

INDICATION AS PRIOR KNOWLEDGE FOR MULTIMODAL DISEASE CLASSIFICATION IN CHEST RADIOGRAPHS WITH TRANSFORMERS

Grzegorz Jacenków¹ Alison Q. O’Neil^{1,2} Sotirios A. Tsafaris^{1,3}

¹ The University of Edinburgh ² Canon Medical Research Europe ³ The Alan Turing Institute

ABSTRACT

When a clinician refers a patient for an imaging exam, they include the reason (e.g. relevant patient history, suspected disease) in the scan request; this appears as the indication field in the radiology report. The interpretation and reporting of the image are substantially influenced by this request text, steering the radiologist to focus on particular aspects of the image. We use the indication field to drive better image classification, by taking a transformer network which is unimodally pre-trained on text (BERT) and fine-tuning it for multimodal classification of a dual image-text input. We evaluate the method on the MIMIC-CXR dataset, and present ablation studies to investigate the effect of the indication field on the classification performance. The experimental results show our approach achieves 87.8 average micro AUROC, outperforming the state-of-the-art methods for unimodal (84.4) and multimodal (86.0) classification. Our code is available at <https://github.com/jacenkow/mmbt>.

Index Terms— Multimodal Learning, Chest X-Ray Classification, Transformers, BERT

1. INTRODUCTION

Chest radiography remains the most common imaging examination for the diagnosis and treatment of a variety of lung conditions such as pneumonia, cancer, and even COVID-19. Automation of X-ray interpretation could considerably improve healthcare systems, lowering costs and addressing the pressing challenge of expert shortage [1]. Yet, current techniques for clinical decision support mostly focus on a single modality (e.g. patient’s X-ray) and do not take into account complementary information which might be already available in a hospital’s database (e.g. patient’s clinical history) [2], [3]. We are particularly interested in providing the indication field, i.e., the motivation for the patient’s screening examination. This field may include the patient’s history, a request to evaluate a particular condition, and other clues which can steer the radiologist’s attention to particular imaging features. The indication field is often the only information provided by the referring physician [4], and can influence the interpretation of the imaging exam [5]. In this paper, we want to design a vision-and-language model that is able to use such text-based

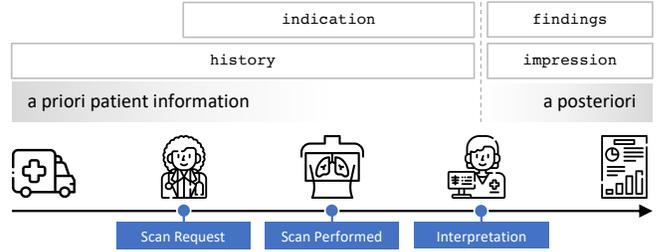


Fig. 1: We consider the problem of classifying chest X-ray images given the patient information in a free-text form. We only use knowledge about the patient collected before the imaging examination and do not require radiologist intervention as opposed to most prior studies.

side information to aid and complement disease classification.

Current state-of-the-art methods for vision-and-language tasks (such as VisualBERT [6]) are mostly based on transformer architectures, which require extensive pre-training. The process typically involves using a dataset with annotated bounding boxes around the objects of interests, such as Conceptual Captions [7], to initialise the weights, which are later fine-tuned to the final task. Unfortunately, the biomedical community lacks domain-specific yet general multimodal datasets which could be used for pre-training large transformer networks. To address this problem, one could leverage existing unimodal models, and fine-tune the models to a multimodal task as proposed in multimodal BERT (MMBT) [8], which we evaluate on a biomedical task. As BERT does not provide the means to process imaging input, MMBT embeds image features from a ResNet-152 [9] classifier.

We evaluate the ability of a unimodally pre-trained BERT model to process biomedical imaging and non-imaging modalities during the fine-tuning step. Specifically, we use chest radiographs and the indication field from associated radiology reports to perform multi-label classification. The network can be pre-trained on unimodal datasets which are more common than multimodal, but it is still capable of learning multimodal interactions during the fine-tuning step.

Contributions: (1) We present a strong baseline for multimodal classification of chest radiographs; (2) We evaluate the model with the prior work achieving the new state-of-the-art

results, and show its robustness to adversarial changes in text.

2. RELATED WORK

Chest X-Ray Classification. Most work for classifying chest radiographs has been based on existing convolutional neural networks (CNNs) with ResNet-50 [9] being the most popular architecture [1]. Several works have proposed to exploit non-imaging data such as patient’s demographics to improve performance. The information is often fused before the final classification layer by concatenating imaging and non-imaging features [1], [10]; this late fusion of modalities limits the methods to model signal-level interactions between imaging and non-imaging information. Moreover, the non-imaging modality has limited expressive power as it only relates to basic demographics and not to the patient’s history. We decide to use the indication field from full-text reports. The free-text input includes relevant information for the imaging procedure, allowing the network to learn more complex interactions between input images and the associated reports. **Learning with Radiology Reports.** TieNet [11] combines image-text pairs to learn a common embedding for classification and report generation. The method uses multi-level attention with CNN and RNN networks for processing radiographs and reports respectively. However, the full report is expected as input, which requires a radiologist to render findings first. Recently, two methods [12], [13] proposed to leverage information available in radiology reports to improve performance of image-only classification. The methods are optimised with a loss encouraging learning a shared representation between two modalities, while keeping the modalities (and the downstream tasks) decoupled. The results show improvement in classification performance, but the methods ignore the additional non-imaging information during inference. Our work follows the same motivation as [14], [15], where the methods only include information available prior to the examination. The first work [15] to include the indication field uses the information only to improve the quality of rendering the diagnosis (impression field) leaving the classification head only dependent on the imaging features. The setup was adapted in [14] to support classification (and impression generation) with both modalities. The authors use an attention layer to merge the output of two feature extractors for image and text, which we term a middle fusion approach. We propose to use a transformer network which is capable of modelling the interactions at the word level, enabling the network to perform more complex fusion. Recently, a study [16] has shown the visual-linguistic BERT models are suitable for processing chest radiographs and the associated radiology reports, outperforming unimodal approaches for text-only. However, the evaluated models use full-text reports making the use of the imaging input negligible and clinically unpractical. By contrast, we propose information only available to the radiologist prior to developing a report to drive better image classifica-

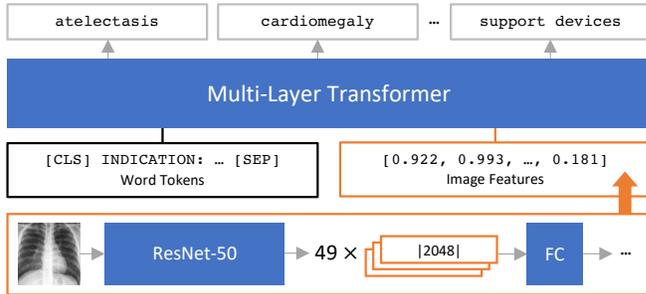


Fig. 2: The overview of method. We extend a multi-layer transformer pre-trained on textual data with imaging input. The images are provided as features extracted from a ResNet-50 network. The features are reshaped to 49 vectors of 2048 dimensions each and combined with two embeddings describing segment (image or text) and position of the token.

tion rather than labelling text reports which already can be effortlessly classified by ruled-based approaches.

3. METHODOLOGY

State-of-the-art methods for modelling vision-language tasks are mostly based on the transformer architecture where the second segment provides visual tokens from an image feature extractor. However, pre-training also requires large, and general multimodal datasets where the visual objects are annotated with bounding boxes, and such datasets are lacking in the biomedical community. We exploit unimodally pre-trained BERT model and fine-tune it to a multimodal task.

Backbone Network: We adapt BERT [17] as our backbone network. We use the Hugging Face implementation of `bert-base-uncased` pre-trained on textual input. As the original model has not been developed for visual-linguistic tasks, we learn a new embedding for the image tokens.

Image Encoder: Our method uses ResNet-50 as the image feature extractor. We first fine-tuned the network pre-trained on ImageNet to classify chest radiographs (also a baseline method) and removed the last pooling layer. The network outputs 2048 feature maps of 7×7 , which we reshape to 49 vectors. Our image tokens are the sum of three embeddings, i.e., the linear projection of the i^{th} vector ($i \in [1, 49]$), the position of the vector i , and the segment indicating the imaging modality. We keep the weights of the image encoder unfrozen during the fine-tuning step of the whole model.

Classification Head: We use the final representation of [CLS] token to fine-tune our model for classification. We apply a multi-layer perceptron $\{768 - 768 - 14\}$ with GELU activation functions and layer normalisation. The last layer applies a *sigmoid* function to each of fourteen nodes.

Loss Function: We optimise a binary cross-entropy loss with class weighting that is inversely proportional to the number of examples in the training set.

Table 1: Quantitative results on the MIMIC-CXR dataset. We report average accuracy, precision, recall, F_1 score, and the area under the ROC (AUROC). The results are reported as average over three runs with standard deviation reported as subscript. The number in **bold** denotes the best performance within the metric.

Method	Modality		Accuracy	Precision		Recall		F_1		AUROC	
	Image	Text		Macro	Micro	Macro	Micro	Macro	Micro	Macro	Micro
CheXpert Labeler	✗	✓	80.6	9.3	13.4	18.8	27.0	8.5	17.9	51.13	53.3
BERT	✗	✓	85.1 \pm 0.2	21.6 \pm 1.0	32.9 \pm 0.9	47.2 \pm 5.8	54.7 \pm 0.9	26.1 \pm 1.1	41.1 \pm 0.9	71.1 \pm 1.2	81.7 \pm 0.8
ResNet-50	✓	✗	86.0 \pm 0.2	26.0 \pm 1.1	43.7 \pm 2.1	34.0 \pm 1.8	57.4 \pm 1.5	27.4 \pm 0.4	49.5 \pm 0.8	73.8 \pm 0.5	84.4 \pm 0.6
Attentive	✓	✓	86.8 \pm 0.1	26.8 \pm 0.5	44.2 \pm 0.7	34.7 \pm 0.2	61.3 \pm 0.6	29.1 \pm 0.4	51.4 \pm 0.3	76.6 \pm 0.3	86.0 \pm 0.2
MMBT	✓	✓	87.7 \pm 0.2	30.8 \pm 0.3	47.8 \pm 0.7	55.4 \pm 1.8	64.7 \pm 0.7	35.0 \pm 0.6	55.0 \pm 0.6	80.6 \pm 0.1	87.8 \pm 0.1

4. EXPERIMENTS

4.1. Dataset

We use the MIMIC-CXR dataset [18]–[20], which consists of 377,110 chest X-ray images associated with 227,835 post-screening reports of 65,379 patients taken at the Beth Israel Deaconess Medical Center Emergency Department. We limit the experiments to examinations with frontal images (AP/PA), and reports with `indication` or `history` fields explicitly. Our final evaluation is based on 210,538 studies following the official splits between training (205,923), validation (1695) and test sets (2920).

Labelling. The original data are not labelled for the classification task. We use the CheXpert Labeler [21] to extract fourteen labels from full radiology reports: atelectasis, cardiomegaly, consolidation, edema, enlarged cardiomeastinum, fracture, lung lesion, lung opacity, no finding, pleural effusion, pleural (other), pneumonia, pneumothorax, and support devices. We set the task as a multilabel problem with positive-vs-rest classification¹.

Pre-processing. The images were taken from the MIMIC-CXR-JPG dataset and resized to 224×224 pixels. We normalise the images to zero mean and unit of standard deviation. The text input has been stripped from special characters (e.g. “_”, “\”) and all characters converted to lower case.

4.2. Baselines

We compare the investigated method to several baselines:

- **CheXpert Labeler [21]:** This is the rules-based method used to extract the original fourteen labels from the full reports. We apply this method to the indication fields.
- **BERT [17]:** We use the unimodal BERT network which is the backbone of the proposed method with no access to the imaging input. We use the same classification head to fine-tune the network for classification.
- **ResNet-50 [9]:** We use the ResNet-50 network pre-trained on ImageNet (image feature extractor in the pro-

posed method), which we fine-tune to classify the chest radiographs.

- **Attentive [14]:** We compare our model to the multimodal approach presented in [14]. The method uses ResNet-50 and BioWordVec [22] with GRU units for feature extraction, with the two branches merged using an attention layer. The original method also generates impression fields (not included in our pipeline).

4.3. Experimental Setup

All baseline methods and the proposed technique were implemented with the multimodal framework (MMF) [23]. We train the models for 14 epochs with a batch size of 128. We use the Adam optimiser with weight decay (0.01). We set the learning rate to 5×10^{-5} with a linear warm-up schedule for the first 2000 steps. We apply the early stopping criterion of multi-label micro F_1 score evaluated on the validation set. We repeat each experiment three times with different seeds to account for variance due to random weight initialisation.

4.4. Results: Classification Performance

We report the performance of the tested methods using label-wise accuracy, precision, and recall metrics where we consider a separate classifier for each of fourteen classes. The overall quantitative results are shown in Table 1. We observe the CheXpert Labeler has the weakest performance across all of the reported metrics. The method is a rule-based approach, so it cannot learn associations between the content of indication fields and the labels, but will pick up only explicit mentions. This problem is mitigated by BERT (text-only) classifier which outperforms the labeler in all metrics (+53.3% improvement in micro AUROC). We further notice the image-only based classifier (ResNet-50) outperforms the BERT in all metrics except recall (macro) with micro AUROC improved by +3.3%. These findings are consistent with our expectation images contain the investigation results requested to help determine a diagnosis, compared to the text modality which describes only the clinician’s suspicion based on patient information prior to imaging. The Attentive [14] baseline, which uses both image and text, outperforms the image-

¹CheXpert Labeler is capable of assigning each label one of four values - positive, negative, uncertain and no mention. We only select the positive instances.

Table 2: The performance of the MMBT to robustness evaluation and manipulation to the indication field. We use the evaluation scheme proposed in [24] and further extend with swapping the indication field (no input, stop words, different patient).

Robustness Evaluation	Accuracy	Precision		Recall		F_1		AUROC	
		Macro	Micro	Macro	Micro	Macro	Micro	Macro	Micro
Baseline	87.7 \pm 0.2	30.8 \pm 0.3	47.8 \pm 0.7	55.4 \pm 1.8	64.7 \pm 0.7	35.0 \pm 0.6	55.0 \pm 0.6	80.6 \pm 0.1	87.8 \pm 0.1
Character Swap	87.0 \pm 0.2	27.4 \pm 0.5	44.7 \pm 0.7	48.3 \pm 8.0	62.1 \pm 0.6	30.4 \pm 0.7	52.0 \pm 0.6	78.0 \pm 0.4	86.5 \pm 0.1
Keyboard Typo	86.9 \pm 0.2	27.6 \pm 0.4	45.3 \pm 0.6	46.6 \pm 3.5	61.8 \pm 1.0	30.4 \pm 0.1	52.3 \pm 0.6	78.2 \pm 0.1	86.4 \pm 0.1
Synonyms	87.2 \pm 0.2	29.1 \pm 0.8	46.1 \pm 0.3	49.0 \pm 1.5	62.6 \pm 0.8	32.7 \pm 0.4	53.1 \pm 0.5	79.3 \pm 0.3	86.9 \pm 0.1
Missing Field	86.2 \pm 0.2	24.4 \pm 1.6	42.0 \pm 1.4	38.3 \pm 3.1	58.9 \pm 1.2	27.0 \pm 1.2	49.0 \pm 0.8	75.1 \pm 0.3	84.6 \pm 0.3
Stop Words Noise	86.2 \pm 0.1	20.2 \pm 0.7	35.0 \pm 1.6	38.6 \pm 4.9	60.7 \pm 0.4	24.0 \pm 0.7	44.4 \pm 1.3	74.7 \pm 0.3	84.3 \pm 0.2
Indication Swap	84.1 \pm 0.1	19.7 \pm 0.8	33.4 \pm 0.6	30.1 \pm 0.6	49.3 \pm 0.5	22.7 \pm 0.8	39.8 \pm 0.6	69.1 \pm 0.4	80.0 \pm 0.4

Table 3: AUROC results per-class for the tested methods.

AUROC	ResNet-50	BERT	Attentive	MMBT
Atelectasis	72.6 \pm 1.4	66.5 \pm 0.6	73.6 \pm 0.3	75.8 \pm 0.2
Cardiomegaly	74.6 \pm 0.9	74.1 \pm 0.4	77.3 \pm 0.6	82.6 \pm 0.2
Consolidation	69.9 \pm 0.9	66.8 \pm 0.9	73.4 \pm 0.6	77.1 \pm 0.3
Edema	81.6 \pm 0.5	70.7 \pm 0.4	81.9 \pm 0.5	84.3 \pm 0.4
Enlarged Card.	63.3 \pm 0.9	69.0 \pm 0.5	66.6 \pm 1.5	74.3 \pm 1.1
Fracture	63.3 \pm 0.9	65.9 \pm 1.5	67.0 \pm 0.7	72.9 \pm 2.0
Lung Lesion	66.5 \pm 2.5	69.9 \pm 1.8	70.9 \pm 0.7	75.9 \pm 1.2
Lung Opacity	68.3 \pm 0.9	62.5 \pm 0.7	69.2 \pm 0.2	71.5 \pm 0.5
No Finding	77.9 \pm 0.8	72.6 \pm 0.9	80.9 \pm 0.2	83.1 \pm 0.3
Ple. Effusion	86.1 \pm 0.6	71.6 \pm 0.4	87.1 \pm 0.1	88.6 \pm 0.1
Pleural Other	79.0 \pm 1.8	72.3 \pm 1.2	78.6 \pm 1.0	86.9 \pm 0.5
Pneumonia	66.1 \pm 2.1	68.9 \pm 0.4	70.9 \pm 0.8	75.2 \pm 0.7
Pneumothorax	78.3 \pm 0.4	86.5 \pm 1.0	85.1 \pm 0.7	88.0 \pm 0.3
Support Devices	85.9 \pm 0.5	88.7 \pm 0.3	90.5 \pm 0.2	92.2 \pm 0.1

and text-only methods in all reported metrics with micro AUROC improved by 1.9% comparing to the best unimodal baseline. Finally, the multimodal BERT outperforms all unimodal and multimodal baselines with 2% margin. The method relies on the early fusion approach (as opposed to middle fusion in Attentive) enabling the network to learn correlation and interactions between the modalities with low-level features. Moreover, we present per-class performance in Table 3, where the investigated method consistently outperforms the baselines in each of the fourteen classes.

4.5. Results: Robustness to Textual Input

Overburdened clinicians may introduce or propagate typographical errors while composing a request for imaging examination. We argue it is essential to evaluate models along with the main performance metrics on robustness to changes of the textual input such as common mistakes and use of synonyms. To achieve this goal, we test the MMBT model to textual changes with an evaluation scheme proposed in [24] which we further extended. We mimic a human operator who commits typographical errors and expresses the original medical terms with synonyms. We only select biomedical terms to proceed with the following word/sentence manipulation:

- **Character Swap:** swapping two consecutive characters at

random, e.g. fever \rightarrow fevre.

- **Keyboard Typo:** selecting a random character and replacing with an adjacent one, e.g. fever \rightarrow f3ver.
- **Synonyms:** selecting a synonym for a given biomedical term using the UMLS database, e.g. fever \rightarrow pyrexia.
- **Missing Field/Stop Words Noise:** replacing the indication field with an empty string or a sentence using only stop words.
- **Indication Swap:** selecting a random indication from another patient ensuring no single positive class is shared between two patients.

The results are presented in Table 2. The tested method is resistant to common typographical errors and capable of processing synonyms affecting the performance at most by -1.7% micro AUROC (keyboard typo). When the method does not have access to the corresponding indication fields, the performance of the multimodal transformer is on par with ResNet-50 (micro AUROC). The experiment has shown the method improves while the patient’s history is provided, yet is still capable of processing only images with no textual input, a common scenario in emergency departments. However, replacing the original indication field with a different patient significantly affects the performance (-16.6 % and -9.8% on macro and micro AUROC, respectively). The test has the most notable effect expected on the method (providing clues conflicting with the imaging input), proving that the model uses both modalities to render a decision.

5. CONCLUSION

We evaluated a unimodally pre-trained BERT model on multimodal chest radiograph classification supported by the indication field. We extended the BERT model with an image feature extractor and show it can successfully learn imaging modality, beating the previous state-of-the-art approaches for this task (+4% and +2% micro AUROC for uni- and multimodal baselines, respectively). These promising results show the model can leverage prior knowledge about the patient for a more accurate image diagnosis. We presented the model as resistant to typographical errors, capable of handling synonyms, and missing text input matching image-only baseline.

6. COMPLIANCE WITH ETHICAL STANDARDS

This research study was conducted retrospectively using human subject data made available in open access. Ethical approval was not required as confirmed by the license attached with the open access data.

Acknowledgements

This work was supported by the Engineering and Physical Sciences Research Council [grant number EP/R513209/1] and Canon Medical Research Europe. S.A. Tsaftaris acknowledges the support of the Royal Academy of Engineering and the Research Chairs and Senior Research Fellowships scheme [grant number RCSR1819\8\25].

References

- [1] I. M. Baltruschat *et al.*, “Comparison of Deep Learning Approaches for Multi-Label Chest X-Ray Classification,” *Scientific Reports*, vol. 9, no. 1, pp. 1–10, 2019.
- [2] G. Jacenków *et al.*, “INSIDE: Steering Spatial Attention with Non-Imaging Information in CNNs,” in *MICCAI*, Springer, 2020, pp. 385–395.
- [3] G. Jacenków *et al.*, “Conditioning Convolutional Segmentation Architectures with Non-Imaging Data,” in *MIDL – Extended Abstract Track*, 2019.
- [4] P. Obara *et al.*, “Evaluating the Referring Physician’s Clinical History and Indication as a Means for Communicating Chronic Conditions That Are Pertinent at the Point of Radiologic Interpretation,” *Journal of Digital Imaging*, vol. 28, no. 3, pp. 272–282, 2015.
- [5] A. Leslie *et al.*, “The influence of clinical information on the reporting of CT by radiologists,” *The British Journal of Radiology*, vol. 73, no. 874, pp. 1052–1055, 2000.
- [6] L. H. Li *et al.*, “VisualBERT: A Simple and Performant Baseline for Vision and Language,” *arXiv preprint arXiv:1908.03557*, 2019.
- [7] P. Sharma *et al.*, “Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning,” in *ACL*, 2018, pp. 2556–2565.
- [8] D. Kiela *et al.*, “Supervised Multimodal Bitransformers for Classifying Images and Text,” *arXiv preprint arXiv:1909.02950*, 2019.
- [9] K. He *et al.*, “Deep Residual Learning for Image Recognition,” in *IEEE CVPR*, 2016, pp. 770–778.
- [10] F. Li *et al.*, “Lesion-aware convolutional neural network for chest radiograph classification,” *Clinical Radiology*, vol. 76, no. 2, pp. 155–e1, 2021.
- [11] X. Wang *et al.*, “TieNet: Text-Image Embedding Network for Common Thorax Disease Classification and Reporting in Chest X-rays,” in *IEEE CVPR*, 2018, pp. 9049–9058.
- [12] G. Chauhan *et al.*, “Joint Modeling of Chest Radiographs and Radiology Reports for Pulmonary Edema Assessment,” in *MICCAI*, Springer, 2020, pp. 529–539.
- [13] T. Sylvain *et al.*, “Cross-Modal Information Maximization for Medical Imaging: CMIM,” *arXiv preprint arXiv:2010.10593*, 2020.
- [14] T. van Sonsbeek *et al.*, “Towards Automated Diagnosis with Attentive Multi-modal Learning Using Electronic Health Records and Chest X-Rays,” in *ML-CDS and CLIP (MICCAI)*, Springer, 2020, pp. 106–114.
- [15] J. Tian *et al.*, “Towards Automatic Diagnosis from Multi-modal Medical Data,” in *iMIMIC and ML-CDS (MICCAI)*, Springer, 2019, pp. 67–74.
- [16] Y. Li *et al.*, “A Comparison of Pre-trained Vision-and-Language Models for Multimodal Representation Learning across Medical Images and Reports,” in *IEEE BIBM*, IEEE, 2020, pp. 1999–2004.
- [17] J. Devlin *et al.*, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” in *NAACL-HLT*, Association for Computational Linguistics, 2019, pp. 4171–4186.
- [18] A. E. Johnson *et al.*, “MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports,” *Scientific Data*, vol. 6, 2019.
- [19] A. E. Johnson *et al.*, “MIMIC-CXR Database (Version 2.0.0),” *PhysioNet*, 2019.
- [20] A. L. Goldberger *et al.*, “PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals,” *Circulation*, vol. 101, no. 23, pp. e215–e220, 2000.
- [21] J. Irvin *et al.*, “CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison,” in *AAAI*, 2019.
- [22] Y. Zhang *et al.*, “BioWordVec, improving biomedical word embeddings with subword information and MeSH,” *Scientific Data*, vol. 6, no. 1, pp. 1–9, 2019.
- [23] A. Singh *et al.*, *MMF: A multimodal framework for vision and language research*, <https://github.com/facebookresearch/mmf>, 2020.
- [24] V. Araujo *et al.*, “On Adversarial Examples for Biomedical NLP Tasks,” *arXiv preprint arXiv:2004.11157*, 2020.