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 vital signs, medicines, 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; this significantly enhances the diagnostic accuracy of generated radiology reports. Our comprehensive evaluation demonstrates the benefits of 018 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

011

014

037

041

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, timely reporting is hindered by a multitude of issues, including high patient volumes and limited availability of radiologists (Bailey et al., 2022).

Automated CXR report generation using multimodal language models is a promising solution (Jones et al., 2021). Potential benefits include enhanced radiologist effectiveness, streamlining report writing, and improved patient outcomes (Shen, 2021; Irmici et al., 2023). Early methods produced a separate report for each image within an exam (Wang et al., 2018). Later methods improved

Patient data



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. Emergencyspecific data, such as reconciled medicines and aperiodic vital signs, are also available for the patient.

065

077

085

on this by considering all images of an exam to generate a single report (Miura et al., 2021; Nicolson et al., 2024a), and incorporating prior exams for a patient (Wu et al., 2022; Nicolson et al., 2024a). Including the reason for the exam (the *indication* section in Figure 1) offered a further improvement (Nguyen et al., 2023). This indicates that CXR report generation benefits from the inclusion of a more comprehensive set of patient data.

Incorporating clinical information, including electronic health record (EHR) data, enhanced the interpretation accuracy, clinical relevance, and reporting confidence of radiologists' findings (Castillo et al., 2021). A growing push to integrate EHR systems into radiology workflows highlights the potential for CXR report generation models to leverage patient data directly (Geeslin and Gaskin, 2016). In this study, we aim to empirically investigate if such data can also improve CXR report generation. To facilitate this, we combine CXR exams from MIMIC-CXR (Johnson et al., 2019) with emergency department (ED) patient records from MIMIC-IV-ED (Johnson et al., 2023). This provides a wide variety of multimodal data per exam, as shown in Figure 1. From MIMIC-CXR, we utilise the images, their metadata, and several sections of the radiology report. Notably, we incorporate the history or comparison section of the report, which has not been investigated previously. From MIMIC-IV-ED, we incorporate triage data, aperiodic vital signs, medicines, and other data to provide a wider clinical context.

We also investigate how to harmonise these heterogeneous data into patient data embeddings to prompt a multimodal language model. In doing so, we develop methods to transform tabular and aperiodic time series data into embeddings that can be used alongside token and image embeddings. We evaluate our model using metrics shown to closely correlate with radiologists' assessments of reporting (Yu et al., 2023). Through our evaluation, we demonstrate that complementary information from different data sources can improve the diagnostic accuracy of CXR report generation. The main contributions of this work are:

- An investigation demonstrating how integrating specific patient data sources, such as medicines and vital signs, enhances CXR report generation and improves diagnostic accuracy.
- Introducing methods to convert numerical, categorical, text, temporal, and image data into embeddings for a multimodal language model.

• A dataset linking MIMIC-CXR exams with MIMIC-IV-ED records, along with the code and Hugging Face model (available at: https://anonymous.4open.science/r/anon-D83E).

095

096

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

2 Background and Related Work

Incorporating more patient data has improved diagnostic accuracy in radiology reporting. Initial improvements came from using multiple images per exam, like EMNLI; CXR exams often include complementary frontal and lateral views of the patient (Miura et al., 2021; Gaber et al., 2005). Methods such as CXRMate enhance diagnostic accuracy by incorporating a patient's prior exams to identify changes over time (Nicolson et al., 2024a; Wu et al., 2022; Kelly, 2012; Bannur et al., 2023; Hou et al., 2023). Including the *indication* section of the radiology report to provide clinical context also provides an improvement (Nguyen et al., 2023). Our investigation focuses on leveraging a more comprehensive set of patient data to improve diagnostic accuracy.

ED records contain a wide range of data, as shown in Figure 1. The reconciled medicines may include furosemide, a diuretic commonly prescribed for managing fluid overload, often associated with conditions such as pulmonary edema or congestive heart failure. Elevated blood pressure and an increased heart rate in a patient's vital signs may correlate with findings such as cardiomegaly or vascular changes. Vital signs such as high temperature, elevated respiratory rate, and low oxygen saturation, along with chief complaints of cough and shortness of breath, may suggest pneumonia. Incorporating such data could complement imaging evidence and provide additional context to support better predictions. Our findings demonstrate that ED patient data can indeed improve CXR report generation.

Recent advancements in integrating multimodal patient data have improved diagnostic and predictive healthcare tasks. A Transformer encoder combining imaging and non-imaging data outperformed single-modality models in diagnosing multiple conditions (Khader et al., 2023b). Similarly, the MeTra architecture, integrating CXRs and clinical parameters, excelled in predicting ICU survival (Khader et al., 2023a), and ETHOS, with zero-shot learning, surpassed single-modality models in predicting mortality, ICU length of stay, and readmission rate (Renc et al., 2024). These stud-

ies underscore the value of multimodal data, and 144 our work demonstrates its benefits for CXR report 145 generation. 146

Multi-task learning has enhanced biomedical 147 models by leveraging shared knowledge across 148 tasks. Med-PaLM M, a generalist biomedical 149 model, excels in classification, question answering, VQA, report summarisation, report generation, and 151 genomic variant calling, using diverse modalities 152 153 like images, text, and genomics, often outperforming specialised models (Tu et al., 2024). Similarly, 154 MIMIC-CXR has been utilised in multi-task learn-155 ing with models like MedXChat, which integrates instruction-tuning and Stable Diffusion for tasks 157 158 like CXR report generation, VQA, and report-to-CXR generation, surpassing other LLM multi-task 159 learners (Yang et al., 2023). RaDialog combines 160 161 visual features and pathology findings to generate accurate radiology reports and enable interactive 162 tasks, improving clinical efficacy. CXR-LLaVA, 163 a multimodal LLM, outperformed models such as 164 GPT-4 Vision and Gemini Pro Vision in CXR re-165 port generation (Lee et al., 2024). 166

Determining the state-of-the-art in CXR report generation is challenging due to model unavailabil-168 ity and limited comparisons with recent methods. The 2024 Shared Task on Large-Scale Radiology Report Generation (RRG24) aimed to address this by benchmarking models on a common leaderboard. The winning model, CXRMate-RRG24 (Nicolson et al., 2024b), a derivative of CXRMate, emerged as a strong contender for state-of-theart. In this work, we compare our model to established models (e.g., EMNLI) and recent benchmarks (e.g., CXRMate-RRG24, CXRMate, CXR-LLaVA, MedXChat, and RaDialog). We ensure a fair comparison by using available code or obtaining generated reports directly from the authors. Our evaluation indicates that our model represents a statistically significant improvement over these.

3 Dataset

167

169

170

171

172

173

174

175

176

178

179

181

183

We construct a dataset of 46 106 patients by linking 185 individual patient information from two separate 186 sources: (1) CXR exams from MIMIC-CXR and 187 (2) emergency records from MIMIC-IV-ED. Thus, we consider MIMIC-CXR exams that occurred dur-189 ing an ED stay from MIMIC-IV-ED. Both datasets 190 are publicly available and originate from the Beth 191 Israel Deaconess Medical Center in Boston, MA. 192 MIMIC-CXR was formed by first extracting pa-193

tient identifiers for exams performed in the ED between 2011–2016, and then extracting all exams for this set of patients from all departments between 2011–2016. Each exam includes a semi-structured free-text radiology report (Figure 1) written by a practising radiologist contemporaneously during routine clinical care. Models are often trained to generate the *findings* and *impression* sections of a radiology report, where the former details the interpretation of a patient's exam and the latter summarises the most important findings. All images and reports were de-identified to protect privacy. Sections from the radiologist reports were extracted using a modification of the official text extraction tool in order to obtain the findings, impression, indication, history, and comparison sections.¹

194

196

197

199

200

202

203

204

205

206

207

208

210

211

212

213

214

215

216

217

218

219

220

221

224

225

226

227

228

229

231

232

233

235

236

237

238

240

MIMIC-IV-ED consists of de-identified data from ED stays between 2011–2019. The data was converted into a denormalised relational database with six primary tables: ED stays, diagnosis, reconciled medicines, administered medicines, triage, and aperiodic vital signs. We do not consider the diagnosis table in this work, as it indicates the outcome of a patient's ED stay. The patients of 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 patient exams with times (formed by the 'StudyDate' and 'StudyTime' columns of the metadata table) that occurred within the 'intime' and 'outtime' of one of the patient's ED stays.² Events during an ED stay that occurred after the exam were removed. 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 45527/343/236 patients, 76398/556/958 exams, and 151 818/1 137/1 812 CXRs. Further details are provided in Appendix A.

Methods 4

We develop a novel approach to transform different sources of patient data from MIMIC-CXR and MIMIC-IV-ED into embeddings; these are then used to prompt a multimodal language model to generate the findings and impression sections of the

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

²Exam 59128861 was removed as it overlapped with two separate ED stays of a patient.



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.



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.

radiology report, as illustrated in Figure 2. Each embedding of the prompt is the summation of a patient data embedding, a source embedding, a position embedding, and a time delta embedding. Source embeddings differentiate the source of the datum, for example, the 'chief complaint' column of the triage table, the indication section, or an image. A time delta embedding represents the time difference between an event and the exam. 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 B and Subsection D.2, respectively.

241

242

245

247

249

250

251

252

4.1 Time, Position, & Source Embeddings

Events from MIMIC-IV-ED, e.g., administered medicines, are timestamped and are more relevant as they occur closer to the exam time (Ben Abacha et al., 2023). Hence, time delta embeddings are used to indicate this to the model. The time delta is the event time subtracted from the exam time, converted to hours and mapped using $D = 1/\sqrt{\Delta + 1}$, emphasising recent events. These mapped time deltas are processed via a feedforward neural network (FNN), $f(DW_1)W_2$, where $W_1 \in \mathbb{R}^{1,2048}$, $W_2 \in \mathbb{R}^{2048,H}$, $f(\cdot)$ is the SiLU activation (Hendrycks and Gimpel, 2016), and *H* is the decoder's hidden size. As shown in Figure 2, these embeddings are applied only to the prompt. 257

258

259

260

261

262

265

267

270

271

321

322

323

324

325

326

327

328

329

The exam time for current and prior exam data was used as the event time for the time delta calculation.

272

273

274

275

276

277

281

285

290

291

292

293

296

297

301

304

305

307

310

312

313

314

315

316

317

320

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 patient data 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 valuecategory columns (described in the next section), prior images, and prior report sections.

4.2 Patient Data Embeddings: Tabular Data

An example table and its conversion to embeddings is shown in Figure 3. The columns of each table were designated as value, category, text, or time columns. Value columns contained numeric data, while category columns contained categorical data. To convert an exam's tabular data to embeddings, data from value and category columns were grouped by their time delta, where each group formed 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 $f(X_i W_1) W_2$ to form the embedding, where $oldsymbol{X}_i \in \mathbb{R}^{|U_C|,|L_i|}$ are the grouped features, U_C is the set of unique time deltas, $W_1 \in \mathbb{R}^{|L_i|,2048}$ and $W_2 \in \mathbb{R}^{2048,H}$, 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 6. 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 C.

4.3 Patient Data Embeddings: Report Sections

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, which are used to capture disease progression. 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.

5 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 four that capture the semantics of radiology reporting - RadGraph-F1 (RG), CheXbert-F1 (CX), CXR-BERT (CB), and GREEN (G) — as well as three natural language generation metrics: BERTScore-F1 (BS), ROUGE-L (R-L), and BLEU-4 (B4). We also propose a metric that measures n-gram repetition rate, namely the absence of repeated *n*-grams (ARN). Comprehensive details on ARN and the other metrics, the model architecture, training procedure, significance testing, and comparison methods are provided in Appendix D.

6 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

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; 958×10 runs). <u>Underlined</u> and <u>dashed</u> <u>underlined</u> 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	СВ	G	BS	R-L	B4	$\overline{ \boldsymbol{\mathcal{E}}[:,0] }$
	In	ages only						
Images	24.54	30.10	59.25	35.16	24.26	25.91	4.75	272.4
Patient emer	gency dep	artment d	ata (MIM	IC-IV-ED)			
Images + ED stays	24.20	29.55	58.37	34.64	24.06	25.77	4.66	273.4
Images + triage	24.59	31.33	<u>62.79</u>	35.78	24.40	25.96	4.76	278.9
Images + vital signs	24.23	30.61	60.61	35.15	24.04	25.86	4.70	274.7
Images + reconciled medicines	25.10	32.05	<u>64.70</u>	36.32	24.71	26.29	4.93	355.6
Images + administered medicines	24.22	30.40	60.13	34.85	23.97	25.61	4.58	273.0
Patier	nt radiolo	gy data (N	AIMIC-C	XR)				
Images + indication	25.01	32.78	<u>65.49</u>	35.88	24.73	26.32	5.15	279.5
Images + history	24.88	31.66	63.91	35.76	24.91	<u>26.70</u>	<u>5.54</u>	277.0
Images + metadata	24.07	30.42	59.75	34.79	23.86	25.59	4.58	273.4
	Pı	rior exams						
Images + $h = 1$	24.71	30.98	<u>62.60</u>	35.81	24.38	26.00	4.82	603.0
Images + $h = 2$	24.56	31.43	62.09	35.43	24.04	25.80	4.84	878.1
Images + $h = 3$	24.50	30.73	59.89	35.21	24.03	25.82	4.70	1134.3
Images + $h = 1$ + comparison	24.92	31.46	<u>62.93</u>	35.84	24.34	26.03	4.89	607.4
Images + $h = 2$ + comparison	24.52	31.01	<u>61.36</u>	34.89	23.90	25.62	4.72	882.6
Images + $h = 3$ + comparison	24.31	30.93	60.10	34.35	23.31	<u>25.39</u>	4.72	1138.8
All effective sources (tri	age, recoi	nciled med	licines, in	dication,	and histor	ry)		
Images + effective sources $(h = 0)$	<u>25.52</u>	32.49	<u>65.93</u>	<u>36.26</u>	<u>25.16</u>	<u>26.81</u>	<u>5.34</u>	373.9
Images + effective sources $(h = 1)$	25.11	31.14	<u>61.19</u>	35.80	<u>24.95</u>	<u>26.45</u>	<u>5.21</u>	704.5
Images + effective sources ($h = 1 + \text{comparison}$)	25.05	30.68	60.99	35.94	<u>24.94</u>	<u>26.48</u>	<u>5.24</u>	709.0
Ablation fro	om Image	s + effecti	ve source.	s(h=0)				
- triage	25.65	32.85	65.38	36.33	25.25	26.75	5.33	367.4
 reconciled medicines 	25.43	32.48	65.63	36.42	25.23	26.86	5.40	290.7
- indication	25.46	32.92	65.69	36.41	25.21	26.79	5.36	366.7
- history	25.41	32.53	65.82	36.65	25.12	26.72	5.37	369.2
- time delta	25.31	33.03	65.72	36.17	25.10	26.75	5.34	373.9

additional data sources improve performance compared to using only images. Significant improvements were observed by incorporating either triage or reconciled medicine data from MIMIC-IV-ED dataset. Notably, this data markedly improved scores on the radiology report metrics (RG, CX, CB, and G). These findings demonstrate that ED patient data can improve the diagnostic accuracy of CXR report generation. Aperiodic vital sign and administered medicine data did not significantly improve the scores overall, likely due to their frequency of occurrence in the exams (62% and 37%, respectively). However, as shown in Table F.1, a significant improvement in performance was attained when evaluated solely on exams that include an aperiodic vital sign table.

371

373

375

377

381

386

390

391

394

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. A history size of one or two significantly increased the scores, which is consistent with previous findings (Wu et al., 2022). However, performance gradually degraded as the history size increased, which contradicts earlier studies. We suspect this is due to the size of the prompt increasing as h grows, combined with the limitations of our model architecture. $|\mathcal{E}[:,0]|$ in Table 1 is the average prompt length over the test set, where $\mathcal{E} = [\mathbf{E}_0, \mathbf{E}_1, \cdots]$. It can be seen that $|\mathcal{E}[:,0]|$ increases substantially as h increases. Since we provide all inputs to the decoder's selfattention, a large input size may cause attention dilution (Qin et al., 2022). 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.

411

412

413

414

415

416

417

418

395

397

Table 2: Benchmark models on the test set described in Section 3 (n = 958). Evaluation is on the **findings** section only. <u>Underlined</u> 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	G	BS	R-L	B4	ARN
EMNLI (Miura et al., 2021)	152 173 exams	29.1	28.9	66.6	41.5	24.4	29.3	4.1	95.1
CMN (Chen et al., 2021)	270 790 images	23.6	24.3	49.4	36.6	19.7	27.8	4.0	99.3
TranSQ (Kong et al., 2022)	368 960 images	28.7	30.4	62.3	38.2	20.4	23.3	4.1	98.5
RGRG (Tanida et al., 2023)	166 512 images	22.9	22.8	37.9	31.1	23.4	22.0	3.7	96.5
CvT2DistilGPT2 (Nicolson et al., 2023)	270 790 images	23.9	29.3	59.8	37.0	24.8	28.6	5.4	99.0
RaDialog (Pellegrini et al., 2023)	276 778 images	24.4	38.4	60.7	34.9	26.2	26.7	4.8	94.4
MedXChat (Yang et al., 2023)	270 790 images	21.0	13.1	21.3	31.4	19.3	23.8	4.0	97.9
CXR-LLaVA-v2 (Lee et al., 2024)	193 513 images	19.4	20.7	44.1	24.0	23.6	21.1	1.7	99.7
CXRMate (Nicolson et al., 2024a)	125 395 exams	26.5	33.9	71.3	40.3	30.5	29.1	7.5	98.2
CXRMate-RRG24 (Nicolson et al., 2024b)	550 395 exams	28.9	31.2	58.2	40.2	31.0	28.7	6.6	97.7
$\frac{1}{\text{Images + effective sources } (h = 0)}$	76 398 exams	25.1	29.6	66.0	36.9	29.4	27.8	5.8	98.5
+ RL (CXR-BERT + BERTScore reward)	76 398 exams	30.4	35.7	<u>79.1</u>	41.6	<u>37.2</u>	<u>31.6</u>	<u>8.7</u>	93.5
+ reward per section	76 398 exams	30.1	33.7	<u>78.3</u>	41.6	<u>37.5</u>	<u>32.2</u>	8.4	94.6
+ ARN reward	76 398 exams	30.2	33.6	<u>78.0</u>	40.7	<u>37.3</u>	<u>31.9</u>	7.6	99.3

	Triage			Cá	ase study	
Se la	heartrate	o2sat	acuity	nain	chiefcomplaint]
		02301	acuity			Indication: New endotracheal tube placement.
A A A	57	90	I	0	ULCER/CHF	
1	Reconciled	medica	tions (nan	nes): Metop	rolol Tartrate, Thiamine H	ICI, Albuterol Sulfate, Provigil, spironolactone,
	Fluoxetine, r	nicotine (polacrilex)	, Imdur, Mul	tivitamin, Ibuprofen, Sand	ctura XR, Metformin, Abilify, Plavix, Furosemide, ProAir
	HFA, Briefs,	Adult-E>	tra Large,	Omeprazole	e, ProFit Precision Scale,	, Senna, Estrace, Lac-Hydrin, triazolam, Lisinopril.
Radiologist findings: The	ere has been	interval	placement	of an endot	racheal tube, which is low	v lying with tip approximately 1.6 cm above the carina.
An esophageal tube is in p	place coursing	inferior	to the dia	phragm; how	ever, tip out of view of th	ne radiograph. Lung volumes remain low with mild
pulmonary edema. No sigi	nificant pleura	al effusio	n or pneur	nothorax is i	dentified. The cardiomed	liastinal silhouette is enlarged, however, unchanged.
Distal tip of ecophageal tu	Interval place	ment of	endotrach	eal tube with	n tip low lying, approxima	itely 1.6 cm above the carina. Mild pulmonary edema.
Distal tip of esophageal tu	De HOL WILLIN					· · · · · · · · · · · · · · · · · · ·
Image + triage + reconci	led medication	ons + in	dication +	(Model: Im	ages + effective sources	(h=0) from Table 1)
Generated findings: An	endotracheal t	ube is in	place with	the tip term	inating at the level of the	thoracic inlet approximately 4 cm above the carina. A
right upper extremity PICC	is in place w	ith the tip	p terminati	ng in the low	SVC. A nasogastric tube	e is seen coursing below the diaphragm and out of view
on this image. A right uppe	er lobe opacity	nas res	solved from	n prior study.	The neart is moderately	enlarged, as before. There is mild pulmonary edema.
No pieural ettusion or phei	umotnorax is	detected	I. I nere is	no tocal con	solidation concerning for	pneumonia.
Generated Impression: E	ndotracheal i	ube tip a	at the level	of the thora	cic inlet 4 cm above the c	carina and nasogastric tube in appropriate position. Mild
pullionary edenia and car	diomegaly.					
Image (Model: Images fro	m Table 1)					``
Generated findings: The	re is an <mark>endot</mark>	racheal	tube in sat	isfactory pos	sition, terminating 3.1 cm	above the level of the carina. A right-sided PICC now
terminates at the cavoatria	al junction. An	enteric	tube is see	en coursing b	pelow the level of the dia	phragm, inferior aspect out of the field of view. There is
moderate cardiomegaly. T	he lungs are	clear wit	hout focal	consolidation	n, pleural effusion or pne	umothorax.
Generated impression: 1	. Appropriate	position	of a right	sided PICC,	terminating at the cavoar	trial junction. 2. No acute cardiopulmonary process.

Figure 4: Case study demonstrating how incorporating auxiliary patient data can aid with report generation.

Next, we combined all effective sources of patient data (those providing a significant improvement). This included 'triage', 'reconciled medicines', 'indication', and 'history'. The best performance was observed with no prior exams (h = 0), indicating that using any prior exams in combination with other sources is detrimental with our model, possibly due to attention dilution. With h = 0, the combination of all effective sources outperformed each individual source. We then conducted an ablation study using 'Images + effective sources (h = 0)', which demonstrated that removing any individual patient data source did not result in a significant change in performance.

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

Following this, we further trained 'Images + effective sources (h = 0)' with reinforcement learning (RL), as described in Subsection 5. Its performance is shown in Table 2; a CXR-BERT and BERTScore composite reward was used, which demonstrates a marked improvement for each metric, except ARN. The low ARN indicates that this reward introduced repetitions. We also propose to calculate the reward separately for the findings and impression section, as described in Appendix E. While this produces similar results for the findings section, as shown in Table 2, this significantly improves the scores on the impression section section, as shown in Table E.1. Finally, we incorporate ARN into the composite reward. This effectively reduces repetitions, as evidenced by the improved ARN, albeit with a slight trade-off in the other metrics. Compared to other benchmark CXR report generation models in the literature that included MIMIC-CXR in their training data, our model significantly outperformed them on multiple metrics in Table 2, despite having substantially fewer training samples. This demonstrates the impact of incorporating auxiliary patient data on CXR report generation.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

Figure 4 demonstrates how auxiliary patient data 457 enhances CXR report generation. Mild pulmonary 458 edema was identified only when this data was incor-459 porated. The patient's low oxygen saturation, chief 460 complaint of congestive heart failure (CHF) — a 461 common cause of pulmonary edema - and recon-462 ciled medicines (Furosemide, Metoprolol Tartrate, 463 Lisinopril, Spironolactone) indicate active manage-464 ment of fluid overload and cardiac dysfunction, all 465 pointing to pulmonary edema. This supplemen-466 tary evidence allowed the model to corroborate the 467 imaging findings and identify pulmonary edema. 468

In Appendix G, we perform an error analysis to assess the influence of auxiliary patient data on the generated reports. Our findings show that incorporating auxiliary patient data increases the AUC for 10 out of the 14 CheXpert labels (Figure E.1), demonstrating its utility across multiple pathologies. Additionally, we analysed its impact on the generated reports for eight exams, with the following key observations:

469

470

471

472

473

474

475

476

477

487

488

489

508

478 **True positives** (n = 2): The model utilised sup-479 portive auxiliary patient data effectively. (See Ap-480 pendix G.2.1 and G.2.2.)

False positives (n = 2): The model was misled by confounding auxiliary patient data. (See Appendix G.2.3 and G.2.4.)

484 **True negatives** (n = 2): The model correctly ig-485 nored confounding auxiliary patient data. (See Ap-486 pendix G.2.5 and G.2.6.)

False negatives (n = 2): The model failed to leverage supportive auxiliary patient data. (See Appendix G.2.7 and G.2.8.)

Auxiliary patient data sources-including the in-490 491 dication and history sections, triage data, and reconciled medicines-collectively contributed to the 492 model's predictions. No single source consistently 493 dominated in providing evidence, with the interplay 494 between these sources frequently complementing 495 one another. A critical challenge for the model lies 496 in its ability to appropriately balance the auxiliary 497 patient data evidence with imaging evidence, partic-498 ularly when conflicting signals are present. To ad-499 dress this limitation, we propose two key improvements: increasing the size of the training dataset, which is currently relatively small, and adopting an LLM-based decoder. LLMs offer advanced reason-504 ing capabilities, enabling them to better synthesise and prioritise evidence from diverse sources. 505

> Table 3 compares different methods for converting value and category columns into embeddings using the triage and reconciled medicines table, as

Table 3: Patient data embedding strategies. <u>Underlined</u> indicates a stat. sig. difference to 'Baseline' (p < 0.05).

Embeddings	RG	CX	СВ	BS
	Images			
Baseline	29.00	25.81	59.04	23.85
Images + triage	+ recond	ciled med	icines	
Grouped embeddings	31.69	<u>26.72</u>	<u>64.01</u>	24.38
Separate embeddings	25.28	25.32	<u>46.29</u>	23.51
Values-to-text, categories- to-embeddings	30.70	<u>26.46</u>	58.62	<u>24.58</u>

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

these 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 learnable 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 was the best representation of heterogeneous patient data for a multimodal language model.

7 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, medicines, and the history section of radiology reports. We present 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.

8 Limitations

548

549

550

555

557

559

567

570

571

573

574

578

579

582

586

590

592

594

598

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.

> This study currently lacks subjective evaluation by radiologists, which is essential for assessing the quality of generated reports. We plan to address this by evaluating with a private dataset and conducting radiologist-led assessments. To facilitate this, we are securing agreements and ethics approval for access to patient data and radiologist time. However, this process is extensive and beyond the scope of this study, and will instead be used to subjectively evaluate future models.

> 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 administered medicines available, 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. It also at times struggles with supportive or confounding evidence from the auxiliary patient data, introducing false positive or false negative predictions. Future work should explore advanced attention mechanisms, hierarchical models, and LLMs to better manage large input sequences and to better balance auxiliary patient data evidence with imaging evidence.

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. 599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

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.

9 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.
- Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fernando Pérez-García, Maximilian Ilse, Daniel C. Castro, Benedikt Boecking, Harshita Sharma, Kenza Bouzid, Anja Thieme, Anton Schwaighofer, Maria Wetscherek, Matthew P. Lungren, Aditya Nori, Javier

- 703 704 705 706 707 708 709 712 713 715 716 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 752
- 753 754 755 756 757

- Alvarez-Valle, and Ozan Oktay. 2023. Learning to Exploit Temporal Structure for Biomedical Vision-Language Processing. In *CVPR*, pages 15016–15027.
- Asma Ben Abacha, Wen-wai Yim, George Michalopoulos, and Thomas Lin. 2023. An investigation of evaluation methods in automatic medical note generation. In *Findings of the ACL*, pages 2575–2588, Toronto, Canada.
- 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.

663

670

671

672

673

674

675

677

683

695

702

- Chelsea Castillo, Tom Steffens, Lawrence Sim, and Liam Caffery. 2021. The effect of clinical information on radiology reporting: A systematic review. *Journal of Medical Radiation Sciences*, 68(1):60–74.
- 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 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. *Frontiers in Medicine*, 8.
- Khalid A Gaber, Clive R McGavin, and Irving P Wells. 2005. Lateral Chest X-Ray for Physicians. *Journal of the Royal Society of Medicine*, 98(7):310–312.
- Matthew G. Geeslin and Cree M. Gaskin. 2016. Electronic Health Record–Driven Workflow for Diagnostic Radiologists. *Journal of the American College of Radiology*, 13(1):45–53.
- Dan Hendrycks and Kevin Gimpel. 2016. Gaussian Error Linear Units (GELUs). arXiv:1606.08415 [cs.LG].
- Wenjun Hou, Yi Cheng, Kaishuai Xu, Wenjie Li, and Jiang Liu. 2023. RECAP: Towards Precise Radiology Report Generation via Dynamic Disease Progression Reasoning. In *Findings of the ACL: EMNLP*, pages 2134–2147.
- Giovanni Irmici, Maurizio Cè, Elena Caloro, Natallia 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? *Diagnostics*, 13(2):216.

- Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Leo Anthony Celi, Roger Mark, and Steven Horng. 2023. MIMIC-IV-ED (version 2.2). PhysioNet.
- 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.
- 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. *Journal of Medical Imaging and Radiation Oncology*, 65(5):538–544.
- Barry Kelly. 2012. The chest radiograph. *The Ulster Medical Journal*, 81(23620614):143–148.
- Firas Khader, Jakob Nikolas Kather, Gustav Müller-Franzes, Tianci Wang, Tianyu Han, Soroosh Tayebi Arasteh, Karim Hamesch, Keno Bressem, 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. *Scientific Reports*, 13(1):10666.
- Firas Khader, Gustav Müller-Franzes, Tianci Wang, Tianyu Han, Soroosh Tayebi Arasteh, Christoph Haarburger, Johannes Stegmaier, Keno Bressem, 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. *Radiology*, 309(1):e230806.
- Ming Kong, Zhengxing Huang, Kun Kuang, Qiang Zhu, and Fei Wu. 2022. TransQ: Transformer-Based Semantic Query for Medical Report Generation. In *MICCAI*, volume 13438, pages 610–620.
- Seowoo Lee, Jiwon Youn, Hyungjin Kim, Mansu Kim, and Soon Ho Yoon. 2024. CXR-LLAVA: a multimodal large language model for interpreting chest X-ray images. arXiv:2310.18341 [cs].
- 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. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12581–12600.
- Chin-Yew Lin and Eduard Hovy. 2003. Automatic evaluation of summaries using N-gram co-occurrence statistics. In *NAACL*, volume 1, pages 71–78.
- Ilya Loshchilov and Frank Hutter. 2022. Decoupled Weight Decay Regularization. In *ICLR*.

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

813

814

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 *NAACL*, pages 5288– 5304.

758

759

762

771

773

776

777

778

779

781

784

787

789

790

795

796

799

801

807

810

811

812

- Dang Nguyen, Chacha Chen, He He, and Chenhao Tan. 2023. Pragmatic Radiology Report Generation. In *Proceedings of the 3rd Machine Learning for Health Symposium*, pages 385–402. PMLR. ISSN: 2640-3498.
- Aaron Nicolson, Jason Dowling, Douglas Anderson, and Bevan Koopman. 2024a. Longitudinal data and a semantic similarity reward for chest X-ray report generation. *Informatics in Medicine Unlocked*, 50:101585.
 - Aaron Nicolson, Jason Dowling, and Bevan Koopman. 2023. Improving chest X-ray report generation by leveraging warm starting. *Artificial Intelligence in Medicine*, 144:102633.
 - 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 G eneration. In *The 23rd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*.
- Sophie Ostmeier, Justin Xu, Zhihong Chen, Maya Varma, Louis Blankemeier, Christian Bluethgen, Arne Edward Michalson Md, Michael Moseley, Curtis Langlotz, Akshay S Chaudhari, and Jean-Benoit Delbrouck. 2024. GREEN: Generative radiology report evaluation and error notation. In *Findings of the ACL: EMNLP*, pages 374–390, Miami, Florida, USA.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2001. BLEU: a method for automatic evaluation of machine translation. In *ACL*, page 311.
- Chantal Pellegrini, Ege Özsoy, Benjamin Busam, Nassir Navab, and Matthias Keicher. 2023. RaDialog: A Large Vision-Language Model for Radiology Report Generation and Conversational Assistance. ArXiv:2311.18681 [cs].
- Zhen Qin, Xiaodong Han, Weixuan Sun, Dongxu Li, Lingpeng Kong, Nick Barnes, and Yiran Zhong. 2022. The Devil in Linear Transformer. In *EMNLP*, pages 7025–7041.
- 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.
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-Critical Sequence Training for Image Captioning. In *CVPR*, pages 1179–1195.
- Dinggang Shen. 2021. Grand Challenges in Radiology. *Frontiers in Radiology*, 1.

- 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 *EMNLP*, pages 1500–1519.
- Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. 2023. Interactive and Explainable Regionguided Radiology Report Generation. In *CVPR*, pages 7433–7442.
- 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].
- 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. NEJM AI, 1(3):AIoa2300138.
- Changhan Wang, Kyunghyun Cho, and Jiatao Gu. 2020. Neural machine translation with byte-level subwords. In *AAAI*, pages 9154–9160.
- 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 *CVPR*.
- Ronald J. Williams and David Zipser. 1989. A Learning Algorithm for Continually Running Fully Recurrent Neural Networks. *Neural Computation*, 1(2):270– 280.
- 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 *ICCL*, pages 2952–2961.
- Ling Yang, Zhanyu Wang, and Luping Zhou. 2023. MedXChat: Bridging CXR Modalities with a Unified Multimodal Large Model. arXiv:2312.02233 [cs].
- 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. *Patterns*, page 100802.

957

958

959

960

961

962

963

916

917

918

919

920

921

922

869 870 871

873

875

877

878

879

884

886

887

888

892

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. In *ICLR*.

A Dataset Details

Each of the exams for the dataset described in Section 3 had one ED stay and triage row; 53% had at least one reconciled medicines 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 administered medicines 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.

B Prior Exam Embeddings

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. Note that all exams from MIMIC-CXR were considered for the priors (train/validation/test 222758/1808/3269 exams), including those that did not occur during an ED stay and those that did not have a findings and/or impression section.

C Table Column Determination

The columns from the tables described in Figure 1 were given the following designations:

- 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.
- 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.

- The column designations for the aperiodic vital signs table were identical to the triage table, except for the rhythm column (e.g., 'Normal Sinus Rhythm'), which was treated as a category column. The aperiodic vital signs table also had no chief complaint column and the 'charttime' column was used as the event time.
- For the reconciled medicines table, the 'intime' from the ED stay table was used as the event time, as it pertains to the patient's medicine 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 descriprion of the etccode column.
- For the administered medicines (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 administered medicines column did not have as high of a cardinality as the name column from the reconciled medicines column, allowing it to be considered as a category column.
- For the metadata table, the 'PerformedProcedureStepDescription', 'ViewPosition', 'ProcedureCodeSequence_CodeMeaning', 'View-CodeSequence_CodeMeaning', and 'PatientOrientationCodeSequence_CodeMeaning' columns were considered, and designated as category columns.

D Experiment Setup

D.1 Metrics

GREEN (Ostmeier et al., 2024), CheXbert-F1 (Smit et al., 2020), RadGraph-F1 2022), (Delbrouck al., BLEU-4 (Papet ineni 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 964 et al., 2023; Ostmeier et al., 2024) and were a 965 part of our evaluation. Additionally, we include CXR-BERT (Boecking et al., 2022; Nicolson 967 et al., 2024a), and ROUGE-L (Lin and Hovy, 2003). GREEN, CheXbert-F1, RadGraph-F1, and 969 CXR-BERT were intended to capture the clinical 970 semantic similarity between the generated and 971 radiologist reports, while BERTscore-F1 was intended to capture general semantic similarity. 973 Finally, ROUGE-L and BLEU-4 were intended 974 to capture the syntactic similarity between the 975 generated and radiologist reports. We also propose 976 a new metric that measures *n*-gram repetition rate, 977 namely the absence of repeated n-grams (ARN). It 978 is calculated as:

$$\operatorname{ARN} = \begin{cases} 1.0 & \text{if } L < n, \\ 1.0 - \frac{\sum_{i=1}^{M} (\operatorname{Count}(g_i) - 1)}{M} & \text{if } L \ge n, \end{cases}$$
(1)

where L is the total number of tokens in the generated report, n is the n-gram size, M = L - n + 1 is the total number of n-grams in the report, g_i is the i^{th} unique n-gram in the report, Count (g_i) is the n-gram frequency in the report. The tokenizer described in Appendix D.2 was used with an n-gram size of three.

983

985

987

991

995

996

997

1001

1002

1004

1005

1006

1007

1008

1009

1010

1012

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,1013Games-Howell post hoc tests were used for pair-1014wise testing. A p-value of 0.05 was used for all1015significance testing. Statistical testing was not per-1016formed for CheXbert, as it is a classification metric.1017

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

D.2 Model

Our model is illustrated in Figure 2; following Nicolson et al. (2024b), we utilised UniFormer as the image encoder (in particular, the 384×384 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 (Elgendi 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 for each epoch.

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 reconciled medicines 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.



Figure D.1: Attention mask for the decoder. Non-causal masking was used for the patient data embeddings and causal masking for the report token embeddings.

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.³ This enabled non-causal masking to be utilised for the prompt and causal masking for the report token embeddings, as shown in Figure D.1. This ensured that the self-attention heads were able to attend to all of the patient data embeddings at each position.

D.3 Training

1064

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1081

1082

1083

1084

1085

1086

1089

1090

1091

1092

1093

Three stages of training were performed. Each stage used *AdamW* (Loshchilov and Hutter, 2022) for mini-batch gradient descent optimisation and gradient clipping with a maximum norm of 1.0 to prevent exploding gradients and maintain training stability. Training and evaluation was performed on a 94GB NVIDIA H100 GPU. The three stages were as follows:

1. Teacher forcing (TF) (Williams and Zipser, 1989) was performed on the MIMIC-CXR dataset with only the images for each exam as input, and exams that contained both a findings and impression section. This gave a training/validation split of 232 855/1 837 images, 125 417/991 exams, and 57 102/436 patients. Training was performed with an initial learn-1094 ing rate of 5e-5, a mini-batch size of 8, a maxi-1095 mum of 32 epochs, and with float16 automatic 1096 mixed precision. All model parameters were 1097 trainable during this stage. The validation 1098 macro-averaged CheXbert-F1 was the mon-1099 itored metric for checkpoint selection. This 1100 stage was necessary, as the language model 1101 struggled to generate reports from multiple 1102 patient data sources without prior learning. 1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1133

1134

1135

1136

1137

1138

- 2. TF was used in the second stage of training, with the MIMIC-CXR & MIMIC-IV-ED 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 se-1114 quence training (SCST) (Rennie et al., 2017) 1115 was performed with the rewards described in 1116 Appendix E in the final stage of training. The 1117 sample report for SCST was generated with 1118 top-k sampling (k = 50). Training was per-1119 formed with an initial learning rate of 5e-6, 1120 a mini-batch size of 32, a maximum of 32 1121 epochs, and with float32 precision. A warmup 1122 phase of 5000 training steps was used for 1123 the learning rate, linearly increasing from 1124 zero. The image encoder's parameters were 1125 frozen during this stage (except for its pro-1126 jection). The validation BERTScore-F1 was 1127 the monitored metric for checkpoint selection. 1128 This stage of training was only applied to the 1129 best model from Table 1, 'Images + effective 1130 sources (h = 0)', with the results presented in 1131 Table 2. 1132

D.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).
 1140

³https://huggingface.co/blog/poedator/4d-masks

TranSQ reports were kindly provided by the authors (Kong et al., 2022).

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1180

- RGRG reports were generated following https://github.com/ttanida/rgrg (Tanida et al., 2023).
 - CvT2DistilGPT2 reports were generated following https://github.com/aehrc/ 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/aehrc/cxrmate (Nicolson et al., 2024a).
 - CXRMate-RRG24 reports were generated following https://huggingface.co/aehrc/ 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.

E Reinforcement Learning Rewards

The separate reward per section was calculated as:

$$r_{s}(\hat{\mathbf{w}}_{f}, \mathbf{w}_{f}, \hat{\mathbf{w}}_{i}, \mathbf{w}_{i}) = \alpha_{1} \cdot r_{f}(\hat{\mathbf{w}}_{f}, \mathbf{w}_{f}) + \alpha_{2} \cdot r_{i}(\hat{\mathbf{w}}_{i}, \mathbf{w}_{i}),$$
(2)

where $r_s(\cdot)$ is the composite reward for the sections 1170 of the report, $r_f(\cdot)$ is the reward for the findings 1171 section, and $r_i(\cdot)$ is the reward for the impression 1172 section, $\hat{\mathbf{w}}_f$ is the generated findings section, \mathbf{w}_f is 1173 the radiologist findings section, $\hat{\mathbf{w}}_i$ is the generated 1174 impression section, \mathbf{w}_i is the radiologist impression 1175 section, and α_1 and α_2 are weights. Normally, 1176 $r_r(\hat{\mathbf{w}}_r, \mathbf{w}_r)$ is calculated, where $\hat{\mathbf{w}}_r$ and \mathbf{w}_r are the 1177 generated and radiologist reports, which include 1178 both the findings and impression sections. 1179

The reward $r_f(\cdot)$, $r_i(\cdot)$, or $r_r(\cdot)$ is calculated as:

$$r(\hat{\mathbf{w}}, \mathbf{w}) = \lambda_1 \cdot \text{CXR-BERT}(\hat{\mathbf{w}}, \mathbf{w}) + \lambda_2 \cdot \text{BERTScore}(\hat{\mathbf{w}}, \mathbf{w}) + \lambda_3 \cdot \text{ARN}(\hat{\mathbf{w}}, \mathbf{w}),$$
(3)

where λ_1 , λ_2 , and λ_3 are weights. For 'Images + 1182 effective source (h = 0) + RL with CXR-BERT 1183 + BERTScore reward', $\lambda_1 = 0.5$, $\lambda_2 = 0.5$, and 1184 $\lambda_3 = 0.0$. For 'Images + effective source (h =1185 0) + RL with CXR-BERT + BERTScore reward 1186 per section', $\alpha_1 = 0.75$, $\alpha_2 = 0.25$, $\lambda_1 = 0.5$, 1187 $\lambda_2 = 0.5$, and $\lambda_3 = 0.0$. A higher weight was 1188 used for the findings section, as it is longer on 1189 average than the impression section. For 'Images 1190 + effective source (h = 0) + RL with CXR-BERT 1191 + BERTScore + ARN reward per section', $\alpha_1 =$ 1192 0.75, $\alpha_2 = 0.25$, $\lambda_1 = 0.45$, $\lambda_2 = 0.45$, and 1193 $\lambda_3 = 0.1$. Only a weak contribution of the ARN 1194 was required to prevent repetitions. 1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1223

1224

1225

1226

1227

1228

1229

1231

The improvement that separating the reward per section has on the findings section is negligible, as seen in Table 2. However, separating the reward per section improves the scores for the impression section, as shown in Table E.1. Separating the reward likely enables the model to better optimise for the concise and summarised nature of the impression section, which was previously overshadowed by the dominance of the findings section's requirement for comprehensive detail when both were jointly considered.

F Ancillary Results

In Figure E.2, the F1-scores for each CheXbert label are shown. The 'Images + effective sources (h = 0)' model from Table 1 attained a higher score than the 'Images' model for 11 of the 14 labels. This suggests that incorporating auxiliary patient data from MIMIC-IV-ED and MIMIC-CXR provides a general improvement, rather than benefiting any specific pathology.

Further improvements can be seen for most labels when reinforcement learning (RL) is used (i.e., our model from Table 2). However, there are performance decreases for 'enlarged cardiomediastinum', 'pneumothorax', and 'fracture'. This might be due to these pathologies being underrepresented in the MIMIC-CXR dataset, leading the model to optimise for more common pathologies during reinforcement learning.

The results for exams that include an aperiodic vital signs table are show in Table F.1. Adding it produced a significant improvement in the scores for CXR-BERT, indicating that it should be considered if available. The results for exams that include an administered medicines table are show in Table F.2. Adding did not produce a significant improve-

Table E.1: Impact of the reward on the impression section of the test set described in Section 3 (n = 9580 exams; 958×10 runs for 'Images + effective sources (h = 0)', n = 1916 exams; 958×3 runs for the remaining models). Evaluation is on the **impression** section only.

Model	RG	CX	СВ	G	BS	R-L	B4	ARN
Images + effective sources $(h = 0)$	20.21	26.81	57.61	28.71	27.90	25.02	4.77	99.59
+ RL (CXR-BERT + BERTScore reward)	23.96	28.07	62.85	30.58	31.58	28.48	7.84	99.89
+ reward per section	24.89	31.08	71.12	30.89	36.27	30.27	6.70	99.33
+ ARN reward	24.87	32.88	71.12	32.14	36.31	30.61	6.84	99.83

ment in the scores, indicating that it is not usefulfor CXR report generation.

G Error analysis

1234

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1264

1265

1266

1268

1269

1270

1271

G.1 Impact of Auxiliary Patient Data on the CheXpert Labels

Figure E.1 demonstrates the impact of incorporating auxiliary patient data for different CheXpert labels. Here, the GREEN score for the 'Images + effective sources (h=0)' model is compared to the 'Images' model from Table 1 for each exam. Note that the generated and radiologist report for each exam will often include findings other than the CheXpert label. Hence, the GREEN scores do not exclusively represent a particular CheXpert label, rather, they represent exams with that label present. The horizontal dashed line where $\Delta = 0$ divides exams where auxiliary patient data improved performance from those where it decreased performance. CheXpert labels with a higher area under the curve (AUC) above the horizontal dashed line suggest that there is a stronger overall benefit from leveraging auxiliary patient data.

Leveraging auxiliary patient data yielded a higher AUC for 10 out of the 14 CheXpert labels, indicating that it is beneficial for many pathologies. For certain CheXpert labels, the influence of auxiliary patient data is less clear, particularly for those associated with smaller sample sizes, such as *enlarged cardiomediastinum* (n = 10), *consolidation* (n = 10), *fracture* (n = 15), *pneumothorax* (n = 5), and *lung lesion* (n = 35). The *no findings* AUC of 6.85 for $\Delta > 0$ being lower than the AUC of 7.72 for $\Delta < 0$ suggests that the auxiliary patient data increases the false positive rate for this model.

G.2 Impact of Auxiliary Patient Data on the Generated Reports

To gain a better understanding of how the auxiliary patient data impacts the generated reports, we analyse multiple case studies where it contributes to either true positive, false positive, true negative, or false negative findings in the generated report: 1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1299

1300

1301

1302

1303

1304

1305

1306

1307

- A true positive is where the model has identified a positive occurrence of a pathology that is also identified as positive in the radiologist's report.
- A false positive is where the model has incorrectly identified a positive occurrence of a pathology that is not identified as positive in the radiologist's report.
- A true negative occurs when a pathology is omitted or absent in the radiologist's report and this is correctly reflected in the generated report, either implicitly through omission or explicitly by stating its absence.
- A false negative is where a pathology is positively identified in the radiologist's report but is not positively identified in the generated report.

Exams with a high Δ from Figure E.1 were selected for true positive and true negative examples, while those with a low Δ were chosen for false positive and false negative examples.⁴ This analysis, though based on only eight exams, exemplifies how auxiliary patient data can both enhance and hinder the CXR report generation process, providing valuable insights into its impact. A more comprehensive analysis would be required to fully characterise the influence of auxiliary patient data across diverse exams and pathologies.

G.2.1 True Positive: Example 1

Table G.1 demonstrates how auxiliary patient data contributed to the true positive detection of increased interstitial markings, which are suggestive of pulmonary fibrosis. The model not using auxiliary patient data failed to detect the interstitial

⁴Out of the 10 training runs, the 'Images + effective sources (h = 0)' and 'Images' models that attained the highest average GREEN score over the test set were selected for the error analysis.

Table F.1: Results for exams that have an aperiodic vital sign table (n = 5250; studies 525×10 runs). Underlined scores indicate a significant difference to the scores of 'Images' (p < 0.05).

Model	RG	CX	СВ	G	BS	R-L	B4
Images	24.73	29.41	58.63	35.11	24.33	25.85	4.89
Images + vital signs	24.55	29.73	<u>60.32</u>	35.21	24.17	25.97	4.87

Table F.2: Results for exams that have a administered medicines table (n = 3520; studies 352×10 runs). Underlined scores indicate a significant difference to the scores of 'Images' (p < 0.05).

Model	RG	CX	CB	G	BS	R-L	B4
Images	25.19	28.29	59.24	36.13	24.81	26.61	5.15
Images + administered medicines	24.70	29.53	59.53	35.82	24.46	26.38	4.85

markings. The patient's triage data included a respiratory rate consistent with tachypnoea and a chief 1309 complaint of dyspnoea, both common features of 1310 pulmonary fibrosis. Additionally, the patient's his-1311 tory of pulmonary fibrosis and worsening shortness 1312 of breath provided further context supporting the observed increase in interstitial markings. In this 1314 case, the inclusion of auxiliary patient data facili-1315 tated a true positive detection. 1316

G.2.2 True Positive: Example 2

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

Table G.2 demonstrates how auxiliary patient data contributed to the true positive detection of pulmonary edema, which was not detected by the model that does not use auxiliary patient data. Recorded in the patient's triage data was a respiratory rate consistent with tachypnoea and a chief complaint of dyspnoea (also documented in the history section), both of which are indicative of pulmonary edema. Additionally, furosemide was listed in the patient's reconciled medicines, which is commonly used to manage pulmonary edema. This example underscores how incorporating auxiliary patient data can enhance true positive detection in CXR report generation.

G.2.3 False Positive: Example 1

Table G.3 provides an example of where the model 1333 leveraging auxiliary patient data introduced a false 1334 positive prediction into the generated report. It incorrectly specifies that there are streaky opacities 1336 in the lung bases, which are reflective of atelectasis. 1337 The model that does not leverage auxiliary patient data did not produce this false positive. Atelectasis 1339 1340 typically presents with symptoms such as dyspnoea, tachypnoea, wheezing, and coughing; how-1341 ever, these symptoms were absent from the indica-1342 tion section or the triage data. Although codeine, 1343 listed among the patient's reconciled medicines, 1344

can contribute to atelectasis in high doses, there 1345 was no evidence of overdose or misuse in this case. 1346 This example suggests that weak or ambiguous ev-1347 idence in the auxiliary data may have influenced 1348 the false positive prediction. Further refinement is 1349 needed to improve the model's ability to appropri-1350 ately weigh auxiliary patient data evidence against 1351 imaging evidence. 1352

1353

G.2.4 False Positive: Example 2

Table G.4 presents a case where the model using 1354 auxiliary patient data generated false positive pre-1355 dictions, identifying mild pulmonary vascular con-1356 gestion and a mildly enlarged cardiac silhouette 1357 (cardiomegaly). In contrast, the model without 1358 auxiliary patient data did not produce these errors. 1359 Shortness of breath, noted in the indication sec-1360 tion and the chief complaint from the triage data, 1361 is a common symptom of both mild pulmonary 1362 vascular congestion and cardiomegaly. The ele-1363 vated respiratory rate recorded in the triage data, 1364 consistent with tachypnoea, may suggest mild pul-1365 monary vascular congestion, while the elevated 1366 systolic blood pressure reflects isolated systolic hy-1367 pertension, a known risk factor for cardiomegaly. 1368 Furosemide, included in the reconciled medicines, 1369 can help manage mild pulmonary vascular congestion associated with fluid overload and conditions 1371 like cardiomegaly. Lisinopril and diltiazem pri-1372 marily treat hypertension, which is a risk factor 1373 for cardiomegaly. This example indicates that there 1374 was evidence in the auxiliary patient data that could 1375 have led the model to multiple false positive pre-1376 dictions. For this example, the model lacked the 1377 ability to correctly balance auxiliary patient data 1378 evidence with imaging evidence. 1379 Table G.1: True positive example for study 51707133. Here, the triage data and the history section provide additional evidence supporting increased interstitial markings. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

				P	atient da	ıta					
Image											
History	year-old fo	emale with	pulmonary	fibrosis	and CHF	F with w	vorsenii	ng shortn	ess of breath.		
Reconciled medicines; name	atorvastatin, az succinate, mul pantoprazole, l	elastine [As tivitamin, gl lidocaine, ko	telin], aspin lucosamine etotifen fur	rin, calciu e sulfate narate.	um carbo [Glucosa	onate-vi mine],	tamin E acetam	93 [Calciu inophen,	ım 500 + D], lorata ferrous sulfate [Fe	adine, metoprolol cosol], torsemide,	
Triago	temperature	heartrate	resprate	o2sat	sbp	dbp	pain	acuity	chiefcomplaint		
mage	99.7	90.0	36.0	100.0	118.0	70.0	0	2.0	Dyspnea		
Findings	AP and lateral Diffusely increase when compare detected.	l views of t eased inters ed to prior.	he chest. titial mark Cardiomed	R Low lun ings are liastinal	adiologi g volum seen thro silhouett	<i>est</i> nes are oughou te is gro	seen co t the lui ossly ur	ompatible ngs, but t ichanged	e with patient's hi hese appear overa . No acute osseou	story of fibrosis. Ill slightly worse is abnormality is	
Impression	Findings comp be excluded an	batible with	pulmonary orrelation i	fibrosis s necessa	with like ary.	ely supe	erimpos	ed edema	a. Please note that	infection cannot	
Findings	Frontal and lat more extensive cardiomediasti at the shoulder	teral views e on the left nal silhouet	of the ches of than on te is stable	t. Low little right the right. No acut	ung volu . There te osseou	imes ar is no e is abnor	e again vidence rmality	= 0.373) noted. In e of over is identif	ncreased interstitia t consolidation no ied. Degenerative	al markings seen or effusion. The changes are seen	
Impression	Increased interstitial markings throughout the lungs which could be due to chronic lung disease and possible chronic lung disease. No definite superimposed acute process, although clinical correlation suggested.										
Findings	Assessment is enlarged. The pulmonary ede or pneumothor	limited due aorta appear ema. Streaky eax is seen.	to patient re rs to be calc y opacities Multilevel	otation and cified. Per in the lun degenera	nd patien erihilar h ng bases itive char	aziness likely r nges are	on. Lun and vas reflect a e noted	g volume scular ind reas of at in the tho	s are low. Heart si listinctness is com electasis. No large pracic spine.	ze appears mildly patible with mild e pleural effusion	
Impression	Mild pulmonat	ry edema an	d bibasilar	atelectas	sis.						

1380

1394

G.2.5 True Negative: Example 1

Table G.5 shows an exam where the model using auxiliary patient data generated a report with true negatives, despite confounding evidence from the auxiliary patient data. The renal failure and upper quadrant pain mentioned in the history section could suggest a condition related to fluid overload, such as pleural effusion. Furosemide and metolazone mentioned in the reconciled medicines are commonly used for fluid management and treating pulmonary edema. Lisinopril and amlodipine, primarily used for cardiovascular conditions such as hypertension, can lead to secondary effects like pulmonary congestion or cardiomegaly, which may be detected radiologically. Despite these confounding factors, the model effectively prioritised the imaging evidence, avoiding false positive predictions. This demonstrates that the model possesses the ability to balance auxiliary patient data evidence with imaging evidence.

1395

1396

1397

1398

1399

1400

G.2.6 True Negative: Example 2

Table G.6 is another exam where the model using 1401 auxiliary patient data generated a report with true 1402 negatives, despite confounding evidence from the 1403 auxiliary patient data. The model utilising auxil-1404 iary patient data accurately identified a dual-lead 1405 pacemaker without introducing any false positive 1406 findings, despite the presence of confounding evi-1407 dence from the auxiliary patient data. The indica-1408 tion section requests evaluation for fluid overload 1409 Table G.2: True positive example for study 52841174. Here, the triage data and reconciled medicines provide additional evidence indicative of pulmonary edema. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

				P	atient da	ta					
Image											
History	year-old w	vith dyspnea	l .								
Reconciled medicines; name	Coumadin, furosemide, metoprolol succinate, Calcarb 600 With Vitamin D, simvastatin, Tylenol Extra Strength, levothyroxine, docusate sodium.										
Triage	temperature	heartrate	resprate	o2sat	sbp	dbp	pain	acuity	chiefcomplaint		
mage	97.0	81.0	22.0	100.0	102.0	58.0	0	2.0	DYSPNEA		
Findings	AP upright and rotated somew appears grossl been no signifi	d lateral vie hat limiting y stable. The cant change	ws of the o the evalua ere are sma from prior	R chest we tion of th Ill layerin r study. H	<i>adiologi</i> re provid le cardio lg bilater Bony stru	<i>st</i> led. Mi mediast al effus ictures	dline s inal sill ions wi are inta	ternotom houette, t th mild in ct.	y wires are again hough cardiomedi nterstitial edema. (noted. Patient is astinal silhouette Overall, there has	
Impression	Mild interstitia	ıl edema, sta	ble cardio	megaly v	vith sma	l bilate	ral effu	sions.			
Findings	AP upright and as a prosthetic edema. Small I structures appe	<i>Ima</i> d lateral view cardiac valv bilateral pleu ear intact. N	ges + effect ws of the cl we. Low lu ral effusio o free air b	<i>ctive sour</i> hest were ng volun ns persist pelow the	rces (h = constructions (h = construction)) = construction (here) = construction (here	= 0) (Gl ed. Mid evaluat s left ba midiapl	REEN line ste tion. Th asilar at hragm.	= 0.375) rnotomy here is hil electasis.	wires and mediast lar congestion and The heart is mildl	inal clips as well mild pulmonary y enlarged. Bony	
Impression	Pulmonary ede	ema, small b	ilateral ple	ural effu Images (sions, le GREEN	ft great $= 0.22$	er than 2)	right.			
Findings	Images (GREEN = 0.222) The patient is status post median sternotomy and CABG. Large hiatal hernia is present. The cardiac silhouette size is mildly enlarged. The aorta is tortuous. Crowding of bronchovascular structures is present with probable mild pulmonary vascular congestion. Small right pleural effusion is present. Patchy opacities in the lung bases may reflect atelectasis. No pneumothorax is demonstrated. There are moderate multilevel degenerative changes seen in the thoracic spine.										
Impression	1. Small right be completely	pleural effus excluded. 2	sion and bi . Mild pulı	basilar o nonary v	pacities ascular o	likely re congest	eflect at ion. 3.	telectasis Moderate	. Infection at the le e cardiomegaly.	ung bases cannot	

1410

1411

- 1417 1418
- 1419

1420

1421

1422

1423

1424

or pneumonia and notes chest pain, which could lead to false positives such as pulmonary edema, pneumonia, pleural effusion, or cardiomegaly. The reconciled medicines, including furosemide and nitroglycerin, suggest the management of conditions such as pulmonary edema or heart failure, which could be associated with pleural effusion or cardiomegaly. Despite these confounding factors, the model effectively prioritised the evidence from the image, avoiding false positive predictions.

G.2.7 False Negative: Example 1

Table G.7 is an example where the model failed to leverage auxiliary patient data to detect trace bilateral pleural effusions and the increased opacity in the right mid-to-lower lung (concerning for pneumonia). The history section notes dyspnea and hy-1425 poxia, which are symptoms associated with pleural 1426 effusion and pneumonia, among other conditions. 1427 It also requests to assess for fluid overload or pneu-1428 monia, both of which should prompt the model 1429 to assess for pleural effusion and opacities. The 1430 significantly reduced oxygen saturation recorded 1431 in the triage data indicates severe hypoxia (also 1432 noted in the history section), which can be caused 1433 by pleural effusion or pneumonia. Despite strong 1434 evidence from the auxiliary patient data to support 1435 pleural effusion and the opacity, the model failed 1436 to combine this with the imaging evidence to make 1437 the correct predictions. 1438

Table G.3: False positive example for study 51274564. This example demonstrates how weak auxiliary patient data evidence may have misled the model. Only the patient data that Images + effective sources (h=0) utilises is shown.

				P	atient a	lata						
Image	PORTABLE											
Indication	Status post new central line placement.											
Reconciled medicines; name	colchicine, Aspirin, nifedipine, blood sugar diagnostic [OneTouch Ultra Test], labetalol, calcitriol, insulin needles (disposable) [BD Insulin Pen Needle UF Mini], fluticasone, codeine-guaifenesin, lisinopril, insulin lispro [Humalog KwikPen], insulin glargine [Lantus Solostar], prednisone, acetaminophen, torsemide, albuterol sulfate [ProAir HFA], mycophenolate mofetil, Multivitamin, tacrolimus, Vitamin E, allopurinol, ferrous sulfate.											
	temperature	heartrate	resprate	o2sat	sbp	dbp	pain	acuity	chiefcomplaint			
Triage	98.1	72.0	16.0	0.0	95.0	46.0	8	2.0	Abnormal labs, Weakness, Diar- rhea			
Findings	A new central there has been	venous cath no significa	eter termin int short-ter	H ates in t rm chan	<i>adiolo</i> he left l ge.	g <i>ist</i> orachio	cephalio	c vein. T	here is no pneumothorax. Otherwise,			
Impression	Status post pla	cement of n	ew left inte	ernal jug	ular cei	ntral ve	nous ca	theter; no	o pneumothorax identified.			
Findings	There is interv SVC. Lung vo Mediastinal an pulmonary ede effusion or pre	al placeme lumes are lo d hilar contr ma. Streak	<i>ges</i> + <i>effec</i> nt of a left ow. This ac ours are un y opacities is demonst	internal contuate changed are note trated. N	<i>rces (h</i> l jugula es the s l. There d in the lild deg	= 0) (0 or centrative of the centrative of t	GREEN al veno he cardi vding of bases, lil ve chan	0 = 0.14 us cathet iac silhou f the bror kely refle ges are n	3) ter with tip terminating in the lower uette which appears mildly enlarged. inchovascular structures without overt ective of atelectasis. No large pleural loted in the thoracic spine.			
Impression	Interval placement of a left internal jugular central venous catheter with tip in the lower SVC. Low lung volumes with streaky bibasilar opacities, likely atelectasis.											
Findings	Images (GREEN $= 0.25$) A PICC line terminates in the mid-to-lower SVC. The cardiomediastinal and hilar contours are within normal limits. The lung fields are clear. There is no pneumothorax, fracture or dislocation. Limited assessment of the abdomen is unremarkable.											
Impression	Left PICC tern	ninates in th	e mid-to-lo	ower SV	C.							

G.2.8 False Negative: Example 2

1439

Table G.8 highlights a false negative example for 1440 the model leveraging auxiliary patient data, where 1441 it failed to identify the right lower lobe opacity 1442 concerning for pneumonia. Despite omitting this 1443 finding, the model received strong evidence sup-1444 porting its presence. Specifically, the indication 1445 1446 section notes a right lower lobe infiltrate, directly pointing to an opacity, alongside dyspnoea, a non-1447 specific symptom of pneumonia. Additionally, the 1448 chief complaint explicitly lists pneumonia, another 1449 strong indicator. The triage data, including the nor-1450 1451 mal temperature and heart rate, might have influenced the model's decision by suggesting a lack of 1452 systemic immune response, which could reduce the 1453 likelihood of pneumonia. However, the reconciled 1454 medicines, including antibiotics like erythromycin 1455

and tobramycin-dexamethasone, support the pos-1456 sibility of an active infection. Despite the align-1457 ment between the auxiliary patient data and the 1458 suspected pneumonia, the model failed to integrate 1459 this evidence with the imaging evidence to make a 1460 correct prediction. This underscores the need for 1461 further model development to better synthesise the 1462 evidence from auxiliary patient data with imaging 1463 evidence. 1464 Table G.4: False positive example for study 54082940. This example demonstrates how the model failed to balance auxiliary patient data evidence with imaging evidence. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

	Patient data										
Image											
Indication	Shortness of breath and wheezing, previously diagnosed with pneumonia or infectious process.										
Reconciled medicines; name	prednisolone acetate, albuterol sulfate [ProAir HFA], gabapentin, Humulin 70/30, cholecalciferol (vitamin D3), sennosides [senna], furosemide, Trusopt, lisinopril, AERO CHAMBER, levobunolol, insulin aspart, insulin aspart [Novolog], fluticasone-salmeterol [Advair Diskus], latanoprost, dorzolamide [Trusopt], aspirin [Enteric Coated Aspirin], diltiazem HCl [DILT-XR], blood sugar diagnostic [FreeStyle Lite Strips], magnesium hydroxide [Milk of Magnesia], Acetaminophen, lidocaine, docusate sodium, brimonidine, calcium carbonate, nebulizer and compressor, olanzapine [Zyprexa].										
	temperature heartrate resprate o2sat sbp dbp pain acuity chiefcomplaint										
Triage	98.0 81.0 24.0 100.0 151.0 66.0 0 2.0 SHORTNESS OF BREATH										
Findings	<i>Radiologist</i> There is no evidence of focal consolidation. There is left lower lobe atelectasis. There is no pleural effusion or pneumothorax. The cardiac and mediastinal contours are normal.										
Impression	No acute cardiopulmonary process.										
Findings	Images + effective sources $(h = 0)$ (GREEN = 0.429) There is mild pulmonary vascular congestion. No definite focal consolidation is seen. No pleural effusion or pneumothorax is seen. Cardiac silhouette is mildly enlarged. The cardiac and mediastinal silhouettes are grossly stable with the cardiac silhouette possibly slightly enlarged compared to prior.										
Impression	Mild pulmonary vascular congestion. Cardiomegaly.										
Findings	Images (GREEN = 0.8) There is no confluent consolidation. No pulmonary edema or pleural effusions are identified. Cardiomediastinal and hilar contours are within normal limits. No pneumothorax is evident.										
Impression	No acute cardiopulmonary process.										



Studies sorted by Δ

Figure E.1: The mean pairwise difference GREEN score for the generated report (findings and impression sections) of each exam from the test set between 10 training runs of the "Images" model and the "Images + effective sources (h=0)" model. This illustrates the performance change (increase or decrease) over the exams resulting from incorporating auxiliary patient data for different CheXpert labels. Δ , m and n are the number of training runs for each model (m = n = 10) and s is the GREEN score for one of the models. The subplots are sorted in descending order based on the ratio of AUC($\Delta > 0$) to AUC($\Delta < 0$).



Figure E.2: F1-score for each CheXbert label. (n = 9580 exams; 958×10 runs for 'Images' and 'Images + effective sources (h = 0)' and n = 2874 exams; 958×3 runs for 'Images + effective sources (h = 0) + RL with CXR-BERT + BERTScore + ARN reward per section'.

Table G.5: True negative example for study 52428322. This demonstrates how the model can avoid false positives despite confounding evidence from the auxiliary patient data. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

				I	Patient	data					
Image		L.									
History	year-old fe	emale with	renal failur	e and rig	ght upp	er qua	drant pa	iin. Hypo	tension.		
Reconciled medicines; name	aspirin, FreeSty perazine malea TTS-1], Humu ferrous sulfate,	aspirin, FreeStyle Lite Strips, metolazone, furosemide, omeprazole, oxycodone-acetaminophen [Endocet], prochlor- perazine maleate, calcitriol, fentanyl, insulin glargine [Lantus], sevelamer HCl [Renagel], clonidine [Catapres- TTS-1], Humulin R, Colace, insulin lispro [Humalog], potassium chloride [Klor-Con M20], FreeStyle Lite Meter, ferrous sulfate, lisinopril, BD Insulin Syringe Ultra-Fine, Glucose Meter, Disp & Strips, Lipitor, amlodipine.									
Triago	temperature	heartrate	resprate	o2sat	sbp	dbp	pain	acuity	chiefcomplaint		
mage	0.0	0.0	0.0	0.0		0.0	None	1.0	GI BLEED		
Findings	Single portable silhouette is no seen below the	view of the ormal. Ossee diaphragm	e chest is co ous and sof	ompared t tissue	to prev structur	<i>ogist</i> vious e res are	exam fro unrema	om T rkable. N	he lungs are clear o visualized free	Cardiomediastinal intraperitoneal air is	
Impression	No acute cardi	opulmonary	process.								
Findings	Images + effective sources $(h = 0)$ (GREEN = 1.0) PA and lateral views of the chest were provided. The lungs are clear bilaterally without focal consolidation, effusion, or pneumothorax. The heart and mediastinal contours are normal. The imaged osseous structures are intact. There is no free air below the right hemidiaphragm.										
Impression	No acute findir	ngs in the cl	nest.	_							
Findings	Images (GREEN = 0.5) Single portable frontal chest radiograph demonstrates unremarkable cardiomediastinal and hilar contours. Lungs are clear. No pleural effusion or pneumothorax evident.										
Impression	No acute intrat	horacic pro	cess.								

Table G.6: True negative example for study 52169517. This demonstrates how the model can avoid false positives despite confounding evidence from the auxiliary patient data. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

Patient data										
Image										
Indication	year-old woman with chest pain. Evaluate for fluid overload or pneumonia.									
Reconciled medica- tions; name	Humalog, atorvastatin, aspirin, gabapentin, nitroglycerin, methylprednisolone, valsartan [Diovan], insulin glargine [Lantus], One Touch Ultra Test, metoprolol tartrate, isosorbide mononitrate, cephalexin, colchicine [Colcrys], furosemide									
Triage	temperature	heartrate	resprate	o2sat	sbp	dbp	pain	acuity	chiefcomplaint	
	0.0	0.0	0.0	0.0 <i>I</i>	0.0 Radiolo	0.0 ogist	None	2.0	CHEST PAIN	
Findings	The lungs are clear. The hilar and cardiomediastinal contours are normal. There is no pneumothorax or pleural effusion. Pulmonary vascularity is normal. A dual-lead pacemaker is present.									
Impression	No acute cardiopulmonary process.									
Findings	Images + effective sources $(h = 0)$ (GREEN = 0.833) The lunge are clear without feed consolidation. No plaural offician or product for a second the condition of the second sec								n The cordiac and	
rindings	mediastinal silhouettes are unremarkable. Dual lead left-sided pacemaker is stable in position. Patient is status post median sternotomy and CABG.									
Impression	No acute cardiopulmonary process.									
T ' 1'	DA 11/ 1	. 6.4	1 4	Images	(GREE	EN = 0).25)	4 11	1	1 1 4 1 4
Findings	PA and lateral views of the chest provided demonstrate a left chest wall pacer device with leads extending to the region of the right atrium and right ventricle. Midline sternotomy wires and mediastinal clips are again noted. Mild left basal atelectasis is noted. No convincing signs of pneumonia, effusion, or pneumothorax. The cardiomediastinal silhouette is stable. Bony structures are intact. No free air below the right hemidiaphragm.									
Impression	No acute findings in the chest.									

Table G.7: False negative example for study 55715754. The model failed to identify the pleural effusions despite evidence from the auxiliary patient data. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

Patient data										
Image										
History	Dyspnea and hypoxia, assess for fluid overload or pneumonia.									
Triage	temperature	heartrate	resprate	o2sat	sbp	dbp	pain	acuity	chiefcomplaint	
	96.4	83.0	20.0	76.0	145.0	70.0	10	1.0	SORE THROAT	
Findings	Semi-upright portable AP view of the chest provided. The heart is massively enlarged. There are trace pleural effusions. Increased opacity in the right mid-to-lower lung is concerning for pneumonia. The left lung appears essentially clear. No pneumothorax. The mediastinal contour appears normal. Bony structures are intact.									
Impression	Massive cardiomegaly with trace bilateral pleural effusions. Opacity within the right mid-to-lower lung is concerning for pneumonia.									
	Images + effective sources $(h = 0)$ (GREEN = 0.2)									
Findings	Single portable radiograph of the chest demonstrates moderate enlargement of the cardiac silhouette, not sig- nificantly changed compared to the prior examination. There is mild pulmonary vascular congestion. No focal consolidation, pleural effusion or pneumothorax is seen. The visualized upper abdomen is unremarkable.									
Impression	Persistent enlargement of the cardiac silhouette, not significantly changed compared to Unchanged mild pulmonary vascular congestion and stable enlargement of the cardiac silhouette.									
Findings	Images (GREEN = 0.333) There is moderate enlargement of the cardiac silhouette. The aorta is unfolded. Mediastinal and hilar contours are otherwise unremarkable. Pulmonary vasculature is not engorged. Hazy opacity in the right lung is compatible with pneumonia. Right midlung linear opacity may be due to atelectasis. No pleural effusion or pneumothorax is identified. No acute osseous abnormalities seen.									
Impression	1. Moderate enlargement of the cardiac silhouette, compatible with pneumonia. 2. Moderate enlargement of the cardiac silhouette. 3. Right lung base opacity, likely scarring. No definite evidence of pneumonia.									

Table G.8: False negative example for study 53964812. Despite strong evidence from the auxiliary patient data supporting pleural effusion, the model failed to detect it. Only the patient data that *Images* + *effective sources* (h=0) utilises is shown.

Patient data									
Image									
History	Asthma, coronary disease, myocardial infarction, diabetes, presenting with dyspnea and right lower lobe infiltrate.								
Reconciled medicines; name	metformin, acetaminophen, erythromycin, fluticasone-salmeterol [Advair Diskus], Boost Diabetic, bupropion HCl, diltiazem HCl, albuterol sulfate, losartan [Cozaar], lorazepam, sennosides [senna], gabapentin, aspirin, tiotropium bromide [Spiriva with HandiHaler], tobramycin-dexamethasone [TobraDex], atorvastatin, tramadol, nut.tx.gluc.intol,lac-free,soy [Glucerna], clopidogrel, ReliOn Ultra Thin Plus Lancets, melatonin, lidocaine, cal- cium carbonate-vitamin D3 [Calcium 600 + D(3)], furosemide, compression socks, medium, sertraline, L.acidoph & sali-B.bif-S.therm [Acidophilus], travoprost [Travatan Z], montelukast [Singulair], Multivitamin with Iron- Mineral, fluticasone.								
Triage	temperature heartrate resprate o2sat sbp dbp pain acuity chiefcomplaint								
	98.2 83.0 18.0 98.0 162.0 100.0 0 2.0 PNEUMONIA								
Findings	Radiologist Tohe cardiac, mediastinal and hilar contours appear stable. There is no pleural effusion or pneumothorax. Since the very recent prior studies, there is a substantial new opacity in the right lower lobe concerning for pneumonia. The bones appear demineralized. There is mild-to-moderate rightward convex curvature again centered along the lower thoracic spine with incompletely characterized lumbar compression deformities. Moderate degenerative changes are again noted along lower thoracic levels.								
Impression	Findings consistent with pneumonia in the right lower lobe. Depending on clinical circumstances, the possibility of aspiration could also be considered.								
Findings	Images + effective sources $(h = 0)$ (GREEN = 0.0) Frontal and lateral views of the chest. Right apical scarring is again seen. The lungs are otherwise clear without consolidation or effusion. Mild cardiomegaly is again noted. Slightly tortuous descending thoracic aorta is similar to prior. No acute osseous abnormality is identified.								
Impression	No acute cardiopulmonary process.								
Findings	Images (GREEN = 0.333) There is bibasilar atelectasis without definite focal consolidation. No pleural effusion or pneumothorax is seen. The cardiac and mediastinal silhouettes are stable. Mild loss of height anteriorly of a lower thoracic vertebral body is unchanged. Evidence of DISH is seen along the spine.								
Impression	No acute cardiopulmonary process. No significant interval change.								