

ChatCAD: Interactive Computer-Aided Diagnosis on Medical Image using Large Language Models

Sheng Wang^{1,2,3} Zihao Zhao¹ Xi Ouyang³ Qian Wang¹ Dinggang Shen^{1,3}
¹ShanghaiTech University ²Shanghai Jiao Tong University ³United Imaging Intelligence

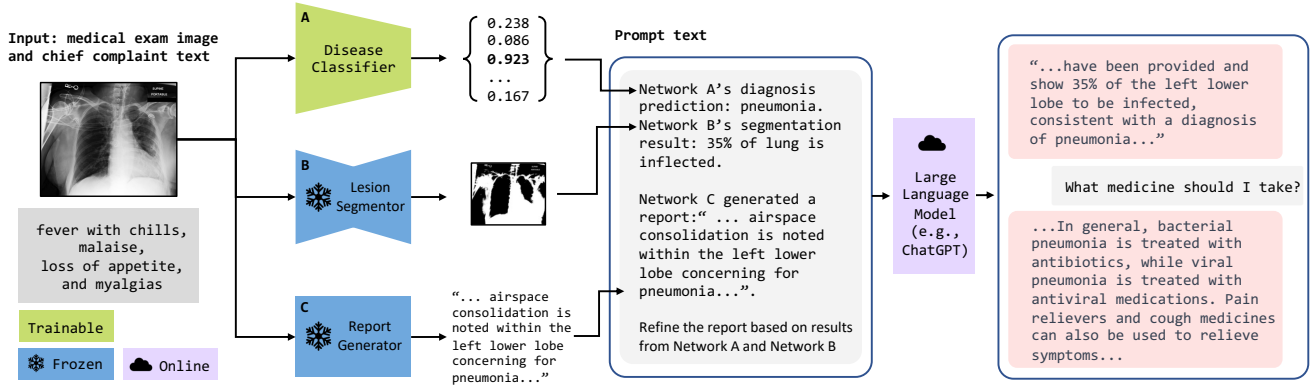


Figure 1. Overview of our proposed strategy. The image is processed by various networks to generate diverse outputs, which are then transformed into text descriptions. The descriptions, served as a link between visual and linguistic information, are combined as inputs to a large language model (LLM). With its ability to reason and its knowledge of medical field, the LLM can provide a condensed report and offer interactive explanations and medical recommendations based on the given image.

Abstract

Large language models (LLMs) have recently demonstrated their potential in clinical applications, providing valuable medical knowledge and advice. For example, a large dialog LLM like ChatGPT has successfully passed part of the US medical licensing exam. However, LLMs currently have difficulty processing images, making it challenging to interpret information from medical images, which are rich in information that supports clinical decisions. On the other hand, computer-aided diagnosis (CAD) networks for medical images have seen significant success in the medical field by using advanced deep-learning algorithms to support clinical decision-making. This paper presents a method for integrating LLMs into medical-image CAD networks. The proposed framework uses LLMs to enhance the output of multiple CAD networks, such as diagnosis networks, lesion segmentation networks, and report generation networks, by summarizing and reorganizing the information presented in natural language text format. The goal is to merge the strengths of LLMs’ medical domain knowledge and logical reasoning with the vision understanding capability of existing medical-image CAD models to cre-

ate a more user-friendly and understandable system for patients compared to conventional CAD systems. In the future, LLM’s medical knowledge can be also used to improve the performance of vision-based medical-image CAD models.

1. Introduction

Large Language Models (LLMs) are advanced artificial intelligence systems that have been trained on vast amounts of text data [5, 22]. These models use deep learning techniques to generate human-like responses, making them useful for a variety of tasks such as language translation, question answering, and text generation. LLMs like OpenAI’s GPT-3 [3] have shown remarkable results in natural language processing and have the potential to revolutionize various industries, including marketing, education, and customer service. The ability of LLMs to process and understand large amounts of data has made them highly sought after for solving complex problems. In the medical domain, LLMs have demonstrated their potential as valuable tools for providing medical knowledge and advice. For instance, a large dialog-based LLM, such as ChatGPT [17], has demonstrated remarkable results in a critical evalua-

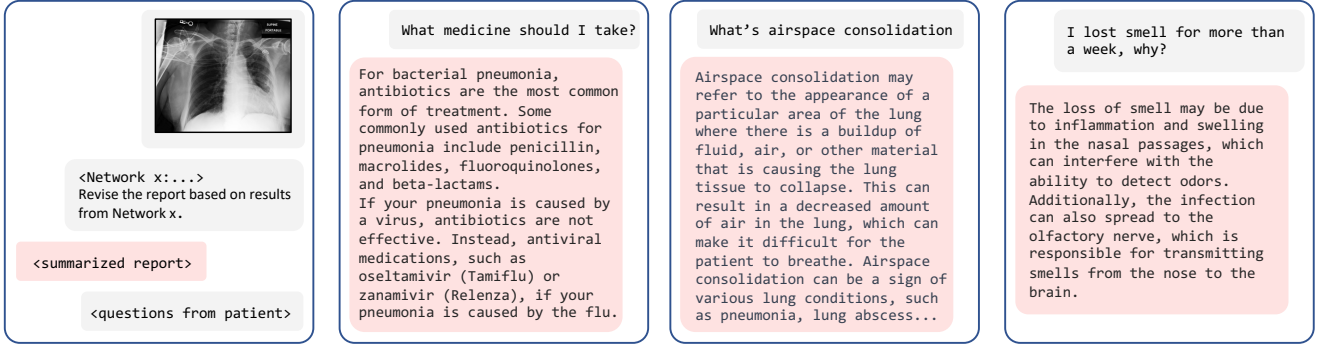


Figure 2. Interactive CAD with LLMs. This example uses the ChatGPT as LLM.

tion of its medical knowledge. ChatGPT has successfully passed part of the US medical licensing exams, showcasing its potential to augment medical professionals in delivering care. Inspired by their remarkable progress in natural language processing, it is an interesting topic to integrate the LLMs to understand visual information in computer vision tasks. Processing images involves understanding the spatial relationships between objects, recognizing patterns and textures, and extracting features that describe the objects in an image. These tasks require a deep understanding of visual information, which is challenging for LLMs that have been primarily trained on text data. This limitation presents a major challenge in the medical field, where images play a crucial role in supporting clinical decisions. Medical images, such as X-rays, CT scans and MRIs, are rich in information that can provide critical insights into a patient’s condition. However, LLMs currently struggle to interpret and extract information from these images, limiting their ability to fully support clinical decision-making processes.

As the “pure” computer vision method, medical-image computer-aided diagnosis (CAD) networks have achieved significant success in supporting clinical decision-making processes in the medical field [24]. These networks leverage advanced deep learning algorithms to analyze medical images and provide valuable insights to support clinical decision-making. CAD networks have been designed specifically to handle the complexities of visual information in medical images, making them well-suited for tasks such as disease diagnosis [30], lesion segmentation [35], and report generation. These networks have been trained on large amounts of medical image data, allowing them to learn to recognize complex patterns and relationships in visual information that are specific to the medical field.

The aim of this paper is to provide a scheme that combines the strength of LLMs and CAD models. In this scheme, namely ChatCAD, the image is first fed into multiple networks, i.e., an image classification network, a lesion segmentation network, and a report generation network as depicted in Figure 1. The results produced by classifica-

tion or segmentation are a vector or a mask, which can not be understood by LLMs. Therefore, we transform these results into the text representation form as shown in the middle panel of Figure 1. These text-form results will then be concatenated together as a prompt “*Revise the report based on results from Network A and Network B*” for the LLM. The LLM then summarizes the results from all the CAD networks. As the example in this figure, the refined report combines the findings from all three networks to provide a clear and concise summary of the patient’s condition, highlighting the presence of pneumonia and the extent of the infection in the left lower lobe. In this way, the LLM could correct errors in the generated report based on the results from CAD networks. Our experiment shows that our scheme could improve the diagnosis performance score of the state-of-the-art report generation methods by 16.42%. A major benefit of our approach is the utilization of LLM’s robust logical reasoning capabilities to combine various decisions from multiple models. This allows us to fine-tune each model individually. For instance, in response to an emergency outbreak such as COVID-19, we can add a pneumonia classification model (differentiating between community-acquired pneumonia and COVID-19 [19]) using very few cases without affecting the other models. Since classifiers are usually less data-hungry than other models, we mark it with “trainable” (green) in Figure 1.

Another advantage of bootstrapping LLMs to CAD models is that their extensive and robust medical knowledge can be leveraged to provide interactive explanations and medical advice as we illustrate on Figure 2. For example, based on an image and generated report, patients can inquire about appropriate treatment options (second panel) or define medical terms such as “airspace consolidation” (third panel). Or with patient’s chief complaint (forth panel), LLMs can explain why such symptom happens. In this manner, patients can gain a deeper understanding of their symptoms, diagnosis, and treatment more efficiently. It can efficiently help patients to reduce consultation costs with clinical experts. As the performances of CAD models and LLMs become

increasingly improved in the future, the proposed scheme has the potential to improve the quality of radiology reports and enhance the feasibility of online healthcare services.

2. Related Works

2.1. Large Language Models

Recent advances in Transformer architecture [28] and computing power have enabled the training of large language models with billions of parameters, leading to a significant improvement in their ability to summarize, translate, predict and generate human-like text [3, 23, 25].

Several domain-specific LLMs have been developed using general-purpose model weight and training schemes. BioBERT [13] and PubMedBERT [8] are examples of BERT [5] models trained on PubMed for biomedical data, while ClinicalBERT [2] was further trained on the MIMIC dataset and outperformed its predecessor. Med-PaLM [25] was developed in late 2022 using curated biomedical corpora and human feedback, and showed promising results, including a 67.6% accuracy on the MedQA exam. ChatGPT, which was not given supplementary medical training, passed all three parts of the USMLE and achieved over 50% accuracy across all exams and surpassed 60% accuracy in the majority of them [12].

2.2. Vision-Language Model

A popular method of converting visual information into language is through image captioning. Deep learning-based image caption models [9, 33] can generate descriptive and coherent captions using large datasets such as Microsoft COCO and Flickr 30K. In medical image analysis, image captioning methods are employed to generate exam image reports. For example, Li et al. [14] implement explicit medical abnormality graph learning for report generation. Zhang et al. [34] utilize a pre-constructed knowledge graph based on disease topics, respectively. Another line of research [4, 29] learns cross-modal patterns using self-attention architecture. The recent emergence of foundation models with more clinical knowledge holds promise as a potential future direction.

Recently, with the increase in model size, advances in the field have shifted towards Vision-Language Pretraining (VLP) and utilizing pre-trained models. CLIP [21] merges visual and language information into a shared feature space, setting new state-of-the-art performance on various downstream tasks. Frozen [27] fine-tunes an image encoder, whose outputs serve as soft prompts for the language model. Flamingo [1] introduces cross-attention layers into the LLM to incorporate visual features, pre-training these new layers on billions of image-text pairs.

3. Method

3.1. Bridge between Image and Text

The key idea is to utilize the powerful logical reasoning capabilities of the LLMs to make more robust disease diagnosis for medical images. Therefore, we need to build a bridge to translate medical images into texts as inputs for the LLM. Our strategy is straightforward: 1) Feed exam images (e.g., X-Ray) into trained CAD models to obtain outputs; 2) Translate these outputs (typically tensors) into natural language; 3) Use language models to summarize the results and make a final conclusion; 4) Based on the results from visual models and pre-trained medical knowledge in the language models, engage in conversation about symptoms, diagnosis, and treatment. In this section, we mainly discuss the details of our proposed scheme.

An example is illustrated in Figure 3, where the output of a disease classifier is a 5-value vector indicating the probabilities of five diseases, i.e., Cardiomegaly, Edema, Consolidation, Atelectasis, and Pleural Effusion. After that, we need to translate this result into a prompt sentence for the LLM. A natural way of prompting is to show all five kinds of pathology and their corresponding scores. We first tell the LLM “Higher disease score means higher possibility of illness” as the basic rule in order to avoid some misconception. Then, we represent the score of each disease as “\${disease}\$ score: \${score}\$” as shown in upper-right panel (Prompt#1). Reports generated using Prompt#1 can be found at second column in Figure 8 and Figure 9. One may notice that the LLMs are heavily influenced by Prompt#1, usually repeating all the numbers in the output. Reports generated from Prompt#1 are very different from radiologist’s reports since the concrete diagnostic scores is not frequently used in clinical settings.

To align with the language commonly used in clinical reports, we propose to transform the concrete scores into descriptions of disease severity as shown in lower-left panel (Prompt#2). Prompt#2 will be designed using a grading system, which will divide the scores into four categories: “No sign” [0.0-0.2), “Small possibility” [0.2-0.5), “Likely” [0.5-0.9), and “Definitely” [0.9 and above). These categories will be used to describe the likelihood of each of the five observations. Prompt#3 is a concise one that reports diseases with diagnosis scores higher than 0.5 in the prompt. If no prediction is made among all five diseases, the prompt will be “No Finding”. Reports generated from Prompt#2 and Prompt#3 are generally acceptable and reasonable in most cases as one can observe in Figure 8 and Figure 9. “Network A” is frequently referenced in the generated reports. Some prompt tricks, e.g., “*Revise the report based on results from Network A but without mentioning Network A*”, can be applied to removing its mention. We do not utilize these tricks in current experiments.

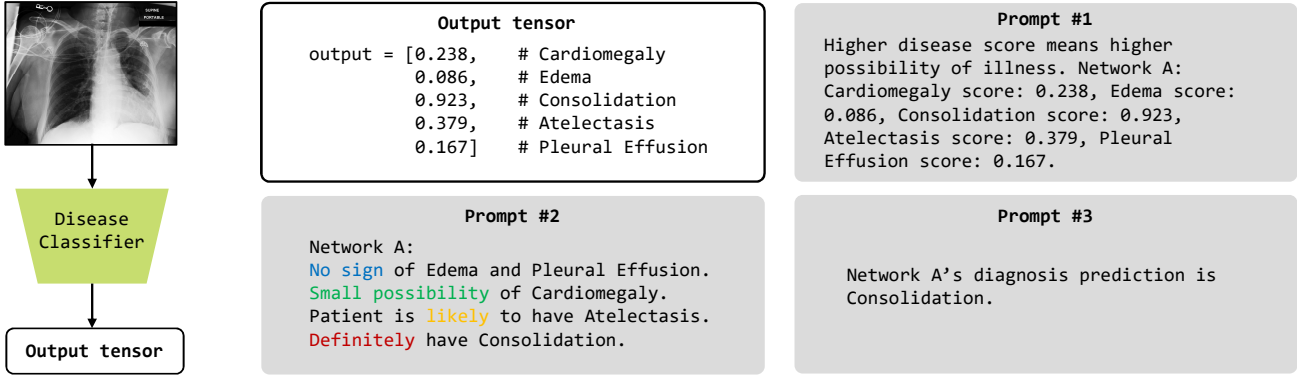


Figure 3. Prompts that bridge between tensor and text. We show three different prompt designs.

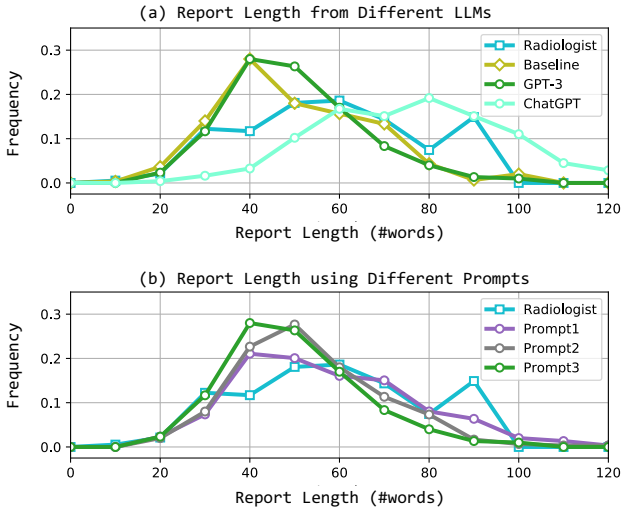


Figure 4. Length comparison of generated reports.

3.2. Dataset and Implementation

In this paper, we evaluate the performance of the combination of a report generation network (R2GenCMN [4]) and a classification network (PCAM [32]). The report generation networks (CvT2DistilGPT2 and R2GenCMN) are trained on the MIMIC-CXR dataset [11]. The MIMIC-CXR dataset is a large-scale public dataset of chest x-ray images with free-text radiology reports. It contains 377,110 images corresponding to 227,835 radiographic studies performed at the Beth Israel Deaconess Medical Center in Boston, MA. At the same time, the classifier is trained on the CheXpert dataset [10]. CheXpert is a large public dataset for chest radiograph interpretation, consisting of 224,316 chest radiographs of 65,240 patients.

The reports from the LLMs are tested on the official test set of the MIMIC dataset. Due to the current limitation of ChatGPT usage (i.e., around 20 requests per hour), we are

unable to test the entire test set of MIMIC-CXR now. Therefore, 300 cases are randomly selected, including 50 cases of Cardiomegaly, 50 cases of Edema, 50 cases of Consolidation, 50 cases of Atelectasis, 50 cases of Pleural effusion, and 50 cases with no findings. During the evaluation process, the text reports were converted to multi-class labels using cheXbert [26].

The LLMs are updating constantly to include more new knowledge and events, leading to the improvement of their reasoning capability. The GPT-3 model we use in this paper is *text-davinci-003* which was released by OpenAI on Feb, 2023 based on InstructGPT [18]. The maxlen of the output is set to 1024 and temperature set to 0.5. The ChatGPT [17] model used is the *Jan-30-2023* version. In the section “Interactive and Understandable CAD”, ChatGPT is used to generate the example. During our test, the GPT-3 can also provide accurate and helpful chat.

4. Report Generation

4.1. Quality Improvement of the Generated Report

In this section, we evaluate the performance of our proposed method with other two report-generation methods, i.e., R2GenCMN [4] and CvT2DistilGPT2 [16]. On the basis of clinical importance and prevalence, we focus on five kinds of observations. Three metrics, including precision (PR), recall (RC), and F1-score (F1), are reported in Table 1.

The strengths of our method are clearly shown in Table 1. It has obvious advantages in RC and F1, and is only weaker than R2GenCMN in the term of PR. Our method has a relatively high Recall and F1-score on MIMIC-CXR dataset. For all five kinds of diseases, both CvT2DistilGPT2 and R2GenCMN show inferior performance to our method concerning RC and F1. Specifically, their performances on Edema and Consolidation are rather low. Their RC values on Edema are 0.468 and 0.252, respectively, while

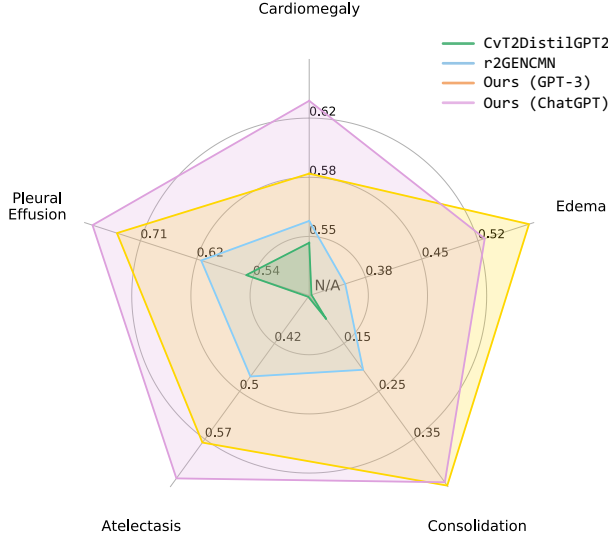


Figure 5. F1-score comparison on 5 observations.

our method achieves the RC value of 0.626 based on GPT-3. The same phenomenon can be observed in Consolidation, where the first two methods hold the values of 0.239 and 0.121 while ours (GPT-3) drastically outperforms them, with the RC value of 0.803. The R2GenCMN has a higher PR value compared to our method on three of five diseases. However, the cost of R2GenCMN’s high performance on Precision is its weakness in the other two metrics, which can lead to biased report generation, e.g., seldomly reporting any potential diseases. At the same time, our method has the highest F1 among all methods, and we believe it can be the most trustworthy report generator.

4.2. How LLMs affect Report Quality

In this section, we compare the performance of different LLMs for report generation. We use Prompt#3 as the default prompt. OpenAI provides four different sizes of GPT-3 models through its publicly accessible API: text-ada-001, text-babbage-001, text-curie-001, and text-davinci-003. The smallest text-ada-001 can not generate meaningful reports and is therefore not included in this experiment. The size of the models has not been officially disclosed. The figures listed in Table 2 are approximate estimates based on the information in [7].

We report the F1-score of all observations in Table 2. It is noteworthy that language models struggle to perform well in clinical tasks when their model size is limited. The diagnostic performances of text-babbage-001 and text-curie-001 is subpar, as demonstrated by their low average F1-scores over five observations compared with the last two models. The improvement in diagnostic performance is evident in text-davinci-003, whose model size is hundreds of times larger than that of text-babbage-001. On aver-

age, text-davinci-003’s F1-score is improved from 0.471 to 0.591. The ChatGPT is slightly better than text-davinci-003, achieving the improvement of 0.014, and their diagnostic abilities are comparable. Overall, the diagnostic capability of language models is proportional to their size, highlighting the critical role of the logistic reasoning capability of LLMs.

In our experiments, it can be observed that more capable models generally produce longer reports as shown in Figure 6. At the same time, nearly 40% of reports generated by text-babbage-001 and nearly 15% of reports generated by text-curie-001 have no meaningful content.

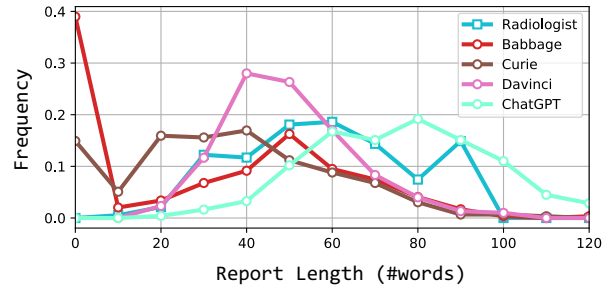


Figure 6. Length of Reports generated by different models. “Babbage”, “Curie” and “Davinci” represent the three GPT-3 models with different model sizes, i.e., text-babbage-001, text-curie-001, and text-davinci-003.

5. Interactive and Understandable CAD

The proposed ChatCAD offers several benefits, including its ability to utilize LLM’s extensive and reliable medical knowledge to provide interactive explanations and advice. As shown in Figure 7, two examples of the interactive CAD are provided, with one chat discussing pleural effusion and the other addressing edema and its relationship to swelling.

Through this approach, patients can gain a clearer understanding of their symptoms, diagnosis, and treatment options, leading to more efficient and cost-effective consultations with medical experts. As language models continue to advance and become more accurate with access to more trustworthy medical training data, ChatCAD has the potential to significantly enhance the quality of online healthcare services.

6. Limitations and Discussion

In this paper, we explore a novel framework, ChatCAD, introducing large language models in CAD. The proposed method, however, still has limitations to be solved.

First, LLM-generated reports are not human like in a certain way. LLM is likely to output sentences like “Network A’s diagnosis prediction is consistent with the findings in

Table 1. Diagnostic accuracy comparison with SOTA label generation methods. Our methods use the Prompt#3 setting.

Observations	CvT2DistilGPT2 [16]			R2GenCMN [4]			Ours (GPT-3)			Ours (ChatGPT)		
	PR	RC	F1	PR	RC	F1	PR	RC	F1	PR	RC	F1
Cardiomegaly	0.512	0.591	0.549	0.590	0.534	0.561	0.606	0.569	0.587	0.663	0.595	0.627
Edema	0.224	0.468	0.303	0.563	0.252	0.348	0.563	0.626	0.593	0.556	0.514	0.534
Consolidation	0.063	0.239	0.099	0.667	0.121	0.205	0.310	0.803	0.447	0.322	0.697	0.440
Atelectasis	0.306	0.388	0.342	0.442	0.504	0.471	0.408	0.991	0.578	0.470	0.981	0.636
Pleural Effusion	0.454	0.692	0.548	0.819	0.500	0.618	0.634	0.916	0.749	0.736	0.845	0.787
Average	0.312	0.476	0.368	0.616	0.382	0.441	0.504	0.781	0.591	0.549	0.726	0.605

Table 2. F1-score comparison of different-size LLMs

Model	Size	Cardiomegaly	Edema	Consolidation	Atelectasis	Pleural Effusion	Average
text-babbage-001	~1.3B	0.350	0.479	0.418	0.471	0.639	0.471
text-curie-001	~6.7B	0.529	0.451	0.369	0.515	0.674	0.508
text-davinci-003	~175B	0.587	0.593	0.447	0.578	0.749	0.591
ChatGPT	~175B	0.627	0.534	0.440	0.636	0.787	0.605

the radiological report” or “The findings from Network A’s diagnosis prediction are supported by the X-ray”. This is reflected on natural language similarity metrics when we compare to our baseline method. ChatCAD improved the diagnosis accuracy but dropped the BLEU score [20]. A promising way to address this issue is to add a module after ChatGPT to filter generated reports. Or add prompt like “please do not mention Network A”.

Additionally, we only design three typical kinds of prompts that are intuitive, and there is room for improvement. LLMs are capable of solving logical reasoning problems without additional computational costs [31]. In current ChatCAD, we did not give the network about patient’s major complaint since there is no such dataset available. We believe the LLMs can process more complex information than what we currently provide. Better datasets and benchmarks are needed.

Our experiments demonstrate significant impact of language model size on diagnostic accuracy. Larger, advanced, and more truthful LLMs such as the upcoming GPT-4 may improve the accuracy and report quality further. However, the role of vision classifiers has not yet been explored, and additional research is necessary to determine if models such as ViT [6] or SwinTransformer [15], which boast larger parameters, can deliver improved results. On the other hand, LLMs can also be used to help the training of vision models, such as correcting outputs of vision models using related medical knowledge learned in LLMs.

In our work, we have only carried out a qualitative analysis of the prompt design instead of a quantitative analysis. Further in-depth investigations will be undertaken once the API for ChatGPT becomes available for use. Moreover, the specifics of this paper have not been discussed with any

clinical professionals, and therefore it still lacks rigor in many places. We will improve it in subsequent versions.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022. 3
- [2] Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*, 2019. 3
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1, 3
- [4] Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, Aug. 2021. 3, 4, 6
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1, 3
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 6

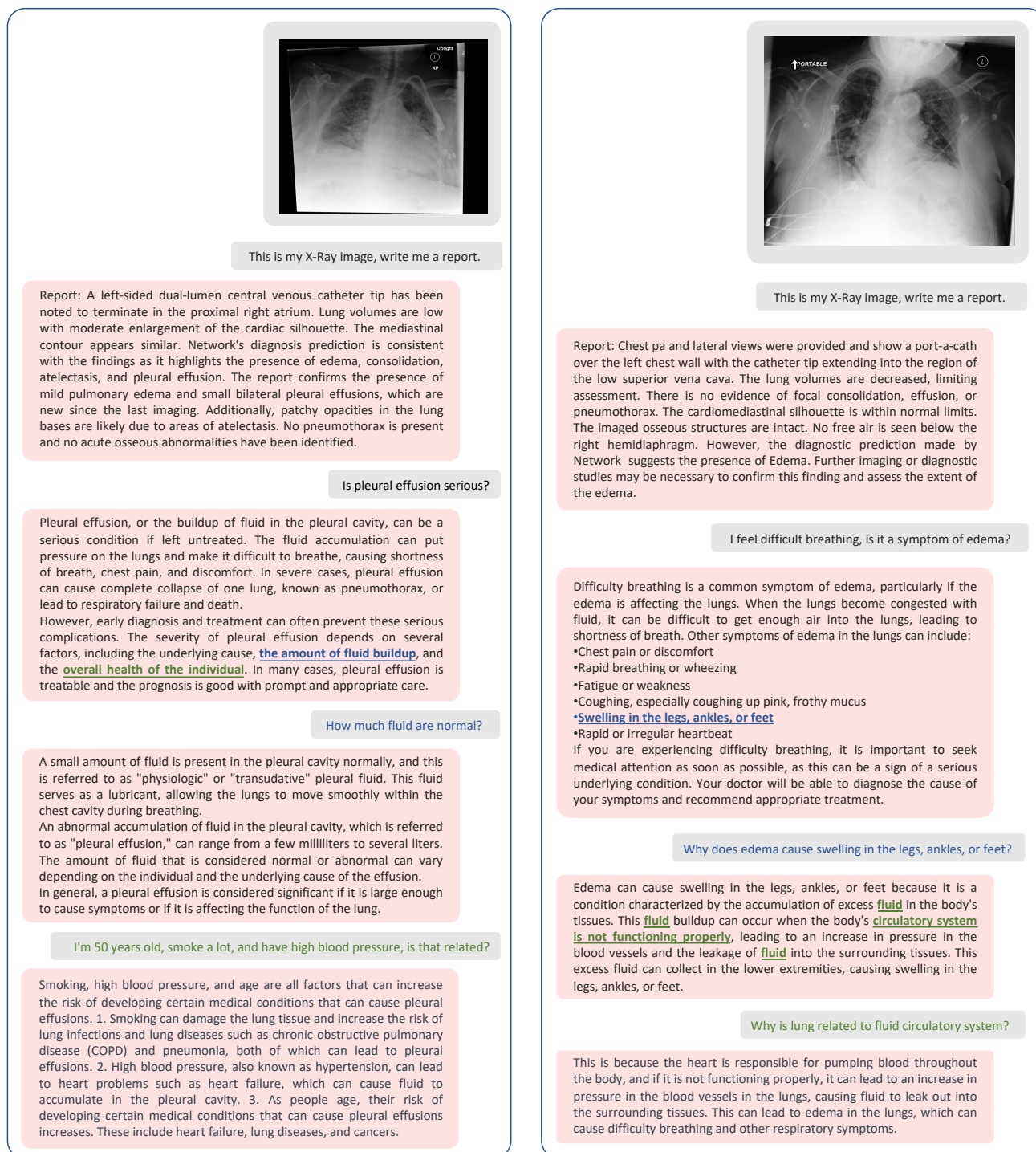


Figure 7. Two examples of interactive CAD using our framework and ChatGPT.

- [7] Leo Gao. On the sizes of openai api models. 5
- [8] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pre-training for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021. 3
- [9] Simao Herdade, Armin Kappeler, Kofi Boakye, and Joao

- Soares. Image captioning: Transforming objects into words. *Advances in neural information processing systems*, 32, 2019. 3
- [10] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019. 4
- [11] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. 4
- [12] Tiffany H Kung, Morgan Cheatham, Arielle Medinilla, ChatGPT, Czarina Sillos, Lorie De Leon, Camille Elepano, Marie Madriaga, Rimel Aggabao, Giezel Diaz-Candido, et al. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. *medRxiv*, pages 2022–12, 2022. 3
- [13] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020. 3
- [14] Guanxiong Liu, Tzu-Ming Harry Hsu, Matthew McDermott, Willie Boag, Wei-Hung Weng, Peter Szolovits, and Marzyeh Ghassemi. Clinically accurate chest x-ray report generation. In *Machine Learning for Healthcare Conference*, pages 249–269. PMLR, 2019. 3
- [15] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 6
- [16] Aaron Nicolson, Jason Dowling, and Bevan Koopman. Improving chest x-ray report generation by leveraging warm-starting. *arXiv preprint arXiv:2201.09405*, 2022. 4, 6
- [17] OpenAI. Chatgpt: Optimizing language models for dialogue, 2023. 1, 4
- [18] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022. 4
- [19] Xi Ouyang, Jiayu Huo, Liming Xia, Fei Shan, Jun Liu, Zhanhao Mo, Fuhua Yan, Zhongxiang Ding, Qi Yang, Bin Song, et al. Dual-sampling attention network for diagnosis of covid-19 from community acquired pneumonia. *IEEE Transactions on Medical Imaging*, 39(8):2595–2605, 2020. 2
- [20] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 6
- [21] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 3
- [22] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. 1
- [23] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 3
- [24] Dinggang Shen, Guorong Wu, and Heung-Il Suk. Deep learning in medical image analysis. *Annual review of biomedical engineering*, 19:221–248, 2017. 2
- [25] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022. 3
- [26] Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y Ng, and Matthew P Lungren. Chexpert: combining automatic labelers and expert annotations for accurate radiology report labeling using bert. *arXiv preprint arXiv:2004.09167*, 2020. 4
- [27] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021. 3
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 3
- [29] Jun Wang, Abhir Bhalerao, and Yulan He. Cross-modal prototype driven network for radiology report generation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*, pages 563–579. Springer, 2022. 3
- [30] Sheng Wang, Xi Ouyang, Tianming Liu, Qian Wang, and Dinggang Shen. Follow my eye: Using gaze to supervise computer-aided diagnosis. *IEEE Transactions on Medical Imaging*, 2022. 2
- [31] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022. 6
- [32] Wenwu Ye, Jin Yao, Hui Xue, and Yi Li. Weakly supervised lesion localization with probabilistic-cam pooling, 2020. 4
- [33] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image captioning with semantic attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4651–4659, 2016. 3
- [34] Yixiao Zhang, Xiaosong Wang, Ziyue Xu, Qihang Yu, Alan Yuille, and Daguang Xu. When radiology report generation meets knowledge graph. In *Proceedings of the AAAI Con-*

ference on Artificial Intelligence, volume 34, pages 12910–12917, 2020. [3](#)

- [35] Xiangyu Zhao, Zengxin Qi, Sheng Wang, Qian Wang, Xuehai Wu, Ying Mao, and Lichi Zhang. Rcps: Rectified contrastive pseudo supervision for semi-supervised medical image segmentation. *arXiv preprint arXiv:2301.05500*, 2023. [2](#)

ChatGPT's Prompt Test


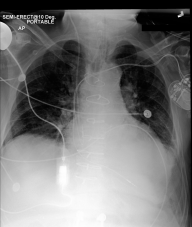
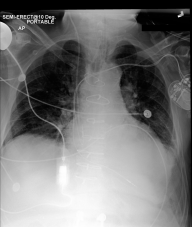

Image	Prompt1	Prompt2	Prompt3
	<p>A left-sided AICD device is noted with leads terminating in the right atrium and right ventricle. Significant enlargement of the cardiac silhouette is demonstrated with a Cardiomegaly score of 0.946. The mediastinal contour is similar. Mild to moderate Pulmonary Edema is present with an Edema score of 0.588, and small bilateral Pleural Effusions are observed with a Pleural Effusion score of 0.849. There are patchy opacities in the lung bases, likely reflecting areas of Atelectasis with a score of 0.708. No pneumothorax is present, and there are no acute osseous abnormalities.</p>	<p>The patient is found to have multiple conditions affecting their chest and heart. A left-sided AICD device is present, with leads terminating in the right atrium and right ventricle. Moderate to severe enlargement of the cardiac silhouette (cardiomegaly) is noted and is confirmed by the results from Network A. The mediastinal contour is unchanged. The presence of mild pulmonary edema is new and is accompanied by small bilateral pleural effusions. Patchy opacities in the lung bases are noted and are likely reflective of areas of atelectasis, as also indicated by the results from Network A. No pneumothorax is present. There are no acute osseous abnormalities seen on imaging. The patient is likely to have edema, consolidation, and atelectasis, along with a pleural effusion, as indicated by the results from Network A.</p>	<p>There is moderate to severe enlargement of the cardiac silhouette, consistent with a diagnosis of cardiomegaly. The contour of the mediastinum is similar. There is evidence of new mild pulmonary edema, with small bilateral pleural effusions and patchy opacities in the lung bases that are likely due to atelectasis. A left-sided AICD device is present with leads terminating in the right atrium and right ventricle. No pneumothorax is identified and there are no acute osseous abnormalities.</p>
	<p>As compared to the previous radiograph, the patient has received a right internal jugular vein catheter. The course of the catheter is unremarkable and the tip of the catheter projects over the mid-superior vena cava. However, the current radiograph shows findings of cardiomegaly with a score of 0.189, edema with a score of 0.713, consolidation with a score of 0.479, atelectasis with a score of 0.803, and pleural effusion with a score of 0.702. Despite these findings, there is still no evidence of complications, notably no pneumothorax. The radiograph remains otherwise unchanged.</p>	<p>The latest radiograph shows that there is no sign of Cardiomegaly. However, there is a small possibility of Consolidation. The patient is likely to have Edema, Atelectasis, or Pleural Effusion. The patient has received a right internal jugular vein catheter, and its course appears to be unremarkable. The tip of the catheter projects over the mid-superior vena cava, and there is no evidence of complications such as pneumothorax. The radiograph is otherwise unchanged from the previous one.</p>	<p>As compared to the previous radiograph, the patient has received a right internal jugular vein catheter and the course of the catheter is unremarkable with the tip of the catheter projecting over the mid SVC. There is no evidence of complications, notably no pneumothorax. However, the diagnosis prediction from Network A suggests the presence of Edema, Atelectasis, and Pleural Effusion, which indicate the patient may have fluid accumulation in the lung and surrounding tissues. Further evaluation and medical intervention may be necessary to address these conditions. The overall radiograph remains otherwise unchanged.</p>
	<p>Ap portable upright view of the chest shows a dialysis catheter with its tip extending into the right atrium. The lung volumes are low, which limits the assessment of the chest. There is a right IJ central venous catheter with its tip in the region of the low SVC. The heart is moderately enlarged with a Cardiomegaly score of 0.178. There is significant pulmonary edema with an Edema score of 0.496 and small bilateral pleural effusions with a Pleural Effusion score of 0.725. There is evidence of consolidation with a Consolidation score of 0.509 and Atelectasis with a score of 0.704. No pneumothorax is present, and the mediastinal contour is stable. The bony structures are intact.</p>	<p>The ap portable upright view of the chest showed the presence of a dialysis catheter with its tip extending into the right atrium and a right ij central venous catheter with its tip in the region of the low svc. The heart remains mildly enlarged with no sign of cardiomegaly, however, there is a small possibility of edema present. The lung volumes are low, limiting assessment, and there is mild pulmonary edema with small bilateral pleural effusions noted. No pneumothorax was detected, and the mediastinal contour is stable. Bony structures were intact.</p>	<p>The ap portable upright view of the chest shows a dialysis catheter with its tip extending into the right atrium. There is also a right internal jugular central venous catheter with its tip in the region of the low superior vena cava. The heart is mildly enlarged and there is evidence of mild pulmonary edema with small bilateral pleural effusions, consistent with the diagnosis of pleural effusion. No pneumothorax is seen, and the mediastinal contour is stable. The bony structures are intact. However, lung volumes are low, limiting the assessment of the chest. Based on the results from Network A, the findings are consistent with a diagnosis of consolidation and atelectasis in addition to pleural effusion.</p>
	<p>Chest X-rays, both ap upright and lateral views, were performed. The results from a deep learning network (Network A) indicate the presence of cardiomegaly with a score of 0.977. There is also mild evidence of pulmonary edema (score: 0.492) and a suggestion of consolidation (score: 0.493) and atelectasis (score: 0.672). The network also indicates the presence of pleural effusion (score: 0.598). Hilar congestion and mediastinal contour stability were noted on the original report. No large effusion or pneumothorax was detected. The bony structures are intact. Although the original report mentioned mild pulmonary edema, the results from Network A suggest that it may be more pronounced. No convincing evidence of pneumonia was found on either the original report or the results from Network A.</p>	<p>The chest x-rays show evidence of cardiomegaly and hilar congestion. Mild pulmonary edema is also noted, but there is no large effusion or pneumothorax. The mediastinal contour appears stable, and the bony structures are intact. The findings are consistent with those from Network A, which suggests a small possibility of edema and consolidation, with a likelihood of atelectasis and pleural effusion. There is no convincing evidence of pneumonia.</p>	<p>The ap upright and lateral views of the chest were reviewed and the results from Network A were taken into consideration. The diagnosis of cardiomegaly is confirmed, with evidence of hilar congestion and mild pulmonary edema. The presence of pleural effusion is also noted. Although no large effusion or pneumothorax was initially identified, the effusion seen on the chest images is consistent with the diagnosis of pleural effusion as identified by Network A. The mediastinal contour remains stable, and the bony structures are intact. No convincing evidence for pneumonia was identified on the initial review.</p>

Figure 8. ChatGPT's generated reports from three different prompt designs.

GPT-3's Prompt Test


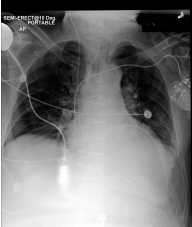


Image	Prompt1	Prompt2	Prompt3
 <p>s55498995</p>	<p>Left-sided AICD device is noted with leads terminating in the right atrium and right ventricle. Moderate to severe cardiomegaly is re-demonstrated. The mediastinal contour is similar. There is mild pulmonary edema with small bilateral pleural effusions and patchy opacities in the lung bases likely reflecting areas of atelectasis and consolidation. No pneumothorax is present. There are no acute osseous abnormalities.</p>	<p>left-sided aicd device is noted with leads terminating in the right atrium and right ventricle . moderate to severe enlargement of the cardiac silhouette is re- demonstrated . the mediastinal contour is similar . there is mild pulmonary edema new in the interval with small bilateral pleural effusions . patchy opacities in the lung bases likely reflect areas of atelectasis . no pneumothorax is present . there are no acute osseous abnormalities .</p>	<p>Left-sided AICD device is noted with leads terminating in the right atrium and right ventricle. Moderate to severe enlargement of the cardiac silhouette is re-demonstrated. The mediastinal contour is similar. There is mild pulmonary edema new in the interval with small bilateral pleural effusions and evidence of edema, consolidation, and atelectasis. No pneumothorax is present. There are no acute osseous abnormalities.</p>
 <p>s59707249</p>	<p>As compared to the previous radiograph, the patient has received a right internal jugular vein catheter. The course of the catheter is unremarkable and the tip of the catheter projects over the mid SVC. There is no evidence of complications, notably no pneumothorax. However, there is an increased cardiomegaly score of 0.189, edema score of 0.713, consolidation score of 0.479, atelectasis score of 0.803, and</p>	<p>as compared to the previous radiograph the patient has received a right internal jugular vein catheter . the course of the catheter is unremarkable the tip of the catheter projects over the mid svc . there is no evidence of complications notably no pneumothorax . otherwise unchanged radiograph</p>	<p>No sign of Cardiomegaly or Consolidation. Patient is likely to have Edema, Atelectasis, Pleural Effusion. There is no evidence of complications notably no pneumothorax. Otherwise unchanged radiograph. The right internal jugular vein catheter is in place with an unremarkable course and the tip of the catheter projects over the mid SVC.</p>
 <p>s53350789</p>	<p>A portable upright view of the chest is seen. A dialysis catheter is again seen with its tip extending into the right atrium. The heart is mildly enlarged with cardiomegaly (score 0.178). There is evidence of pulmonary edema (score 0.496) and mild bilateral pleural effusions (score 0.725). There is also evidence of consolidation (score 0.509) and atelectasis (score 0.704). No pneumothorax is seen</p>	<p>ap portable upright view of the chest . dialysis catheter is again seen with its tip extending into the right atrium . lung volumes are low limiting assessment . there is a right ij central venous catheter with its tip in the region of the low svc . the heart remains mildly enlarged . there is mild pulmonary edema with small bilateral pleural effusions . no pneumothorax . mediastinal contour is stable . bony structures are intact .</p>	<p>Ap portable upright view of the chest. Dialysis catheter is again seen with its tip extending into the right atrium. Lung volumes are low limiting assessment. There is a right IJ central venous catheter with its tip in the region of the low SVC. The heart remains mildly enlarged. There is no sign of Cardiomegaly. There is small possibility of Edema. Patient is likely to have Consolidation, Atelectasis, and Pleural Effusion. No</p>
 <p>s59379638</p>	<p>Ap upright and lateral views of the chest provided. Cardiomegaly is again noted with mild pulmonary edema and mild pleural effusion. Hilar congestion is present. No large effusion or pneumothorax. No convincing evidence for consolidation or atelectasis. Mediastinal contour is stable. Bony structures are intact.</p>	<p>ap upright and lateral views of the chest provided . cardiomegaly is again noted with hilar congestion and mild pulmonary edema . no large effusion or pneumothorax . no convincing evidence for pneumonia . mediastinal contour is stable . bony structures are intact</p>	<p>Ap upright and lateral views of the chest provided. Cardiomegaly is again noted with hilar congestion and mild pulmonary edema. No large effusion or pneumothorax. No convincing evidence for pneumonia. Atelectasis and pleural effusion are present. Mediastinal contour is stable. Bony structures are intact.</p>

Figure 9. GPT-3's generated reports from three different prompt designs.