

# The Impact of Auxiliary Patient Data on Automated Chest X-Ray Report Generation and How to Incorporate It

Anonymous ACL submission

## Abstract

This study investigates the integration of diverse patient data sources into multimodal language models for automated chest X-ray (CXR) report generation. Traditionally, CXR report generation relies solely on CXR images and limited radiology data, overlooking valuable information from patient health records, particularly from emergency departments. Utilising the MIMIC-CXR and MIMIC-IV-ED datasets, we incorporate detailed patient information such as aperiodic vital signs, medications, and clinical history to enhance diagnostic accuracy. We introduce a novel approach to transform these heterogeneous data sources into embeddings that prompt a multimodal language model, significantly enhancing the diagnostic accuracy of generated radiology reports. Our comprehensive evaluation demonstrates the benefits of using a broader set of patient data, underscoring the potential for enhanced diagnostic capabilities and better patient outcomes through the integration of multimodal data in CXR report generation.

## 1 Introduction

Chest X-ray (CXR) exams, which consist of multiple images captured during an imaging session, are essential for diagnosing and managing a wide range of conditions, playing a significant role in patient care. Radiologists interpret these exams and produce a written report with their findings. However, prompt reporting is hindered by a multitude of issues, including high patient volumes and limited availability of radiologists (Bailey et al., 2022).

Machine learning for automated CXR report generation is a promising solution that has garnered significant attention in the literature (Jones et al., 2021). By leveraging multimodal language models, exams can be rapidly interpreted and reported, potentially providing quick and reliable diagnostic insights crucial for decision-making, such as triaging patients. Models are often trained to generate

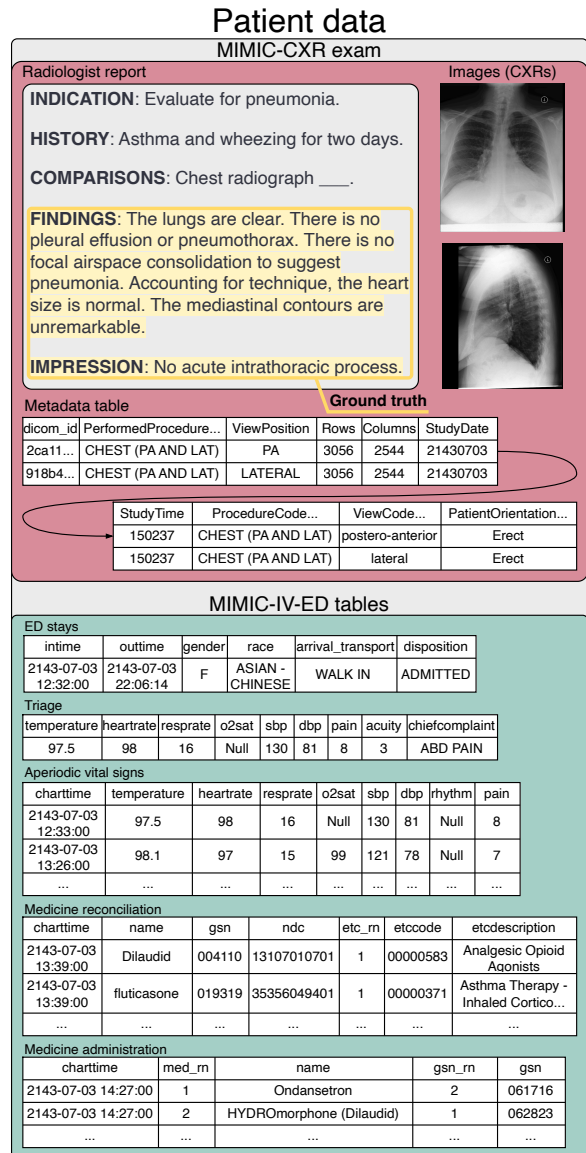


Figure 1: The patient data from MIMIC-IV-ED associated with a CXR exam from MIMIC-CXR. This includes the exam’s images, the corresponding radiology report, and the associated image metadata. The findings and impression sections of the radiology report form the ground truth for CXR report generation. ED-specific data, such as medicine reconciliation and aperiodic vital signs, is also available for the patient.

042 the *findings* and *impression* sections of a radiol- 094  
043 ogy report (Figure 1), where the former details 095  
044 the interpretation of a patient’s exam and the latter 096  
045 summarises the most important findings. Potential 097  
046 benefits include enhanced radiologist effectiveness, 098  
047 a reduced workload, alleviation of the burden of re- 099  
048 port writing, and improved patient outcomes (Shen, 100  
049 2021; Irmici et al., 2023). 101

050 Early methods for CXR report generation pro- 102  
051 duced a separate report for each image within an 103  
052 exam (Wang et al., 2018). Later methods improved 104  
053 on this by considering all images of an exam to gener- 105  
054 ate a single report (Miura et al., 2021; Nicolson 106  
055 et al., 2024a), and incorporating prior exams for 107  
056 a patient (Wu et al., 2022; Nicolson et al., 2024a). 108  
057 Additionally, including the reason for conducting 109  
058 the exam (the *indication* section in Figure 1) of- 110  
059 fered a further improvement (Nguyen et al., 2023). 111  
060 This indicates that CXR report generation could 112  
061 benefit from the inclusion of a more comprehen- 113  
062 sive set of patient data. 114

063 Modern patient record systems are another rich 115  
064 source of patient data, containing detailed informa- 116  
065 tion that may be valuable for CXR report genera- 117  
066 tion. However, (1) the utility of this data has not 118  
067 been empirically investigated, and (2) it is unclear 119  
068 how to harmonise this heterogeneous data into a 120  
069 unified multimodal language model. This paper 121  
070 aims to address these two points. To achieve this, 122  
071 we combine CXR exams from MIMIC-CXR (John- 123  
072 son et al., 2019) with emergency department (ED) 124  
073 patient records from MIMIC-IV-ED (Johnson et al., 125  
074 2023). This means that for a single exam, a wide 126  
075 variety of multimodal data is available, as shown 127  
076 in Figure 1. From MIMIC-CXR, we utilise the 128  
077 images, their metadata, and several sections of the 129  
078 radiology report. Notably, incorporating the com- 130  
079 parison or history section is a novel approach in 131  
080 the literature. From MIMIC-IV-ED, we investigate 132  
081 triage information, aperiodic vital signs, medica- 133  
082 tions, and other data to provide a wider clinical 134  
083 context. 135

084 We explore combining these sources of patient 136  
085 data as patient embeddings to prompt a multimodal 137  
086 language model. We demonstrate that complemen- 138  
087 tary information from different data sources can 139  
088 improve the diagnostic accuracy of CXR report 140  
089 generation. To achieve this, we develop methods 141  
090 to transform tabular and aperiodic time series data 142  
091 into embeddings that can be used alongside token 143  
092 and image embeddings. We evaluate our model on 144  
093 MIMIC-CXR exams with accompanying patient

data from MIMIC-IV-ED, using metrics shown to 094  
closely correlate with radiologists’ assessments of 095  
reporting (Yu et al., 2023). The main contributions 096  
of this work are: 097

- An investigation into how patient data impacts 098  
CXR report generation, focusing on the effects 099  
of specific data sources, such as medications and 100  
vital signs. 101
- An empirical evaluation demonstrating that using 102  
multiple patient data sources — from a patient’s 103  
CXR exams and their ED record — significantly 104  
improves diagnostic accuracy. 105
- Introducing methods to convert multimodal pa- 106  
tient data into embeddings for a language model, 107  
including numerical, categorical, free text, tem- 108  
poral, and image data. 109
- A release of dataset splits based on MIMIC- 110  
CXR and MIMIC-IV-ED, linking patient exams 111  
with their associated ED records (available as 112  
a Hugging Face dataset). This, along with our 113  
code repository and Hugging Face checkpoint 114  
can be found at: [https://anonymous.4open. 115  
science/r/anon-D83E](https://anonymous.4open.science/r/anon-D83E), enabling others to ex- 116  
periment with new methods for multimodal pa- 117  
tient data. 118

## 2 Background and Related Work 119

120 There is evidence to suggest that incorporating 121  
122 more patient data improves diagnostic accuracy 123  
124 in radiology reporting. Initial improvements came 125  
126 from using multiple images per exam, like EMNLI, 127  
128 which often includes complementary frontal and 129  
130 lateral views of the patient (Miura et al., 2021; 131  
132 Gaber et al., 2005). Methods such as CXRMate 133  
134 enhance diagnostic accuracy by incorporating a pa- 135  
136 tient’s prior exams to identify changes over time 137  
138 (Nicolson et al., 2024a; Wu et al., 2022; Kelly, 139  
140 2012). Including the *indication* section of the ra- 141  
142 diology report to provide clinical context also pro- 143  
144 vides an improvement (Nguyen et al., 2023). This 144  
trend indicates that providing more comprehensive 145  
patient data improves diagnostic accuracy, which 146  
we investigate in this work. 147

148 ED records contain a myriad of data, including 149  
150 vital signs such as respiratory rate, temperature, 151  
152 and blood pressure, which can aid in the identifica- 153  
154 tion of various diseases. A high respiratory rate and 155  
156 low blood oxygen saturation are indicative of condi- 157  
158 tions that compromise pulmonary function, such as 159  
160 pulmonary embolism. Similarly, an elevated body 161  
162 temperature is suggestive of an infectious process, 163

144 such as pneumonia or tuberculosis. Incorporating  
145 such data into a CXR report generator could  
146 help corroborate subtle radiographic signs typical  
147 of these infections. Our findings demonstrate that  
148 patient data from the ED can indeed enhance CXR  
149 report generation.

150 Recent advancements in integrating multimodal  
151 patient data have enhanced diagnostic and predic-  
152 tive healthcare capabilities. A study showed that a  
153 Transformer encoder combining imaging and non-  
154 imaging data outperformed single-modality mod-  
155 els, diagnosing up to 25 conditions with higher  
156 AUC scores (Khader et al., 2023b). Similarly, the  
157 MeTra architecture, which integrates CXRs and  
158 clinical parameters, demonstrated superior perfor-  
159 mance in predicting ICU patient survival compared  
160 to using either CXRs or clinical data alone (Khader  
161 et al., 2023a). ETHOS, using a zero-shot learn-  
162 ing approach, outperformed single-modality mod-  
163 els in predicting inpatient mortality, ICU length  
164 of stay, and readmission rates (Renc et al., 2024).  
165 These studies highlight the importance of multi-  
166 modal data for improved healthcare analytics. Our  
167 work demonstrates that incorporating a compre-  
168 hensive set of multimodal patient data enhances CXR  
169 report generation.

170 Recent advancements in multi-task learning have  
171 significantly improved biomedical models by lever-  
172 aging shared knowledge. Med-PaLM M, a gen-  
173 eralist biomedical model, excels in multiple tasks  
174 including classification, question answering, visual  
175 question answering (VQA), report summarisation,  
176 report generation, and genomic variant calling, us-  
177 ing diverse input modalities like images, text, and  
178 genomics. It often outperforms specialised models,  
179 demonstrating superior performance and generalis-  
180 ation (Tu et al., 2024).

181 Similarly, MIMIC-CXR has been leveraged for  
182 multi-task learning with models like MedXChat,  
183 which integrates instruction-tuning and Stable Dif-  
184 fusion to perform CXR report generation, VQA,  
185 and report-to-CXR generation, outperforming other  
186 LLM multi-task learners (Yang et al., 2023). RaDi-  
187 alog, another LLM-based method, combines visual  
188 features and pathology findings to generate accu-  
189 rate radiology reports and support interactive tasks,  
190 significantly improving clinical efficacy. CXR-  
191 LLaVA, a multimodal LLM integrating a vision  
192 transformer with a language model, outperformed  
193 models like GPT-4 Vision and Gemini Pro Vision  
194 in CXR report generation (Lee et al., 2024).

195 Determining the state-of-the-art CXR report gen-

196 eration model can be challenging due to the un-  
197 availability of some models and the lack of com-  
198 parison to recent methods. The 2024 Shared  
199 Task on Large-Scale Radiology Report Generation  
200 (RRG24) aimed to address this by benchmarking  
201 models on a common leaderboard. The winning  
202 model, CXRMate-RRG24 (Nicolson et al., 2024b),  
203 a derivative of CXRMate, emerged as a strong  
204 contender for state-of-the-art. In this work, we  
205 compare our model to established models (e.g.,  
206 EMNLI) and recent benchmarks (e.g., CXRMate-  
207 RRG24, CXRMate, CXR-LLaVA, MedXChat, and  
208 RaDialog). We ensure a fair comparison by us-  
209 ing available code or obtaining generated reports  
210 directly from the authors. Our findings indicate  
211 our model produces significantly better results than  
212 these models.

### 213 3 Dataset

214 We construct a dataset of 46 106 patients by linking  
215 individual patient information from two separate  
216 sources: (1) CXR exams from MIMIC-CXR and  
217 (2) emergency records from MIMIC-IV-ED. Thus  
218 we consider MIMIC-CXR exams that occurred dur-  
219 ing an ED stay from MIMIC-IV-ED. Both datasets  
220 are publicly available and originate from the Beth  
221 Israel Deaconess Medical Center in Boston, MA.

222 MIMIC-CXR was formed by first extracting pa-  
223 tient identifiers for exams performed in the ED  
224 between 2011–2016, and then extracting all exams  
225 for this set of patients from all departments between  
226 2011–2016. Each exam includes a semi-structured  
227 free-text radiology report (Figure 1) that describes  
228 the radiological findings of the images, written by  
229 a practising radiologist contemporaneously during  
230 routine clinical care. All images and reports were  
231 de-identified to protect privacy. Sections from the  
232 radiologist reports were extracted using a modifica-  
233 tion<sup>1</sup> of the official text extraction tool<sup>2</sup> in order to  
234 obtain the findings, impression, indication, history,  
235 and comparison sections.

236 MIMIC-IV-ED consists of de-identified data  
237 from ED stays between 2011–2019. The data was  
238 converted into a denormalised relational database  
239 with six primary tables: ED stays, diagnosis,  
240 medicine reconciliation, medicine administration,  
241 triage, and aperiodic vital signs. We do not con-  
242 sider the diagnosis table in this work, as it indicates  
243 the outcome of a patient’s ED stay. The patients of

<sup>1</sup><https://anonymous.4open.science/r/anon-D83E>

<sup>2</sup><https://github.com/MIT-LCP/mimic-cxr/tree/master/txt>

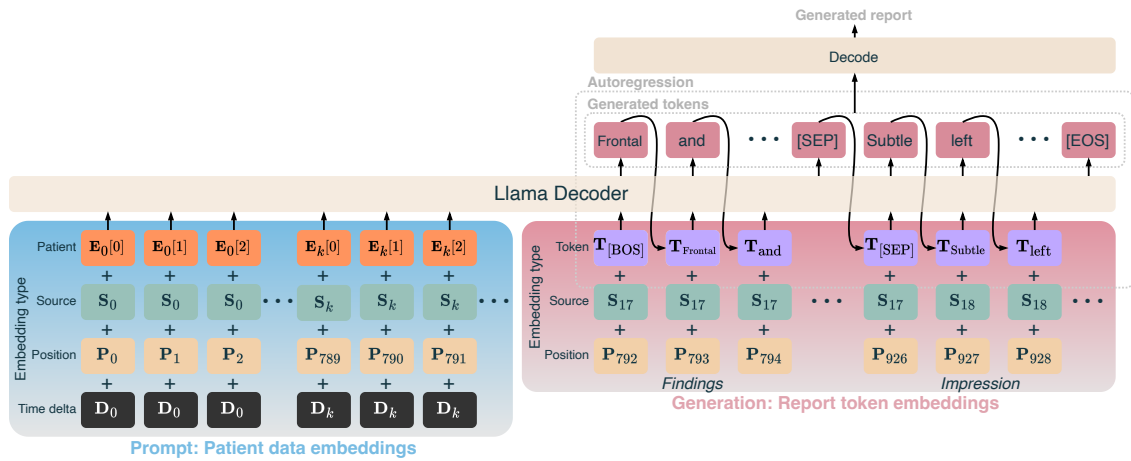


Figure 2: Multimodal language model for CXR report generation. The patient data embeddings prompt the decoder to generate the findings and impression sections of a radiology report.

MIMIC-CXR can be linked to MIMIC-IV-ED via an identifier, allowing an ED specific dataset to be formed.

Example tables for a patient’s exam are shown in Figure 1. The dataset was formed by extracting a patient’s exams whose times (formed by the ‘StudyDate’ and ‘StudyTime’ columns of the metadata table) occurred within the ‘intime’ and ‘outtime’ of one of their ED stays.<sup>3</sup> Exams with either a missing findings or impression section were not considered. Using the official splits of MIMIC-CXR, this gave a train/validation/test split of 45 527/343/236 patients, 76 398/556/958 exams, and 151 818/1 137/1 812 CXRs. Each of these exams had one ED stay and triage row; 53% had at least one medicine reconciliation row with up to 106 rows; 62% had at least one vital signs row with up to 69 rows; and 37% had at least one medication administration row with up to 52 rows. Exams had an indication section 66% of the time with a maximum of 75 words, a history section 34% of the time with a maximum of 74 words, and a comparison section 97% of the time with a maximum of 129 words. Only one exam had both an indication and a history section.

## 4 Methods

The patient data from MIMIC-CXR and MIMIC-IV-ED for an exam are transformed into embeddings, which are used to prompt a multimodal language model to generate the findings and impression sections of the radiology report, as illustrated in Figure 2. Additionally, ‘Source’ embeddings differentiate the source of the data (e.g., the ‘chief complaint’ column from the triage table, the indi-

<sup>3</sup>Exam 59128861 was removed as it overlapped with two separate ED stays for the patient.

cation section, etc.), and time delta embeddings represent the time difference between an event and the exam. Standard embeddings, such as position and token embeddings, are also included. The patient data embeddings originate from three main groups: the tables of MIMIC-IV-ED; the report, images, and metadata of the current exam from MIMIC-CXR; and the patient’s prior exams (also originating from MIMIC-CXR). The prior exam and image embeddings are described in Section A and Subsection C.2, respectively.

### 4.1 Time, Position, & Source Embeddings

The ED information from MIMIC-IV-ED is typically recorded as discrete events, such as medications administered or vital signs measured, each with a specific timestamp. Events that occur closer to the time of the patient’s exam are generally more relevant for diagnostic purposes. To capture this, a time delta is calculated by subtracting the time of an event from the time of the exam. The exam time originates from MIMIC-CXR’s metadata table (Figure 3), whereas most of the MIMIC-IV-ED tables have event times for each row. As shown in Figure 3, the time delta is first converted to hours and then mapped using  $1/\sqrt{\Delta + 1}$ , assigning higher weights to events that occurred closer to the exam. The mapped time deltas are then passed through a feedforward neural network (FNN) defined as  $f(\Delta \mathbf{W}_1) \mathbf{W}_2$ , where  $\mathbf{W}_1 \in \mathbb{R}^{1,2048}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{2048,H}$ ,  $f(\cdot)$  is the sigmoid linear unit (SiLU) activation function (Hendrycks and Gimpel, 2016), and  $H$  is the hidden size of the decoder. This process generates the time delta embeddings, which are subsequently added to the embeddings of their respective sources. As shown in Figure 2, time delta embeddings are only applied

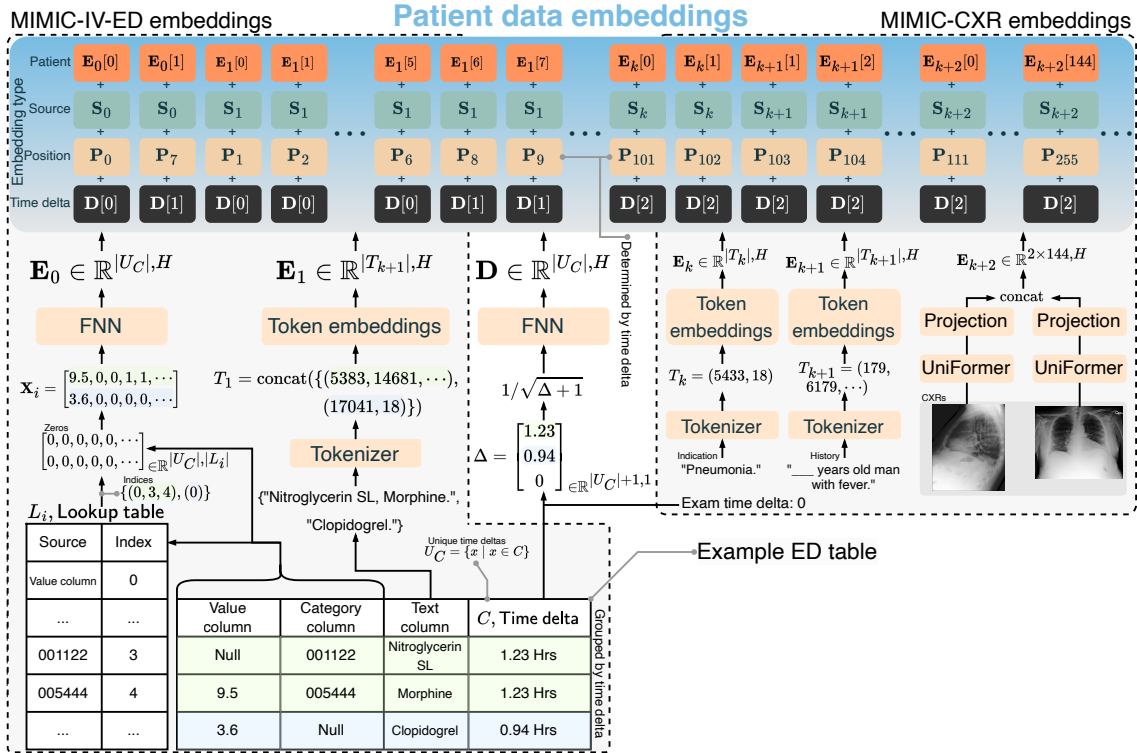


Figure 3: Proposed patient data embeddings from the multiple heterogeneous data types taken from MIMIC-IV-ED and MIMIC-CXR. The embeddings are formed from numerical, categorical, textual, temporal, and image data.

to the prompt. Patient data from the current exam, such as the images, have a time delta of zero, while data from prior exams have a positive time delta.

The position embeddings are ordered by the time delta (Figure 3). This is due to the rotary position embeddings of the decoder; tokens that are closer together are given more importance. Hence, the smaller the time delta, the closer the embedding’s position is to the report token embeddings. Following Nicolson et al. (2024a), each unique patient data source is given its own source embedding. This includes the images, each report section, each table’s text column and value-category columns (described in the next section), and prior images and report sections.

## 4.2 Tabular Data

An example table and its conversion to embeddings is shown in Figure 3. To convert an exam’s tabular data to embeddings, columns were designated as value, category, text, or time columns. Value columns contained numeric data, while category columns contained categorical data. Datum from value and category columns were grouped by their time delta, with each group forming a feature vector. The feature vector initially consisted of zeros. Values and categories from the group were then used to set its values based on indices determined

by a lookup table. For value columns, the lookup table determined the index where the numeric value was placed. For category columns, it determined which indices were activated (set to 1).

Next, the feature vector was passed through an FNN  $f(\mathbf{X}_i \mathbf{W}_1) \mathbf{W}_2$  to form the embedding, where  $\mathbf{X}_i \in \mathbb{R}^{|U_C|, |L_i|}$  are the grouped features,  $\mathbf{W}_1 \in \mathbb{R}^{|L_i|, 2048}$  and  $\mathbf{W}_2 \in \mathbb{R}^{2048, H}$ ,  $L_i$  is a lookup table, and  $i$  designates the table. Each table has a unique FNN and lookup table. Rows for a value column always had a unique time, preventing multiple values from the same column in a group. We investigated alternatives to form the value-category embeddings in Section 5. The described framework was found to be the most efficient. Columns with a high cardinality were set as text columns. Text embeddings were formed via the decoder’s tokenizer and token embeddings. Text embeddings were given the time delta embedding from their respective row. The column designation for each table in Figure 1 is described in the Appendix B.

## 4.3 Report Section Embeddings

Here, we consider five sections of the radiology report: the findings, impression, indication, history, and comparison sections. The findings and impression sections serve as the ground truth to be

generated. The remainder form part of the patient data embeddings. The indication section explains the reason for the exam, such as symptoms or suspected conditions. The history section provides relevant medical history, such as past conditions and treatments. The comparison section mentions any prior exams used to identify changes over time. These sections provide context that guides the interpretation of the exam, influencing the content of the findings and impression sections. The embeddings were formed via the decoder’s tokenizer and token embeddings. Of these, the history and comparison sections have not been investigated for CXR report generation. The comparison section was used only when prior exams were considered.

#### 4.4 Experiment Setup

Our multimodal language model, illustrated in Figure 2, is based on CXRMate-RRG24; it features a Llama decoder and the UniFormer as the image encoder. The training procedure for our model involved three stages: (1) initial training on the MIMIC-CXR training set using only images as input with Teacher Forcing (TF) (Williams and Zipser, 1989), (2) further training on the dataset described in Section 1 with the inputs detailed in Table 1, again using TF, and (3) reinforcement learning on the same dataset through self-critical sequence training (SCST) (Rennie et al., 2017) (only for Table 2). Our evaluation metrics included three that capture the semantics of radiology reporting — RadGraph-F1 (RG), CheXbert-F1 (CX), and CXR-BERT (CB) — as well as five natural language generation metrics: BERTScore-F1 (BS), CIDEr (C), METEOR (M), ROUGE-L (R-L), and BLEU-4 (B4). Comprehensive details on the model architecture, training procedure, significance testing, and comparison methods are provided in Appendix C.

### 5 Results & Discussion

The impact of different patient data sources on the performance of CXR report generation is summarised in Table 1. This analysis identifies which additional data sources enhance performance compared to using only images.

Significant improvements were observed by incorporating either the ED stays, triage, medicine reconciliation, or vital signs data from MIMIC-IV-ED dataset. Notably, the ED data markedly improved scores on the radiology report metrics (RG, CX, and CB). The medicine administration

table did not significantly improve the scores overall, likely due to its infrequent occurrence in the exams (37%). (However, as shown in Table 4, it significantly improves performance when evaluated solely on exams that include a medicine administration table.) These findings demonstrate that ED patient data can enhance the diagnostic accuracy of CXR report generation.

Incorporating the indication or history section led to significant score improvements. This demonstrates the substantial influence these sections have on the findings and impression sections. Conversely, adding the metadata table did not result in significant score improvements, indicating it lacks valuable information for CXR report generation. While previous studies have established that the indication section boosts CXR report generation (Nguyen et al., 2023), our findings demonstrate that the history section is equally important.

When examining the impact of prior exams, we considered a maximum history size  $h$  of up to three, incorporating the findings and impression sections, and images from prior exams. Any history size significantly increases the scores compared to using solely the images, consistent with previous findings (Wu et al., 2022). However, performance gradually degrades as the history size increases, which contradicts earlier studies. Additionally, the comparison section appears to slightly degrade performance. We suspect this is due to the increasing number of inputs as  $h$  grows, combined with the limitations of our model architecture.  $|\overline{\mathcal{E}[:, 0]}|$  in Table 1 is the average prompt length over the test set, where  $\mathcal{E} = [\mathbf{E}_0, \mathbf{E}_1, \dots]$ . It can be seen that  $|\overline{\mathcal{E}[:, 0]}|$  increases substantially as  $h$  increases. Since we provide all inputs to the decoder’s self-attention, a large input size may cause *attention dilution*. With more inputs, the attention weights must be distributed across a larger number of inputs, resulting in each input receiving a smaller share of the attention, making it harder for the model to focus on the most relevant inputs (Qin et al., 2022).

We then combined all the effective sources of patient data (those providing a significant improvement). This excluded ‘medicine administration’, ‘metadata’, and ‘comparison’. The best performance was observed with no prior exams ( $h = 0$ ), indicating that using any prior exams in combination with other sources is detrimental due to attention dilution. With  $h = 0$ , the combination of all effective sources outperformed each individual source. We then performed an ablation study

Table 1: Results of the various patient data sources on the test set described in Section 3. Results were calculated over ten training runs ( $n = 9580$  exams;  $958 \times 10$  runs). Underlined and Dashed underlined scores indicate a significant difference to the scores of ‘Images’ and ‘Images + effective sources ( $h = 0$ )’, respectively ( $p < 0.05$ ). Evaluation is performed on both the **findings** and **impression** sections.

| Patient data sources   | RG           | CX           | CB           | BS           | C            | M            | R-L          | B4          | $ \mathcal{E}[:, 0] $ |
|--|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|-----------------------|
| <i>Images only</i>   |              |              |              |              |              |              |              |             |                       |
| Images   | 26.00        | 29.24        | 58.87        | 24.10        | 12.24        | 14.35        | 24.34        | 6.33        | 272.4                 |
| <i>Patient Emergency Department (ED) data (MIMIC-IV-ED)</i>                        |              |              |              |              |              |              |              |             |                       |
| Images + ED stays  | 26.10        | 29.47        | <u>60.65</u> | 24.17        | 12.39        | 14.52        | 24.50        | 6.36        | 273.4                 |
| Images + triage  | <u>26.46</u> | 31.27        | <u>63.06</u> | 24.29        | 12.32        | <u>14.66</u> | 24.58        | 6.44        | 278.9                 |
| Images + vital signs   | <u>26.47</u> | 31.72        | <u>63.39</u> | 24.32        | 13.16        | <u>14.61</u> | <u>24.74</u> | 6.47        | 274.7                 |
| Images + medicine reconciliation   | <u>26.86</u> | 31.37        | <u>63.98</u> | 24.52        | 12.77        | <u>14.90</u> | <u>24.85</u> | 6.60        | 343.5                 |
| Images + medicine administration   | 26.15        | 29.47        | 59.21        | 24.25        | 12.30        | 14.44        | 24.47        | 6.38        | 273.0                 |
| <i>Patient additional radiology data (MIMIC-CXR)</i>                               |              |              |              |              |              |              |              |             |                       |
| Images + indication  | <u>26.94</u> | <u>32.13</u> | <u>65.43</u> | <u>24.74</u> | <u>14.16</u> | <u>15.19</u> | <u>25.16</u> | <u>7.02</u> | 279.5                 |
| Images + history   | <u>27.00</u> | 31.88        | <u>65.06</u> | <u>25.05</u> | <u>14.32</u> | <u>15.30</u> | <u>25.48</u> | <u>7.33</u> | 277.0                 |
| Images + metadata  | 26.34        | 29.63        | 59.55        | 24.37        | 12.40        | 14.55        | 24.50        | 6.43        | 273.4                 |
| <i>Prior exams</i>   |              |              |              |              |              |              |              |             |                       |
| Images + $h = 1$   | <u>26.98</u> | 31.42        | <u>63.98</u> | <u>24.65</u> | 12.65        | <u>15.11</u> | <u>25.03</u> | 6.78        | 558.9                 |
| Images + $h = 1$ + comparison  | 26.76        | 31.55        | <u>64.20</u> | 24.42        | 13.36        | <u>15.03</u> | <u>24.82</u> | 6.74        | 563.4                 |
| Images + $h = 2$   | <u>26.67</u> | 30.48        | <u>61.27</u> | 24.53        | 13.60        | <u>14.94</u> | <u>24.85</u> | <u>6.72</u> | 810.6                 |
| Images + $h = 2$ + comparison  | 26.20        | 30.19        | <u>61.24</u> | 24.05        | 12.43        | <u>14.80</u> | 24.55        | 6.58        | 815.0                 |
| Images + $h = 3$   | 26.47        | 29.96        | 59.95        | 24.14        | 12.90        | <u>14.94</u> | 24.66        | 6.65        | 1037.1                |
| Images + $h = 3$ + comparison  | 26.14        | 30.09        | <u>60.51</u> | 23.90        | 13.22        | <u>14.87</u> | 24.56        | 6.64        | <b>1041.5</b>         |
| <i>All effective sources (no medicine administration, metadata, or comparison)</i> |              |              |              |              |              |              |              |             |                       |
| Images + effective sources ( $h = 0$ )   | <u>27.11</u> | <u>32.23</u> | <u>64.80</u> | <u>25.07</u> | <u>14.48</u> | <u>15.15</u> | <u>25.40</u> | <u>7.07</u> | 365.0                 |
| Images + effective sources ( $h = 1$ )   | <u>26.78</u> | 31.83        | <u>63.85</u> | <u>24.75</u> | <u>14.10</u> | <u>15.15</u> | <u>25.25</u> | <u>7.01</u> | 651.7                 |
| <i>Ablation from Images + effective sources (<math>h = 0</math>)</i>               |              |              |              |              |              |              |              |             |                       |
| - medicine reconciliation  | 26.78        | <b>32.81</b> | <b>65.60</b> | 24.84        | 14.44        | 15.21        | 25.33        | 7.19        | 293.9                 |
| - ED stays   | 26.94        | 31.56        | <u>64.87</u> | 25.02        | 14.08        | 15.14        | 25.37        | 7.09        | 364.0                 |
| - triage   | 27.15        | <u>32.45</u> | 65.18        | 25.15        | <b>14.80</b> | 15.27        | <b>25.54</b> | <u>7.25</u> | 358.5                 |
| - vital signs  | <b>27.27</b> | 31.78        | 65.44        | 25.14        | 14.07        | <b>15.35</b> | 25.49        | <u>7.22</u> | 362.6                 |
| - indication   | 26.89        | 31.25        | 64.65        | 24.99        | 13.87        | 15.07        | 25.39        | 7.00        | 357.9                 |
| - history  | 26.96        | 31.87        | 64.02        | 24.86        | 14.60        | 15.10        | 25.24        | 7.04        | 360.3                 |
| - time delta   | 27.17        | 32.11        | 65.10        | <b>25.18</b> | 14.64        | 15.24        | <b>25.54</b> | 7.16        | 365.0                 |

using ‘CXRs + effective sources ( $h = 0$ )’. Removing ‘medicine reconciliation’ significantly increased performance, specifically for CXR-BERT. This improvement was also likely due to attention dilution, as removing medicine reconciliation substantially decreased  $|\mathcal{E}[:, 0]|$ .

Next, we further trained ‘Images + effective sources ( $h = 0$ ) - medicine reconciliation’ with reinforcement learning, as described in Subsection 4.4. This model, denoted as ‘Ours’ in Table 2, was compared to other benchmark CXR report generation models in the literature that included MIMIC-CXR in their training data. Despite having substantially fewer training samples than the other models, our model significantly outperformed them on CXR-BERT, BERTScore-F1, METEOR, ROUGE-L, and BLEU-4. This demonstrates the impact of incorporating a more comprehensive set of patient data on CXR report generation.


A case study is presented in Figure 4 demonstrating how a diverse set of patient data can impact report generation. Here, the first model is given the image only, and fails to identify key findings that

the radiologist noted in their report. The second model is given the additional patient data available for this exam; the indication section and triage data. Hypoxia, as indicated by the low oxygen saturation (‘o2sat’), along with the elevated respiratory rate (‘resprate’) and systolic blood pressure (‘SBP’), are consistent with the physiological responses to pulmonary edema. Given this, the second model was able to identify the moderate pulmonary edema, echoing the radiologist’s findings.

Table 3 compares different methods for converting value and category columns into embeddings. This evaluation includes images, the triage table, and the medicine reconciliation table, as these tables contain multiple value and category columns. The aforementioned method of producing embeddings by grouping data from value and category columns (‘Grouped embeddings’) is compared to two other methods. The first is separate embeddings for each datum, where each value column datum is separately transformed using the previously described FNN, while each category column datum is converted to an embedding using a learn-

Table 2: Benchmark models on the test set described in Section 3 ( $n = 958$ ). Evaluation is on the **findings** section only. Underlined indicates statistical significance between the top two scores ( $p < 0.05$ ). In the ‘Train samples’ column, ‘images’ means the model generates reports per image, while ‘exams’ means a report generated per exam.

| Model                                  | Train samples  | RG          | CX          | CB          | BS          | C           | M           | R-L         | B4          |
|--|----------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| EMNLI (Miura et al., 2021)             | 152 173 exams  | <u>32.8</u> | 28.9        | 66.6        | 24.4        | 19.4        | 17.1        | 28.1        | 8.9         |
| CMN (Chen et al., 2021)                | 270 790 images | 25.3        | 24.3        | 49.4        | 19.7        | 16.9        | 15.1        | 26.4        | 7.6         |
| TranSQ (Kong et al., 2022)             | 368 960 images | 29.8        | 30.4        | 62.3        | 20.4        | 14.9        | 17.6        | 22.6        | 7.9         |
| RGRG (Tanida et al., 2023)             | 166 512 images | 23.2        | 22.8        | 37.9        | 23.4        | 7.6         | 12.4        | 21.1        | 5.4         |
| CvT2DistilGPT2 (Nicolson et al., 2023) | 270 790 images | 25.8        | 29.3        | 59.8        | 24.8        | 20.9        | 16.0        | 27.3        | 8.8         |
| RaDialog (Pellegrini et al., 2023)     | 276 778 images | 26.8        | <u>38.4</u> | 60.7        | 26.2        | 14.6        | 14.7        | 25.4        | 6.9         |
| MedXChat (Yang et al., 2023)           | 270 790 images | 22.6        | 13.1        | 21.3        | 19.3        | 9.8         | 14.3        | 23.2        | 7.0         |
| CXR-LLaVA-v2 (Lee et al., 2024)        | 193 513 images | 20.7        | 20.7        | 44.1        | 23.6        | 5.2         | 11.3        | 19.9        | 2.7         |
| CXRMate (Nicolson et al., 2024a)       | 125 395 exams  | 28.8        | 33.9        | 71.3        | 30.5        | 22.4        | 17.7        | 28.1        | 9.7         |
| CXRMate-RRG24 (Nicolson et al., 2024b) | 550 395 exams  | 30.4        | 31.2        | 58.2        | 31.0        | 20.6        | 16.7        | 27.5        | 9.1         |
| <b>Ours</b>                            | 76,398 exams   | <u>33.7</u> | 35.1        | <u>79.1</u> | <u>35.8</u> | <u>24.1</u> | <u>19.1</u> | <u>30.6</u> | <u>11.9</u> |



**Indication:** Hypoxia.

**Case study**

**Radiologist findings:** A portable frontal chest radiograph demonstrate an unchanged cardiomeastinal silhouette, which is top-normal in size. Bilateral opacities are consistent with moderate pulmonary edema. No definite focal consolidation or pneumothorax is identified. There are likely trace bilateral pleural effusions.

**Radiologist impression:** Moderate pulmonary edema.

**Triage:**

| temperature | heartrate | resprate | o2sat | sbp | dbp | pain | acuity | chiefcomplaint |
|-------------|-----------|----------|-------|-----|-----|------|--------|----------------|
| 100.3       | 93        | 24       | 83    | 175 | 74  | Null | 1      | ILI, Fever     |

---

**Image** (Model: Images from Table 1)  
**Generated findings:** Cardiomeastinal silhouette is normal. There is no focal consolidation. There is no pneumothorax or pleural effusion. There is no significant pleural effusion.  
**Generated impression:** No acute cardiopulmonary process.

---

**Image + Indication + Triage** (Model: Images + effective sources ( $h=0$ ) - medicine reconciliation from Table 1)  
**Generated findings:** There is moderate pulmonary edema. No definite focal consolidation is identified. There are probable small bilateral pleural effusions. The cardiac silhouette is mildly enlarged. There is no pneumothorax.  
**Generated impression:** Moderate pulmonary edema and small bilateral pleural effusions.

Figure 4: Case study demonstrating how incorporating a diverse set of patient data can aid with report generation.

able weight matrix, akin to how token embeddings are produced (‘Separate embeddings’). The second method modifies ‘Separate embeddings’ by instead converting the value column data to text and using the decoder’s tokenizer and token embeddings (‘Values-to-text, categories-to-tokens’). The results indicate that the grouped embeddings method generally works best and is useful for encoding heterogeneous patient data for multimodal models.

## 6 Conclusion

This paper demonstrates the value of incorporating diverse patient data into automated CXR report generation. By integrating patient data from the MIMIC-CXR and MIMIC-IV-ED datasets, we have shown significant improvements in the diagnostic accuracy of generated radiology reports. Our empirical evaluation uncovers new sources of patient information that enhance CXR report generation, including data from ED stays, triaging information, aperiodic vital signs, medications, and the history section of radiology reports. We present

Table 3: Formatting strategies for the value-category columns. Four training runs were used ( $n = 3832$ ; exams  $958 \times 4$  runs). Underlined indicates a stat. sig. difference to ‘Baseline’ ( $p < 0.05$ ).

| Embeddings                                       | CX           | RG           | CB           | BS           |
|--|--------------|--------------|--------------|--------------|
| <i>Images</i>                                    |              |              |              |              |
| Baseline   | 25.81        | 29.00        | 59.04        | 23.85        |
| <i>Images + triage + medicine reconciliation</i> |              |              |              |              |
| Grouped embeddings                               | <u>26.72</u> | <u>31.69</u> | <u>64.01</u> | 24.38        |
| Separate embeddings                              | 25.32        | 25.28        | 46.29        | 23.51        |
| Values-to-text, categories-to-embeddings         | <u>26.46</u> | 30.70        | 58.62        | <u>24.58</u> |

specific methods to convert multimodal patient data into embeddings for a language model, encompassing numerical, categorical, textual, temporal, and image data. We encourage further research and experimentation using our released dataset splits, code, and model checkpoints to explore innovative methods for multimodal patient data integration, with the ultimate goal of enhancing diagnostic accuracy and patient care.



## 7 Limitations

Despite the promising results demonstrated in this study, several limitations must be acknowledged. Firstly, the generalisability of our findings may be constrained by the datasets utilised, specifically MIMIC-CXR and MIMIC-IV-ED, which are derived from a single institution, the Beth Israel Deaconess Medical Center. This could introduce biases unique to the demographic and clinical practices of this institution, potentially limiting the applicability of our model to other healthcare settings with different patient populations or clinical workflows. Our reliance on these datasets is due to the fact that they are the only publicly available sources that link CXR exams with ED records.

Another limitation pertains to the completeness and quality of the patient data. Despite incorporating a wide range of data sources, the datasets still contain missing or incomplete information, which can affect model performance. For example, not all exams include a history section, and not all ED patient records have medicine administration details, leading to potential gaps in the data that the model can utilise. However, this reflects the nature of real patient records where issues of data quality and completeness are to be expected.

Our model’s architecture, while effective, has certain limitations. It struggles with large input sizes, especially when incorporating multiple prior exams, likely due to attention dilution. Future work should explore advanced attention mechanisms or hierarchical models to better manage large input sequences.

The interpretability of the model also poses a challenge. While our model shows improved diagnostic accuracy, the decision-making process within the multimodal language model remains a black box. Developing methods to enhance the interpretability and explainability of the model’s outputs would be beneficial, especially in clinical settings where understanding the rationale behind a diagnosis is critical.

Finally, while we provide a comprehensive set of metrics to evaluate our model’s performance, these metrics focus primarily on the diagnostic accuracy and quality of the generated reports. Broader evaluations considering clinical outcomes, such as the impact on patient management or reduction in radiologist workload, would offer a more holistic view of the benefits and limitations of CXR report generation models in general. Conducting such

assessments could help to better understand the practical implications of deploying these models in a clinical setting.

In summary, while our study provides valuable insights into the integration of multimodal patient data for CXR report generation, addressing these limitations will be crucial for further advancements and broader adoption of such models in clinical practice. Future research should explore alternative architectures and training strategies, find alternative datasets to evaluate generalisability, improve model interpretability, and comprehensively assess the practical impact on patient care and radiologist workflow.

## 8 Ethical Considerations

In this research, we used real-world patient data from the MIMIC-CXR and MIMIC-IV-ED datasets. Since these datasets are de-identified, we consider privacy leakage risks to be minimal. Our method employs a language model to generate medical reports from patient data. However, we acknowledge that language models can exhibit bias and produce hallucinations, which may result in incorrect content in the generated reports.

## References

- Christopher R. Bailey, Allison M. Bailey, Anna Sophia McKenney, and Clifford R. Weiss. 2022. Understanding and Appreciating Burnout in Radiologists. *RadioGraphics*, 42(5):E137–E139.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *ACL Workshops*, pages 65–72.
- Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, Hoifung Poon, and Ozan Oktay. 2022. Making the Most of Text Semantics to Improve Biomedical Vision–Language Processing. In *ECCV*, pages 1–21.
- Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. 2021. Cross-modal Memory Networks for Radiology Report Generation. In *IJCNLP*, pages 5904–5914.
- Jean-Benoit Delbrouck, Pierre Chambon, Christian Bluethgen, Emily Tsai, Omar Almusa, and Curtis Langlotz. 2022. Improving the Factual Correctness of Radiology Report Generation with Semantic Rewards. In *EMNLP*, pages 4348–4360.
- Mohamed Elgendi, Muhammad Umer Nasir, Qunfeng Tang, David Smith, John-Paul Grenier, Catherine

|     |  |   |  |
|-----|--|---|--|
| 647 | Batte, Bradley Spieler, William Donald Leslie, Carlo Menon, Richard Ribbon Fletcher, Newton Howard, Rabab Ward, William Parker, and Savvas Nicolaou. 2021. The Effectiveness of Image Augmentation in Deep Learning Networks for Detecting COVID-19: A Geometric Transformation Perspective. <i>Frontiers in Medicine</i> , 8.   | Seowoo Lee, Jiwon Youn, Hyungjin Kim, Mansu Kim, and Soon Ho Yoon. 2024. CXR-LLAVA: a multi-modal large language model for interpreting chest X-ray images. ArXiv:2310.18341 [cs].  | 703<br>704<br>705<br>706               |
| 654 | Khalid A Gaber, Clive R McGavin, and Irving P Wells. 2005. Lateral Chest X-Ray for Physicians. <i>Journal of the Royal Society of Medicine</i> , 98(7):310–312.  | Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. 2023. UniFormer: Unifying Convolution and Self-Attention for Visual Recognition. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(10):12581–12600.                     | 707<br>708<br>709<br>710<br>711<br>712 |
| 657 | Dan Hendrycks and Kevin Gimpel. 2016. Gaussian Error Linear Units (GELUs). <i>arXiv:1606.08415 [cs.LG]</i> .   | Chin-Yew Lin and Eduard Hovy. 2003. Automatic evaluation of summaries using N-gram co-occurrence statistics. In <i>NAACL</i> , volume 1, pages 71–78.   | 713<br>714<br>715                      |
| 660 | Giovanni Irmici, Maurizio Cè, Elena Caloro, Natalia Khenkina, Gianmarco Della Pepa, Velio Ascenti, Carlo Martinenghi, Sergio Papa, Giancarlo Oliva, and Michaela Cellina. 2023. Chest X-ray in Emergency Radiology: What Artificial Intelligence Applications Are Available? <i>Diagnostics</i> , 13(2):216.   | Ilya Loshchilov and Frank Hutter. 2022. Decoupled Weight Decay Regularization. In <i>ICLR</i> .   | 716<br>717                             |
| 666 | Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Leo Anthony Celi, Roger Mark, and Steven Horng. 2023. MIMIC-IV-ED (version 2.2). PhysioNet.   | Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. 2021. Improving Factual Completeness and Consistency of Image-to-Text Radiology Report Generation. In <i>NAACL</i> , pages 5288–5304.   | 718<br>719<br>720<br>721<br>722        |
| 669 | Alistair E. W. Johnson, Tom J. Pollard, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G. Mark, Seth J. Berkowitz, and Steven Horng. 2019. MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. PhysioNet.  | Dang Nguyen, Chacha Chen, He He, and Chenhao Tan. 2023. <b>Pragmatic Radiology Report Generation</b> . In <i>Proceedings of the 3rd Machine Learning for Health Symposium</i> , pages 385–402. PMLR. ISSN: 2640-3498.   | 723<br>724<br>725<br>726<br>727        |
| 675 | Catherine M Jones, Quinlan D Buchlak, Luke Oakden-Rayner, Michael Milne, Jarrel Seah, Nazanin Esmaili, and Ben Hachey. 2021. Chest radiographs and machine learning – Past, present and future. <i>Journal of Medical Imaging and Radiation Oncology</i> , 65(5):538–544.  | Aaron Nicolson, Jason Dowling, and Bevan Koopman. 2023. Improving chest X-ray report generation by leveraging warm starting. <i>Artificial Intelligence in Medicine</i> , 144:102633.   | 728<br>729<br>730<br>731               |
| 681 | Barry Kelly. 2012. The chest radiograph. <i>The Ulster Medical Journal</i> , 81(23620614):143–148.   | Aaron Nicolson, Jason Dowling, and Bevan Koopman. 2024a. Longitudinal Data and a Semantic Similarity Reward for Chest X-Ray Report Generation. arXiv:2307.09758 [cs].   | 732<br>733<br>734<br>735               |
| 683 | Firas Khader, Jakob Nikolas Kather, Gustav Müller-Franzes, Tianci Wang, Tianyu Han, Soroosh Tayebi Arasteh, Karim Hamesch, Keno Bresslem, Christoph Haarburger, Johannes Stegmaier, Christiane Kuhl, Sven Nebelung, and Daniel Truhn. 2023a. Medical transformer for multimodal survival prediction in intensive care: integration of imaging and non-imaging data. <i>Scientific Reports</i> , 13(1):10666. | Aaron Nicolson, Jinghui Liu, Jason Dowling, Anthony Nguyen, and Bevan Koopman. 2024b. e-Health CSIRO at RRG24: Entropy-Augmented Self-Critical Sequence Training for Radiology Report Generation. In <i>The 23rd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks</i> . | 736<br>737<br>738<br>739<br>740<br>741 |
| 691 | Firas Khader, Gustav Müller-Franzes, Tianci Wang, Tianyu Han, Soroosh Tayebi Arasteh, Christoph Haarburger, Johannes Stegmaier, Keno Bresslem, Christiane Kuhl, Sven Nebelung, Jakob Nikolas Kather, and Daniel Truhn. 2023b. Multimodal Deep Learning for Integrating Chest Radiographs and Clinical Parameters: A Case for Transformers. <i>Radiology</i> , 309(1):e230806.                                | Kishore Papineni, Salim Roukos, Todd Ward, and Weijing Zhu. 2001. BLEU: a method for automatic evaluation of machine translation. In <i>ACL</i> , page 311.   | 742<br>743<br>744                      |
| 699 | Ming Kong, Zhengxing Huang, Kun Kuang, Qiang Zhu, and Fei Wu. 2022. TransSQ: Transformer-Based Semantic Query for Medical Report Generation. In <i>MICCAI</i> , volume 13438, pages 610–620.   | Chantal Pellegrini, Ege Özsoy, Benjamin Busam, Nasir Navab, and Matthias Keicher. 2023. RaDialog: A Large Vision-Language Model for Radiology Report Generation and Conversational Assistance. ArXiv:2311.18681 [cs].   | 745<br>746<br>747<br>748<br>749        |
| 702 |  | Zhen Qin, Xiaodong Han, Weixuan Sun, Dongxu Li, Lingpeng Kong, Nick Barnes, and Yiran Zhong. 2022. The Devil in Linear Transformer. In <i>EMNLP</i> , pages 7025–7041.  | 750<br>751<br>752<br>753               |
|     |  | Pawel Renc, Yugang Jia, Anthony E. Samir, Jaroslaw Was, Quanzheng Li, David W. Bates, and Arkadiusz Sitek. 2024. A Transformer-Based Model for Zero-Shot Health Trajectory Prediction. MedRxiv.   | 754<br>755<br>756<br>757               |

|     |  |  |  |
|-----|--|--|--|
| 758 | Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-Critical Sequence Training for Image Captioning. In <i>CVPR</i> , pages 1179–1195.   |  |  |
| 759 |  |  |  |
| 760 |  |  |  |
| 761 |  |  |  |
| 762 | Dinggang Shen. 2021. Grand Challenges in Radiology. <i>Frontiers in Radiology</i> , 1.   |  |  |
| 763 |  |  |  |
| 764 | Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Ng, and Matthew Lungren. 2020. Combining Automatic Labelers and Expert Annotations for Accurate Radiology Report Labeling Using BERT. In <i>EMNLP</i> , pages 1500–1519.   |  |  |
| 765 |  |  |  |
| 766 |  |  |  |
| 767 |  |  |  |
| 768 |  |  |  |
| 769 | Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. 2023. Interactive and Explainable Region-guided Radiology Report Generation. In <i>CVPR</i> , pages 7433–7442.   |  |  |
| 770 |  |  |  |
| 771 |  |  |  |
| 772 |  |  |  |
| 773 | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. arXiv:2302.13971 [cs].   |  |  |
| 774 |  |  |  |
| 775 |  |  |  |
| 776 |  |  |  |
| 777 |  |  |  |
| 778 |  |  |  |
| 779 |  |  |  |
| 780 | Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaekermann, Mohamed Amin, Pi-Chuan Chang, Andrew Carroll, Charles Lau, Ryutaro Tanno, Ira Ktena, Anil Palepu, Basil Mustafa, Aakanksha Chowdhery, Yun Liu, Simon Kornblith, David Fleet, Philip Mansfield, Sushant Prakash, Renee Wong, Sunny Virmani, Christopher Semturs, S. Sara Mahdavi, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Karan Singhal, Pete Florence, Alan Karthikesalingam, and Vivek Natarajan. 2024. Towards Generalist Biomedical AI. <i>NEJM AI</i> , 1(3):A10a2300138.  |  |  |
| 781 |  |  |  |
| 782 |  |  |  |
| 783 |  |  |  |
| 784 |  |  |  |
| 785 |  |  |  |
| 786 |  |  |  |
| 787 |  |  |  |
| 788 |  |  |  |
| 789 |  |  |  |
| 790 |  |  |  |
| 791 |  |  |  |
| 792 | Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. CIDEr: Consensus-Based Image Description Evaluation. In <i>CVPR</i> , pages 4566–4575.   |  |  |
| 793 |  |  |  |
| 794 |  |  |  |
| 795 | Changhan Wang, Kyunghyun Cho, and Jiatao Gu. 2020. Neural machine translation with byte-level subwords. In <i>AAAI</i> , pages 9154–9160.  |  |  |
| 796 |  |  |  |
| 797 |  |  |  |
| 798 | Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, and Ronald M. Summers. 2018. TieNet: Text-Image Embedding Network for Common Thorax Disease Classification and Reporting in Chest X-Rays. In <i>CVPR</i> .   |  |  |
| 799 |  |  |  |
| 800 |  |  |  |
| 801 |  |  |  |
| 802 |  |  |  |
| 803 | Ronald J. Williams and David Zipser. 1989. A Learning Algorithm for Continually Running Fully Recurrent Neural Networks. <i>Neural Computation</i> , 1(2):270–280.   |  |  |
| 804 |  |  |  |
| 805 |  |  |  |
| 806 |  |  |  |
| 807 | Xian Wu, Shuxin Yang, Zhaopeng Qiu, Shen Ge, Yangtian Yan, Xingwang Wu, Yefeng Zheng, S. Kevin Zhou, and Li Xiao. 2022. DeltaNet: Conditional Medical Report Generation for COVID-19 Diagnosis. In <i>ICCL</i> , pages 2952–2961.  |  |  |
| 808 |  |  |  |
| 809 |  |  |  |
| 810 |  |  |  |
| 811 |  |  |  |
|     | Ling Yang, Zhanyu Wang, and Luping Zhou. 2023. MedXChat: Bridging CXR Modalities with a Unified Multimodal Large Model. ArXiv:2312.02233 [cs].   |  | 812<br>813<br>814  |
|     | Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser Ururahy Nunes Fonseca, Henrique Min Ho Lee, Zahra Shakeri Hossein Abad, Andrew Y. Ng, Curtis P. Langlotz, Vasantha Kumar Venugopal, and Pranav Rajpurkar. 2023. Evaluating progress in automatic chest X-ray radiology report generation. <i>Patterns</i> , page 100802.  |  | 815<br>816<br>817<br>818<br>819<br>820<br>821<br>822   |
|     | Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. In <i>ICLR</i> .  |  | 823<br>824<br>825  |
|     | <b>A Prior exam embeddings</b>   |  | 826  |
|     | The images, findings section, and impression section from previous exams were considered. For prior exams, the time delta was positive, calculated by subtracting the time of the prior exam from the current exam. The images, findings section, and impression section from prior exams were given distinct source embeddings, separate from the current exam, to enhance differentiation. The comparison section from the current exam was also investigated, anticipating that references to prior exams in this section would prompt the decoder to reflect this in the generated report. We explored prior exams with a history size $h$ of up to three.   |  | 827<br>828<br>829<br>830<br>831<br>832<br>833<br>834<br>835<br>836<br>837<br>838<br>839                                    |
|     | <b>B Table column determination</b>  |  | 840  |
|     | The columns from the tables described in Figure 1 were given the following designations:   |  | 841<br>842   |
|     | <ul style="list-style-type: none"> <li>• For the ED stay table, the patients ‘intime’ was used as the event time. Gender (e.g., ‘F’), race (e.g., ‘HISPANIC OR LATINO’), and arrival transport (e.g., ‘AMBULANCE’) were designated as category columns. The disposition column was not considered.</li> <li>• For the triage table, the ‘intime’ from the ED stay table was used. Temperature (e.g., ‘100.6’), heart rate (e.g., ‘93’), respiratory rate (e.g., ‘16’), O2 saturation (e.g., ‘94’), systolic blood pressure (SBP) (e.g., ‘110’), diastolic blood pressure (DBP) (e.g., ‘56’), and acuity (e.g., ‘2’) were designated as value columns. Pain (e.g., ‘6-9’ and ‘yes.’) and the chief complaint (e.g., ‘BILATERAL FOOT PAIN’) were designated as text columns.</li> <li>• The column designations for the vital sign table were identical to the triage table, except</li> </ul> |  | 843<br>844<br>845<br>846<br>847<br>848<br>849<br>850<br>851<br>852<br>853<br>854<br>855<br>856<br>857<br>858<br>859<br>860 |

for the rhythm column (e.g., ‘Normal Sinus Rhythm’), which was treated as a category column. The vital signs table also had no chief complaint column and the ‘charttime’ column was used as the event time.

- For the medicine reconciliation table, the ‘in-time’ from the ED stay table was used as the event time, as it pertains to the patient’s medication history prior to the ED stay. The name column was designated as a text column, while the gsn, ndc, etc\_rn, and etccode columns were designated as category columns. The etcdescription column was not considered, as it is a description of the etccode column.
- For the medicine administration (pyxis) table, ‘charttime’ was used as the event time. The med\_rn, name, gsn\_rn, and gsn columns were all treated as category columns. The name column for the medicine reconciliation column did not have as high of a cardinality as the name column from the medicine reconciliation column, allowing it to be considered as a category column.
- For the metadata table, the ‘PerformedProcedureStepDescription’, ‘ViewPosition’, ‘ProcedureCodeSequence\_CodeMeaning’, ‘ViewCodeSequence\_CodeMeaning’, and ‘PatientOrientationCodeSequence\_CodeMeaning’ columns were considered, and designated as category columns.

## C Experiment setup

### C.1 Metrics

CheXbert-F1 (Smit et al., 2020), RadGraph-F1 (Delbrouck et al., 2022), BLEU-4 (Papineni et al., 2001), and BERTScore-F1 (roberta-large\_L17\_no-idf\_rescaled) (Zhang et al., 2020) have been found to correlate with radiologists’ assessment of reporting (Yu et al., 2023) and were a part of our evaluation. Additionally, we include CXR-BERT (Boecking et al., 2022; Nicolson et al., 2024a), CIDEr (Vedantam et al., 2015), METEOR (Banerjee and Lavie, 2005), and ROUGE-L (Lin and Hovy, 2003) as part of our evaluation. CheXbert-F1, RadGraph-F1, and CXR-BERT were intended to capture the clinical semantic similarity between the generated and radiologist reports, while

BERTscore-F1 was intended to capture general semantic similarity. Finally, CIDEr, METEOR ROUGE-L, and BLEU-4 were intended to capture the syntactic similarity between the generated and radiologist reports.

For the models in Table 2 that generate a report for each image in an exam, the average score was taken across all reports for an exam. Following this, the final average score was computed across all exams for both models that generate a report per image and those that generate a report per exam.

For CheXbert, the macro-averaged F1 was computed between the 14 CheXbert observations extracted from the generated and radiologist reports. “No mention”, “negative”, and “uncertain” were considered negative, while “positive” was considered positive. Here, the true positives, false positives, and false negatives were averaged over the reports of each exam for the models that generate a report per image.

We also perform statistical testing; first, a Levene’s test was conducted to reveal if the variances across model scores was homogeneous or not. If the assumption of equal variances was upheld, a one-way ANOVA was conducted to determine if there was a significant difference between models. Finally, pairwise Tukey-HSD post-hoc tests were used for pairwise testing. If the assumption of equal variances was violated, a one-way Welch’s ANOVA was conducted to determine if there was a significant difference between models. Finally, Games-Howell post hoc tests were used for pairwise testing. A  $p$ -value of 0.05 was used for all significance testing. Statistical testing was not performed for CheXbert, as it is a classification metric.

### C.2 Model

Our model is illustrated in Figure 2; following (Nicolson et al., 2024b), we utilised UniFormer as the image encoder (in particular, the  $384 \times 384$  base model warm started with its token labelling fine-tuned checkpoint) (Li et al., 2023). The image embeddings are formed by processing each image in the exam separately with the image encoder and then projecting its last hidden state to match the decoder’s hidden size using a learnable weight matrix. Each image was resized using bicubic interpolation so that its smallest side had a length of 384 and its largest side maintained the aspect ratio. Next, the resized image was cropped to a size of  $\mathbb{R}^{3 \times 384 \times 384}$ . The crop location was random during training and centred during testing. Following (El-

gendi et al., 2021), the image was rotated around its centre during training, where the angle of rotation was sampled from  $\mathcal{U}[-5^\circ, 5^\circ]$ . Finally, the image was standardised using the statistics provided with the UniFormer checkpoint. A maximum of five images per exam were used during training. If more were available, five were randomly sampled uniformly without replacement from the exam.

Again following (Nicolson et al., 2024b), we employed the Llama architecture for the decoder, which is notable for features such as its rotary positional encoding (RoPE), root mean square normalisation (RMSNorm), and SwiGLU activation function (Touvron et al., 2023). A byte-level byte pair encoding tokenizer (Wang et al., 2020) was trained with a vocabulary size of 30 000. It was trained on the findings, impression, indication, and history sections (not the comparison section) of the entire MIMIC-CXR training set, as well as the ‘pain’ and ‘chiefcomplaint’ columns from the triage table, the ‘name’ column of the medicine reconciliation table, and the ‘pain’ column from the vital signs table (from the entire MIMIC-IV-ED dataset). Newline, tab, repeated whitespaces, and leading and trailing whitespaces were removed from any text before tokenization.

The hyperparameters of the Llama decoder were six hidden layers, a hidden size of 768, 12 attention heads per layer, and an intermediate size of 3 072. The maximum number of position embeddings was set to 2048 to accommodate all the patient data embeddings and the report tokens. The maximum number of tokens that could be generated was set to 256, which was also the limit for the radiologist reports during training. During testing, a beam size of four was utilised. The Llama decoder allows a custom attention mask to be provided in current implementations.<sup>4</sup> This enabled non-causal masking to be utilised for the prompt and causal masking for the report token embeddings, as shown in Figure 5. This ensured that the self-attention heads were able to attend to all of the patient data embeddings at each position.

### C.3 Training

Three stages of training were performed. Each stage used *AdamW* (Loshchilov and Hutter, 2022) for mini-batch gradient descent optimisation, where training and evaluation was performed on a 94GB NVIDIA H100 GPU. The three stages were

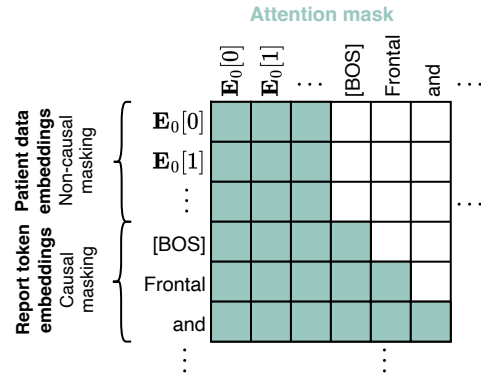


Figure 5: Attention mask for the decoder. Non-causal masking was used for the patient data embeddings and causal masking for the report token embeddings.

as follows:

1. Teacher forcing (TF) (Williams and Zipser, 1989) was performed on the MIMIC-CXR dataset with only the images for an exam as input, and exams that contained both a findings and impression section. This gave a training/validation split of 232 853/1 837 images, 125 416/991 exams, and 57 101/436 patients. Training was performed with an initial learning rate of  $5e-5$ , a mini-batch size of 8, a maximum of 32 epochs, and with float16 automatic mixed precision. All model parameters were trainable during this stage. The validation macro-averaged CheXbert-F1 was the monitored metric for checkpoint selection. This stage was necessary, as the language model struggled to generate reports from multiple sources without prior learning.
2. TF on the dataset described in Section 3 with the inputs described in Table 1. The training strategy was identical to the previous stage, except that a maximum of 16 epochs was performed, and the image encoder’s parameters were frozen (except for its projection). The models featured in Table 1 were trained using only the first two stages.
3. Reinforcement learning using self-critical sequence training (SCST) (Rennie et al., 2017) with CXR-BERT and BERTScore as the reward (each weighted with 0.5) was performed in the final stage of training. The sample report for SCST was generated with top- $k$  sam-

<sup>4</sup><https://huggingface.co/blog/poedator/4d-masks>

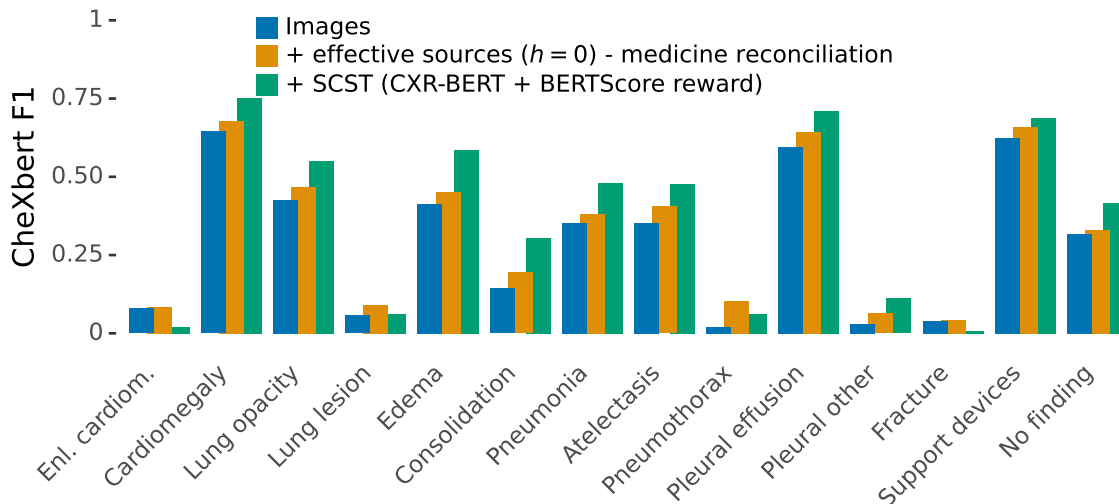


Figure 6: F1-score for each CheXbert label. ( $n = 9580$  exams;  $958 \times 10$  runs for ‘Images’ and ‘Images + effective sources ( $h = 0$ ) - medicine reconciliation’ and  $n = 2874$  exams;  $958 \times 3$  runs for ‘Images + effective sources ( $h = 0$ ) - medicine reconciliation + SCST (CXR-BERT + BERTScore reward)’.)

pling ( $k = 50$ ). Training was performed with an initial learning rate of  $5e-6$ , a mini-batch size of 32, a maximum of 24 epochs, and with float32 precision. The image encoder’s parameters were frozen during this stage (except for its projection). The validation BERTScore-F1 was the monitored metric for checkpoint selection, as it helped to select checkpoints less prone to repetitions. This stage of training was only applied to the best model from Table 1, ‘Images + effective sources ( $h = 0$ ) - medicine reconciliation’, with the results presented in Table 2. This model had 161 185 728 parameters.

#### C.4 Comparison Models

The generated reports for the models in Table 2 were attained as follows:

- EMNLI reports were generated following <https://github.com/ysmiura/ifcc> (Miura et al., 2021).
- CMN reports were generated following <https://github.com/zhjohnchan/R2GenCMN> (Chen et al., 2021).
- TranSQ reports were kindly provided by the authors (Kong et al., 2022).
- RGRG reports were generated following <https://github.com/ttanida/rgrg> (Tanida et al., 2023).

- CvT2DistilGPT2 reports were generated following <https://github.com/aeherc/cvt2distilgpt2> (Nicolson et al., 2023).
- RaDialog reports were kindly provided by the authors (Pellegrini et al., 2023).
- MedXChat reports were kindly provided by the authors (Yang et al., 2023).
- CXR-LLaVA-v2 reports were generated following <https://huggingface.co/ECOFRI/CXR-LLAVA-v2> (Lee et al., 2024).
- CXRMate reports were generated following <https://huggingface.co/aeherc/cxrmate> (Nicolson et al., 2024a).
- CXRMate-RRG24 reports were generated following <https://huggingface.co/aeherc/cxrmate-rrg24> (Nicolson et al., 2024b).

CXRMate-RRG24 was trained on five datasets, including MIMIC-CXR. RGRG was trained on the ImaGenome dataset derived from MIMIC-CXR — which may have some overlap with our test set.

#### D Ancillary results

In Figure 6, the F1-scores for each CheXbert label are shown. The ‘Images + effective sources ( $h = 0$ ) - medicine reconciliation’ model from Table 1 improves performance across all labels compared to the ‘Images’ model. This suggests that

Table 4: Results for exams that have a medicine administration table ( $n = 3520$ ; studies  $352 \times 10$  runs). Underlined scores indicate a significant difference to the scores of ‘Images’ ( $p < 0.05$ ).

| Inputs                           | RG           | CX           | CB           | BS           |
|----------------------------------|--------------|--------------|--------------|--------------|
| Images                           | 26.24        | 28.36        | 57.17        | 24.33        |
| Images + medicine administration | <u>26.95</u> | <u>28.53</u> | <u>58.94</u> | <u>24.93</u> |

1095 incorporating ancillary data from MIMIC-IV-ED  
 1096 and MIMIC-CXR provides a general improvement,  
 1097 rather than benefiting any specific pathology.

1098 Further improvements are seen when training the  
 1099 ‘Images + effective sources ( $h = 0$ ) - medicine rec-  
 1100 onciliation’ model with SCST (i.e., our model from  
 1101 Table 2) for most pathologies. However, there are  
 1102 performance decreases for ‘enlarged cardiome-  
 1103 diastinum’, ‘lung lesion’, ‘pneumothorax’, and ‘frac-  
 1104 ture’. This might be due to these pathologies be-  
 1105 ing underrepresented in the MIMIC-CXR dataset,  
 1106 leading the model to optimise for more common  
 1107 pathologies during SCST.

1108 The results for exams that include a medicine  
 1109 administration table are show in Table 4. Adding  
 1110 the medicine administration table produced a sig-  
 1111 nificant improvement in the scores, indicating that  
 1112 it should be considered if available.