

# Generative Data Augmentation using LLMs improves Distributional Robustness in Question Answering

**Arijit Ghosh Chowdhury**

University of Illinois Urbana Champaign

arijit10@gmail.com

**Aman Chadha**

Stanford University

Amazon GenAI\*

hi@aman.ai

## Abstract

Robustness in Natural Language Processing continues to be a pertinent issue, where state of the art models under-perform under naturally shifted distributions. In the context of Question Answering, work on domain adaptation methods continues to be a growing body of research. However, very little attention has been given to the notion of domain generalization under natural distribution shifts, where the target domain is unknown. With drastic improvements in the quality of and access to generative models, we answer the question: How do generated datasets influence the performance of QA models under natural distribution shifts? We perform experiments on 4 different datasets under varying amounts of distribution shift, and analyze how "in-the-wild" generation can help achieve domain generalization. We take a two-step generation approach, generating both contexts and QA pairs to augment existing datasets. Through our experiments, we demonstrate how augmenting reading comprehension datasets with generated data leads to better robustness towards natural distribution shifts.

## 1 Introduction

In this work, we perform a systematic study of how "in-the-wild" generation can affect the distributional robustness of question-answering models trained on the popular Stanford Question Answering Dataset (SQUAD) (Rajpurkar et al., 2016). Synthetic data generation is a widely adopted method for domain adaptation in QA systems (Shakeri et al., 2020) (Yue et al., 2021) (Yue et al., 2022). However, domain adaptation methods have access to unlabelled/labelled data belonging to the target domain, and do not account for unseen natural distribution shifts. Our work studies the effect of generated data on distribution shifts where the target domain is unseen.

The conception of a dataset has undergone significant evolution in recent times. This transformation has been catalyzed by the advent of generative models trained 'in-the-wild', such as those described in (Brown et al., 2020), (Bubeck et al., 2023), and (Touvron et al., 2023). These models, which use vast and diverse datasets across a range of domains, have facilitated the infusion of the web with synthesized data of high calibre, applicable to an extensive array of conceptual topics. Interestingly, these models are not merely confined to generation based on a pre-established distribution; they possess the capacity for repeated prompting, resulting in the creation of markedly diverse data. In the context of this emerging model paradigm, our research investigates the following query: How do generated datasets affect the distributional robustness of Question Answering models? Specifically, **natural distribution shifts** in NLP can arise due to differences in the text genre and style, text topics and vocabulary, demographics of the authors, medium of the text (written vs spoken), and other attributes (Wang et al., 2022). A key challenge is that NLP models trained on one data distribution often fail to generalize well to these naturally occurring shifts. For instance, (Miller et al., 2020) found that question answering models experienced average F1 score drops of 3.8 points on news articles, 14 points on Reddit posts, and 17.4 points on Amazon reviews compared to Wikipedia articles. This reveals brittleness of NLP models to natural distribution shifts.

We present an overview of our generation setup in Figure 1. For generating data, use GPT-3.5 (Brown et al., 2020), and create a question-answering dataset using questions provided in the SQUAD (Rajpurkar et al., 2016) dataset. We use a dual generation approach, by first prompting the language model to generate a context for a question given in the SQUAD dataset, and then generating question-answer pairs for the newly generated con-

---

\*Work does not relate to position at Amazon.

text.

Recent surveys, such as (Ramponi and Plank, 2020), discuss domain adaptation in NLP and divide approaches into *data centric* and *model centric*. We take a data-centric approach, as highlighted by findings from (Wang et al., 2022) that demonstrate overlap in test-train data for QA models. The scarcity of research on generalization in QA models, especially with natural distribution shifts, is a motivation for our work, backed by observations from (Arora et al., 2021) on out-of-distribution data in NLP.

Initial experiments like (Longpre et al., 2019) ventured into domain-agnostic question answering using data augmentation. New datasets introduced by (Miller et al., 2020), sourced from various platforms, emphasize the effect of natural distribution shifts on QA models. While these studies provide extensive evaluations, our work builds on them by focusing on the impact of large language model (LLM)-generated datasets for QA tasks and further leveraging these datasets for our data augmentation method.

The benefits of generated data have been explored by (Gowal et al., 2021), showing its potential in adversarial robustness. (Bartolo et al., 2021) and (Mekala et al., 2022) use synthetic and context-generated data respectively for QA and text classification. Our method uses a GPT-3.5 model, as described by (Wei et al., 2022), to generate context for questions. With similar motivations, (Bansal and Grover, 2023) demonstrates the application of Stable Diffusion in diverse dataset creation for image tasks.

#### OUR CONTRIBUTIONS

- ▣ We propose a framework to improve the distributional robustness of reading comprehension models in the presence of natural distribution shifts.
- ▣ Through a thorough quantitative evaluation, we evaluate the capabilities of LLMs to generate high quality synthetic data for question answering tasks.

## 2 Methodology

### 2.1 Context Generation

We first generate contexts by conditioning it on a question present in the SQUAD dataset. This allows the language model to generate a paragraph that can be used to generate question-answer pairs. Since the paragraph is generated using an existing question, the generated context is consistent

with the informative trivia format of SQUAD-like datasets. We also ensure that the generated contexts are diverse yet complimentary to the original dataset, as highlighted by (Gowal et al., 2021). To maintain further consistency, the generated context is clipped to be within 250 words, based on the average context length present in the SQUAD dataset. We prompt GPT 3.5 (gpt-3.5-turbo)<sup>1</sup> in the following manner: *Generate a paragraph which answers the following question: (question)* . Here the question is sampled from the SQUAD dataset. Figure 1 demonstrates the generation process. Additionally, the Appendix A contains examples from the generation process.

### 2.2 Question Answer Generation

After the context is created, the generated paragraph is used to create question-answer pairs. This is done by using a T5 based question generation model (Lopez et al., 2020) that is trained on the SQUAD dataset, which takes a paragraph as an input and returns a question-answer pair. We use the open source<sup>2</sup> implementation for this model. Additionally we also filter out QA pairs based on round-trip consistency (Alberti et al., 2019).

## 3 Experiments

### 3.1 Setup

We train an extractive reading comprehension model using SQUAD V1.1, using the RoBERTA-Base model across all our experiments. We use a learning rate of  $3e - 5$ , a batch size of 16 and run our experiments for 3 epochs each. We use the implementation provided by HuggingFace, and run our models on a stand-alone Nvidia A100 GPU provided by Google Colab. We do not use GPT-3.5 as a baseline since the purpose of this study is to specifically measure the performance by smaller models.

For all our experiments, we measure F1 and Exact Match scores to quantify performance on Natural Distribution Shift (NDS) datasets.

### 3.2 Datasets

We use the following datasets created by (Miller et al., 2020) to set up our testbed:

The **New Wikipedia** dataset contains newer QA pairs from wikipedia articles used by the SQUAD V1.1 dataset. Contains 7,938 test samples from

<sup>1</sup><https://platform.openai.com/docs/models>

<sup>2</sup><https://github.com/patil-suraj/question-generation>

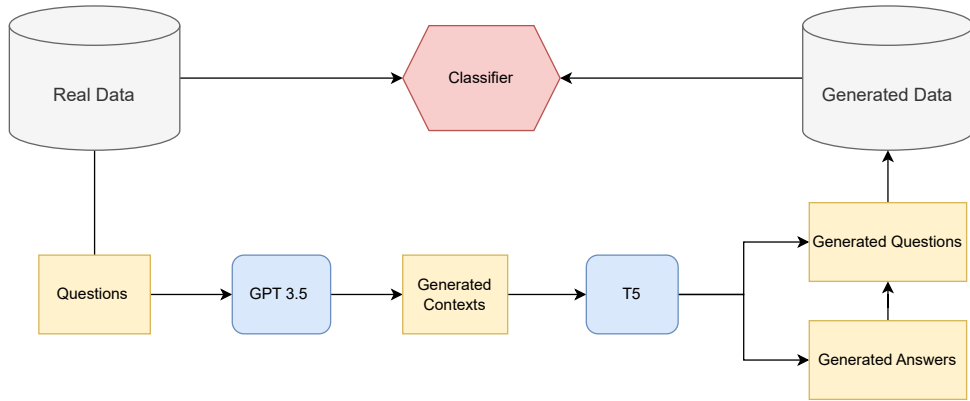


Figure 1: Overview of the generation system. Our method creates a generated dataset which is then augmented with the real dataset to train a question answering model.

Dataset	SQUAD		NewWiki		NYT		Amazon		Reddit	
	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM
Real data	90.4	83.0	89.4	79.2	86.4	76.1	79.9	66.4	80.1	67.1
Generated data	79.5	64.6	80.1	65.3	76.5	63.2	72.4	59.5	72.7	60.2
Real + Wiki-samples	93.4	85.2	89.3	77.3	79.4	78.1	76.4	66.6	78.8	63.2
Real + Generated data	<b>92.7</b>	<b>84.7</b>	<b>91.1</b>	<b>80.4</b>	<b>88.9</b>	<b>79.3</b>	<b>80.3</b>	<b>67.1</b>	<b>81.7</b>	<b>68.7</b>

Table 1: Generated datasets demonstrate robustness to natural distribution shifts.

48 contexts. The **New York Times** dataset contains articles from New York times which are then used to annotate QA pairs in the same format as SQUAD. It is ensured that the passage length statistics stay the same. Contains 10,065 test samples from 46 articles. **Reddit** dataset contains articles from Reddit where the authors concatenated each post’s title with its body. This dataset contains 9,803 test samples from 1,969 posts. The **Amazon Product Reviews** dataset contains user generated product reviews from the "Home and Kitchen" category on Amazon. This data contains 9,885 test samples from 1,909 reviews.

## 4 Results

### 4.1 Does generated data help with distributional robustness?

We evaluate the F1 and Exact Match scores of models trained with different datasets on natural distribution shifts (NDS) benchmarks. We note the average EM and F1 numbers across three random seeds in Table 1. The models are trained on an equal amount of real and generated data.

We find that the model, when trained on SQUAD, when subjected to natural distribution shift datasets, the model’s performance significantly deteriorates. A noteworthy observation was that exclusive training on the generated data resulted in substandard

performance on both the SQUAD and its Natural Distribution Shift (NDS) datasets. The inferior absolute performance could be potentially attributed to the distribution disparity between the source and the generated training datasets. Interestingly, we observe that for the model trained on the generated data, the performance gaps on the real validation dataset and its NDS datasets are low, which might be attributed to the benefits of training on diverse generated data. This highlights the contributions of the generated data in improving robustness, as opposed to simply generating more data for training.

We also sample paragraphs from Wikipedia and generate questions from those paragraphs, instead of letting GPT3.5 generate the paragraphs. This improves in-domain performance on SQUAD, but leads to drops in performance across out of domain datasets, further emphasizing on the effectiveness of the in-the-wild context generation on distribution shifts.

Finally, we expose our model to an evenly-distributed blend of real and generated datasets, with the goal of investigating the impact of generative augmentations. Our results reveal that the absolute performance of the model, when trained with a combination of real and generated data, either parallels or exceeds the performance of models trained exclusively on either real or generated

Dataset	SQUAD		NewWiki		NYT		Amazon		Reddit	
Metrics	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM
Real + 50% Generated data	91.4	81.1	90.4	82.2	87.4	77.1	79.7	65.4	80.3	67.4
Real + 100% Generated data	92.7	84.7	91.1	80.4	<b>88.9</b>	<b>79.3</b>	80.3	67.1	<b>81.7</b>	<b>68.7</b>
Real + 200% Generated data	<b>92.9</b>	<b>84.8</b>	<b>91.3</b>	<b>80.7</b>	88.5	79.1	<b>80.9</b>	<b>67.3</b>	80.8	68.1

Table 2: Performance on varying amounts of data. Using equal measures of real and generated data is essential.

Dataset	SQUAD		NYT		Amazon	
Metrics	F1	EM	F1	EM	F1	EM
Real data	90.4	83.0	86.4	76.1	79.9	66.4
Real + Generated data (Questions Only)	91.5	82.7	85.7	75.6	77.4	63.5
Real + Generated data (Contexts + Questions)	<b>92.7</b>	<b>84.7</b>	<b>88.9</b>	<b>79.3</b>	<b>80.9</b>	<b>67.3</b>

Table 3: Ablation Study demonstrating how context generation is key to robustness.

datasets, across all naturally distributed datasets. This observation suggests that the incorporation of real data into the training process is indeed essential for attaining superior absolute performance.

To summarize, while using solely generated data improves robustness at the expense of absolute performance, a blend of real and artificially generated data presents the ideal balance for robust and precise training.

#### 4.2 How much generated data is needed?

Here, we investigate how different combinations of the generated dataset can help the classifiers take advantage of the complementary strengths of the two data sources (Table 2).

To do so, we assessed the average performance of models trained with three different input mixing combinations created by using 50%, 100%, and 200% of the generated dataset. We observed an increase in performance on shifted datasets as the size of the generated data increases while keeping the amount of real data fixed. However, when the proportion of the generated data increases twofold while keeping the proportion of the real data fixed, we observe that the performance gains are only marginal. Additionally, we note that using only half of the generated data does not provide enough meaningful signal in terms of diversity and does not lead to major performance improvements compared to training on real data.

Overall, we found that the ideal split between real and generated data is a 50-50 split where the two datasets are able to compliment each other, in terms of providing both diversity and in-domain samples at the same time.

#### 4.3 Is context generation needed?

Table 3 demonstrates the importance of generating both contexts and questions for improving model robustness to distribution shifts. When only questions are generated for existing contexts, performance on the original SQuAD dataset improves slightly, while performance degrades substantially on the out-of-distribution NYT and Amazon datasets. This indicates that generating questions alone overfits models to the SQuAD distribution, reducing robustness. In contrast, generating both contexts and questions leads to consistent improvements in performance across all datasets. The dual generation approach enhances model robustness by exposing the model to more diversity during training, leading to better generalization. The results demonstrate that generating varied contexts in addition to targeted question generation is crucial for improving robustness to natural distribution shifts, rather than question generation alone.

### 5 Conclusion and Future Avenues

We created a framework that enhances the performance of reading comprehension models by supplementing real datasets with a diverse dataset generated by contemporary, real-world generative models. Our findings indicate that this training method yields superior results on test datasets and those with natural distribution shifts, due to the added robustness from training on the generated data as opposed to traditional methods. In the future, we want to explore a more extensive comparison against question generation methods and how this paradigm fits into fine-tuning larger models.



## References

- Chris Alberti, Daniel Andor, Emily Pitler, Jacob Devlin, and Michael Collins. 2019. [Synthetic qa corpora generation with roundtrip consistency](#).
- Udit Arora, William Huang, and He He. 2021. [Types of out-of-distribution texts and how to detect them](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10687–10701, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hritik Bansal and Aditya Grover. 2023. [Leaving reality to imagination: Robust classification via generated datasets](#).
- Max Bartolo, Tristan Thrush, Robin Jia, Sebastian Riedel, Pontus Stenetorp, and Douwe Kiela. 2021. [Improving question answering model robustness with synthetic adversarial data generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8830–8848, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4](#).
- Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and Timothy Mann. 2021. [Improving robustness using generated data](#).
- Shayne Longpre, Yi Lu, Zhucheng Tu, and Chris DuBois. 2019. [An exploration of data augmentation and sampling techniques for domain-agnostic question answering](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 220–227, Hong Kong, China. Association for Computational Linguistics.
- Luis Enrico Lopez, Diane Kathryn Cruz, Jan Christian Blaise Cruz, and Charibeth Ko Cheng. 2020. [Transformer-based end-to-end question generation](#). *ArXiv*, abs/2005.01107.
- Dheeraj Mekala, Tu Vu, Timo Schick, and Jingbo Shang. 2022. [Leveraging QA datasets to improve generative data augmentation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9737–9750, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- John Miller, Karl Krauth, Benjamin Recht, and Ludwig Schmidt. 2020. [The effect of natural distribution shift on question answering models](#).
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Alan Ramponi and Barbara Plank. 2020. [Neural unsupervised domain adaptation in NLP—A survey](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6838–6855, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Siamak Shakeri, Cicero Nogueira dos Santos, Henry Zhu, Patrick Ng, Feng Nan, Zhiguo Wang, Ramesh Nallapati, and Bing Xiang. 2020. [End-to-end synthetic data generation for domain adaptation of question answering systems](#). *arXiv preprint arXiv:2010.06028*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#).
- Xuezhi Wang, Haohan Wang, and Diyi Yang. 2022. [Measure and improve robustness in NLP models: A survey](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4569–4586, Seattle, United States. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Y. 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](#).
- Zhenrui Yue, Bernhard Kratzwald, and Stefan Feuerriegel. 2021. [Contrastive domain adaptation for question answering using limited text corpora](#). *arXiv preprint arXiv:2108.13854*.
- Zhenrui Yue, Huimin Zeng, Ziyi Kou, Lanyu Shang, and Dong Wang. 2022. [Domain adaptation for question answering via question classification](#). *arXiv preprint arXiv:2209.04998*.

## **Frequently Asked Questions (FAQs)**

### **\* How are we sampling questions to generate paragraphs?**

▮ One question is sampled per context in the original SQUAD dataset to condition the paragraph generation.

### **\* Why don't we generate new contexts from scratch?**

▮ This is done to create topically consistent datasets, and run a controlled experiment where only determine whether LLM generated contexts provide linguistic and stylistic diversity.

### **\* Why don't we prompt GPT-3/3.5/4 to generate QA pairs as well?**

▮ In our attempts to do so, we realised that GPT prompts do not generate outputs that stay faithful to the extractive nature of the SQUAD dataset. Furthermore, we wanted to limit our API usage.

## A Appendix: Example Generations

This section provides supplementary material in the form of additional examples to bolster the reader’s understanding of the concepts presented in this work.

<b>Original Context</b>	The American bison is a large bovid which inhabited much of western North America prior to the 1800s, living on the prairies in large herds. However, the vast herds of bison attracted market hunters, who killed dozens of bison for their hides only, leaving the rest to rot. Thousands of these hunters quickly eliminated the bison herds, bringing the population from several million in the early 1800s to a few hundred by the 1880s. Conservation efforts have allowed the population to increase, but the bison remains near-threatened.
<b>Sampled Question</b>	"What type of creature is the American bison?"
<b>Prompt</b>	Generate a paragraph that answers the following question: "What type of creature is the American bison?"
<b>Generated Context</b>	The American Bison, often colloquially referred to as buffalo, is a North American species of bison that once roamed across the vast expanses of the continent in large herds. These colossal creatures are mammals, belonging to the family Bovidae, which includes other large, hooved mammals such as domestic cattle and goats. Adult bison are imposing creatures; they are known to stand up to 6.5 feet tall at the shoulder and can weigh over a ton, making them the heaviest land animals in North America. They have a robust body structure with a large, muscular hump over the shoulders and a shaggy, dark brown winter coat which gets lighter in the summer. Their large heads sport a pair of short, sharp horns, used both for defense and in fights for status within the herd. Despite their bulk, bison are surprisingly agile and fast, capable of reaching speeds up to 35 mph. The American Bison plays a crucial role in maintaining the health and biodiversity of the prairie ecosystem.
<b>Generated Questions</b>	<b>Q:</b> To which species does the term "buffalo" colloquially refer in North America? <b>A:</b> American Bison, <b>Q:</b> How does the American Bison use its horns? <b>A:</b> for defense and in fights for status within the herd

Table 4: Example #1 – Context, Question, and Answer Generation.

<b>Original Context</b>	Exhibitions and annual horse shows in all districts and a national horse and cattle show at Lahore are held with the official patronage. The national horse and cattle show at Lahore is the biggest festival where sports, exhibitions, and livestock competitions are held. It not only encourages and patronises agricultural products and livestock through the exhibitions of agricultural products and cattle but is also a colourful documentary on the rich cultural heritage of the province with its strong rural roots.
<b>Sampled Question</b>	"What is Punjab’s biggest festival?"
<b>Prompt</b>	Generate a paragraph that answers the following question: "What is Punjab’s biggest festival?"
<b>Generated Context</b>	Punjab’s biggest festival is Baisakhi, also known as Vaisakhi, celebrated with great fervor and enthusiasm in the region. It’s a significant festival for the Punjabi people and is observed annually on the 13th or 14th of April. Baisakhi is not only a harvest festival, marking the time when farmers reap their rabi (winter) crops, but it also has deep religious significance, particularly for Sikhs. It commemorates the formation of the Khalsa Panth, a Sikh political and military group, by Guru Gobind Singh Ji in 1699. The day is marked by visiting gurdwaras, participating in processions, dancing the traditional Bhangra and Gidda, and sharing meals in the community. The joyous festival is a celebration of culture, heritage, and communal unity, encapsulating the vibrant spirit of Punjab.
<b>Generated Questions</b>	<b>Q:</b> When is Baisakhi observed? <b>A:</b> annually on the 13th or 14th of April, <b>Q:</b> Why is Baisakhi significant for farmers? <b>A:</b> a harvest festival, marking the time when farmers reap their rabi (winter) crops

Table 5: Example #2 – Context, Question, and Answer Generation.

<b>Original Context</b>	In Canada, the Royal Newfoundland Constabulary was founded in 1729, making it the first police force in present-day Canada. It was followed in 1834 by the Toronto Police, and in 1838 by police forces in Montreal and Quebec City. A national force, the Dominion Police, was founded in 1868. Initially the Dominion Police provided security for parliament, but its responsibilities quickly grew. The famous Royal Northwest Mounted Police was founded in 1873. The merger of these two police forces in 1920 formed the world-famous Royal Canadian Mounted Police.
<b>Sampled Question</b>	"When did Canada get its first police?"
<b>Prompt</b>	Generate a paragraph that answers the following question: "When did Canada get its first police?"
<b>Generated Context</b>	The first policing institution in Canada was established in the early 19th century. The Royal Newfoundland Constabulary, founded in St. John's, Newfoundland in 1729, is considered the oldest civil law enforcement body in North America. However, it primarily served as a paramilitary unit, dealing with the city's disorderly sailors and soldiers. Later, the North-West Mounted Police (NWMP), which is often seen as the direct predecessor to the current Royal Canadian Mounted Police (RCMP), was formed on May 23, 1873. The establishment of the NWMP marked a significant shift in the approach to law enforcement in Canada, transitioning from a military-style model to a community-focused approach. The NWMP, and later the RCMP, played a crucial role in establishing the rule of law across Canada, particularly in the west and north during the periods of rapid expansion and development in the late 19th and early 20th centuries.
<b>Generated Questions</b>	<b>Q:</b> What is the oldest civil law enforcement body in North America? <b>A:</b> The Royal Newfoundland Constabulary, <b>Q:</b> Which organization is seen as the direct predecessor to the Royal Canadian Mounted Police? <b>A:</b> the North-West Mounted Police

Table 6: Example #3 – Context, Question, and Answer Generation.

<b>Original Context</b>	The Earth of the early Archean (4,000 to 2,500 million years ago) may have had a different tectonic style. During this time, the Earth's crust cooled enough that rocks and continental plates began to form. Some scientists think because the Earth was hotter, that plate tectonic activity was more vigorous than it is today, resulting in a much greater rate of recycling of crustal material. This may have prevented cratonisation and continent formation until the mantle cooled and convection slowed down. Others argue that the subcontinental lithospheric mantle is too buoyant to subduct and that the lack of Archean rocks is a function of erosion and subsequent tectonic events.
<b>Sampled Question</b>	"During what time period was the Archean era?"
<b>Prompt</b>	Generate a paragraph that answers the following question: "During what time period was the Archean era?"
<b>Generated Context</b>	The Archean Eon, one of the four principal eons of Earth's geologic time scale, occurred approximately 4.0 to 2.5 billion years ago. It follows the Hadean Eon and precedes the Proterozoic Eon. The Archean era is divided into four geologic time periods, known as eras: the Eoarchean, Paleoarchean, Mesoarchean, and Neoarchean. This time period was crucial in Earth's history as it saw the formation of the planet's crust and the development of the earliest forms of life. Despite the extreme conditions, with high temperatures and a lack of oxygen, early life forms such as cyanobacteria began to emerge, which eventually played a role in shaping the Earth's atmosphere by producing oxygen through photosynthesis.
<b>Generated Questions</b>	<b>Q:</b> What eon follows the Hadean Eon and precedes the Proterozoic Eon? <b>A:</b> The Archean Eon, <b>Q:</b> Despite what conditions did early life forms begin to emerge during the Archean Eon? <b>A:</b> the extreme conditions, with high temperatures and a lack of oxygen

Table 7: Example #4 – Context, Question, and Answer Generation.