright$  Apply the heuristic to  $QB_e$ 
9:          $QB_e \leftarrow$ PostCondition( $QB_e$ )                          $\triangleright$  Check for syntax errors after applying the heuristic application
10:      else
11:         $QB_e$  is unchanged
12:      end if
13:    end for
14:  end for
15: end procedure

```

Algorithm 2 In transforming QB clues into NQ-like questions, we split the clues via conjunction and construct two independent clauses by splitting them. We give question ‘q’ as input and the algorithm returns two separate questions (first question, second question) (split by conjunctions if applicable)

```

1: procedure POS(word)
2:   Return parts of speech of ‘word’
3: end procedure
4: procedure DEP(word)
5:   Return dependency of ‘word’ in the parse tree
6: end procedure
7: procedure POSITION(word)
8:   Return position of ‘word’ in the question q
9: end procedure
10: procedure PARSE(question)
11:   Return parse tree of question
12: end procedure
13: if question contains conjunctions then
14:   Parse(q)  $\leftarrow$  parse tree for the question
15:   root verb  $\leftarrow$  [ $x \in$  Parse(q) | PoS(x)= "VERB" and x has no ancestors in Parse(q)]
16:   verbs = [ $x \in$  Parse(x) | PoS(x)= "VERB" and x.head  $\in$  root verb]
17:   verb conj  $\leftarrow$  []                                              $\triangleright$  Initialize an empty list for verb-conjunction pairs
18:   for verb  $\in$  verbs do
19:     for child  $\in$  verb.children do
20:       if PoS(child) = coordinating conjunction then
21:         verb conj.add((verb, child))
22:       end if
23:     end for
24:   end for
25:   for (verb, conj)  $\in$  verb conj do                                $\triangleright$  Check to see if this is the second verb and if it has no ancestors
26:     if Position(verb) > Position(verbs[0]) and verb has no ancestor in Parse(q) then  $\triangleright$  Two independent clauses found,
yield parts around the conjunction
27:       First question  $\leftarrow$  [ $x \in$  Parse(q) and Position(x) < Position(conj)]
28:       Second question  $\leftarrow$  [ $x \in$  Parse(q) and Position(x) > Position(conj)]
29:     else if Position(verb) < Position(verbs[-1]) and Dep(verbs[-1]) = "conj" then  $\triangleright$  Two sentences with the same subject,
get what is before the verb that does not modify it
30:       left tokens  $\leftarrow$  [ $x \in$  Parse(q) | Position(x) < Position(verb) and not (head(x) == verb and (PoS(x) = "ADVERB" or
"AUX"))]
31:       first verb  $\leftarrow$  [ $x \in$  Parse(q) | Position(x) < Position(conj) and  $x \notin$  left tokens]
32:       second verb  $\leftarrow$  [ $x \in$  Parse(q) | Position(x) > Position(conj) and  $x \notin$  left tokens]
33:       First question  $\leftarrow$  left tokens + first verb
34:       second question  $\leftarrow$  left tokens + second verb
35:     end if
36:   end for
37: end if

```

Algorithm 3 No Wh-words: In converting question with for No Wh-words we need to introduce wh-words. We determine the appropriate transformation and modify the question accordingly.

```

1: Flag ← Check if question has no wh-words
2: if Flag is True then                                     ▷ No wh-words found in the question
3:   answer type ← Find the canonical type of the answer for the question
4:   if question contains “this” then
5:     final question ← Replace “this” with ”which” in the question
6:   else if If the subject of the question is pronoun then
7:     final question ← Replace the subject of the question with “which” + answer type
8:   else
9:     final question ← Add “which” + answer type at the beginning of the question
10:  end if
11: end if

```

Algorithm 4 Heuristics for Imperative to Interrogative: If the question starts with verbs like “name,” “give,” or “identify”, it converts it to standardized imperative question form.

```

1: procedure PARSE(question)
2:   Return parse tree of question
3: end procedure
4: procedure INTERROGATIVE(question)
5:   Patterns ← {(ftp | FTP | Ftp) (give | identify | name) (this | these) }, {(For | for) (ten | 10 | 20 | 5 | 15) (Points | points | points)} (give | identify | name) (this | these)}
6:   for x do ∈ Patterns such that isSubstring(x,q)
7:     verb position ← find the minimum position of verbs [“name”, “give”, “identify”] in Parse(q)
8:     head = the head of the verb using verb position in Parse(q)                                     ▷ Get the first noun after the verb
9:     if There is a relative clause in the children for the head in the dependency for the parse tree then
10:      relative head ← relative clause’s head from the parse tree                                     ▷ Find the relative clause head
11:      relative head ← first element in relative head list
12:      continuation ← concatenate text from Parse(q)[relative head’s left edge + 1 : relative head’s right edge + 1]
13:      final question ← “Which” + answer type + continuation
14:     else if length of parse tree > head’s index + 1 AND parse [head’s index + 1] is comma then
15:      continuation ← concatenate text from parse(q)[head’s index + 2:]
16:      final question ← answer type + “is” + continuation
17:     else
18:      reduced ← question after cutting off the “For 10 ... points [name/identify]”
19:      final question ← “Which is the” + reduced
20:     end if
21:   end for
22: end procedure

```

Algorithm 5 Find Canonical Answer Type. In rewriting elicitations into questions, we need to replace uncommon, odd answer mentions (e.g., “this polity”) with more traditional ones (e.g., “this country”). Thus, we count all mentions used to refer to an answer a , then store the most frequent in M . This becomes the canonical mention we will always use for rewriting questions.

```

1: Mention count  $C := |a| \times |m|$  zero array
2: for Elicitation  $e$ , Answer  $a$  in Dataset do
3:   for Noun Phrase  $n \in \text{Parse}(e)$  do                                     ▷ The mention could be any noun phrase.
4:     if Yield( $n$ )[0] ∈ { this, these, ... } then                                     ▷ Mentions start with specific determiners.
5:       Mention  $m \leftarrow \text{Yield}(n)[1 : ]$ 
6:        $C[a][m] \leftarrow C[a][m] + 1$                                      ▷ Record all mentions of this answer
7:     end if
8:   end for
9: end for
10: Canonical Mention  $M := a \mapsto m$ 
11: for Answer  $a \in C$  do
12:    $M[a] \leftarrow \arg \max_m C[a][m]$                                      ▷ The canonical mention is the most frequent
13: end for

```

Heuristic	Purpose	Example before Heuristic	Example after Heuristic
substitute non answer pronouns	Substitute non answer pronouns to noun+possession.	she founded Carthage and reigned as its queen from 814-759 BC	she founded Carthage and reigned as carthage's queen from 814-759 BC
clean marker	Remove punctuation patterns at the beginning and the end of the question.	which german philosopher is this philosopher wrote a work , . "	which german philosopher also wrote glowing reviews of which german philosopher's own works in ecce homo
drop after semicolon	Remove contents after semicolon in NQlike.	which molecule is this compound 's presence can be quantified in spectrophotometry by observing an intense absorption peak at 255 nanometers ; that peak is the	which molecule 's presence can be quantified in spectrophotometry by observing an intense absorption peak at 255 nanometers
convert continuous to present	Add verb if elicitation has verbal	which particle consisting of a charm quark and an anti - charm quark	which particle consists of a charm quark and an anti - charm quark
fix no wh words	Convert "this" to "which"+answer_type when there's no "wh-" words.	this play begins with the protagonist arriving at the elysian fields to see her sister stella	which play begins with the protagonist arriving at the elysian fields to see her sister stella
replace this is	Replace "this" to "which"+answer_type within "this is" pattern.	this is the first party name , followed by kraemer , in that supreme court case , which held that racially restrictive covenants are unconstitutional	which name the first party name , followed by kraemer , in that supreme court case , which held that racially restrictive covenants are unconstitutional
replace which with that	Convert "which" to "that" and check if no "which" present anymore, if so, convert "this" to "which".	michael green is a current professor at this university , which is where watson and crick discovered dna 's structure	michael green a current professor at which university , that is where watson and crick discovered dna 's structure
add question word	Adding "which"+answer_type when no "wh-" words present.	a chamberlain named cleaner was killed on the orders of marcia , a mistress of this man who was involved in the plot that eventually assassinated him and replaced him with pertinax	a chamberlain named cleaner killed on the orders of marcia , a mistress of which man who was involved in the plot that eventually assassinated him and replaced him with pertinax
add subject	Add "which"+answer_type at the beginning when question starting with VERB/AUX and missing the subject.	were refused real employment because of " logical discrimination , " an excuse which belied the employers ' fear of their " death taint	which se people were refused real employment because of " logical discrimination , " an excuse which belied the employers ' fear of their " death taint
fix what is which	Remove "what is" from "what is which".	what is which desert lying mostly in northern china and mongolia	which desert lying mostly in northern china and mongolia
remove end BE verbs	Remove "is/are" at the end of NQlike questions.	which jewish holiday is that hymn is	which jewish holiday is that hymn
remove extra AUX	Remove extra auxiliary words.	which number is it is the base for solutions to the differential equation	which number is the base for solutions to the differential equation
remove patterns	Remove bad patterns in NQlike.	This country is home to the author (*) of Miss Brill, Bliss	Which country is home to the author (*) of Miss Brill, Bliss
remove rep subject	remove repetition of the subject "is this".	which goddess is this goddess is considered a daughter of ra	which goddess is considered a daughter of ra
remove BE determiner	Change is his/is her/is its to 's.	which greek goddess's is her wedding night lasted three hundred years	which greek goddess's wedding night lasted three hundred years
remove repeated pronoun	Removes repeated pronouns like "which character who is", "is who is".	which character who is the character who never appears to linus in a peanuts halloween special	which character never appears to linus in a peanuts halloween special

Table 8: List of Heuristics

Heuristic	Purpose	Example before Heuristic	Example after Heuristic
fix no verb	Ensure there's at least one verb per question.	which greek god wielding chief greek god	which greek god is wielding chief greek god
add space before punctuation	Add space before punctuation because in NQ there's space before all types of punctuation	which greek goddess's wedding night lasted three hundred years	which greek goddess 's wedding night lasted three hundred years
rejoin whose	replace "who's" with "whose"	which wife who 's kidnaping by paris began the trojan war	which wife whose kidnaping by paris began the trojan war

Table 9: List of Heuristics.

Original QB	Answer	QB-TRANS	LLAMA2
Performing "electrodeoxidation" on an oxide of this metal may be able to improve on the current method of producing it and is called the Fray-Farthing-Chen Process.	Titanium	Performing "electrodeoxidation" on an oxide of which metal may be able to improve on the current method of producing it and is called the Fray-Farthing-Chen Process?	What metal can be produced through "electrodeoxidation"?
His government also endured the Dreyfus Affair.	French Third Republic	Whose government also endured the Dreyfus Affair?	What government was established in France after the fall of Napoleon III in 1870?
This country is home to the author of "Miss Brill," "Bliss," and a story in which Mr. Scott's death fails to derail the title event, despite Laura Sheridan's objections.	New Zealand	which country is home to the author of "Miss Brill," "Bliss," and a story in which Mr. Scott's death fails to derail the title event, despite Laura Sheridan's objections?	What country is home to Katherine Mansfield?
For 10 points, name this homeland of the authors of The Bone People and "The Garden Party," Keri Hulme and Katherine Mansfield.	New Zealand	What is the homeland of the authors of The Bone People and "The Garden Party," Keri Hulme and Katherine Mansfield?	What homeland did Keri Hulme and Katherine Mansfield come from?
Communist historiography blames the loss of this battle on the government's brutal suppression of the peasant rebellion of Gregory Dozsa twelve years earlier.	Battle of Mohacs	Communist historiography blames the loss of which battle on the government's brutal suppression of the peasant rebellion of Gregory Dozsa twelve years earlier?	What battle did the Ottoman Empire lose in 1526, which led to the death of the emperor and the loss of much of its territory?
A later battle fought at this site saw troops under Charles of Lorraine claim victory in 1687.	Battle of Mohacs	A later battle fought at which site saw troops under Charles of Lorraine claim victory in 1687?	What site was the location of a significant battle fought between the Ottoman Empire and the Holy Roman Empire in 1526?
This period ended after the failure of the Dyle Plan.	French Third Republic	Which period ended after the failure of the Dyle Plan?	What period in French history saw the rise of the Third Republic?

Table 10: LLAMA2 converted question can introduce error in question. For example, conversion of QB clues using rule-based and large language model-based techniques (LLAMA2). In the case of LLAMA2, words were removed from the question, and additional random words were added or the entity of the question was changed, implying LLAMA2 transformation is worse than that of QB-TRANS.

Original QB	Answer	QB-TRANS	GPT
For 10 points, name this tongue spoken by the members of O-Zone and Nicolae Ceaufescu, an Eastern Romance language spoken in Bucharest.	Daco-Romanian	What language is spoken by the members of O-Zone and Nicolae Ceaufescu, an Eastern Romance language spoken in Bucharest?	What are the distinguishing features of the language spoken in Bucharest?
This state’s largest city endured an 1855 Rum Riot.	Maine	Which state’s largest city endured an 1855 Rum Riot?	What was the 1855 Rum Riot in the largest city of the state?
One of this state’s Native American tribes won back 300,000 acres of land from the state government in Passamaquoddy v. Morton.	Maine	One of which state’s Native American tribes won back 300,000 acres of land from the state government in Passamaquoddy v. Morton?	What was the outcome of Passamaquoddy v. Morton and which Native American tribes in Maine were involved in the case?
In addition to the De Laceys, another character in this work educated himself by reading about the Northwest Passage.	Frankenstein, or the Modern Prometheus	In addition to the De Laceys, a character in which work educated himself by reading about the Northwest Passage?	Who is the character in Frankenstein who educated himself by reading about the Northwest Passage?
MCTs are needed to pump this substance out of the cell.	Lactic acid	MCTs are needed to pump which substance out of the cell?	What is the role of MCTs in removing a byproduct of anaerobic respiration from cells?
This substance moves to the liver where it is turned into glucose in the Cori cycle.	Lactic acid	Which substance moves to the liver where it is turned into glucose in the Cori cycle?	What happens to the byproduct of anaerobic respiration in muscles and how is it converted into glucose in the body?
Performing “electrodeoxidation” on an oxide of this metal may be able to improve on the current method of producing it and is called the Fray-Farthing-Chen Process.	Titanium	Performing “electrodeoxidation” on an oxide of which metal may be able to improve on the current method of producing it and is called the Fray-Farthing-Chen Process.?	What is the Fray-Farthing-Chen Process and how does it relate to improving the production of a certain metal oxide?

Table 11: GPT converted question can introduce error and hallucinations in question. In the case of GPT like LLAMA2, words were removed from the question, and additional random words were added or the entity of the question was changed, implying GPT transformation is also worse than that of QB-TRANS.

Original clue in QB	GPT converted question	Llama converted question
For 10 points, name this native of Rokovoko and savage companion of Ishmael in Moby-Dick.	Who is Ishmael’s savage companion in Moby-Dick from Rokovoko?	What native of Rokovoko and savage companion of Ishmael in Moby-Dick?
This state’s largest city endured an 1855 Rum Riot, put down with the help of 1880 Prohibitionist Candidate for President Neal S. Dow.	What state’s largest city experienced the 1855 Rum Riot and was aided by Neal S. Dow in 1880?	What state’s largest city endured an 1855 Rum Riot, put down with the help of 1880 Prohibitionist Candidate for President Neal S. Dow?

Table 12: Conversion of QB elicitation using two large language model-based technique (GPT and Llama). In the case of both models, we can see, the converted questions are similar.

Question	A fortification overlooking which city was renamed “narin qala” or “little fortress” by mongolinvaders in the 13th century.
Answer	Tbilisi
Positive context	City in the Caucasus, with its at least 50,000 inhabitants and thriving commerce. Several intellectuals born or living in Tbilisi, bearing the nisba al-Tiflisi were known across the Muslim world. The Abbasid Caliphate weakened after the Abbasid civil war in the 810s, and caliphal power was challenged by secessionist tendencies among peripheral rulers, including those of Tbilisi . At the same time, the emirate became a target of the resurgent Georgian Bagrationi dynasty who were expanding their territory from Tao-Klarjeti across Georgian lands. The Emirate of Tbilisi grew in relative strength under Ishaq ibn Isma’il, who was powerful enough to
Negative context	near the shores of Kasagh River, during the reign of king Orontes I Sakavakyats of Armenia (570 ² 013560 BC). However, in his first book “Wars of Justinian”, the Byzantine historian Procopius has cited to the city as “Valashabad” (Balashabad), named after king “Valash” (Balash) of Armenia. The name evolved into its later form by the shift in the medial “L” into a “Gh”, which is common in the Armenian language. Movses Khorenatsi mentioned that the Town of Vardges was entirely rebuilt and fenced by king Vagharsh I to become known as “Noarakaghak” (“New City”) and later “Vagharshapat”. The territory of

Table 13: We have a QB question: *A fortification overlooking which city was renamed “narin qala” or “little fortress” by mongolinvaders in the 13th century.* with answer *Tbilisi*. Now, for the positive context of the DPR training we have used those passage which contain the answer string and the rest of the passages are selected as negative context. One of the examples of positive contexts and negative contexts for this question is shown here.

Dataset	Size	Wrong	Examples of Error in Original Dataset	Comment
Trivia QA	138384	859(0.620%)	There are around 60.000 miles of veins, arteries and capillaries in the human body. True or false? We all knew him as Radar, but was the actual first name of the pride of Ottumwa, Iowa, Corporal O’Reilly on the TV series MASH?	There are some true/false questions in TriviaQA. In our heuristics of “no wh-words”, it is wrongly transformed.
Jeopardy	216930	35(0.016%)	Hits hard 1 of the 2 born in Vermont	No words to generate the question
AI King	22335	155(0.693%)	Is Ichiro a right-handed or left-handed batter in the major leagues? In horse racing, a “10,000 horse racing ticket” refers to a horse racing ticket with multiple odds? Will the 2020 Olympics in Tokyo be the Summer Olympics or the Winter Olympics?	There are some yes/no and either/or questions in the dataset. We have no heuristics to handle those clues.
Hotpot QA	90447	21(0.023%)	Are Patrick White and Katherine Anne Porter both writers? Did both Carl Boese and Franco Zeffirelli direct and produce film? Are Pam Veasey and Jon Jost both American?	There are some yes/no questions in the dataset. We have no heuristics to handle those clues.

Table 14: Error analysis of four clue-based datasets after applying our heuristics. We can see from the above analysis, is that our heuristics mostly fail to convert questions when there is an error in the question or the question is specific to the context of the game.

Original Question	Heuristic Applied from List in 3.1	Syntactic Transformed Question
Dataset Name: Jeopardy		
For the last 8 years of his life, Galileo was under house arrest for espousing this man’s theory	No wh-words	For the last 8 years of his life, Galileo was under house arrest for espousing which man’s theory
The city of Yuma in this state has a record average of 4,055 hours of sunshine each year	No wh-words	The city of Yuma in which state has a record average of 4,055 hours of sunshine each year
In 1963, live on "The Art Linkletter Show", this company served its billionth burger		In 1963, live on "The Art Linkletter Show", which company served its billionth burger
Signer of the Dec. of Indep., framer of the Constitution of Mass., second President of the United States’		Who is Signer of the Dec. of Indep., framer of the Constitution of Mass., second President of the United States’
In the title of an Aesop fable, this insect shared billing with a grasshopper		In the title of an Aesop fable, which insect shared billing with a grasshopper
In the winter of 1971-72, a record 1,122 inches of snow fell at Rainier Paradise Ranger Station in this state		In the winter of 1971-72, a record 1,122 inches of snow fell at Rainier Paradise Ranger Station in which state
This housewares store was named for the packaging its merchandise came in & was first displayed on Cows regurgitate this from the first stomach to the mouth & chew it again		Which housewares store was named for the packaging its merchandise came in & was first displayed on Cows regurgitate this from the first stomach to the mouth & chew it again
In 1000 Rajaraja I of the Cholas battled to take this Indian Ocean island now known for its tea		In 1000 Rajaraja I of the Cholas battled to take which Indian Ocean island now known for its tea
Dataset Name: TriviaQA		
Name the 1980’s hit sung by Tina Turner and Rod Stewart?	Imperative to Interrogative	What is the 1980’s hit sung by Tina Turner and Rod Stewart?
Name the two tiles with the highest score in Scrabble?		What is the two tiles with the highest score in Scrabble?
Name the Dick Francis mount that collapsed approaching the finishing line in the 1956 ‘Grand National’?		What is the Dick Francis mount that collapsed approaching the finishing line in the 1956 ‘Grand National’?
Name the 1972 musical starring David Essex as Jesus Christ?		What is the 1972 musical starring David Essex as Jesus Christ?
Name the male lead in the 1946 film The Big Sleep?		Who is the male lead in the 1946 film The Big Sleep?
Name the stretch of water separating Anglesey from the Welsh mainland?		What is the stretch of water separating Anglesey from the Welsh mainland?
For a point each, name the characters in a bottle of Flintstones Chewable Vitamins.		What is the characters in a bottle of Flintstones Chewable Vitamins.
For a point each, name the state(s) bordering Maine		What is the state(s) bordering Maine
Name the year: NAFTA is ratified, Nancy Kerrigan gets clubbed, Kurt Cobain eats his shotgun, OJ Simpson offs his ex wife and her friend.		What is the year: NAFTA is ratified, Nancy Kerrigan gets clubbed, Kurt Cobain eats his shotgun, OJ Simpson offs his ex wife and her friend.

Table 15: To show the generalization of our dataset, we applied the heuristics from Section 3.1 to different domain datasets. At first, heuristics are applied to two similar clue-based datasets– *Jeopardy!* and *TriviaQA*. We can see, for similar clue-like questions’ datasets like QB, our heuristics convert them into NQ-like questions successfully.

Original Question	Heuristic Applied from List in 3.1	Syntactic Transformed Question
Dataset Name: AI King official distribution dataset		
In 1960, while studying abroad from Nankai, he achieved a record of 5 wins, 1 loss, and 9 seasons in his one year on the job, and was promoted to the San Francisco Giants, becoming the first Japanese major leaguer.	Split Conjunction and No wh words	In 1960, while studying abroad from Nankai, who achieved a record of 5 wins, 1 loss, and 9 seasons in his one year on the job, Who was promoted to the San Francisco Giants, becoming the first Japanese major leaguer. In 1960, while studying abroad from Nankai, who achieved a record of 5 wins, 1 loss, and 9 seasons in his one year on the job, and was promoted to the San Francisco Giants, becoming the first Japanese major leaguer.
It is Germany’s second largest trading port after Hamburg, and is also featured in the Grimm fairy tales that feature musical bands.		What is Germany’s second largest trading port after Hamburg, and is also featured in the Grimm fairy tales that feature musical bands? What is Germany’s second largest trading port after Hamburg? What is featured in the Grimm fairy tales that feature musical bands?
This fish is said to have gotten its name from the fact that it eats by cutting its body into two?		Which fish is said to have gotten its name from the fact that it eats by cutting its body into two, but why are its ovaries called “herring roe”?
On July 16th of this year, Katsura Saegusa will become the 6th generation of the famous Kamigata Rakugo story.		On July 16th of which year, Katsura Saegusa will become the 6th generation of the famous Kamigata Rakugo story.
Dataset Name: Hotpot QA		
This is the place of fish and is the capital city of Frobisher Bay south?	Split conjunction and No wh words	1. Which is the place of fish and is the capital city of Frobisher Bay south? 2. Which is the place of fish? 3. Which is the capital city of Frobisher Bay south?
This Ghanaian footballer was a notable graduate of SC Bastia Reserves and Academy?		Which Ghanaian footballer was a notable graduate of SC Bastia Reserves and Academy?
Name one comedy series that stars the younger brother of Arthur White ?		Which comedy series that stars the younger brother of Arthur White ?
Bottom Points railway station is on a heritage railway system that is situated near this town?		Bottom Points railway station is on a heritage railway system that is situated near which town?
Barry Moltz taught entrepreneurship as an adjunct professor in this city?		Barry Moltz taught entrepreneurship as an adjunct professor in which city?
Adebayo Akinfenwa was a star in the 2006 Football League Trophy Final, but know plays for this team?		Adebayo Akinfenwa was a star in the 2006 Football League Trophy Final, but know plays for which team?
Topics covered by this author include corporate control of government, the harshness of war, gender polarities and sexual identity.		Topics covered by which author include corporate control of government, the harshness of war, gender polarities and sexual identity.

Table 16: To show the generalization of our dataset, we applied the heuristics from Section 3.1 to different domain datasets. At first, heuristics are applied to a different lingual dataset (Japanese). Secondly, it is applied to a multi-hop dataset HotpotQA. We can see, for similar clue-like questions’ datasets like QB, our heuristics convert them into NQ-like questions successfully.