

# LLM Ability to Answer University Student Questions in CS Trained Exclusively on Auto-Generated QA

Anonymous ACL submission

## Abstract

With the advent of generative AI as commonplace to daily-life within the past couple of years, there is a strong need for extensive research as we apply this technology to educational contexts. This study supports that body of research as we explore two driving questions: **(1) Can we use LLMs to create synthetic student-like question-answer datasets? and (2) Can we train an LLM to embody an instructor in answering real student questions?** In this paper, we explore both of these questions, grounded by prior works and approaches. Ultimately, our findings suggest that synthetic QA data can be generated, but still requires significant improvement to aptly represent the range of questions asked by real students. Additionally, while LLMs can be trained to give reasonable answers, the generated responses often struggle with alignment to instructional intent and semantic accuracy, requiring further fine-tuning and advanced evaluation frameworks.

## 1 Introduction

The release of ChatGPT in November 2022 brought with it a new wave of AI and the potential of its applications in all aspects of modern life. It showcased to the public the power of LLMs, and their potential to reduce human workloads. Based on its extensive and ever growing training data, ChatGPT and other generative AI models can write papers, read and understand code, source information across the internet to answer niche questions, and so much more. Companies tend to hold a shared view of AI as a way to optimize their workforce, opening greater potential for cost savings. In the realm of software development, coders can use LLMs to optimize code, debug, write pseudocode and boilerplate functions, etc., potentially streamlining and simplifying the development process.

Depending on the context, these uses are immensely powerful and time-saving, but with un-

known effects and unregulated usage also carry strong potential for negative ramifications. The integration of this technology into educational contexts is naturally a lot more careful and cautious. Instead of accepting AI in full force, education also has to consider problems with students cheating and AI hallucinating as potential threats to learning. Unfortunately, given the opportunity, many students will have the tendency to offload portions of their work to generative models, causing them to miss out on critical learning. Enforcing learning in a world where this technology is readily, openly available is a difficult question, and one that has inspired innumerable research projects in the past couple of years. Instructors are seeking how to safely integrate generative technologies into their curriculum while also GPT-proofing assignments, bringing exams back to being written in-person, and taking a number of other measures to ensure student learning. Ultimately, we see gen-AI as an inevitability; we must turn our focus to improving/structuring interactions such that it can enhance the learning process.

As we will touch on in our related works, the accessibility of 1-1 teaching (ie. instructors and course staff) is a weakness of conventional teaching environments. As students aren't able to ask questions of their teachers, they turn to question-answering technologies as an alternative. Unfortunately, AI is often not trained to fully understand the context/assignment students are asking about and often yields incorrect or incomplete responses. To some end, this can actually also result in students garnering incorrect understandings based on their interactions with hallucinative/un-informed AI systems. The accessibility issue is an ongoing one, even prior to the advent of generative AI, and brings up an important question: can we fine-tune/train generative AI systems to act like instructors when they aren't available? Specifically, can we train them to understand the requirements of

084 specific assignments/projects such that they can  
085 provide useful feedback without also giving away  
086 answers to students? Finally, to what extent is prior  
087 data (rather than synthetic data) necessary to train  
088 such a system to sensibly answer real student ques-  
089 tions?

090 Answering these questions is complex. In  
091 this paper, we explore how well modern general-  
092 purpose fine-tuned LLMs can generate synthetic  
093 QA data; how well can they generate the kinds of  
094 questions we expect from students with regard to a  
095 project specification? Next, we compare the use of  
096 this synthetic data with real student QA to train a  
097 model to answer student questions how instructors  
098 would. To contextualize the importance of the syn-  
099 thetic data, we acknowledge that new courses have  
100 non-existent prior student data, and even existing  
101 courses may have gaps in datasets. So, we will  
102 examine both the process of generating synthetic  
103 QA data and using that data to inform a model in  
104 answering student questions.

## 105 2 Related Works

106 There are a number of related, but fundamentally  
107 different works that inspired our pursuit of this  
108 project. To contextualize these works, we note  
109 that the progress of LLMs in the past five years  
110 has been especially significant. Integrating these  
111 technologies into educational contexts holds much  
112 promise, but also has to be done carefully. Edu-  
113 cational research has been around for a long time;  
114 a hallmark of this space comes from **Bloom’s 2**  
115 **Sigma Challenge**, published in 1984 (BLOOM,  
1984).

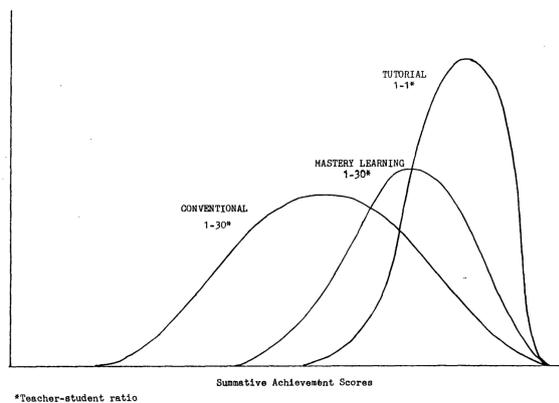


Figure 1: Curves comparing scores for students in varied learning environments.

116 Figure 1, pulled from this study, highlights the  
117 differences between conventional (1-30) teaching  
118

119 environments, personalized (1-1) tutoring settings,  
120 and enhanced (1-30) “*mastery learning*” settings.  
121 As expected, average student performance is the  
122 highest in 1-1 tutoring settings; while such teach-  
123 ing settings are most desirable, they are also highly  
124 impractical and expensive. Mastery learning serves  
125 as a generalization for settings in which students re-  
126 ceive a sense of personalized learning, often aided  
127 by technology (ie. cognitive tutors) in an otherwise  
128 conventional (1-30) teaching setting. There are a  
129 number of technologies that have been developed  
130 in support of mastery learning, and in our paper  
131 we discuss those related to answering student ques-  
132 tions through formal course forums (ie. Piazza).

133 In our paper, we turn to Question-Answer gener-  
134 ative systems. While round-the-clock instructor  
135 QA is the most desirable scenario, this is again far  
136 too expensive to be practical. Instead, we aim to  
137 observe to what extent generative language mod-  
138 els can support course instructors in correctly an-  
139 swering student questions. While some courses  
140 are likely to have extensive prior data (from past  
141 semesters), we further seek to answer the question  
142 if **synthetic data**, QA generated by an LLM based  
143 on extracting data from an assignment specifica-  
144 tion, would provide sufficient training for such a  
145 model?

146 Prior work has explored QA generation in a num-  
147 ber of ways. First [Basu et al.](#) aims to reduce work-  
148 load for instructors by generating multiple choice  
149 questions from text inputs. Next, [Riza et al.](#) looks  
150 at generating reading comprehension based short-  
151 answer questions via KNN techniques. Finally, [Vi-  
152 rani et al.](#) looks at how QA generation systems can  
153 be designed to support generation of a variety of  
154 question types. Ultimately, our work differentiates  
155 itself in two ways. First, we narrow the problem  
156 scope to the context of Computer Science education  
157 as we will seek to analyze synthesis of questions  
158 from coding project specifications; can we train an  
159 LLM to generate a range of sensible questions for  
160 a CS assignment based on its specification? Sec-  
161 ondly, we want to look at using generated QA as  
162 synthetic training data for a final student-question  
163 answering LLM; can synthetic QA train a system to  
164 answer real student questions with high accuracy?  
165 Ultimately, the prior work we have touched on so  
166 far highlights how our project fits into the broader  
167 space of QA generation for educational purposes.

168 Another interestingly related work is the  
169 CodeAid system produced by [Kazemitabaar et al.](#)  
170 which offers as a custom “coding assistant” as a

direct alternative to larger-purpose LLMs (ie. chat-GPT) that helps students learn to code with strictly no-code responses. Their approach focuses more on few-shot learning in prompt engineering for OpenAI API calls rather than training/fine-tuning an LLM, and focuses more on the usefulness of an LLM answering coding-specific questions (ie. fixing, writing, and understanding code). In contrast, our project aims to build two models working in tandem: one that can generate synthetic student QA data, and another that can answer hyper-specific project/assignment-based questions. While this work provides a vastly different approach, it does address a problem similar to the broader one (around-the-clock support for student questions) we are exploring in this paper.

Next, we will turn to some of the research that informed our technical approach to this project. Majority of the works we've mentioned utilize a T5 model as a baseline encoder-decoder tool; we do the same in building our LMs (Raffel et al., 2023). SQuAD (Stanford Question Answering Dataset) is a very popular QA dataset used across works in the space, and consists of a combination of answerable and unanswerable questions (Rajpurkar et al., 2018). For example, one of the T5 models we use to generate our synthetic QA data has been fine-tuned extensively on SQuAD (Manakul et al., 2023). The other model we use fine-tunes on SQuAD as well as CoQA and MSMARCO, two other large general-purpose QA datasets (Reddy et al., 2019; Bajaj et al., 2018). The SQuAD dataset is sourced from Wikipedia articles, CoQA has a focus on conversational QA, and MSMARCO is sourced from questions asked on Bing, such that the data from each of these datasets is expectedly quite different from the types of questions CS students may be asking about project specifications. In an ideal world, we would be able to use a model that has been fine-tuned on coding-related questions. Regardless the SQuAD (and other large dataset) fine-tuning aids the baseline T5 models in generating more sensible synthetic QA.

A significant problem in the realm of generative AI is model hallucination. Specifically in the QA space, large LLMs have been trained on a lot of information, and have the tendency to carry answers beyond the scope of the question being asked. Simultaneously, fine-tuning a model to concepts within a limited domain is a challenging, resource-intensive process. Retrieval Augmented Generation (RAG) is a popular technique for combating

these issues by basing model responses in a set knowledge base (Meyur et al., 2024; Barron et al., 2024). In our final question answering model we attempt a model configuration that emphasizes RAG techniques, utilizing a combination of synthetic and existing CS QA data to ground the model.

For evaluation, we utilized a combination of metrics to capture different aspects of quality in natural language generation. BERTScore (Zhang et al., 2020) assesses semantic similarity using contextual embeddings from pre-trained transformers, offering a nuanced understanding beyond token-level overlap. BLEU (Post, 2018), a standard metric for machine translation, measures n-gram overlap to evaluate lexical precision, while METEOR (Banerjee and Lavie, 2005) accounts for linguistic variations such as stemming and synonymy, aligning more closely with human judgment. ROUGE-L (Lin, 2004) evaluates the longest common sub-sequence between generated and reference texts, emphasizing fluency and recall. Together, these metrics provide a comprehensive evaluation framework, balancing semantic, lexical, and structural quality.

### 3 Implementation

As we have introduced, our project seeks to answer two key research questions:

1. Can synthetically made datasets effectively train LLMs?
2. Can an LLM be trained to answer real student questions like an instructor?

In order to answer these questions, we devised a research project that uses two LLMs. The first LLM takes in a project or assignment specification and outputs a set of corresponding question and answer pairs. The second LLM utilizes the synthetic data (question-answer pairs) generated by the first LLM as training data to ultimately take real student questions as input and output instructor-like responses. Throughout our report, we will divide analysis into two, corresponding to each of these two LLMs.

#### 3.1 LLM #1

Curating the first LLM had three core steps:

1. Generate Question Contexts
2. Generate Question/Answer (QA) Pairs
3. Evaluate the Question/Answer (QA) Pairs

### 3.1.1 Input Data

We selected the University of Michigan EECS 280 Project 3 spec as our input for this LLM due to our access to real student Piazza question-answer data. Additionally, our familiarity with the project enabled a level of manual analysis of the quality of generated questions and answers.

### 3.1.2 Generate Question Contexts

Upon first examination of the EECS 280 P3 (Eu-chre) specification, we had to consider how best to consolidate the combination of code snippets, images, tables, and raw textual paragraphs into a format that would be readable by a model. While the visual aspects of a specification are certainly important and insightful, we opted to extract just the text from the project assignments. Future work might attempt to better integrate different information formats into model inputs.

After paring specifications down to just the text, we moved on to consideration of how much information to give the model at a time. A **question context** is a snippet of a larger project specification or assignment. The pre-existing models we use in this project require contexts as input to generate QA pairs. Further, context sizes are limited; we cannot pass an entire specification into a model and ask it to generate questions (this is too large), nor would we want to if we could (questions would likely be too generic). Instead, we have to feed the model chunks at a time and ask it to generate QA pairs for each context.

Our initial approach to this problem was splitting an assignment by sentences and/or paragraphs. Intuitively, this would be a natural way to break apart text, and would ensure that each context passed to the model was as sensible as possible. Unfortunately, this approach introduced several complexities. Sometimes ideas are split between sentences and/or paragraphs wherein creating a good context for good, fully-informed questions would require including multiple sentences in a single context. However, the presence of long sentences, sentence groupings or paragraphs introduces the issue of how they should be split when they are too long. Should sentences be split in half? Should splits be included with the previous or next context? Should long sentences be moved to their own context altogether? Unfortunately, it doesn't feel as though there is a singular correct answer to this question; the answer is context dependent and we would have to opt for a solution that is good enough for the

general case. Again, future work might seek to use a more sophisticated approach for context creation; deeming which sentences should be grouped together, and which could be separated, whilst considering the size limitations of a context.

Instead, we opted for a slightly different approach: splitting our specification with a **sliding window algorithm**. We converted our specifica-

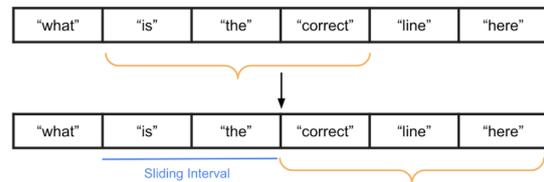


Figure 2: Visualization of sliding window algorithm concept.

tion into an array of words and used a fixed window and fixed sliding interval. Starting from the beginning of the assignment, we insert the first `<window_size>` words into a single context. Then, the window shifts by the sliding interval to generate a new context, and so on through the remainder of the document. Figure 2 below gives a visual representation of this algorithm. The downfall of this approach is the lack of intentionality in ensuring particular words/sentences/paragraphs stay together. But, the upsides are that we don't have to try to decide which window/context a given set of words is included with, as all text gets included in multiple contexts. This also provides us with a fixed/set context size, and addresses all the ambiguities/complexities we discussed in a sentence-/paragraph splitting approach.

### 3.1.3 Generate Question/Answer (QA) Pairs

For generating our QA pairs, we tested two models identified in prior works:

- [potsawee/t5-large-generation-squad-QuestionAnswer](#)
- [iarfmoose/t5-base-question-generator](#)

The first model we tested was `potsawee/t5-large-generation-squad-QuestionAnswer`. The input and output format for this model is included below. Note: `<sep>` serves as a separator token.

Input: context

Output: question `<sep>` answer

We did most of our parameter fine-tuning on this model; we discuss our results more extensively in the Evaluation section of the report, but we used this model to identify the best combination of hyper-parameters which we then carried on when evaluating the second model.

The second model evaluated was the [iarfmoose/t5-base-question-generator](#). Again, the input and output format for this model is described below.

Input: <answer> answer <context> context

Output: questions

The difference in input and output formats between these two models is quite significant. Namely, the first model only requires a context and actually produces a QA pair while the second model requires a context and answer and aims to produce a corresponding question. For the second model, its documentation specifies that for short answer questions the context could also be passed as the answer. Since we limit our input for this model to just the project specification (and don't have answers readily available, nor would we expect this of new courses), we opt for this approach. As we will discuss later, we hypothesize that the quality of this second model performed worse than the first partly because of this limitations; our contexts weren't exactly answers (they just contained answers), meaning the model sought to curate questions based on strangely formatted "answers". For both models, we generated multiple questions per context. This was a parameter that we varied to find the best results.

### 3.1.4 Evaluate the Question/Answer (QA) Pairs

The final step for LLM #1 is to evaluate the quality of the QA pairs generated by each model. We utilized a combination of quantitative and qualitative methods to do so:

- Quantitative:  
[iarfmoose/bert-base-cased-qa-evaluator](#)  
numerical scoring
- Qualitative: manual evaluation

The first method we applied was an evaluator LLM provided by HuggingFace: [iarfmoose/bert-base-cased-qa-evaluator](#). The evaluator accepts QA pairs and outputs

their corresponding "scores" wherein higher scores are indicative of "better" QA pairs. A known/acknowledged limitation of this model is that it only evaluates QA pairs based on "if they are semantically related," and not on the validity/correctness of the information nor the relevance of the question. Because of the sheer volume of questions generated we used the evaluator to set a threshold score of questions that "passed" and those that "failed"; questions below a set score threshold were deemed "failed." Moreover, we hoped that sorting questions by score would provide some semblance of a ranking of the "best" and "worst" questions generated so manual evaluation would be slightly easier.

While we could utilize the evaluator for high-level semantic assessments of the generated QA pairs, we determined that manual inspection was the best way to assess their quality. For this, we randomly inspected questions from both passed/failed datasets (for each model configuration) to validate accuracy and analyze the quality of generated QAs. This process allowed us to determine our optimal set of hyper-parameters as well as which model produced stronger QA pairs. We needed our evaluators to have a strong understanding of the project specification used as input for this evaluation to be effective. Two of our group members are veteran EECS 280 GSIs which we felt to be apt experience in comparing the questions generated by these models to those they would field in their instruction while helping students in OH and on Piazza for this project.

## 3.2 LLM #2

In this task, we explored approaches to develop a question-answering system tailored for a course with limited (or non-existent) historical Piazza data. The system aims to take a question as input and generate an answer in the style of course instructors. To achieve this, we harness the question-answering capabilities of large language models (LLMs) and experiment with various prompting techniques to address two key research questions:

1. Is synthetic data generated by our first LLM sufficiently reliable as a reference for building a generative QA system?
2. What challenges arise in building a system that can fully replicate a course instructor's question response style?

### 3.2.1 Experimental Setup

Our experiments are run on a server from Chameleon Cloud’s TACC cluster, running Ubuntu 22.04, equipped with 2 AMD EPYC 7763 64-Core CPUs, 252 GB RAM, and a single NVIDIA A100 40GB GPU. The experiments were ran on torch 2.5.1, transformers 4.46.0.

### 3.2.2 Models

To evaluate different approaches fairly and consistently, while keeping the computing budget affordable, we use the google/gemma-2-2b-it variant of Gemma (Team et al., 2024b), a lightweight, state-of-the-art open model that can be viewed as the open source version of the Google Gemini (Team et al., 2024a). For generation of the question embeddings, we use the all-MiniLM-L6-v2 model.

### 3.2.3 Datasets

We collected 4468 question-and-answer pairs on EECS280’s Piazza Forum. Each pair contains

1. **Question Subject:** A student-written title for the question.
2. **Question Body:** The content of the student question submitted.
3. **Instructor Answer:** An official instructor response to the question asked.
4. (Optional) **Student Follow-ups:** Back-and-forth discourse in response to the instructor answer, if present.

The question subject and body are concatenated to form what we consider the “final question,” while the instructor answer and follow-ups are concatenated to form what we consider the “ground truth answer.” These real student questions and instructor answer pairs form our validation data. In addition, we use synthetic data generated during the optimal run of LLM #1 as QA pairs for training.

### 3.2.4 Methodology

To address the challenges of building a course-specific question-answering system with limited historical Piazza data, we employed two complementary methodologies:

- Zero-shot Question Answering
- Few-shot In-context Learning (RAG) Question Answering

Both approaches leverage the inherent capabilities of large language models (LLMs) to generate and refine responses tailored to the course’s style and requirements.

### Zero-Shot

In this approach, we directly input course-specific questions into an LLM without providing additional contextual examples. The synthetic data generated by our first LLM, and any other collected QA data are not used. The purpose of a zero-shot approach is to test the **inherent knowledge** of the LLM; how well can it naturally generate relevant, contextually appropriate answers to student questions based solely on its pre-trained abilities? We analyze an LLMs raw capacity to respond like real instructors without any supplemental or context-specific training. Because collecting existing data is not always possible, and generating synthetic data can be expensive, we want to observe the quality of responses when not informed by data. Finally, the results of this approach serve as a sort of baseline for how much and in what ways other techniques we apply improve generated responses.

### Few-Shot In-Context

With this approach, we supply the LLM with a small number of QA pairs that are representative of the types of questions it should expect to receive and the corresponding types of answers we want it to generate. We do this prior to “asking” our LLM any question, and select the examples as a combination of synthetic data (collected from LLM #1) and validation data (from EECS 280 Piazza). The aim of this approach is to evaluate how much our model can be improved by providing it with examples. Is it better able to answer questions with an instructional tone, in the style of course instructors? Do examples enhance the factual relevance of responses?

### 3.2.5 Evaluation

In order to evaluate the results of this LLM, that requires comparing responses against a ground truth (actual instructor responses), to assess accuracy, relevance, and stylistic alignment. More broadly, we want to examine:

- Are generated answers semantically correct?
- Do generated answers align with expected instructor responses?

549 While semantic correctness is a relatively obvious 598  
550 need, recall that instructor answers have their 599  
551 own stylistic/conceptual/structural tendencies as 600  
552 a means of best enabling learning for students. 601  
553 For example, we would expect instructors not 602  
554 to readily give away answers, but rather lead 603  
555 students towards a better way to think about and/or 604  
556 approach problems. They may list steps or hints 605  
557 that help the student explore and discover the 606  
558 answer on their own. For these reasons (as with 607  
559 LLM #1) we apply a combination of quantitative 608  
560 evaluation through automated text-similarity based  
561 metrics, and qualitative evaluation through manual  
562 human inspection. Here, we discuss the three  
563 automated metrics we utilized as well as how we  
564 went about manual evaluation.  
565

### 566 **Quantitative Metric #1: BLEU and ROUGE**

567 These metrics measure the overlap between the  
568 generated answers and the ground truth at the word  
569 and phrase levels, providing insights into lexical  
570 and structural similarities.  
571

### 572 **Quantitative Metric #2: BERTScore**

573 This metric evaluates semantic similarity by  
574 computing embeddings of the generated and  
575 reference answers, offering a deeper understanding  
576 of how well the generated responses capture the  
577 meaning and intent of the instructor's answers.  
578

### 579 **Quantitative Metric #3: Meteor**

580 Meteor extends traditional n-gram-based  
581 evaluation metrics by incorporating more advanced  
582 linguistic matching techniques. It calculates  
583 similarity between generated and reference  
584 texts by considering synonyms, stemming, and  
585 paraphrasing, thus providing a more nuanced  
586 assessment of translation and text generation  
587 quality beyond exact word matching.  
588

### 589 **Qualitative Metric: Human Evaluation**

590 We employ a comprehensive human review pro-  
591 cess to assess the quality of generated responses.  
592 We evaluate a subset of answers across three critical  
593 dimensions:

- 594 • Style and tone matching instructional ap-  
595 proach
- 596 • Factual accuracy aligned with course content
- 597 • Clarity in addressing specific questions

598 When significant deviations are identified, we  
599 conduct a systematic error analysis to uncover  
600 potential issues such as question misinterpretation,  
601 training data limitations, or response overgeneral-  
602 ization.  
603

604 By combining these evaluation methods, we aim  
605 to gain a holistic understanding of LLM #2's perfor-  
606 mance, identify its strengths and limitations, and  
607 inform future improvements to better align with  
608 course instructor expectations.

### 609 **3.3 Implementation of the In-Context** 610 **Learning Approach**

611 To enable our LLM to better answer student ques-  
612 tions, we feed it sample QA pairs. The intuitions  
613 here are that (1) real student questions pulled from  
614 Piazza may show underlying patterns over time,  
615 and (2) the existing answers can act as "role mod-  
616 els" from whom LLMs can learn about the appropri-  
617 ate tone and explicitness. The specific implemen-  
618 tation is shown in Figure 3, on the following page.  
619 We first apply the all-MiniLM-L6-v2 to generate  
620 embeddings for the questions. It maps sentences  
621 to a 384 dimensional dense vector space and can  
622 be used for tasks like clustering or semantic search.  
623 When answering new questions, the embedding  
624 of the new question will be used to search in the  
625 embedding space for which questions are similar  
626 to this one. We fetch the top-k similar question  
627 answer pairs and use them to construct the final  
628 prompt to the system. Finally, the system generates  
629 the answer for the new question.

## 630 **4 Evaluation**

### 631 **4.1 LLM #1**

632 For our first LLM, we ran a total of 13  
633 different configurations (via altering hyper-  
634 parameters and models). The first 12 configu-  
635 rations were all performed with our first model  
636 ([potsawee/t5-large-generation-squad-  
637 QuestionAnswer](#)) and the final configuration was  
638 performed using the optimal hyper parameters  
639 identified for the first model on our second one  
640 ([iarfmoose/t5-base-question-generator](#)). In  
641 the following section, we will discuss the options  
642 we considered and decisions we made in selecting  
643 values for each hyper parameter, supported by QA  
644 examples. The approach we took in testing hyper  
645 parameters was first setting each of them to a sensi-  
646 ble default value. We then progress through the list,

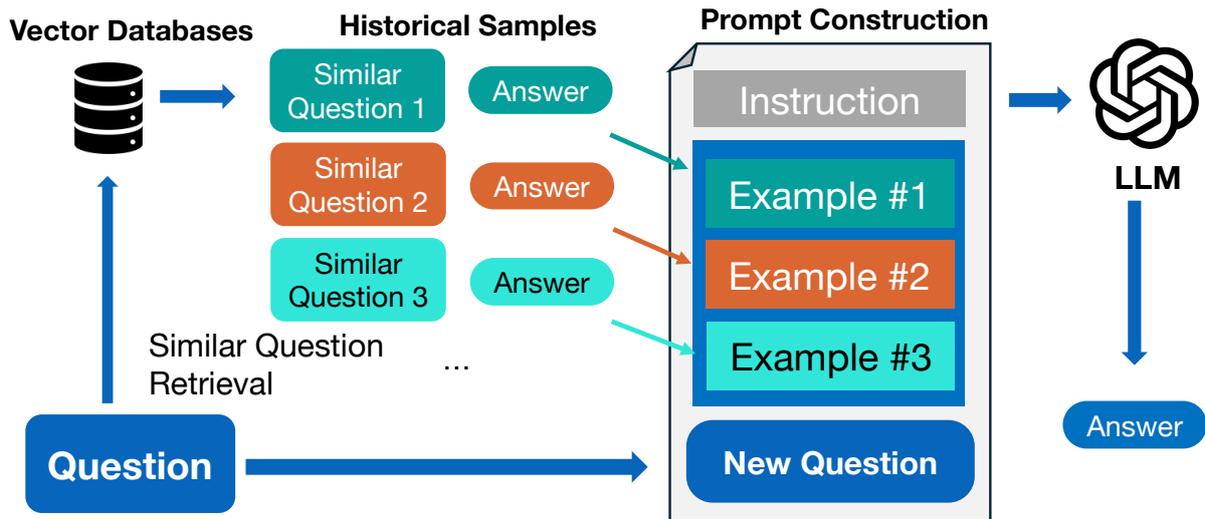


Figure 3: Workflow of In-Context Learning Approach for Question Answering

647 modifying one at a time, determining the optimal  
 648 value for that parameter and then holding it at that  
 649 value for the remainder of our configurations. First,  
 650 here's a list of the parameters we varied:

- 651 • Window Step Size
- 652 • Context Window Size
- 653 • Number of Questions per Context
- 654 • Score Threshold
- 655 • Model

#### 656 4.1.1 Window Step Size

657 The window step size refers to how much our con-  
 658 text window was shifted forward in the input text  
 659 between each iteration (context generated). The  
 660 essence of this parameter is balancing a step size  
 661 that is too small wherein QA generated across dif-  
 662 ferent contexts are overly redundant, against one  
 663 that is too large wherein information does not get  
 664 properly represented by any of the contexts it is  
 665 included in.

666 For the larger window step sizes (5, 10), it felt  
 667 like questions either (a) didn't have enough con-  
 668 text and/or (b) were missing relevant context such  
 669 that both passed and failed questions weren't super  
 670 sensible. We noticed several questions that were  
 671 actually just sentences, sometimes pulled from the  
 672 specification. For a window step size of 1, we  
 673 saw noticeably better performance; questions made  
 674 more sense, but were slightly more redundant. The  
 675 contrast between some of the top-scored QA pairs

676 for a larger vs. smaller window step size was quite  
 677 apparent.

678 The top scoring QA pair and score for a window  
 679 size of 10 was:

680 **Question:** "Why might your code be eas-  
 681 er to test and debug?"  
 682 **Answer:** "May make for easier testing  
 683 and debugging"  
 684 **Score:** 3.76

685 The top scoring QA pair and score for a window  
 686 size of 1 was:

687 **Question:** "What does a user enter in  
 688 order to discard a card?"  
 689 **Answer:** "The user will then enter the  
 690 number corresponding to the card they  
 691 want to discard"  
 692 **Score:** 3.79

693 While the questions produced by all configura-  
 694 tions had some weaknesses, the one utilizing a win-  
 695 dow size of 1 was remarkably better. This was con-  
 696 sistent throughout the QA pairs for these configura-  
 697 tions. Not only were questions more grammatically  
 698 correct, they also made more sense, had better cov-  
 699 erage of information contained in the assignment,  
 700 and were more similar to the types of questions  
 701 we would expect from real students. From these  
 702 observations, we moved forward with a window  
 703 step size of 1 for the remaining configurations.

#### 704 4.1.2 Context Window Size

705 The context window size refers to how much text  
 706 (how large a context) we opted to send to the gen-

erative QA model. This parameter aids in the processes of determining the optimal way to send contexts to our model within a sliding window approach. Specifically, it balances a context window that is too small wherein contexts cannot hold enough information to construct sensible questions despite a small window step size, against a window that is too large wherein the model fails to understand contexts in a modular way.

For smaller context windows (40, 50, 60), the QA pairs produced were not noticeably different in any meaningful way. A context window size of 70 produced some really good questions, but definitely still produced questions that don't make sense and are slightly repetitive. For the largest window size (100), there wasn't a huge qualitative difference in the questions generated (compared to 70). It seems that at a certain point, expanding the window size doesn't improve question quality but limits the diversity of questions asked.

While we observed that expanding the window size did have a positive effect on the quality of QA pairs produced by the model, varying this parameter also demonstrated the limitations of having a model that cannot hold onto the knowledge conveyed in neighboring contexts. Regardless, we saw that smaller context window sizes limited the sensibility of questions generated while context windows that are too large cause the model to be too unfocused. For these reasons, we moved forward with a context window size of 70 in the remaining configurations.

**4.1.3 Number of Questions per Context**

The models we used allows for multiple QA pairs to be generated for each context sent to the model, so this variable specifies that value. This parameter balances generating too few QA pairs per context where the model fails to generate questions covering the entire context space, against too many QA pairs per context where questions become overly redundant. We tested generating 1, 2, and 4 questions per context. There was already some redundancy in questions in prior configurations, and the 4 QA pair generation configuration aligned with this observation as it failed to produce questions that were noticeably unique. At the same time, we did notice some of the QA pairs generated by the 2-question configuration were meaningfully unique. For example, here's a student-like question we didn't see in the 1-question configuration:

**Question:** "Does a Simple Player con-

sider whether their opponent is the dealer?"

**Answer:** "Do not consider whether they are the dealer and could gain an additional trump by picking up the upcard"

**Score:** 3.73

Again, the QA generated here is by no means semantically perfect. But, it does capture meaningful content from the project specification. We moved forward generating 2 questions per context in the remaining configurations.

**4.1.4 Score Threshold**

As we mentioned in our Methods section for LLM #1, we utilized a QA evaluator to numerically score each of the questions we generated (for semantic correctness), and sorted them in descending order. Thus far, we have separated passed/failed questions with a threshold of 1.5. In this step, we examined questions with scores surrounding this threshold, specifically values of 1.5, 2, 2.5, and 3. Ultimately, we observed several questions with scores in the high 2s that were still relevant, but felt that a lower threshold than that was likely to include more bad questions than good ones. Here's an example:

**Question:** "What two players are in the game?"

**Answer:** "A simple AI player and a human-controlled player that reads instructions from standard input cin"

**Score:** 2.72

Ultimately, we moved forward with a score threshold of 2.5 in the remaining configurations.

**4.1.5 Model**

The model refers to whether QA pairs were generated with [potsawee/t5-large-generation-squad-QuestionAnswer](#) or [iarfmoose/t5-base-question-generator](#). So far, all examples we have observed were generated from the first model. In running our second model, questions were significantly less semantically correct, didn't take the form of actual questions, and were long/rambly. Here's one such example:

**Question:** "<pad> True ID: a vector is added to a card and no member variable can be used to type it. When using a const function, sort result in your swap function showing a weird error.</s>"

**Answer:** "in Card.hpp. Use the STL to

sort a vector<Card> hand: Pitfall: Using sort on a member variable in a const member function leads to a confusing error, no matching function for call to 'swap'. Instead, call sort when adding a”  
**Score: 3.81**

Despite somehow having a higher evaluation score than many of the questions produced by our first model, because of the significantly lower quality of the questions produced by this model, we carried forward with the results of our first model.

#### 4.1.6 Summary

Ultimately, the combination of quantitative and qualitative metrics used to evaluate the QA pairs generated by our tested model in this section led us to the following final parameters:

- Window Step Size = 1
- Context Window Size = 70
- Number of Questions per Context = 2
- Score Threshold = 2.5
- Model = `potsawee/t5-large-generation-squad-QuestionAnswer`

This run resulted in a decent combination of sensible, student-like project questions and some attempted questions which still held room for improvement. Regardless, this was the best synthetic QA data we were able to produce, and what we carried forward as partial training data for LLM #2.

## 4.2 LLM #2

Our evaluations of LLM #2 serve as answers to the key questions guiding this project:

1. Can synthetic data serve as credible samples for effective curation of question answer bots?
2. Is the RAG-based in-context learning approach effective for answering previously unseen questions?

Finally, based on our findings with regards to each of these questions, we identify what challenges there are to building a good question-answering system to replicate 1-1 instructor-student tutoring.

### 4.2.1 Main Results

The evaluations/comparisons made in the following sections are an analysis of the data presented in Table 1.

	BLEU	ROUGE-L	BERTScore	Meteor
zero-shot	0	0.013	0.792	0.002
rag	0.009	<b>0.133</b>	<b>0.801</b>	<b>0.167</b>
rag simulation	<b>0.013</b>	0.128	0.775	0.148

Table 1: Evaluation Results of LLM #2. *zero-shot*: our baseline, no additional context given to the model, *rag*: our in-context learning approach that uses the synthetic data as the reference, *rag simulation*: our in-context learning approach that uses the validation data capped up to the date of the current question as the reference.

### Synthetic Data as Credible Samples

The experimental results suggest that synthetic data generated by LLM #1 can serve as credible samples for question-answering tasks. Specifically, the ‘RAG’ approach, which uses synthetic data as a reference for in-context learning, achieves substantial improvements over the ‘zero-shot’ baseline (which has no training). For instance, the **ROUGE-L** score improves from 0.013 to 0.133, and the **Meteor** score increases significantly from 0.002 to 0.167. These improvements indicate that the synthetic data provide valuable contextual information that aids the model in generating more sensible, accurate answers. Interestingly, the ‘rag simulation’ approach, which uses validation data (real Piazza QA) capped to the question date as the reference, performs slightly worse than ‘rag’ (synthetic data) in terms of **ROUGE-L** (0.128 vs. 0.133) and **Meteor** (0.148 vs. 0.167). This outcome highlights that synthetic data might better align with the model’s pre-training distribution compared to real-world validation data, making it easier for the model to generalize. However, the marginal gap also suggests that real-world validation data still plays a vital role as it also still improves on the zero-shot context.

### Effectiveness of RAG-Based In-Context Learning

The results strongly support the effectiveness of the RAG-based in-context learning approach for answering previously unseen questions. When compared to the ‘zero-shot’ baseline, ‘rag’ demonstrates significant improvements across all evaluation metrics except for **BLEU**. Notably, the **Meteor** score sees an 83.5% increase from 0.002 to 0.167, and the **ROUGE-L** score improves by an order of magnitude. These gains emphasize the ability of RAG to leverage context effectively, thereby bridging the gap between training data and unseen questions. It is worth noting that the **BERTScore**

metric for 'rag' (0.801) shows only a marginal improvement over the baseline (0.792). This observation suggests that while RAG aids in structuring answers more coherently, the semantic similarity to the ground truth is still limited. This limitation could arise from inconsistencies in synthetic or reference data or the inherent complexity of the domain-specific questions.

### 4.3 Challenges in Building a Question-Answering System

Despite the positive results, the study reveals several challenges in building an effective question-answering system:

First, all approaches yield relatively low **BLEU** scores, with the highest being 0.013 for 'rag simulation'. This indicates that lexical overlap between generated and reference answers remains a significant bottleneck. Such low scores highlight the difficulty in achieving exact word-level matches, especially in domain-specific and highly nuanced questions. Simultaneously, we acknowledge that lexical similarity and exact overlaps in words/phrases are not necessarily the most important metrics when analyzing answer quality.

Second, the improvement in **BERTScore** across approaches is modest, with the highest being 0.801 for 'rag'. This suggests that while the generated answers improve in structure, their semantic alignment with the ground truth is still limited. A deeper analysis of the quality of both the synthetic and real validation data is required to address this issue.

Third, the dependence on synthetic data quality is evident. The superior performance of 'rag' over 'rag simulation' indicates that the synthetic data generated by LLM #1 might better align with the pre-training distribution of the model, thereby enabling more effective generalization. This observation underscores the importance of high-quality and diverse synthetic data for training robust systems.

Lastly, the experiments highlight the limitations of in-context learning for complex question-answering tasks. While RAG significantly outperforms the baseline, the absolute metric values remain low, suggesting that additional techniques, such as fine-tuning or enhanced retrieval mechanisms, might be necessary to further improve performance to a deployable level.

#### 4.3.1 Challenges: Evaluation Fairness and Style Alignment

```
*Actual Instructor Answer*:
If you follow step 1 in @6, instructors
will be able to see your unqiqname
even
if you post anonymously (i.e., if you
post anonymously other students won't
see your unqiqname). If you want to
understand what you missed on the
private
test cases, I recommend going to
Proffice hours. You can attend any
Proffice
hours even if you don't attend that
Professor's lectures.

*Model Answer*:
It's a good question, and it highlights
a common misconception. While you
can
certainly add getter functions to your
`player` class in `euchre.cpp`, you
**shouldn't**. Let's break it down:

1. Visibility Within the Class
- Getter functions are intended to
allow clients or other parts of
your code
to safely access an object's
private data.
- By defining a getter function
within your `player` class, you
are
potentially making the hand data
accessible from anywhere.
```

Figure 4: Example Generated Answer and the Ground Truth

One possible explanation for the extremely low scores for all of our approaches might be related to the fairness of our evaluation metrics and style alignment. The primary challenge in evaluating educational QA systems lies in ensuring fairness and aligning evaluation metrics with the desired style and intent of responses. Lexical-based metrics, such as BLEU and ROUGE, often penalize responses that diverge in phrasing or structure from the ground truth, even if the generated answers are semantically relevant or meaningful. For example, as shown in Figure 4, the model-generated answer focuses on coding practices and getter functions, which is tangential to the actual query about private tests but may still provide valid insights. However, such responses are unfairly penalized due to their deviation from the expected response format. Furthermore, educational contexts often require adherence to an instructive and supportive tone, which can lead to further mismatches if models produce technical or overly casual responses. These chal-

962	lenges underscore the need for evaluation metrics	1012
963	that balance semantic correctness, relevance, and	
964	stylistic alignment to ensure fair and contextually	1013
965	appropriate assessments of model performance. Ul-	1014
966	timately, it is a very complex problem to determine	1015
967	how to form an evaluation metric that can fairly as-	1016
968	sess the quality of these kinds of question answers.	1017
969	To an extent, real instructor answers are highly sub-	1018
970	jective such that there is no ultimately “correct”	1019
971	(truth) answer to these kinds of questions. Even	1020
972	within real instruction, there is a lot of ambigu-	1021
973	ity in how to balance vagueness in responses with	1022
974	giving away answers to students, and this issue is	1023
975	only exacerbated when we try to assign answersxw	
976	numerical scores that assess their quality.	
977	<b>5 Discussion of Results</b>	1024
978	The results of this project demonstrate both the	1025
979	promise and the limitations of utilizing synthetic	1026
980	data and retrieval-augmented generation (RAG)-	1027
981	based in-context learning for educational QA sys-	1028
982	tems. Synthetic data generated by LLM #1 proved	1029
983	to be credible and effective for training purposes, as	1030
984	evidenced by significant improvements in metrics	1031
985	such as ROUGE-L and Meteor when compared to	1032
986	the zero-shot baseline. This highlights the poten-	1033
987	tial of synthetic QA data to fill gaps in historical	1034
988	datasets, particularly for new or evolving courses.	1035
989	However, the study also reveals the inherent chal-	1036
990	lenges in aligning generated answers with both the	1037
991	semantic and stylistic expectations of instructional	1038
992	responses. For instance, while RAG-based methods	1039
993	significantly improved contextual understanding	1040
994	and coherence, metrics like BLEU and BERTScore	1041
995	indicate limitations in achieving precise lexical and	
996	semantic alignment. This is further complicated by	
997	the evaluative metrics themselves, which often pe-	
998	nalize stylistic or structural deviations, even when	
999	the content of the generated answer is relevant and	
1000	insightful. The discrepancy between the gener-	
1001	ated answers and the ground truth highlights the	
1002	need for alternative metrics that account for stylis-	
1003	tic alignment and the instructive tone essential in	
1004	educational contexts.	
1005	Additionally, while synthetic data demonstrated	
1006	a closer alignment with the pretraining distribu-	
1007	tion of the models, real-world validation data offered	
1008	complementary benefits, emphasizing the impor-	
1009	tance of a hybrid approach. The reliance on high-	
1010	quality synthetic data underscores the necessity of	
1011	refining QA generation techniques, including im-	
	proved context creation and filtering strategies.	
	Overall, while RAG-based in-context learning	
	has shown considerable potential, the relatively low	
	absolute metric values indicate that further refine-	
	ments—such as domain-specific fine-tuning, en-	
	hanced retrieval mechanisms, and the integration of	
	diverse data sources—are required to build robust	
	and reliable question-answering systems. These	
	findings underscore the complexities of creating	
	AI systems capable of emulating the nuanced re-	
	sponses of human instructors while providing ac-	
	tionable insights for future research.	
	<b>6 Conclusion</b>	1024
	This research project explores the potential of	1025
	LLMs to improve computer science education. It	1026
	shows that while using LLMs to generate prac-	1027
	tice questions and offer student support is possible,	1028
	there are still obstacles to overcome. Improving the	1029
	quality of the synthetically generated data is crucial	1030
	for creating a reliable system that can accurately	1031
	answer student questions. This involves exploring	1032
	better techniques for creating context and retaining	1033
	knowledge from previous interactions. The evalu-	1034
	ation of the question-answering LLM suggests fu-	1035
	ture research in having more accurate benchmarks	1036
	for complex question answering systems, and how	1037
	LLMs can align with the actual instructors to give	1038
	inspiring yet not explicit answers. These findings	1039
	highlight both the promise and the complexities of	1040
	integrating LLMs into educational settings.	1041
	<b>7 Division of Work</b>	1042
	Our team consists of four members. For the actual	1043
	implementations of each LLM, we split into two	1044
	subteams (one for each LLM). The team members	1045
	with teaching experience for EECS 280 were se-	1046
	lected to work on LLM #1 while the other team	1047
	members worked on LLM #2. The team worked	1048
	collaboratively on the project proposal, presenta-	1049
	tion, and final paper to ensure all findings were	1050
	well represented and to provide complete, insight-	1051
	ful final results.	1052
	<b>8 Codebase URL</b>	1053
	<a href="https://github.com/stoneann/AutoQATrainedLLM/tree/main">https://github.com/stoneann/</a>	1054
	<a href="https://github.com/stoneann/AutoQATrainedLLM/tree/main">AutoQATrainedLLM/tree/main</a>	1055

## References

- 1056
- 1057 Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng,  
1058 Jianfeng Gao, Xiaodong Liu, Rangan Majumder,  
1059 Andrew McNamara, Bhaskar Mitra, Tri Nguyen,  
1060 Mir Rosenberg, Xia Song, Alina Stoica, Saurabh  
1061 Tiwary, and Tong Wang. 2018. *Ms marco: A human  
1062 generated machine reading comprehension dataset.*  
1063 *Preprint*, arXiv:1611.09268.
- 1064 Satanjeev Banerjee and Alon Lavie. 2005. *METEOR:*  
1065 *An automatic metric for MT evaluation with im-*  
1066 *proved correlation with human judgments.* In *Pro-*  
1067 *ceedings of the ACL Workshop on Intrinsic and Ex-*  
1068 *trinsic Evaluation Measures for Machine Transla-*  
1069 *tion and/or Summarization*, pages 65–72, Ann Arbor,  
1070 Michigan. Association for Computational Linguistics.  
1071
- 1072 Ryan C. Barron, Ves Grantcharov, Selma Wanna, Mak-  
1073 sim E. Eren, Manish Bhattarai, Nicholas Solovyev,  
1074 George Tompkins, Charles Nicholas, Kim Ø. Ras-  
1075 mussen, Cynthia Matuszek, and Boian S. Alexandrov.  
1076 2024. *Domain-specific retrieval-augmented genera-*  
1077 *tion using vector stores, knowledge graphs, and ten-*  
1078 *sor factorization.* *Preprint*, arXiv:2410.02721.
- 1079 Rahul Basu, Dolasaha, Chiranjib Dutta, and Ananjan  
1080 Maiti. 2023. *Automatic Transformation and Gen-*  
1081 *eration of Question Papers in Education with NLP*  
1082 *Techniques.*
- 1083 BENJAMIN S. BLOOM. 1984. *The 2 sigma problem:*  
1084 *The search for methods of group instruction as effec-*  
1085 *tive as one-to-one tutoring.* *Educational Researcher*,  
1086 13(6):4–16.
- 1087 Majeed Kazemitabaar, Runlong Ye, Xiaoning Wang,  
1088 Austin Zachary Henley, Paul Denny, Michelle Craig,  
1089 and Tovi Grossman. 2024. *Codeaid: Evaluating a*  
1090 *classroom deployment of an llm-based programming*  
1091 *assistant that balances student and educator needs.* In  
1092 *Proceedings of the 2024 CHI Conference on Human*  
1093 *Factors in Computing Systems, CHI '24*, New York,  
1094 NY, USA. Association for Computing Machinery.
- 1095 Chin-Yew Lin. 2004. *ROUGE: A package for auto-*  
1096 *matic evaluation of summaries.* In *Text Summariza-*  
1097 *tion Branches Out*, pages 74–81, Barcelona, Spain.  
1098 Association for Computational Linguistics.
- 1099 Potsawee Manakul, Adian Liusie, and Mark JF Gales.  
1100 2023. *Mqag: Multiple-choice question answering*  
1101 *and generation for assessing information consistency*  
1102 *in summarization.* *arXiv preprint arXiv:2301.12307.*
- 1103 Rounak Meyur, Hung Phan, Sridevi Wagle, Jan Strube,  
1104 Mahantesh Halappanavar, Sameera Horawalavithana,  
1105 Anurag Acharya, and Sai Munikoti. 2024. *Weqa:*  
1106 *A benchmark for retrieval augmented generation in*  
1107 *wind energy domain.* *Preprint*, arXiv:2408.11800.
- 1108 Matt Post. 2018. *A call for clarity in reporting bleu*  
1109 *scores.* *Preprint*, arXiv:1804.08771.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine  
Lee, Sharan Narang, Michael Matena, Yanqi Zhou,  
Wei Li, and Peter J. Liu. 2023. *Exploring the limits*  
of transfer learning with a unified text-to-text trans-  
former. *Preprint*, arXiv:1910.10683.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.  
Know what you don’t know: Unanswerable questions  
for squad. *Preprint*, arXiv:1806.03822.
- Siva Reddy, Danqi Chen, and Christopher D. Manning.  
2019. *Coqa: A conversational question answering*  
challenge. *Preprint*, arXiv:1808.07042.
- Lala Septem Riza, Yahya Firdaus, Rosa Ariani Sukamto,  
Wahyudin, and Khyrina Airin Fariza Abu Samah.  
2023. *Automatic generation of short-answer ques-*  
tions in reading comprehension using nlp and knn.  
*Multimedia Tools and Applications*, 82(27):41913–  
41940.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, and  
Jean-Baptiste Alayrac et al. 2024a. *Gemini: A fam-*  
ily of highly capable multimodal models. *Preprint*,  
arXiv:2312.11805.
- Gemma Team, Thomas Mesnard, Cassidy Hardin,  
Robert Dadashi, Surya Bhupatiraju, and  
Shreya Pathak et al. 2024b. *Gemma: Open*  
models based on gemini research and technology.  
*Preprint*, arXiv:2403.08295.
- Altaj Virani, Rakesh Yadav, Prachi Sonawane, and  
Smita Jawale. 2023. *Automatic question answer*  
generation using t5 and nlp. In *2023 International*  
Conference on Sustainable Computing and Smart  
Systems (ICSCSS), pages 1667–1673.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.  
Weinberger, and Yoav Artzi. 2020. *Bertscore:*  
Evaluating text generation with bert. *Preprint*,  
arXiv:1904.09675.