CREPE: Rapid Chest X-ray Report Evaluation by Predicting Multi-category Error Counts

Anonymous ACL submission

Abstract

We introduce CREPE (Rapid Chest Xray Report Evaluation by Predicting Multicategory Error Counts), a rapid, interpretable, and clinically grounded metric for automated chest X-ray report generation. CREPE uses a domain-specific BERT model fine-tuned with a multi-head regression architecture to predict error counts across six clinically meaningful categories. Trained on a large-scale synthetic dataset of 32,000 annotated report pairs, CREPE demonstrates strong generalization and interpretability. On the expert-annotated ReX-Val dataset, CREPE achieves a Kendall's τ cor-015 relation of 0.786 with radiologist error counts, outperforming traditional and recent metrics. CREPE achieves these results with an inference speed approximately 280 times faster than large language model (LLM)-based approaches, enabling rapid and fine-grained evaluation for scalable development of chest X-ray report generation models.

1 Introduction

011

012

017

037

041

The automatic generation of radiology reports from chest X-ray images has seen rapid progress with the advent of advanced generative AI technologies. Such systems hold substantial promise to reduce radiologists' workload and enhance clinical workflow efficiency by enabling the automated production of clinically meaningful reports at scale (Chen et al., 2024; Zambrano Chaves et al., 2025).

A central challenge that remains is the reliable evaluation of these generated reports. Current metrics for evaluating chest X-ray report generation fall into four main categories: overlapbased, embedding-based, named entity recognition (NER)-based, and large language model (LLM)based approaches. Overlap-based metrics, like BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004) measure surface-level lexical similarity but miss clinical meaning. Embedding-based metrics,

such as BERTScore (Zhang et al., 2020) and SembScore (Smit et al., 2020) capture semantic alignment but not factual or relational accuracy. NERbased metrics, like F1 RadGraph (Jain et al., 2021; Delbrouck et al., 2024) and RaTEScore (Zhao et al., 2024) focus on extracted medical entities and relations, but their reliability depends on the underlying NER system and they can miss nuanced contextual errors. LLM-based metrics, such as GREEN (Ostmeier et al., 2024) and FineRadScore (Huang et al., 2024) achieve strong alignment with human judgment, yet are limited by high computational cost and slow inference. These trade-offs motivate the need for a fast, clinically appropriate, and robust evaluation metric.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

081

In this work, we introduce CREPE (Rapid Chest X-ray Report Evaluation by Predicting Multicategory Error Counts), a novel evaluation metric designed for both rapid inference and clinically interpretable assessment. CREPE leverages a medical domain-specific BERT model with multiple regression heads to predict continuous error counts across six clinically meaningful categories. By explicitly modeling these error categories, CREPE outputs both a fast, interpretable overall score and a detailed error breakdown, thus providing actionable, clinically relevant feedback beyond conventional summary metrics.

Our primary contributions are as follows:

- Fast and Fine-Grained Evaluation: We propose CREPE, a rapid and fine-grained evaluation metric that predicts clinically interpretable multi-category error counts via regression, achieving performance comparable to or exceeding state-of-the-art LLM-based metrics in alignment with radiologist assessments.
- Large-Scale Synthetic Training Data: We construct a large-scale synthetic dataset comprising 32,000 report pairs with detailed,



Figure 1: **Overview of CREPE and Comparative Performance.** (Left) Schematic illustration of CREPE's evaluation pipeline: a ground truth and generated chest X-ray report are concatenated and processed by the CREPE model, which predicts continuous error counts across six clinically defined categories; the total error count is used as the final CREPE score. (Right) Comparative analysis of evaluation metrics on the ReXVal dataset, plotting Kendall's τ correlation with radiologist error counts (x-axis) against average computation time per report pair (y-axis, log scale). CREPE achieves superior correlation with expert assessments while maintaining low computational cost, enabling fast and clinically aligned report evaluation.

category-specific error annotations, automatically generated using LLMs under strict data governance.

- **Robust Generalization**: We demonstrate strong generalization and robustness of CREPE across diverse public evaluation benchmarks, including challenging out-of-distribution datasets.
- Efficiency and Scalability: CREPE delivers substantial reductions in computational cost and inference time compared to LLM-based evaluation methods, enabling practical and scalable assessment for both research and realworld development pipelines.

An overview of the CREPE pipeline and its comparative performance is illustrated in Figure 1.

2 Related Work

2.1 General Text Evaluation Metrics

Traditional natural language generation (NLG) metrics are widely adopted for evaluating generated
text in machine translation and summarization.
Overlap-based metrics such as BLEU (Papineni
et al., 2002) and ROUGE (Lin, 2004) measure

n-gram precision and recall, respectively, providing straightforward and interpretable scores but offering limited insight into the semantic or factual correctness of generated content. Embeddingbased metrics like BERTScore (Zhang et al., 2020) address some of these limitations by comparing contextualized token embeddings from pretrained language models, thus better capturing semantic similarity. 105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

More recently, LLM-based evaluators, including Prometheus (Kim et al., 2023) and G-Eval (Liu et al., 2023) have been proposed to deliver holistic and dimension-specific quality assessments through instruction-tuned LLMs. Although these approaches achieve stronger alignment with human judgments, their high computational cost and inference latency limit their practical application in large-scale or rapid evaluation settings.

2.2 Radiology Report Evaluation Metrics

Metrics tailored for clinical text generation aim to124address the unique demands of the medical domain,125where factual accuracy and clinical interpretability126are critical. NER-based metrics such as F1 Rad-127Graph (Jain et al., 2021; Delbrouck et al., 2024)128focus on the extraction and matching of medical129entities and their relations to evaluate factual con-130

sistency, while SembScore (CheXbert vector similarity) (Smit et al., 2020) represents findings as structured vectors for semantic comparison.

131

132

133

134

136

137

140

141

142

143

144

145

146

147

148

149 150

151

152

153

155

156

157

158

159

160

161

162

163

165

166

167

168

170

171

172

173

174

176

177

178

179

Composite metrics such as RadCliQ (Yu et al., 2023a) employ linear regression to combine the outputs of BLEU, BERTScore, SembScore, and F1 RadGraph into a single score. RaTEScore (Zhao et al., 2024) uses medical NER model to emphasize entity-level factual precision. LLM-based metrics such as GREEN (Ostmeier et al., 2024) and FineRadScore (Huang et al., 2024) leverage the reasoning capabilities of language models to provide fine-grained, expert-aligned scoring. While each of these methods advances clinical alignment or interpretability, they often introduce complexity, require specialized domain resources, or incur substantial computational cost, which can hinder scalability and routine deployment.

2.3 Regression-Based Text Evaluation

Regression-based approaches aim to bridge the gap between metric predictions and human judgments by training models to directly regress to quality scores or error counts from text features. Notable examples include COMET (Rei et al., 2020) and BLEURT (Sellam et al., 2020), which use pretrained language model encoders with regression heads to approximate human evaluation. While these methods can offer improved correlation with human ratings, they typically operate as black-box predictors and may not provide detailed, clinically interpretable feedback.

Overall, despite substantial progress in the development of evaluation metrics, a fast, interpretable, and clinically grounded method for chest X-ray report generation remains an open challenge.

3 CREPE

Our proposed metric, CREPE, evaluates the quality of chest X-ray reports by directly predicting error counts across clinically relevant categories. By providing structured, category-specific error estimates, CREPE enables both quantitative assessment and interpretable feedback, thereby facilitating robust comparison and targeted diagnostic improvement in automated report generation.

3.1 Problem Formulation

Given a reference report R_{ref} and a generated report R_{cand} , the goal is to predict a vector of error counts,

$$\mathbf{n} = [n_A, n_B, \dots, n_F]$$



Figure 2: CREPE Model Architecture. Tokenized reference and candidate reports are encoded by a medical domain-specific BERT model, and the pooled representation is used by the error counting module to output category-wise error counts.

207

where each element corresponds to one of six clini- cally defined error categories:	180 181
(A) False prediction of finding	182
(B) Omission of finding	183
(C) Incorrect location or position of finding	184
(D) Incorrect severity of finding	185
(E) Mention of comparison not present in the ref- erence impression	186 187
(F) Omission of comparison describing a change from a previous study	188 189
as defined in the annotation protocol used for the ReXVal dataset (Yu et al., 2023b). The overall error count score, S , is given by the sum of the predicted error counts:	190 191 192 193
$S = \sum_{c \in \mathcal{E}} n_c, \tag{1}$	194
where $\mathcal{E} = \{A, B, C, D, E, F\}$. Lower values of <i>S</i> indicate higher report quality, with category-level predictions providing actionable insights for further model refinement.	195 196 197 198
3.2 Model Architecture	199
The CREPE model is built upon a pretrained BERT model specifically tailored for biomedical text, which we fine-tune for regression-based error pre- diction. For each error category, the model jointly	200 201 202 203
estimates both the expected error count and the presence of any error, thereby enabling detailed	203 204 205

and clinically interpretable evaluation. The overall

architecture of the model is shown in Figure 2.

208 209

211

214

215

216

217

218

219

224

227

228

236

239

241

244

245

247

248

249

Encoder and Feature Extraction. Let x denote the tokenized input sequence formed by concatenating R_{ref} and R_{cand} , separated by special tokens. 210 This sequence is processed by the BERT encoder, yielding a pooled representation $\mathbf{h} \in \mathbb{R}^d$:

$$\mathbf{h} = BERT(\mathbf{x})_{[CLS]} \tag{2}$$

A dropout layer is applied to h to prevent overfitting. This pooled output serves as the shared feature for all downstream prediction heads.

Error Regression Heads. For each error category $c \in \mathcal{E}$, an independent regression head predicts the error count:

$$\hat{n}_c = f_c(\mathbf{h}) \tag{3}$$

where $f_c(\cdot)$ is a category-specific feedforward layer outputting a predicted error counts \hat{n}_c . The collection of outputs, $\hat{\mathbf{n}} = [\hat{n}_A, \hat{n}_B, \dots, \hat{n}_F]$, constitutes the predicted error vector.

Error Detection Heads. To address the challenge of class imbalance, where certain error categories are infrequently represented in the training data as shown in Figure 4, the model incorporates auxiliary error detection heads during training. These heads are designed to predict the presence or absence of each error type, and their outputs are used exclusively for loss calculation to enhance learning for rare categories. Specifically, for each category c, the model produces a logit:

$$\hat{p}_c = g_c(\mathbf{h}) \tag{4}$$

where $q_c(\cdot)$ is a category-specific feedforward layer. During training, a sigmoid activation is applied to obtain a probability estimate for error presence, but these predictions are not used at inference time.

Loss Function. As shown in Figure 3, the CREPE model is trained with a dual-objective loss that captures both the count and presence of clinically meaningful errors. For each error category $c \in \mathcal{E}$, the model predicts a continuous error count \hat{n}_c as well as an auxiliary presence logit \hat{p}_c . Let $\mathbf{n} = [n_A, \ldots, n_F]$ denote the ground-truth error counts, and $\mathbf{p} = [p_A, \dots, p_F]$ denote the binary presence indicators, where $p_c = \mathbb{I}[n_c > 0]$.

The regression loss encourages accurate estimation of error counts and is computed using mean squared error (MSE):

252
$$\mathcal{L}_{\text{reg}} = \frac{1}{|\mathcal{E}|} \sum_{c \in \mathcal{E}} \text{MSE}(\hat{n}_c, n_c)$$
(5)



Figure 3: Prediction Heads Architecture. The pooled [CLS] embedding is passed through a dropout layer and then fed into two parallel sets of output heads for each error category: error regression heads (left) predict continuous error counts, while auxiliary error detection heads (right) predict the presence or absence of errors during training.

The presence loss penalizes incorrect predictions of whether an error of a given type is present and is computed using binary cross-entropy (BCE):

$$\mathcal{L}_{\text{pres}} = \frac{1}{|\mathcal{E}|} \sum_{c \in \mathcal{E}} \text{BCE}(\hat{p}_c, p_c)$$
(6)

The total loss used for training is the average of the regression and presence losses:

$$\mathcal{L} = \frac{\mathcal{L}_{\text{reg}} + \mathcal{L}_{\text{pres}}}{2} \tag{7}$$

254

256

257

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

278

This combined objective allows the model to learn both precise error count estimation and improved sensitivity to rare or underrepresented error categories, resulting in more accurate and interpretable evaluation of chest X-ray reports.

Inference and Scoring. As shown in Figure 1, the CREPE model outputs predicted error counts \hat{n}_c for each of the six error categories given a pair of reference and generated reports. The CREPE metric is defined as the sum of predicted error counts:

$$CREPE = \sum_{c \in \mathcal{E}} \hat{n}_c \tag{8}$$

A lower CREPE score indicates higher report quality, and this aggregate value serves as the primary metric for evaluating overall performance through correlation-based analysis. In addition, categorylevel error predictions can be analyzed individually, for example using mean absolute error (MAE), to provide more granular diagnostic feedback or to target specific clinical priorities.

Metric	ReXVal		ReFiSco-v0		RadEvalX		RaTE-Eval		RaTE-Eval [†]	
Metric	au	ho	au	ρ	au	ho	au	ρ	au	ho
BLEU-4	0.383	0.516	0.489	0.616	0.074	0.096	0.197	0.247	0.134	0.166
ROUGE-L	0.570	0.748	0.524	0.662	0.257	0.356	0.200	0.281	0.220	0.302
METEOR	0.484	0.653	0.468	0.617	0.201	0.284	0.174	0.245	0.248	0.338
BERTScore	0.521	0.694	0.541	0.689	0.326	0.452	0.224	0.315	0.256	0.351
F1 RadGraph	0.585	0.765	0.475	0.609	0.171	0.243	0.306	0.393	0.258	0.328
SembScore	0.495	0.666	0.461	0.605	0.318	0.434	0.198	0.280	0.245	0.336
RaTEScore	0.520	0.697	0.433	0.571	0.316	0.438	<u>0.339</u>	0.460	0.310	0.419
RadCliQ-v1	0.623	0.809	0.510	0.656	0.326	0.449	0.299	0.415	0.304	0.414
GREEN	0.626	0.798	0.592	0.709	0.411	0.539	0.374	<u>0.457</u>	<u>0.409</u>	0.494
GREEN EC	<u>0.775</u>	<u>0.899</u>	0.723	<u>0.811</u>	0.448	<u>0.577</u>	0.252	0.315	0.432	0.517
CREPE	0.786	0.933	<u>0.697</u>	0.825	0.580	0.745	0.267	0.375	0.407	0.541

Table 1: Correlation with Radiologist Error Counts Across Datasets. Kendall's τ and Spearman's ρ correlation coefficients for each evaluation metric on five benchmark datasets. Bold indicates the best and underline the second-best value for each metric and dataset. 95% confidence intervals are reported in Table 5.

Synthetic Training Data Generation 3.3

279

281

288

290

291

Obtaining large-scale, expert-annotated error counts for chest X-ray reports is logistically and financially challenging. To address this, we constructed a synthetic training dataset through an automated pipeline, designed to maximize clinical relevance while adhering strictly to responsible data use policies for sensitive medical data.

Report Pair Sampling. We randomly sampled 32,000 image-report pairs from the MIMIC-CXR (Johnson et al., 2019) training set to serve as the basis for synthetic data generation.

Synthetic Report Generation. For each selected image, a candidate report was generated using CheXagent (Chen et al., 2024), a vision-language foundation model specialized for chest X-ray interpretation.

Error Analysis. To identify and classify errors 296 between reference and candidate reports, we used 297 Gemini 2.5 Pro 03-25 preview version deployed via Vertex AI on Google Cloud (Google, 2025). 299 This ensured that all analysis was performed in a secure environment, fully compliant with the MIMIC-301 CXR data use agreement, which prohibits sending 302 data to external services such as public LLM APIs.

Label Extraction. We adapted the prompt from 304 GREEN (Ostmeier et al., 2024) to obtain, for each reference-candidate report pair, error counts across 306 six predefined categories, separately for clinically significant and clinically insignificant errors. For each category, we summed the significant and in-309



Figure 4: Distribution of per-category error counts in the synthetic training dataset. The highly skewed distributions highlight substantial class imbalance, motivating the use of a dual-objective loss to improve learning of rare error categories.

significant error counts to create the final label vector $[n_A, n_B, \ldots, n_F]$ used for CREPE model training. The distribution of error counts for each category is illustrated in Figure 4.

310

311

312

313

314

315

316

317

318

319

320

321

322

324

325

329

This data generation pipeline enabled us to create a diverse and representative training dataset, supporting robust model learning across a wide range of clinically meaningful error types while strictly adhering to data governance requirements. For full details on data generation, including prompts and hyperparameters used for the generation models, see Appendix A.1.

4 **Experiments**

Experimental Setup 4.1

We designed a comprehensive experimental framework to evaluate CREPE's performance across multiple aspects of chest X-ray report evaluation. Be-326 low, we describe the datasets used for benchmark-327 ing, the baseline metrics for comparison, and the 328 main implementation and optimization details for



Figure 5: Kendall's τ Correlation with Radiologist Error Counts Across Datasets. Bar plots show Kendall's τ values (with 95% confidence intervals) for each evaluation metric on multiple public benchmarks. CREPE consistently achieves the highest or near-highest correlation with expert error counts across all datasets.

reproducibility.

330

Evaluation Datasets. We evaluated CREPE on four publicly available datasets containing radiologist-annotated errors: ReXVal (Yu et al., 2023b), ReFiSco-v0 (Tian et al., 2023), RadEvalX (Calamida et al., 2024), and RaTE-Eval (Zhao et al., 2024). ReXVal consists of expert-labeled report-level errors from the MIMIC-CXR dataset, and we additionally use a filtered variant, ReX-Val*, with identical report pairs removed to reduce class imbalance. ReFiSco-v0 provides line-level severity annotations, which we map to binary error labels and aggregate at the report level. RadEvalX comprises 100 IU-Xray (Pavlopoulos et al., 2019) reports with expert annotations for six standard error categories, plus two uncertainty-related categories. For RaTE-Eval, we use only the sentencelevel human rating task, which covers nine imaging modalities and 22 anatomies and includes normalization of error counts by potential error opportunities. We also report raw counts for comparison. Detailed descriptions of dataset construction, annotation protocols, and pre-processing are provided in Appendix A.2.

Baselines. To contextualize CREPE's performance, we compared it against general text evaluation metrics, including BLEU-4 (Papineni et al., 2002), ROUGE-L (Lin, 2004), METEOR (Banerjee and Lavie, 2005), and BERTScore (Zhang et al., 2020), as well as radiology report evaluation metrics, such as F1 RadGraph (Jain et al., 2021), SembScore (Smit et al., 2020), RaTEScore (Zhao et al., 2024), RadCliQ (Yu et al., 2023a), GREEN (Ostmeier et al., 2024), and GREEN error count (GREEN EC).

365 Implementation Details. Our experiments were
366 conducted using BiomedBERT (Gu et al., 2021), a

medical domain-specific BERT model, as the foundation for the CREPE model. The entire BERT encoder and all regression and presence heads were fine-tuned on the synthetic dataset for 10 epochs, with a validation split of 0.1 and a batch size of 64. Sequences were truncated or padded to a maximum length of 512 tokens, and mixed-precision training (FP16) was employed to accelerate computation. All experiments were performed on a single NVIDIA A6000 GPU. The source code will be made publicly available upon acceptance of the paper. 367

369

370

371

372

373

375

376

377

378

379

381

382

383

384

386

388

390

391

394

395

396

397

398

400

401

402

Hyperparameter Optimization. Key hyperparameters, including learning rate, weight decay, and warmup ratio, were optimized via automated search using Optuna (Akiba et al., 2019). The optimal configuration was selected based on validation set performance, as measured by Kendall's τ correlation.

4.2 Results

We comprehensively evaluate CREPE across four key dimensions: correlation with human expert judgments (Sec. 4.3), robustness to class imbalance (Sec. 4.4), computational efficiency (Sec. 4.5), and absolute error-prediction accuracy (Sec. 4.6).

4.3 Correlation with Human Judgments

To assess clinical validity, we evaluate the correlation between each metric's predictions and radiologist-annotated error counts across five public benchmarks with varying annotation protocols and levels of difficulty. Table 1 reports Kendall's τ and Spearman's ρ for all metrics. CREPE achieves the highest or second-highest correlation on most datasets, demonstrating strong alignment with human expert ratings. On ReXVal, CREPE attains a Kendall's τ of 0.786 and a Spearman's ρ of



Figure 6: Correlation Between Metric Scores and Radiologist-Identified Errors (ReXVal). Scatter plots show the relationship between each evaluation metric's score and the total number of radiologist-identified errors on the ReXVal dataset. Each subplot includes a regression line and 95% confidence intervals. CREPE achieves the highest correlation with expert error counts. Kendall's τ and Spearman's ρ are shown in the subplot titles.

0.933, while maintaining competitive performance on line-level (ReFiSco-v0), category-extended (RadEvalX), and out-of-distribution multi-modality (RaTE-Eval) datasets.

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431 432

433

434

435

436

Figure 5 summarizes Kendall's τ (with 95% confidence intervals) across all benchmarks, illustrating that CREPE is consistently among the topperforming metrics in terms of agreement with radiologist judgments.

For a more detailed view, Figure 6 visualizes the relationship between metric scores and total radiologist-identified errors on ReXVal. The scatter plots show that CREPE exhibits a monotonic association with expert error counts, in line with its correlation statistics.

Taken together, these results suggest that CREPE provides reliable alignment with expert judgment across a range of evaluation settings and dataset characteristics.

4.4 Robustness to Class Imbalance

To examine the impact of class imbalance, particularly the prevalence of zero-error pairs in ReX-Val, we constructed a filtered version, ReXVal*, by removing report pairs where the reference and candidate reports were identical. As shown in Table 2, this filtering results in a notable decline in the performance of conventional metrics; the average decrease in Kendall's τ is approximately 0.12. In contrast, CREPE's correlation with human judgments decreases by only 0.033, corresponding to a 4.2 percent relative drop. These results suggest that CREPE remains robust and continues to reliably differentiate clinically meaningful errors even in more challenging evaluation settings.

Metric	ReXVal	ReXVal*	$\Delta \tau$
BLEU-4	0.383	0.215	-0.168
ROUGE-L	0.570	0.459	-0.111
METEOR	0.484	0.355	-0.129
BERTScore	0.521	0.404	-0.117
F1 RadGraph	0.585	0.484	-0.101
SembScore	0.495	0.368	-0.127
RaTEScore	0.520	0.408	-0.113
RadCliQ-v1	0.623	0.500	-0.123
GREEN	0.626	0.500	-0.126
GREEN EC	0.775	0.729	-0.046
CREPE	0.786	0.753	-0.033

Table 2: **Robustness to Class Imbalance.** Kendall's τ for each metric on the original ReXVal dataset and the filtered set ReXVal*, after removing 26 identical report pairs. $\Delta \tau$ shows the change in correlation.

4.5 Computational Efficiency

We compare the inference time per sample for all metrics using the ReXVal dataset. As shown in Figure 1 and Table 3, CREPE processes a report in 9.5 milliseconds, which is comparable to the fastest neural baselines such as BERTScore. In contrast, LLM-based methods such as GREEN require over 2,600 milliseconds per sample, making them approximately 280 times slower. This substantial speed advantage demonstrates that CREPE is wellsuited for rapid and large-scale evaluation, without sacrificing alignment with expert judgment. 437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

4.6 Error Score Accuracy

Among public benchmarks, RaTE-Eval uniquely provides category-level error counts, enabling direct assessment of absolute prediction accuracy.

Metric	Time (ms)
BLEU-4	4.29
ROUGE-L	125.99
METEOR	24.83
BERTScore	9.65
F1 RadGraph	36.59
SembScore	18.86
RaTEScore	160.28
RadCliQ-v1	370.58
GREEN	2642.37
GREEN EC	2642.37
CREPE	9.53

Table 3: **Speed Comparison.** Average inference time per sample (milliseconds) for each evaluation metric on the ReXVal dataset.

Head Counts	au
1	0.723
2	0.730
6 (ours)	0.753
Loss Function	au
MSE only	0.718
Poisson NLL only	0.738
MSE + BCE (ours)	0.753
Encoder Backbone	au
ClinicalBERT	0.693
Bio_ClinicalBERT	0.716
BiomedBERT (ours)	0.753

Table 4: Ablation on Model Components. Impact of regression head granularity, loss function, and encoder backbone on Kendall's τ for ReXVal*.

On this dataset, CREPE achieves a mean absolute error (MAE) of 0.739 ± 0.627 , which is 33% lower than the MAE of GREEN EC (1.102 ± 0.632). Both metrics offer category-specific predictions, but CREPE's lower error demonstrates the effectiveness of regression-based modeling when detailed supervision is available.

5 Ablation Studies

453

454

455

456

457

458

459

460

We conducted ablation studies on the balanced
ReXVal* dataset to quantify the contributions of
key architectural and modeling choices. Unless otherwise specified, all variants share the same hyperparameters as the full model (six regression heads,
dual-objective MSE + BCE loss, and BiomedBERT
backbone).

5.1 Number of Regression Heads

To assess the impact of output granularity, we compared three regression head configurations: (1) a single head predicting the total error count, (2) two heads for clinically significant and insignificant errors, and (3) six heads for category-specific errors, which is the default setting. As shown in Table 4, increasing the number of regression heads leads to consistently higher agreement with radiologist judgments, with the six-head model achieving the best Kendall's τ . 468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

5.2 Loss Function

We evaluated three loss formulations: MSE only, Poisson negative log likelihood (NLL) only, and a dual-objective MSE+BCE loss that incorporates an auxiliary error presence term. The results in Table 4 demonstrate that the dual-objective loss yields the highest correlation with expert annotations, indicating the value of explicitly modeling both error counts and the presence of rare error types.

5.3 Encoder Backbone

To determine the importance of domain-specific pre-training, we replaced BiomedBERT with two clinical BERT models: ClinicalBERT (Wang et al., 2023: Liu et al., 2025) and Bio_ClinicalBERT (Alsentzer et al., 2019), keeping all other settings fixed. As reported in Table 4, BiomedBERT yields the highest Kendall's τ . The reason for BiomedBERT's superior performance may be related to its pre-training on a broader biomedical literature corpus, potentially offering better coverage of radiology terminology than models trained solely on clinical notes. Details about encoder backbones are provided in Appendix A.3.

6 Conclusion

We presented **CREPE**, an efficient evaluation metric for automated chest X-ray report generation that predicts clinically meaningful error counts using a domain-specific BERT model with a multi-head regression architecture. CREPE provides both an overall score and interpretable category-level feedback, demonstrating strong correlation with expert judgments, robustness to class imbalance, and fast inference compared to existing evaluation methods. Future work includes extending this approach to additional medical imaging domains and further investigating its generalizability.

516 Limitations

Despite its advantages, CREPE has several limita-517 tions. First, the model requires a GPU for efficient 518 inference, which may limit accessibility in low-519 resource environments compared to purely rule-520 based or string-matching metrics. Second, although 521 inference is fast, the generation of synthetic train-522 ing data involves significant computational cost, including the use of large foundation models and commercial LLM APIs, which could restrict repro-525 ducibility or scalability in some settings. Third, 526 while CREPE demonstrates robust performance 527 on several public benchmarks, its accuracy on out-528 of-distribution modalities and reporting styles is not guaranteed and may degrade when applied to 530 clinical scenarios substantially different from those 531 in the training data. Additionally, the reliance on synthetic error annotations, rather than large-scale real-world expert labeling, may introduce biases or affect generalizability. Addressing these challenges will be important for future extensions and 536 real-world deployment. 537

References

538

539

542

543

544

545

546

547

548

549

550 551

552

553

554

555

557

560

563

564

565

566

567

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A nextgeneration hyperparameter optimization framework. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '19, page 2623–2631, New York, NY, USA. Association for Computing Machinery.
 - Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. Publicly Available Clinical BERT Embeddings. In Proceedings of the 2nd Clinical Natural Language Processing Workshop, pages 72–78, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
 - Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
 - Amos Rubin Calamida, Farhad Nooralahzadeh, Morteza Rohanian, Mizuho Nishio, Koji Fujimoto, and Michael Krauthammer. 2024. Radiology Report Generation Models Evaluation Dataset For Chest X-rays (RadEvalX).
 - Zhihong Chen, Maya Varma, Justin Xu, Magdalini Paschali, Dave Van Veen, Andrew Johnston, Alaa

Youssef, Louis Blankemeier, Christian Bluethgen, Stephan Altmayer, Jeya Maria Jose Valanarasu, Mohamed Siddig Eltayeb Muneer, Eduardo Pontes Reis, Joseph Paul Cohen, Cameron Olsen, Tanishq Mathew Abraham, Emily B. Tsai, Christopher F. Beaulieu, Jenia Jitsev, and 4 others. 2024. A Vision-Language Foundation Model to Enhance Efficiency of Chest X-ray Interpretation. *arXiv preprint*. ArXiv:2401.12208 [cs]. 568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

585

586

587

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

- Jean-Benoit Delbrouck, Pierre Chambon, Zhihong Chen, Maya Varma, Andrew Johnston, Louis Blankemeier, Dave Van Veen, Tan Bui, Steven Truong, and Curtis Langlotz. 2024. RadGraph-XL: A Large-Scale Expert-Annotated Dataset for Entity and Relation Extraction from Radiology Reports. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 12902–12915, Bangkok, Thailand. Association for Computational Linguistics.
- Google. 2025. Gemini 2.5: Our most intelligent AI model.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing. *ACM Trans. Comput. Healthcare*, 3(1):2:1–2:23.
- Alyssa Huang, Oishi Banerjee, Kay Wu, Eduardo Pontes Reis, and Pranav Rajpurkar. 2024. FineRadScore: A Radiology Report Line-by-Line Evaluation Technique Generating Corrections with Severity Scores. *arXiv preprint*. ArXiv:2405.20613 [cs].
- Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Truong, Du Nguyen Duong N. Duong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew Lungren, Andrew Ng, Curtis Langlotz, and Pranav Rajpurkar. 2021. RadGraph: Extracting Clinical Entities and Relations from Radiology Reports. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, 1.
- Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Roger G. Mark, and Steven Horng. 2019. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(1):317.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. Prometheus: Inducing Fine-Grained Evaluation Capability in Language Models. In 12th International Conference on Learning Representations (ICLR 2024).
- Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

- 625 626 627
- 6
- 633 634
- 6
- 638
- 639 640 641
- 6
- 6
- 6
- 647
- 6 6
- 651
- 6

653 654

- 65
- 65

65

66

662 663 664

66

- 667
- 66
- 670 671
- 672
- 673 674 675

676 677

678 679

- Xiaohong Liu, Hao Liu, Guoxing Yang, Zeyu Jiang, Shuguang Cui, Zhaoze Zhang, Huan Wang, Liyuan Tao, Yongchang Sun, Zhu Song, Tianpei Hong, Jin Yang, Tianrun Gao, Jiangjiang Zhang, Xiaohu Li, Jing Zhang, Ye Sang, Zhao Yang, Kanmin Xue, and 5 others. 2025. A generalist medical language model for disease diagnosis assistance. *Nature Medicine*, 31(3):932–942. Publisher: Nature Publishing Group.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG Evaluation using Gpt-4 with Better Human Alignment. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 2511–2522, Singapore. Association for Computational Linguistics.
- 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 Association for Computational Linguistics: EMNLP 2024, pages 374–390, Miami, Florida, USA. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- John Pavlopoulos, Vasiliki Kougia, and Ion Androutsopoulos. 2019. A survey on biomedical image captioning. In *Proceedings of the Second Workshop on Shortcomings in Vision and Language*, pages 26–36, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A Neural Framework for MT Evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning Robust Metrics for Text Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.
- 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 Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1500–1519, Online. Association for Computational Linguistics.

Katherine Tian, Sina J Hartung, Andrew A Li, Jaehwan Jeong, Fardad Behzadi, Juan Calle-Toro, Subathra Adithan, Michael Pohlen, David Osayande, and Pranav Rajpurkar. 2023. ReFiSco: Report Fix and Score Dataset for Radiology Report Generation. 681

682

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

- Guangyu Wang, Xiaohong Liu, Zhen Ying, Guoxing Yang, Zhiwei Chen, Zhiwen Liu, Min Zhang, Hongmei Yan, Yuxing Lu, Yuanxu Gao, Kanmin Xue, Xiaoying Li, and Ying Chen. 2023. Optimized glycemic control of type 2 diabetes with reinforcement learning: a proof-of-concept trial. *Nature Medicine*, 29(10):2633–2642. Publisher: Nature Publishing Group.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. HuggingFace's Transformers: State-of-theart Natural Language Processing. *arXiv preprint*. ArXiv:1910.03771 [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. 2023a. Evaluating progress in automatic chest X-ray radiology report generation. *Patterns*, 4(9):100802.
- Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser Ururahy Nunes Fonseca, Henrique Lee, Zahra Shakeri, Andrew Ng, Curtis Langlotz, Vasantha Kumar Venugopal, and Pranav Rajpurkar. 2023b. Radiology Report Expert Evaluation (ReXVal) Dataset.
- Juan Manuel Zambrano Chaves, Shih-Cheng Huang, Yanbo Xu, Hanwen Xu, Naoto Usuyama, Sheng Zhang, Fei Wang, Yujia Xie, Mahmoud Khademi, Ziyi Yang, Hany Awadalla, Julia Gong, Houdong Hu, Jianwei Yang, Chunyuan Li, Jianfeng Gao, Yu Gu, Cliff Wong, Mu Wei, and 8 others. 2025. A clinically accessible small multimodal radiology model and evaluation metric for chest X-ray findings. *Nature Communications*, 16(1):3108. Publisher: Nature Publishing Group.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. In 8th International Conference on Learning Representations (ICLR 2020).
- Weike Zhao, Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2024. RaTEScore: A Metric for Radiology Report Generation. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 15004–15019, Miami, Florida, USA. Association for Computational Linguistics.

A Appendix	739
A.1 Details for Synthetic Data Generation	740
We provide here the details and prompts used to generate the synthetic training data for the CREPE model.	741
A.1.1 Candidate Report Generation	742
To generate candidate reports, we utilized StanfordAIMI/CheXagent-8b from the Hugging Face	743
Hub (Wolf et al., 2020) with default inference settings. From the MIMIC-CXR training split, we	744
randomly sampled 32,000 examples. For each sample, the 'FINDINGS' section was used as the reference	745
report; if this section was unavailable, we instead used the 'IMPRESSION' section. The candidate report	746
was generated by providing the associated chest X-ray image along with a standardized prompt.	747

{image} Provide a radiological report for the following image. ASSISTANT:

A.1.2 Error Count Data Generation

To obtain error counts, we used gemini-2.5-pro-preview-03-25 accessed via Vertex AI on Google Cloud, again with default parameters (temperature 1.0, topP 0.95, candidateCount 1). For each reference–candidate pair, the model was prompted to compare the candidate report against the reference, following explicit instructions to assess both clinically significant and clinically insignificant errors across six predefined categories. The prompt also required the model to return a structured output, detailing error counts and explanations for each category, as well as matched findings between the reports.

Objective: Evaluate the accuracy of a candidate radiology report in comparison to a reference radiology report composed by expert radiologists. Process Overview: You will be presented with: 1. The criteria for making a judgment. 2. The reference radiology report. 3. The candidate radiology report. 4. The desired format for your assessment. 1. Criteria for Judgment: For each candidate report, determine: - The count of clinically significant errors. - The count of clinically insignificant errors. Errors can fall into one of these categories: a) False report of a finding in the candidate. b) Missing a finding present in the reference. c) Misidentification of a finding's anatomic location/position. d) Misassessment of the severity of a finding. e) Mentioning a comparison that isn't in the reference. f) Omitting a comparison detailing a change from a prior study. Note: Concentrate on the clinical findings rather than the report's writing style. Evaluate only the findings that appear in both reports. 2. Reference Report: {reference_report} 3. Candidate Report: {candidate_report} 4. Reporting Your Assessment: Follow this specific format for your output, even if no errors are found: Do NOT abbreviate the name of the error type. [Explanation]: <Explanation> [Clinically Significant Errors]: (a) <Error Type>: <The number of errors>. <Error 1>; <Error 2>; ...; <Error n> (f) <Error Type>: <The number of errors>. <Error 1>; <Error 2>; ...; <Error n> [Clinically Insignificant Errors]: (a) <Error Type>: <The number of errors>. <Error 1>; <Error 2>; ...; <Error n> (f) <Error Type>: <The number of errors>. <Error 1>; <Error 2>; ...; <Error n> [Matched Findings]: <The number of matched findings>. <Finding 1>; <Finding 2>; ...; <Finding n>

A.2 Dataset Details and Processing

We provide detailed descriptions of each benchmark used in our evaluation. Figure 7 visualizes the distribution of total error counts across all datasets.



Figure 7: Total error count distributions for all evaluation datasets.

A.2.1 ReXVal

ReXVal contains radiologist-annotated errors for generated reports compared to ground-truth reports from the MIMIC-CXR dataset. Six radiologists independently evaluated both clinically significant and insignificant errors, assigning counts for each of six predefined error categories. For each ground-truth report, four candidate reports were generated by different automated methods. To address class imbalance, we also evaluate on ReXVal*, a variant where pairs with identical reference and candidate reports are removed.

A.2.2 ReFiSco-v0

ReFiSco-v0 provides line-level radiologist error annotations for MIMIC-CXR reports, categorizing each line as 'No error', 'Not actionable', 'Actionable nonurgent error', 'Urgent error', or 'Emergent error.' We map 'No error' to zero and all other categories to one, then aggregate the binary error labels across lines to produce report-level error counts.

A.2.3 RadEvalX

RadEvalX consists of 100 reports sampled from the IU-Xray dataset, selected to balance normal and abnormal findings. Each report and its generated counterpart were annotated by experts for six standard error categories, consistent with ReXVal, as well as two additional categories related to uncertainty (mention or omission of uncertainty). For our evaluation, we sum error counts across all categories.

A.2.4 RaTE-Eval

For RaTE-Eval, we use the sentence-level human rating benchmark, which spans nine imaging modalities and 22 anatomical regions, representing a multi-modality and out-of-distribution test case. Annotators counted errors per sentence, and scores are normalized by the number of potential error opportunities. For consistency, we also report results using the raw, unnormalized error counts, denoted as RaTE-Eval[†]

A.2.5 Additional Notes

All datasets were used in accordance with their respective data use agreements and ethical guidelines. 781 Where necessary, we standardized error category definitions across datasets for consistency in evaluation. 782

755 756

757

758

759

760

761

764

765

766

767

768

770

771

772

773

776

778

779

A.3 Details for BERT Encoder Backbones

784 We summarize here the key characteristics and pretraining corpora for each encoder backbone used in our785 experiments.

A.3.1 medicalai/ClinicalBERT

787

788

791

793

ClinicalBERT is a domain-adapted BERT model initially trained on general English text and then further pre-trained on a large corpus of de-identified clinical notes. Its training corpus encompasses approximately 1.2 billion words of clinical narratives, including diverse disease phenotypes and a wide range of free-text observations from electronic health records. The model is designed for masked language modeling and fine-tuned for downstream clinical NLP tasks such as information extraction, symptom detection, and temporal representation of patient trajectories. ClinicalBERT has demonstrated high performance in extracting clinically relevant information from free-text notes, with an average F1 score of 94.5% in symptom extraction tasks on annotated samples.

A.3.2 emilyalsentzer/Bio_ClinicalBERT

Bio_ClinicalBERT is based on the BERT architecture, initialized from BioBERT, and further pre-trained on approximately two million clinical notes from the MIMIC-III database. This corpus includes a comprehensive range of note types, including both discharge summaries and general clinical narratives. The pre-training process follows standard masked language modeling objectives and leverages the full MIMIC corpus for greater domain coverage. Bio_ClinicalBERT embeddings are intended as a community resource for downstream medical NLP tasks. Pre-training required significant computational resources, estimated at 17–18 days on a high-end GPU.

A.3.3 microsoft/BiomedNLP-BiomedBERT-base-uncased-abstract

BiomedBERT, previously named PubMedBERT, used as the primary encoder backbone for CREPE, is pretrained on large-scale biomedical text corpora, including PubMed abstracts and clinical literature. The model employs the standard BERT-base architecture and is optimized for masked language modeling to capture domain-specific biomedical terminology and semantics. For our experiments, BiomedBERT was further fine-tuned on the CREPE training set for the specific task of clinical error count regression.

Motrio	ReX	Wal	ReX	Val*	ReFiSco-v0		
wietric	au	ho	au	ho	au	ho	
BLEU-4	0.383 [0.274, 0.482]	0.516 [0.371, 0.633]	0.215 [0.104, 0.333]	0.294 [0.144, 0.447]	0.489 [0.408, 0.561]	0.616 [0.520, 0.692]	
ROUGE-L	0.570 [0.490, 0.636]	0.748 [0.661, 0.809]	0.459 [0.373, 0.537]	0.633 [0.523, 0.716]	0.524 [0.443, 0.602]	0.662 [0.569, 0.742]	
METEOR	0.484 [0.400, 0.558]	0.653 [0.546, 0.738]	0.355 [0.260, 0.444]	0.497 [0.369, 0.606]	0.468 [0.388, 0.538]	0.617 [0.514, 0.694]	
BERTScore	0.521 [0.441, 0.598]	0.694 [0.597, 0.776]	0.404 [0.313, 0.497]	0.558 [0.438, 0.668]	0.541 [0.467, 0.606]	0.689 [0.604, 0.759]	
F1 RadGraph	0.585 [0.512, 0.650]	0.765 [0.686, 0.823]	0.484 [0.401, 0.563]	0.661 [0.562, 0.742]	0.475 [0.392, 0.547]	0.609 [0.509, 0.690]	
SembScore	0.495 [0.416, 0.574]	0.666 [0.571, 0.751]	0.368 [0.271, 0.455]	0.513 [0.388, 0.621]	0.461 [0.387, 0.538]	0.605 [0.512, 0.691]	
RaTEScore	0.520 [0.439, 0.589]	0.697 [0.601, 0.770]	0.408 [0.319, 0.495]	0.564 [0.448, 0.667]	0.433 [0.356, 0.504]	0.571 [0.472, 0.654]	
RadCliQ-v1	0.623 [0.566, 0.676]	0.809 [0.749, 0.855]	0.540 [0.475, 0.602]	0.730 [0.654, 0.791]	0.510 [0.433, 0.576]	0.656 [0.564, 0.729]	
GREEN	0.626 [0.555, 0.685]	0.798 [0.729, 0.843]	0.541 [0.459, 0.614]	0.713 [0.617, 0.789]	0.592 [0.518, 0.663]	0.709 [0.632, 0.781]	
GREEN EC	0.775 [0.728, 0.814]	0.899 [0.861, 0.924]	0.729 [0.667, 0.776]	0.864 [0.807, 0.900]	0.723 [0.660, 0.780]	0.811 [0.744, 0.866]	
CREPE	0.786 [0.749, 0.816]	0.933 [0.907, 0.949]	0.753 [0.703, 0.794]	0.911 [0.871, 0.937]	0.697 [0.640, 0.747]	0.825 [0.761, 0.873]	
Motrio	RadEvalX		RaT	E-Eval	RaTE-Eval†		
Metric	au	ho	au	ho	au	ho	
BLEU-4	0.074 [-0.092, 0.231]	0.096 [-0.120, 0.301]	0.197 [0.119, 0.270]	0.247 [0.151, 0.339]	0.134 [0.060, 0.208]	0.166 [0.075, 0.255]	
ROUGE-L	0.257 [0.111, 0.382]	0.356 [0.155, 0.519]	0.200 [0.136, 0.260]	0.281 [0.192, 0.362]	0.220 [0.146, 0.284]	0.302 [0.201, 0.387]	
METEOR	0.201 [0.065, 0.334]	0.284 [0.100, 0.458]	0.174 [0.105, 0.244]	0.245 [0.150, 0.340]	0.248 [0.183, 0.310]	0.338 [0.251, 0.420]	
BERTScore	0.326 [0.195, 0.442]	0.452 [0.274, 0.596]	0.224 [0.160, 0.284]	0.315 [0.226, 0.395]	0.256 [0.196, 0.316]	0.351 [0.270, 0.432]	
F1 RadGraph	0.171 [0.053, 0.294]	0.243 [0.070, 0.409]	0.306 [0.235, 0.374]	0.393 [0.304, 0.476]	0.258 [0.185, 0.329]	0.328 [0.235, 0.418]	
SembScore	0.318 [0.165, 0.447]	0.434 [0.228, 0.593]	0.198 [0.134, 0.258]	0.280 [0.190, 0.360]	0.245 [0.186, 0.308]	0.336 [0.256, 0.421]	
RaTEScore	0.316 [0.192, 0.435]	0.438 [0.270, 0.583]	0.339 [0.280, 0.396]	0.460 [0.379, 0.534]	0.310 [0.250, 0.369]	0.419 [0.340, 0.494]	
RadCliQ-v1	0.326 [0.206, 0.459]	0.449 [0.288, 0.609]	0.299 [0.238, 0.357]	0.415 [0.332, 0.492]	0.304 [0.247, 0.363]	0.414 [0.338, 0.488]	
GREEN	0.411 [0.308, 0.517]	0.539 [0.408, 0.661]	0.374 [0.303, 0.439]	0.457 [0.371, 0.532]	0.409 [0.343, 0.470]	0.494 [0.419, 0.562]	
GREEN EC	0.448 [0.341, 0.546]	0.577 [0.441, 0.686]	0.252 [0.184, 0.320]	0.315 [0.229, 0.396]	0.432 [0.369, 0.495]	0.517 [0.443, 0.585]	
CDEDE	0.590	0 745 10 (25 0 02)	0.267 10.202 0.2011	0 375 10 200 0 4571	0.407 10.240 0.4661	0 541 10 460 0 6111	

Table 5: Correlation with Radiologist Error Counts and Confidence Intervals. Kendall's τ and Spearman's ρ correlation coefficients, with 95% confidence intervals, for each evaluation metric on six datasets. Bold indicates the best and underline the second-best values for each metric and dataset.