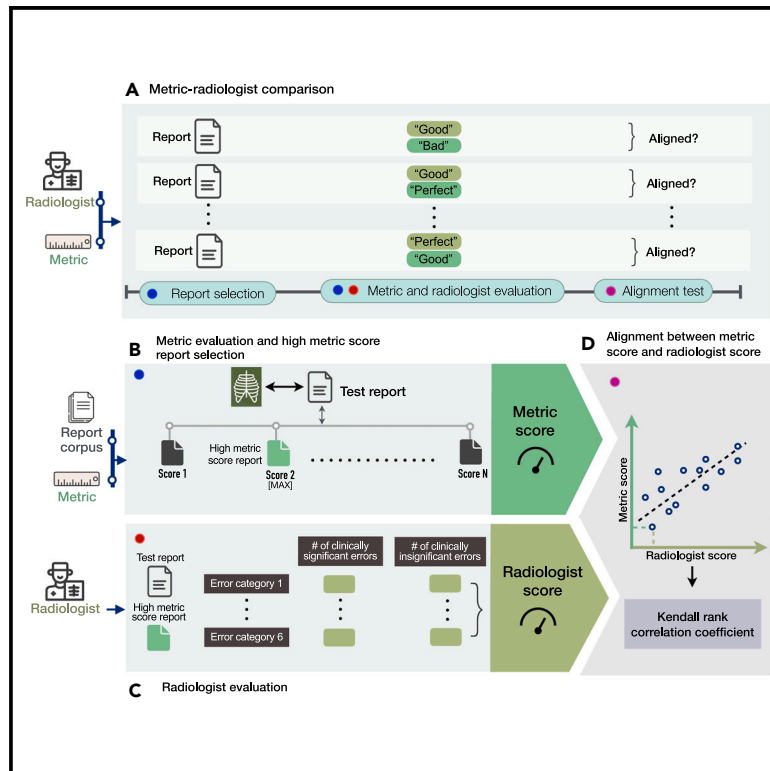


# Patterns

## Evaluating progress in automatic chest X-ray radiology report generation

### Graphical abstract



### Authors

Feiyang Yu, Mark Endo,  
Rayan Krishnan, ..., Curtis P. Langlotz,  
Vasantha Kumar Venugopal,  
Pranav Rajpurkar

### Correspondence

pranav\_rajpurkar@hms.harvard.edu

### In brief

Yu et al. quantitatively examine the correlation between automated metrics and the scoring of radiology reports by radiologists to understand how to meaningfully measure progress on automatic report generation. They propose RadGraph F1, a metric based on overlap in clinical entities and relations, and RadCliQ, a composite metric that combines individual metrics and aligns better with radiologists. They analyze the types of information metrics fail to capture to further understand metric usefulness and evaluate state-of-the-art report generation approaches.

### Highlights

- Examined correlation between automated metrics and scoring of reports by radiologists
- Proposed metric based on overlap in clinical entities and relations named RadGraph F1
- Proposed composite metric RadCliQ with better alignment with radiologists
- Analyzed failure modes of automated metrics



## Article

# Evaluating progress in automatic chest X-ray radiology report generation

Feiyang Yu,<sup>1,9</sup> Mark Endo,<sup>1,9</sup> Rayan Krishnan,<sup>1,9</sup> Ian Pan,<sup>2</sup> Andy Tsai,<sup>3</sup> Eduardo Pontes Reis,<sup>4</sup> Eduardo Kaiser Ururahy Nunes Fonseca,<sup>4</sup> Henrique Min Ho Lee,<sup>4</sup> Zahra Shakeri Hossein Abad,<sup>5</sup> Andrew Y. Ng,<sup>1</sup> Curtis P. Langlotz,<sup>6</sup> Vasantha Kumar Venugopal,<sup>7</sup> and Pranav Rajpurkar<sup>8,10,\*</sup>

<sup>1</sup>Department of Computer Science, Stanford University, Stanford, CA 94305, USA

<sup>2</sup>Department of Radiology, Brigham and Women's Hospital, Boston, MA 02115, USA

<sup>3</sup>Department of Radiology, Boston Children's Hospital, Harvard Medical School, Boston, MA 02115, USA

<sup>4</sup>Cardiothoracic Radiology Group, Hospital Israelita Albert Einstein, São Paulo, São Paulo 05652, Brazil

<sup>5</sup>Dalla Lana School of Public Health, University of Toronto, Toronto, ON M5T 3M7, Canada

<sup>6</sup>AIMI Center, Stanford University, Stanford, CA 94304, USA

<sup>7</sup>CARPL.ai, New Delhi, Delhi 110016, India

<sup>8</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115, USA

<sup>9</sup>These authors contributed equally

<sup>10</sup>Lead contact

\*Correspondence: [pranav\\_raipurkar@hms.harvard.edu](mailto:pranav_raipurkar@hms.harvard.edu)

<https://doi.org/10.1016/j.patter.2023.100802>

**THE BIGGER PICTURE** Artificial intelligence (AI) has made formidable progress in the interpretation of medical images, but its application has largely been limited to the identification of a handful of individual pathologies. In contrast, the generation of complete narrative radiology reports more closely matches how radiologists communicate diagnostic information. While recent progress on vision-language models has enabled the possibility of generating radiology reports, the task remains far from solved. Our work aims to tackle one of the most important bottlenecks for progress: the limited ability to meaningfully measure progress on the report generation task. We quantitatively examine the correlation between automated metrics and the scoring of reports by radiologists and investigate the failure modes of metrics. We also propose a metric based on overlap in clinical entities and relations extracted from reports and a composite metric, called RadCliQ, that is a combination of individual metrics.



**Proof-of-Concept:** Data science output has been formulated, implemented, and tested for one domain/problem

## SUMMARY

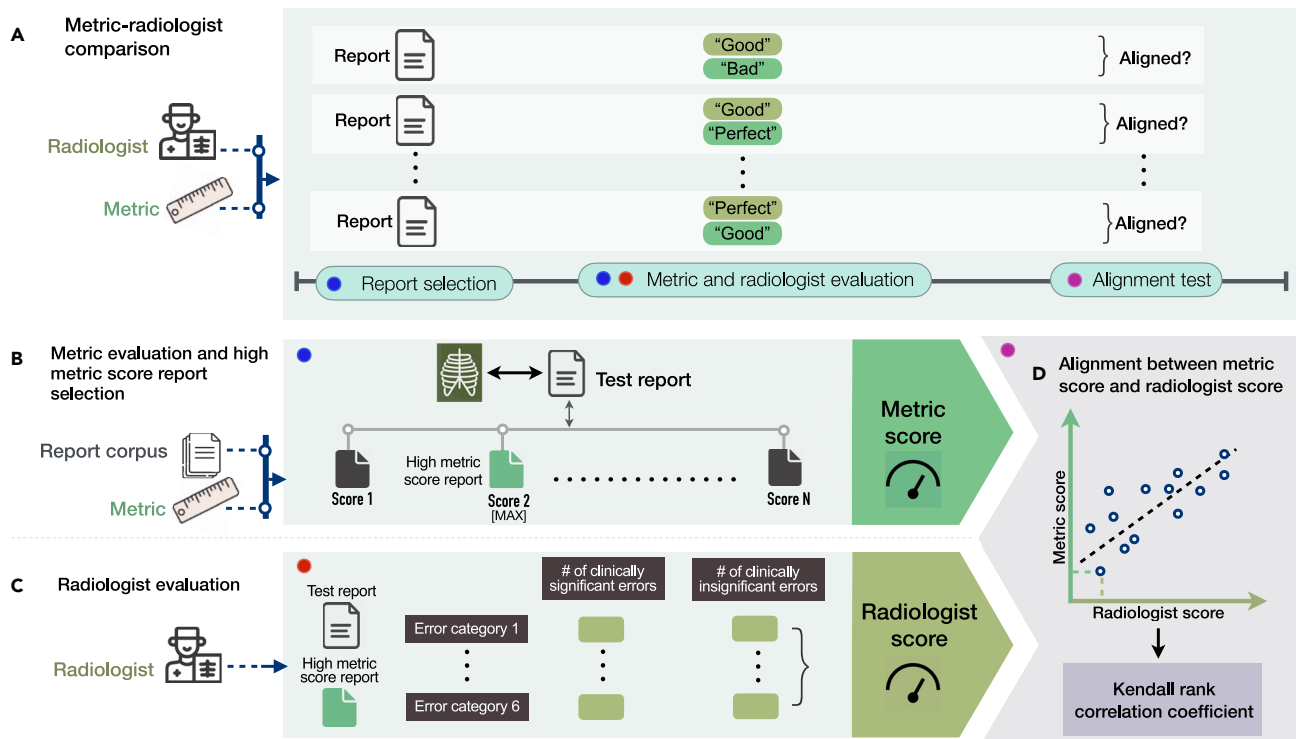
Artificial intelligence (AI) models for automatic generation of narrative radiology reports from images have the potential to enhance efficiency and reduce the workload of radiologists. However, evaluating the correctness of these reports requires metrics that can capture clinically pertinent differences. In this study, we investigate the alignment between automated metrics and radiologists' scoring of errors in report generation. We address the limitations of existing metrics by proposing new metrics, RadGraph F1 and RadCliQ, which demonstrate stronger correlation with radiologists' evaluations. In addition, we analyze the failure modes of the metrics to understand their limitations and provide guidance for metric selection and interpretation. This study establishes RadGraph F1 and RadCliQ as meaningful metrics for guiding future research in radiology report generation.

## INTRODUCTION

Artificial intelligence (AI) has been making great strides in tasks that require expert knowledge,<sup>1-4</sup> including the interpretation of medi-

cal images.<sup>5</sup> In recent years, medical AI models have been demonstrated to achieve expert-level performance,<sup>6</sup> generalize to hospitals beyond which they were trained,<sup>3</sup> and assist specialists in their interpretation.<sup>7</sup> However, the application of AI to image





**Figure 1. Method overview**

(A) Experimental design for selecting radiology reports and comparing metrics and radiologists in evaluating reports. (B) Given a test report, selecting the report with the highest metric score from the training report corpus with respect to the test report and a particular metric. (C) Conducting radiologist evaluation on the high metric score report relative to the test report, where radiologists identify the number of clinically significant and insignificant errors in the high metric score report across six error categories. (D) Determining the alignment between metric scores and radiologist scores assigned to the same reports using the Kendall rank correlation coefficient.

interpretation tasks has often been limited to the identification of a handful of individual pathologies,<sup>8–10</sup> representing an over-simplification of the image interpretation task. In contrast, the generation of complete narrative radiology reports<sup>11–21</sup> moves past that simplification and is consistent with how radiologists communicate diagnostic information: the narrative report allows for highly diverse and nuanced findings, including association of findings with anatomic location, and expressions of uncertainty. Although the generation of radiology reports from medical images in their full complexity would signify a tremendous achievement for AI, the task remains far from solved. Our work aims to tackle one of the most important bottlenecks for progress: the limited ability to meaningfully measure progress on the report generation task.

Automatically measuring the quality of generated radiology reports is challenging. Most prior works have relied on a set of metrics inspired by similar setups in natural language generation, where radiology report text is treated as generic text.<sup>22</sup> However, unlike generic text, radiology reports involve complex, domain-specific knowledge and critically depend on factual correctness. Even metrics that were designed to evaluate the correctness of radiology information by capturing domain-specific concepts do not align with radiologists.<sup>23</sup> Therefore, improvement on existing metrics may not produce clinically meaningful progress or indicate the direction for further progress. This fundamental bottleneck hinders understanding of the quality of report generation methods thereby impeding work toward improvement of ex-

isting methods. We seek to remove this bottleneck by developing meaningful measures of progress in radiology report generation. The answer to this question is imperative to understanding which metrics can guide us toward generating reports that are clinically indistinguishable from those generated by radiologists.

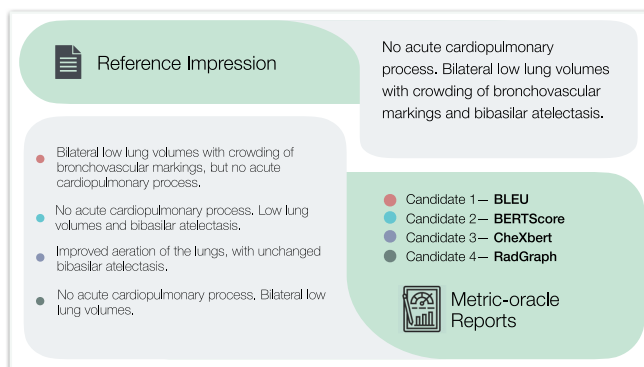
In this study, we quantitatively examine the correlation between automated metrics and the scoring of reports by radiologists. We propose a new automatic metric that computes the overlap in clinical entities and relations between a machine-generated report and a radiologist-generated report, called RadGraph<sup>24</sup> F1. We develop a methodology to predict a radiologist-determined error score from a combination of automated metrics, called RadCliQ. We analyze failure modes of the metrics, namely the types of information the metrics do not capture, to understand when to choose particular metrics and how to interpret metric scores. Finally, we measure the performance of state-of-the-art report generation models using the investigated metrics. The result is a quantitative understanding of radiology report generation metrics and clear guidance for metric selection to guide future research on radiology report generation.

## RESULTS

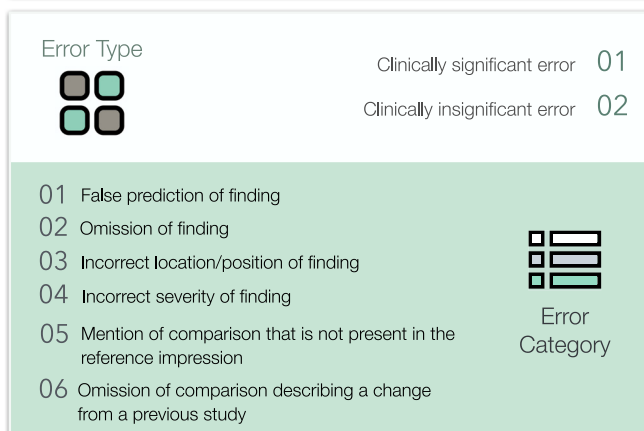
### Alignment between automated metrics and radiologists

We study whether there is alignment between automated metric and radiologist scores assigned to radiology reports. An

**A**  
Example study (test report and metric-oracle reports)



**B**  
Error types and error categories



overview of our methodology is shown in Figure 1. Figure 1A shows the experimental design for determining alignment. Given a test report from MIMIC-CXR,<sup>25–27</sup> we select a series of candidate reports from the MIMIC-CXR training set that score highly according to various metrics, including BLEU,<sup>28</sup> BERTScore,<sup>29</sup> CheXbert vector similarity ( $s_{emb}$ ),<sup>9</sup> and a novel metric RadGraph<sup>24</sup> F1. Specifically, we select a candidate report by finding the test report’s *metric-oracle*: the highest-scoring report from the MIMIC-CXR training set with respect to a particular metric (Figure 1B). We choose this set of reasonably accurate reports so we can study their quality with more precision. An example study with a reference report and candidate metric-oracle reports is shown in Figure 2A.

Next, we have six board-certified radiologists score how well the candidates match the test report (Figure 1C). Radiologists scored the number of errors that various candidate reports make compared with the test report, and errors are categorized as clinically significant or insignificant. Radiologists subtyped every error into the following six categories: (1) false prediction of finding (i.e., false positive), (2) omission of finding (i.e., false negative), (3) incorrect location/position of finding, (4) incorrect severity of finding, (5) mention of comparison that is not present in the reference impression, and (6) omission of comparison describing a change from a previous study. The error types and error categories are summarized in Figure 2B. The instructions and interface presented to radiologists can be seen in Figures S1 and S2. The radiologist error scores on the 50 studies are shown in Figure S3.

We quantify metric-radiologist alignment using the Kendall rank correlation coefficient ( $\tau_b$ ) between metric scores and

**Figure 2. Example study of reports, and error types and categories**

(A) Example study of a test report and four metric-oracle reports corresponding to BLEU, BERTScore, CheXbert vector similarity, and RadGraph F1 that radiologists evaluate to identify errors. (B) Two error types and six error categories that radiologists identify for each pair of test report and metric-oracle report.

number of radiologist-reported errors in the reports (Figure 1D). We determine the metric-radiologist alignment from metric-oracle generations from 50 chosen studies on both a total error and significant error level. The coverage of pathologies, as determined by the CheXpert<sup>8</sup> labels in MIMIC-CXR, for the 50 randomly sampled reports is shown in Table S1. The per-radiologist Kendall rank correlation coefficients are listed in Table S2.

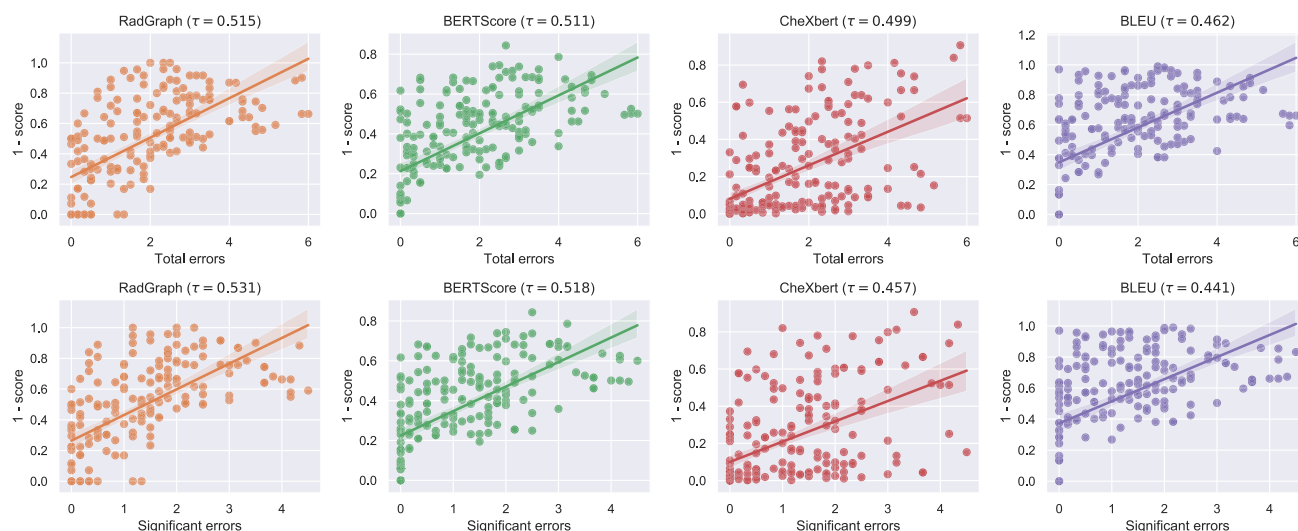
We find that RadGraph F1 and BERTScore are the metrics with the two highest alignments with radiologists. Specifically, RadGraph has a tau value of 0.515 (95% CI, 0.449 0.578) for total number of clinically significant and insignificant errors and 0.531 (95% CI, 0.465 0.584) for significant errors. BERTScore has a tau

value of 0.511 (95% CI, 0.429 0.584) for total number of clinically significant and insignificant errors and 0.518 (95% CI, 0.440 0.586) for significant errors. We find that CheXbert vector similarity is the third best metric under this evaluation with a 0.499 (95% CI, 0.417 0.576) tau value for total number of clinically significant and insignificant errors and 0.457 (95% CI, 0.370 0.538) for significant errors. Finally, BLEU has the worst alignment with a tau value of 0.462 (95% CI, 0.368 0.539) for total number of clinically significant and insignificant errors and 0.441 (95% CI, 0.350 0.521) for significant errors. From these results, we see that RadGraph and BERTScore are the metrics with closest alignment to radiologists. For the total number of clinically significant and insignificant errors, BERTScore has a significantly higher alignment than BLEU. Looking at significant errors, BERTScore and RadGraph have a significantly higher alignment than BLEU and, additionally, RadGraph has a significantly higher alignment than CheXbert. CheXbert, and BLEU have alignment with radiologists but are less concordant than the other two metrics. The metric-radiologist alignment graphs are shown in Figure 3.

**Failure modes of metrics**

In addition to evaluating the clinical relevance of metrics in terms of the total number of clinically significant and insignificant errors, we also examine the particular error categories of metric-oracles to develop a granular understanding of the failure modes of different metrics, as shown in Figure 4. We use the following six error categories as described earlier:

1. false prediction of finding



**Figure 3. Correlations between metric scores and radiologist scores**

Scatterplots and correlations between metric scores and radiologist scores of four metric-oracle generations from 50 studies, where radiologist scores are represented by the total number of clinically significant and insignificant errors (top row) and number of clinically significant errors (bottom row) identified by the radiologists. The translucent bands around the regression line represent 95% confidence intervals.

2. omission of finding
3. incorrect location/position of finding
4. incorrect severity of finding
5. mention of comparison that is not present in the reference impression
6. omission of comparison describing a change from a previous study and analyze the total number of errors and the number of clinically significant errors within each error category

BLEU exhibits a prominent failure mode in identifying false predictions of finding in reports. Metric-oracle reports with respect to BLEU produce more false predictions of finding than BERTScore and RadGraph in terms of both the total number of clinically significant and insignificant errors (0.807 average number of errors per report versus 0.477 and 0.427 for BERTScore and RadGraph) and the number of clinically significant errors (0.607 average number of errors per report versus 0.363 and 0.300 for BERTScore and RadGraph). BLEU exhibits a less prominent failure mode in identifying incorrect locations/positions of finding compared with CheXbert vector similarity. Metric-oracle reports with respect to BLEU have fewer incorrect locations/positions of finding than CheXbert in terms of both the total number of clinically significant and insignificant errors (0.113 average number of errors per report versus 0.227 for CheXbert) and the number of clinically significant errors (0.087 average number of errors per report versus 0.193 for CheXbert). These differences are statistically significant after accounting for multiple-hypothesis testing. Metric-oracle reports of the four metrics exhibit similar behavior in the other error categories, as the differences in number of errors are not statistically significant. The raw error counts and the statistics testing results for two-sample t tests and the Benjamini-Hochberg procedure

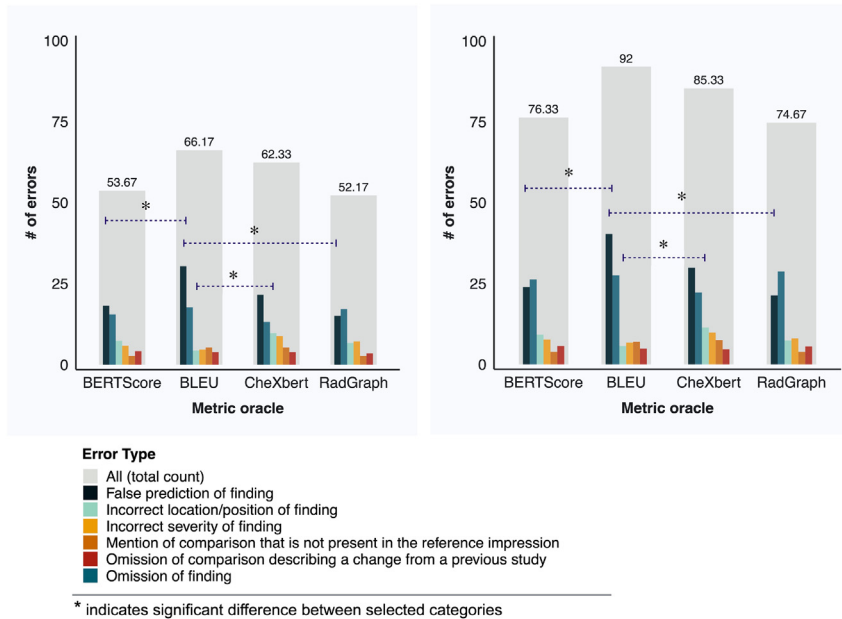
for accounting for multiple-hypothesis testing are shown in [Tables S3](#) and [S4](#).

### Measuring progress of prior methods in report generation

Using the four metrics, we evaluated the following state-of-the-art radiology report generation methods: M<sup>2</sup> Trans,<sup>11</sup> R2Gen,<sup>12</sup> CXR-RePaiR,<sup>13</sup> WCL,<sup>14</sup> and CvT2DistilGPT2.<sup>15</sup> As a baseline, we also implemented a random radiology report generation model, which retrieves a random report from the training set for each test report. The prior methods were trained to generate different sections of radiology reports: CXR-RePaiR generates the impression section, M<sup>2</sup> Trans the findings section, and R2Gen, WCL, and CvT2DistilGPT2 jointly the findings and impression sections. For each method, we compute metric values using the corresponding section(s) of radiology reports it generates as the ground-truth report to ensure accurate evaluation of the method. We also generated three versions of random baselines that retrieved different sections of the reports and compared each method with its corresponding random baseline. Because the impression section of radiology reports is an interpretation of the findings section, we can assume that both sections use the same medical vocabulary and style, which the report metrics evaluate. Conclusions about the report metrics drawn from the radiologist experiment and associated analyses, which used the impression section, can carry over to the evaluation of different report sections.

The performances of metric-oracle selection models, prior models, and random retrieval baselines on the MIMIC-CXR test set are shown in [Tables 1](#), [2](#), and [3](#), grouped by the report sections they generate. Note that the results are comparable within each table, but not across. With respect to the most radiologist-aligned metric RadGraph F1, among





**Figure 4. Distribution of errors across error categories for metric-oracle reports**

Distribution of errors across six error categories for metric-oracle reports corresponding to BERTScore, BLEU, CheXbert vector similarity, and RadGraph F1, in terms of the number of clinically significant errors (left) and the total number of clinically significant and insignificant errors (right). Statistical significance is determined using the Benjamini-Hochberg procedure with a false discovery rate (FDR) of 1% to correct for multiple-hypothesis testing.

impression-generating models, metric-oracle models significantly outperform real report generation models, achieving a maximum score of 0.677. Among findings-generating models, M<sup>2</sup> Trans performs the best (0.244). Among models that jointly generate the findings and impression sections, CvT2DistilGPT2 performs the best (0.154).

### Composite metric RadCliQ

To improve upon individual metrics, we propose a novel composite metric RadCliQ (radiology report clinical quality) that combines the four investigated metrics. We trained a model to predict the total number of clinically significant and insignificant errors that radiologists would assign to a report. The model input consisted of the four metric scores computed for each report. We applied zero-mean unit-variance normalization on the scores of each type of metric before passing the scores as model input. Prediction of the trained model therefore combines evaluations of BLEU, BERTScore, CheXbert vector similarity, and RadGraph F1.

We had 200 metric-oracle reports that were evaluated by radiologists, containing 50 metric-oracle reports corresponding to each of the four investigated metrics. These training data correspond to a subset of 50 studies from the MIMIC-CXR test set. We split our dataset by 8:2 into a development set (160 data points) and a test set (40 data points). On the development set, we conducted a cross-validation of 10-fold with 16 data points per validation fold to experiment with different model formulations for RadCliQ and build a fair comparison between RadCliQ and existing metrics.

Specifically, for each cross-validation setup, we built a normalizer with zero-mean and unit-variance and a linear regression model that took in the normalized metric values, on the cross-validation training set (144 data points). We then used the normalizer and regression model to normalize and generate predictions on the held-out validation set (16 data points). Finally, we computed the Kendall tau b correlation on the held-out set predictions with respect to the held-out set ground-truth radiologist total number

of errors. We also computed the Kendall tau b correlation for each existing metric. Across the 10 cross-validation setups, we computed the mean Kendall tau b correlations for the composite metric and existing metrics, and verified that the composite metric had stronger alignment with radiologists.

Our proposed model builds upon the standard linear regression by introducing constraints that improve its performance. Specifically, we require the negative of the coefficients to be non-negative and sum up to 1, resulting in a well-defined and interpretable convex function. Thus, we ensure that, when one metric score increases, while the others remain constant, the predicted number of errors will decrease. This property makes our model more sensitive to changes in individual metrics and thus more accurate in predicting error rates. Furthermore, the constraint makes the coefficients interpretable as weights, providing insights into the relative importance of each metric in predicting errors. To obtain the constrained coefficients, we use the convex optimization solver CVXPY, which guarantees global optimality and fast convergence. With this approach, we can effectively balance the trade-off between accuracy and interpretability, and obtain a robust and reliable model for error prediction.

After finalizing the model formulation, we fit the normalizer and composite metric model on the full development set and obtain RadCliQ. The coefficients were 0.000 for BLEU,  $-0.370$  for BERTScore,  $-0.253$  for CheXbert, and  $-0.377$  for RadGraph F1. The intercept value for the regression model was 0.000. Finally, on the held-out test set, the composite metric (RadCliQ) has higher Kendall tau b correlations than the other metrics, as shown in Table 4. This result indicates that RadCliQ has the strongest alignment with radiologists than any individual metric.

We used RadCliQ to evaluate all generations of metric-oracle models, prior models, and random retrieval baselines for the MIMIC-CXR test set. The metric scores are shown in Tables 5, 6, and 7. Among impression-generating models, the BERTScore metric-oracle model performs the best ( $-0.095$ ). CXR-RePair (1.642) outperforms the random retrieval baseline (1.755). Among findings-generating models, M<sup>2</sup> Trans performs the best (1.059). Among models that jointly generate findings and impression sections, CvT2DistilGPT2 performs the best (1.463).

**Table 1. Metric scores of impression-generating models, including metric-oracle models, CXR-RePaiR, and the random retrieval baseline model**

	BLEU	BERTScore	CheXbert vector similarity	RadGraph F1
BLEU metric-oracle	0.557*	0.661	0.689	0.476
BERTScore metric-oracle	0.491	0.721*	0.738	0.498
CheXbert metric-oracle	0.381	0.573	0.954*	0.403
RadGraph metric-oracle	0.366	0.541	0.739	0.677*
CXR-RePaiR	0.055	0.193	0.379	0.090
Random retrieval of impression	0.048	0.222	0.269	0.050

The 95% confidence interval and range of metric scores are available in [Table S5](#).

\*indicates the best-performing model.

## DISCUSSION

The purpose of this study was to investigate how to meaningfully measure progress in radiology report generation. We studied popular existing automated metrics and designed novel metrics, the RadGraph graph overlap metric and the composite metric RadCliQ, for report evaluation. We quantitatively determined the alignment of metrics with clinical radiologists and the reliability of metrics against specific failure modes, clarifying whether metrics meaningfully evaluate radiology reports and therefore can guide future research in report generation. We also showed that selecting the best-match report from a large corpus performs better on most metrics than the current state-of-the-art radiology report generation methods. Although the best-match method is unlikely to be clinically viable, it served as a useful tool to derive the RadCliQ composite metric developed in this study and could serve as a useful benchmark against which to evaluate report generation algorithms developed in the future.

The design of automated evaluation metrics that are aligned with manual expert evaluation has been a challenge for research in radiology report generation as well as medical report generation as a whole. Prior works have used metrics designed to improve upon n-gram matching<sup>28–32</sup> or include clinical awareness,<sup>8,9,11,13,24</sup> such as with BLEU<sup>28</sup> and CheXpert labels.<sup>8</sup> However, these evaluations nevertheless poorly approximate radiologists' evaluation of reports. The expressivity of prior metrics is often restricted to a curated set of medical conditions. Therefore, the quantitative investigation of metric-radiologist alignment conducted in this study is necessary for understanding whether these metrics meaningfully evaluate reports. Prior works have investigated the alignment between metrics and human judgment.<sup>23,33</sup> However, to the best of our knowledge, these works pose one of two limitations for radiology report evaluation: (1) they study metric alignment with humans for general image captioning, which does not involve radiology-specific terminology, a high prevalence of negation, or expert human evaluators,

and (2) they do not create a leveled comparison between metrics and radiologists, where metrics and radiologists assign scores to reports in identical experimental settings, or a granular understanding of metric behavior beyond the overall metric score. Our work builds a fair comparison between general natural language and clinically aware metrics and radiologists by providing them with the same set of information that is the reports and goes beyond metric scores to examine six granular failure modes of each metric. In addition, our work proposes a novel composite metric, RadCliQ, that aligns more strongly than any individual metric. We also show that current radiology report generation algorithms exhibit relatively low performance by all of these metrics.

To study metric-radiologist alignment, we designed *metric-oracles*: the reports selected from a large corpus with the highest metric score with respect to test reports. We had metrics and radiologists assign scores to the metric-oracles based on how well the metric-oracles match their respective test reports, and computed the alignment between metric and radiologist scores on the same reports. Pairing metric-oracles with test reports produces a narrower distribution of scores than using random reports. However, metric-oracles are necessary because comparisons with test reports are only reliable when the differences are small. If a random report, rather than a high-scoring report, was paired with the test report, the two reports could diverge to the extent that they were difficult to compare directly. In contrast, metric-oracles are comparable with test reports and therefore allow a meaningful evaluation of errors.

To generate metric-oracles, any report generation model is theoretically feasible. There are three main categories: the first generates free text based on semantics extracted from input chest X-ray images,<sup>16,34,35</sup> the second retrieves existing text that best matches input images from a report corpus,<sup>13,36</sup> and the third selects curated templates corresponding to a predefined set of abnormalities.<sup>10,37</sup> We chose to use retrieval-based models to generate metric-oracles because retrieval from a training report corpus produces a controlled output space,

**Table 2. Metric scores of findings-generating models, including M<sup>2</sup> Trans, and the random retrieval baseline model**

	BLEU	BERTScore	CheXbert vector similarity	RadGraph F1
M <sup>2</sup> Trans	0.220*	0.386*	0.452*	0.244*
Random retrieval of findings	0.123	0.323	0.235	0.105

The 95% confidence interval and range of metric scores are available in [Table S6](#).

\*indicates the best-performing model.

**Table 3. Metric scores of models that jointly generate findings and impression sections, including R2Gen, WCL, CvT2DistilGPT2, and the random retrieval baseline model**

	BLEU	BERTScore	CheXbert vector similarity	RadGraph F1
R2Gen	0.137	0.271	0.286	0.134
WCL	0.144*	0.275	0.309	0.143
CvT2DistilGPT2	0.143	0.280*	0.335*	0.154*
Random retrieval of jointly the findings and impression	0.100	0.256	0.190	0.090

The 95% confidence interval and range of metric scores are available in Table S7.

\*indicates the best-performing model.

instead of an unpredictable one produced by models that generate free text. Retrieval-based models also improve upon templating-based models in terms of flexibility and generalizability because the report corpus better captures real-world occurring conditions, combinations of conditions, and uncertainty. Furthermore, retrieval-based metric-oracle models outperformed existing report generation methods by a large margin.

By investigating the different categories of errors that radiologists identified in metric-oracle reports, we also uncovered specific metric failure modes that valuably inform the choice of metrics and interpretation of metric scores for evaluating generated reports. We find that BLEU performs worse than BERTScore and RadGraph in evaluating false prediction of finding. Yet, BLEU performs better than CheXbert vector similarity in evaluating incorrect position/location of finding. Therefore, RadGraph and BERTScore, which offer the strongest radiologist-alignment, also have better overall reliability against failure modes.

Using the individual metrics and RadCliQ, we also measured the progress of prior state-of-the-art models. Among impression-generating models, we find a significant performance gap between real report generation models and metric-oracle models, which represent the theoretical performance ceiling of retrieval-based methods on MIMIC-CXR for a given metric. This gap suggests that prior models in report generation still have significant room for improvement in creating high-quality reports that are useful to radiologists. We identify M<sup>2</sup> Trans to be the best findings-generating model and CvT2DistilGPT2 to be the best model that jointly generates findings and impression sections. Overall, RadGraph is the best individual metric to use for its strong alignment with radiologists and reliability across failure modes. RadCliQ, a composite metric, offers the strongest alignment with radiologists.

This study has several important limitations. A main limitation is the inter-observer variability in radiologist evaluation.

**Table 4. Kendall tau b correlations of individual metrics and the composite metric (RadCliQ) on the held-out test set of 40 data points with radiologist error annotations**

	Kendall tau b correlation
BLEU	0.414 (95% CI, 0.156 0.635)
BERTScore	0.505 (95% CI, 0.273 0.671)
CheXbert vector similarity	0.537 (95% CI, 0.330 0.717)
RadGraph F1	0.528 (95% CI, 0.357 0.687)
Composite metric (RadCliQ)	0.615 (95% CI, 0.450 0.749)*

\*indicates the best-aligned metric.

Although the evaluation scheme—the separation of clinically significant and insignificant errors, and the six error categories—was designed to be objective and consistent across radiologist evaluation, the same report often received varying scores between radiologists, a common occurrence in experiments that employ subjective ratings from clinicians. This suggests a potential limitation of the evaluation scheme used, but may also present an intrinsic problem with objective evaluation of radiology reports. Another limitation is the coverage of metrics. Although a variety of general and clinical natural language metrics are investigated, there exist other metrics in these two categories that may have different behaviors than the four investigated metrics. For instance, other text overlap-based metrics are commonly used in natural language generation beyond BLEU, such as CIDEr,<sup>31</sup> METEOR,<sup>30</sup> and ROUGE,<sup>32</sup> which may have better or worse radiologist-alignment and reliability than BLEU in report generation.

In this study, we determined that the novel metrics RadGraph F1 and RadCliQ meaningfully measure progress in radiology report generation and hence can guide future report generation models in becoming clinically indistinguishable from radiologists. We have open-sourced the code for computing the individual metrics and RadCliQ on reports in the hope of facilitating future research in radiology report generation.

## EXPERIMENTAL PROCEDURES

### Resource availability

#### Lead contact

The lead contact for this work is Pranav Rajpurkar ([pranav\\_rajpurkar@hms.harvard.edu](mailto:pranav_rajpurkar@hms.harvard.edu)).

#### Materials availability

Does not apply.

#### Data and code availability

- Original data for the radiologist error annotations have been deposited to the Radiology Report Expert Evaluation (ReXVal) Dataset<sup>38</sup> with credentialed access at <https://physionet.org/content/rexval-dataset/1.0.0/> (<https://doi.org/10.13026/2fp8-qr71>). The radiology report data used in the study are available with credentialed access at: <https://physionet.org/content/mimic-cxr-jpg/2.0.0/> (<https://doi.org/10.13026/8360-t248>). Credentialed access can be obtained via an application to PhysioNet.
- The code for computing the composite metric RadCliQ and individual metrics is made publicly available at: <https://doi.org/10.5281/zenodo.7579952>.<sup>39</sup>

### Datasets

We used the MIMIC-CXR dataset to conduct our study. The MIMIC-CXR dataset<sup>25–27</sup> is a de-identified and publicly available dataset containing chest X-ray



**Table 5. RadCliQ scores of impression-generating models, including metric-oracle models, CXR-RePaiR, and the random retrieval baseline model**

	RadCliQ
BLEU metric-oracle	0.081
BERTScore metric-oracle	-0.095*
CheXbert metric-oracle	0.052
RadGraph metric-oracle	-0.020
CXR-RePaiR	1.642
Random retrieval of impression	1.755

Lower is better. The 95% confidence interval and range of metric scores are available in Table S5.

\*indicates the best-performing model.

images and semi-structured radiology reports from the Beth Israel Deaconess Medical Center Emergency Department. There are 227,835 studies with 177,110 images conducted on 65,379 patients. We used the recommended train/validation/test split. We pooled the train and validation splits as the training report corpus from which metric-oracles are retrieved and used the test split as the set of ground-truth reports. We preprocessed the reports by filtering nan reports and extracting the impression and findings sections of reports, which contain key observations and conclusions drawn by radiologists. We follow the section extraction code provided in the MIMIC-CXR repository. In the training set, 187,383 impression reports, 153,415 findings reports, and 214,344 findings and impression joint reports are available. In the test set, 2,191 impression reports, 1,597 findings reports, and 2,192 findings and impression joint reports are available. We refer to the impression section when discussing reports for the metric-oracle reports and failure modes. When evaluating prior models, we use either the impression section, the findings section, or jointly the findings and impression sections based on what the prior model generates.

### Metric-oracle reports

We constructed metric-oracle reports for four metrics. These include BLEU,<sup>28</sup> BERTScore,<sup>29</sup> CheXbert vector similarity (s\_emb),<sup>9</sup> and a novel metric RadGraph<sup>24</sup> F1. BLEU and BERTScore are general natural language metrics for measuring the similarity between machine-generated and human-generated texts. BLEU computes n-gram overlap and is representative for the family of text overlap-based natural language generation metrics such as CIDEr,<sup>31</sup> METEOR,<sup>30</sup> and ROUGE.<sup>32</sup> BERTScore has been proposed for capturing contextual similarity beyond exact textual matches. CheXbert vector similarity and RadGraph F1 are metrics designed to measure the correctness of clinical information. CheXbert vector similarity computes the cosine similarity between the CheXbert model embeddings for machine-generated and human-generated radiology reports. The CheXbert model is designed to evaluate radiology-specific information but its training supervision was limited to 14 pathologies. To address this limitation, we propose the use of the knowledge graph of the report to represent arbitrarily diverse radiology-specific information. We design a novel metric, RadGraph F1, that computes the overlap in clinical entities and relations that RadGraph extracts from machine- and human-generated reports. The four metrics are detailed in the “textual based

**Table 6. RadCliQ scores of findings-generating models, including M<sup>2</sup> Trans and the random retrieval baseline model**

	RadCliQ
M <sup>2</sup> Trans	1.059*
Random retrieval of findings	1.553

Lower is better. The 95% confidence interval and range of metric scores are available in Table S6.

\*indicates the best-performing model.

and natural language generation performance metrics” subsection and the “clinically aware performance metrics” subsection.

For every test report, we generated the matching metric-oracle report by selecting the highest scoring report according to each of the four investigated metrics from the training set. We specifically used the impression section of the report. As an example of our setup, for the test report of “No acute cardiopulmonary process. Bilateral low lung volumes with crowding of bronchovascular markings and bibasilar atelectasis,” the metric-oracle retrieved with respect to BERTScore was: “No acute cardiopulmonary process. Low lung volumes and bibasilar atelectasis,” while the metric-oracle retrieved with respect to RadGraph F1 was: “No acute cardiopulmonary process. Bilateral low lung volumes.”

Using metric-oracles as the candidate reports as opposed to using other strategies such as randomly sampling reports offers two primary advantages: (1) metric-oracles are sufficiently accurate for radiologists to pinpoint specific errors and not be bogged down by candidate reports that are not remotely similar to the test reports, and (2) metric-oracles allow us to analyze where certain metrics fail since the reports are the hypothetical top retrievals.

### Radiologist scoring criteria

In this work, we develop a scoring system for radiologists to evaluate the quality of candidate reports. The goals of our scoring system are to be objective, limit radiologist bias, and change linearly with report quality. To this end, scores are determined by counting the number of errors that candidate reports make where types of errors are broken down into six different categories. By explicitly defining each error category, we clarify what should be classified as an error. Following ACR’s RADPEER<sup>40</sup> program for peer review, we differentiate between clinically significant and clinically insignificant errors. The detailed scoring criteria allow us to analyze report quality based on the accuracy of its findings and the clinical impact of its mistakes.

### Textual-based and natural language generation performance metrics

In this study we make use of two natural language generation metrics: BLEU and BERTScore. The BLEU scores were computed as BLEU-2 bigrams with the fast\_bleu library for parallel scoring. BERTScore uses the contextual embeddings from a BERT model to compute similarity of two text sequences. We used the bert\_score library directly and used the “distilroberta-base” version of the model. We used the unscaled scores for metric-oracle retrieval and the baseline-scaled scores for all other analyses.

### Clinically aware performance metrics

In addition to traditional natural language generation metrics, we also investigated metrics that were designed to capture clinical information in radiology reports. Since radiology reports are a special form of structured text that communicate diagnostics information, their quality depends highly on the correctness of clinical objects and descriptions, which is not a focus of traditional natural language metrics. To address this gap, the CheXbert labeler (which is improved from the CheXpert labeler<sup>5,9</sup> and RadGraph,<sup>24</sup> were developed to parse radiology reports. We investigated whether they could be used as clinically aware metrics. We defined a metric as the cosine similarity between CheXbert model embeddings of the generated report and test report. We extracted the CLS token output embeddings before the final dropout layer and prediction heads. We used the implementation here: <https://github.com/stanfordmlgroup/CheXbert>. In prior literature, a common way of comparing generated reports against ground-truth reports is to compute the micro- and/or macro-F1/precision/recall scores averaged over 14 observation labels outputted by CheXpert/CheXbert. For instance, CXR-RePaiR computes the macro-average F1 over 14 observations to evaluate generations. Positive observation labels are treated as positive, while other labels, including negative, uncertain, and blank labels, are treated as negative. However, this approach limits the evaluation of generated reports to 14 observations and discrete outputs. Because radiology reports can reference diverse observations beyond the 14 and contain more nuanced semantics about the observations, we decided to use the CheXbert model embedding before the final classifiers, which produce 14 outputs to capture a more accurate representation of the report. Our design is supported by a prior work that uses the same CheXbert model embeddings as deep representations of radiology reports

**Table 7. RadCliQ scores of models that jointly generate findings and impression sections, including R2Gen, WCL, CvT2DistilGPT2, and the random retrieval baseline model**

	RadCliQ
R2Gen	1.552
WCL	1.511
CvT2DistilGPT2	1.463*
Random retrieval of jointly the findings and impression	1.726

Lower is better. The 95% confidence interval and range of metric scores are available in Table S7.

\*indicates the best-performing model.

for heart failure patient mortality prediction.<sup>41</sup> In their experiments, they also found that these hidden features led to better prediction performance than the features of 14 observations extracted by CheXpert. This suggests that the model embeddings may preserve more information about the reports than the final model output of observation labels. CXR-RePaiR also adopts the same formulation of CheXbert vector similarity as a report evaluation metric. We propose a novel metric as the overlap in parsed RadGraph graph structures: the RadGraph entity and relation F1 score. RadGraph is an approach for parsing radiology reports into knowledge graphs containing entities (nodes) and relations (edges), which can capture radiology concept dependencies and semantic meaning. We used the model checkpoint as provided here: <https://physionet.org/content/radgraph/1.0.0/>,<sup>27</sup> and inference code as provided here: <https://github.com/dwadden/dygiepp>,<sup>42</sup> to generate RadGraph entities and relations on generated and test reports.

#### Retrieval-based metric-oracle models

To generate metric-oracle reports, the most immediate attempt is to adopt methods akin to those for multi-label classification tasks. Namely, we can curate a set of medical conditions and obtain radiologist annotations for each condition over a training set of reports. Then, we can train a classifier that outputs the likelihood of having each condition given an X-ray image, and proceed to select the corresponding report templates for conditions with high likelihood.<sup>10</sup> Some more nuanced approaches paraphrase the curated templates after selection.<sup>37</sup> The attempt at templating for report generation is well-grounded in abundant experience in multi-label image classification as well as its highly controlled output space. However, its flaw is also prominent, in that it is restricted to a manually curated predefined set of medical conditions and report templates. It does not generalize to unseen or complex conditions, express combinations of conditions, or capture uncertainty in diagnoses. The CheXbert labeler, for instance, can classify 13 conditions and the no-finding observation.<sup>9</sup> This set is representative of common medical observations but not comprehensive. Therefore, while we may define a larger set of conditions with the help of radiologists, manual curation and templating are nevertheless too inflexible for optimizing with respect to automated metrics. To generate reports of higher quality, we consider matching reports more closely onto test reports. We can do so by either generating new text from scratch or retrieving free text from an existing corpus of reports written by radiologists, given an X-ray image.<sup>34,36</sup> Out of the two approaches, retrieval-based methods have the advantage of a controlled output space that is the set of training report corpus. Therefore, in this study, we use retrieval-based methods to generate metric-oracle reports.

#### RadGraph metric-oracle model entities and relations match

The RadGraph F1 metric-oracle model retrieves reports with the highest F1 score match in terms of entities and relations. Specifically, we treat two entities as matched if their tokens (words in the original report) and labels (entity type) match. We treat two relations as matched if their start and end entities match and the relation type matches. These criteria are consistent with what the RadGraph authors have done. For combining entities and relations, we take the average of F1 score of entity match and relation match, respectively. We generated RadGraph entities and relations for each report in the training and

test corpora. We implemented the metric-oracle model by finding, for each report in the test set, which report in the training set is the best match based on the average of entity and relation F1 scores. For reports without nonzero F1 score matches, we used the most frequent report in the training set, “No acute cardiopulmonary process,” as the metric-oracle report in the radiologist experiment.

#### Statistical analysis

##### Metric-radiologist alignment

The alignment of metrics with radiologists’ scoring was determined using the Kendall tau b correlation coefficient. We construct 95% bootstrap confidence intervals by creating 1,000 resamples with replacement where each resample size is the number of studies (50). In this calculation, the number of errors is the mean number across all raters. We additionally test for the difference in correlations between two metrics by counting the number of positive correlation differences computed on 1,000 resamples of metric scores. The fraction of positive correlation differences indicate the p value for the null hypothesis that there is zero difference in correlation between the metrics.

##### Metric failure modes

We conduct one-sided two-sample t tests on pairs of metrics’ error counts for total number of clinically significant and insignificant errors and clinically significant errors within each of the six error categories. We assume equal population variances for the t tests. We take the error count of one radiologist and one study as one data point. Because there are 6 radiologists and 50 studies, we have 300 data points per metric for either total number of clinically significant and insignificant errors or clinically significant errors and for 1 error category. With 4 metrics, there are 12 unique pairs of 2 different metrics for one-sided two-sample t tests with  $(300 + 300 - 2 = 598)$  degrees of freedom. We use the Benjamini-Hochberg procedure with a false discovery rate of 1% to account for multiple-hypothesis testing on 12 tests within an error type and an error category, and determine the significance of a metric having a more-/less-prominent failure mode compared with other metrics.

##### Prior models evaluation

To evaluate performance of metric-oracle models and prior state-of-the-art models, we construct 95% bootstrap confidence intervals by taking 5,000 resamples with replacement of metric scores assigned to generated reports.

##### Composite metric RadCliQ

The composite metric model used to predict the total number of errors was evaluated using the Kendall tau b statistical test. This test produces a tau value correlation coefficient and a corresponding p value, which was used to determine the significance of the result ( $p < 0.01$ ). The same statistical comparison procedure described in [metric-radiologist alignment](#) with 5,000 resamples was used to compare the correlation of RadCliQ with that of other metrics.

The analyses were performed using statsmodels, scikit-learn, and SciPy packages in Python.

#### Implementation of prior report generation methods

We used the following implementations of prior methods in radiology report generation: M<sup>2</sup> Trans, <https://github.com/ymsiura/ifcc><sup>11,42</sup>; R2Gen, <https://github.com/cuhksz-nlp/R2Gen><sup>12</sup>; CXR-RePaiR, <https://github.com/rajpurkarlab/CXR-RePaiR><sup>13</sup>; WCL, <https://github.com/zxslp/WCL><sup>14</sup>; CvT2DistilGPT2, <https://github.com/aehrc/cvt2distilgpt2><sup>15</sup>. CXR-RePaiR was trained to generate the impression section through retrieval. M<sup>2</sup> Trans and CvT2DistilGPT2 were trained to generate the findings section, with maximum sequence lengths of 128 and 60, respectively. R2Gen and WCL were trained to jointly generate the findings and impression sections, with maximum sequence lengths of 60 and 100, respectively. We did not shorten or cut off any part of the actual reports when evaluating our report generation method to avoid creating a problem in our evaluation process. If we had shortened the reports, it could have allowed a generation method to be trained to only produce very short reports that lack important information, but still receive good evaluation scores. Thus, to ensure accurate evaluation, we did not truncate the ground-truth reports. We used these prior methods to generate reports for all available studies in the test set. For impression only and findings only generations, there are fewer test reports and metric outputs. This is considered acceptable, because there are still sufficiently large numbers of reports for a reliable estimate of model performance. For each study ID, if the model generated multiple reports corresponding to different X-ray images for the same study, we used the generated report corresponding to the

anterior-posterior or posterior-anterior views if any were present. If both were present, we randomly chose a report out of the two. If neither was present, we randomly chose a report out of the available reports corresponding to other views. Among variations of CXR-RePaIR, we chose CXR-RePaIR-2 to be consistent with their original study.<sup>13</sup>

### SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.patter.2023.100802>.

### ACKNOWLEDGMENTS

We thank M.A. Endo MD for helpful review and feedback on the radiologist evaluation survey design and the manuscript. Support for this work was provided in part by the Medical Imaging Data Resource Center (MIDRC) under contracts 75N92020C00008 and 75N92020C00021 from the National Institute of Biomedical Imaging and Bioengineering (NIBIB) of the National Institutes of Health.

### AUTHOR CONTRIBUTIONS

P.R. conceived the study. F.Y., M.E., and R.K. contributed to the design, implementation, and analyses of all aspects of this study. I.P., A.T., E.P.R., E.K.U.N.F., H.M.H.L., and V.K.V. provided suggestions on the setup of the radiologist evaluation survey and provided annotations in the radiologist evaluation process. Z.S.H.A. contributed to the design of the illustrations and figures. A.Y.N., C.P.L., and V.K.V. provided guidance on the study. P.R. supervised the study. All authors approved the final version.

### DECLARATION OF INTERESTS

The authors declare no competing non-financial interests but the following competing financial interests: I.P. is a consultant for MD.ai and Diagnosticos da America (Dasa). C.P.L. serves on the board of directors and is a shareholder of Bunkerhill Health. He is an advisor and option holder for GalileoCDS, Sirona Medical, Adra, and Kheiron. He is an advisor to Sixth Street and an option holder in whiterabbit.ai. His research program has received grant or gift support from Carestream, Clarity, GE Healthcare, Google Cloud, IBM, IDEXX, Hospital Israelita Albert Einstein, Kheiron, Lambda, Lunit, Microsoft, Nightingale Open Science, Nines, Philips, Subtle Medical, VinBrain, Whiterabbit.ai, the Paustentbach Fund, the Lowenstein Foundation, and the Gordon and Betty Moore Foundation.

### INCLUSION AND DIVERSITY

We support inclusive, diverse, and equitable conduct of research.

Received: November 8, 2022

Revised: February 3, 2023

Accepted: June 29, 2023

Published: August 3, 2023

### REFERENCES

- Rajpurkar, P., Chen, E., Banerjee, O., and Topol, E.J. (2022). AI in health and medicine. *Nat. Med.* 28, 31–38.
- Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Židek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589.
- Rajpurkar, P., Joshi, A., Pareek, A., Ng, A.Y., and Lungren, M.P. (2021). CheXternal: Generalization of Deep Learning Models for Chest X-ray Interpretation to Photos of Chest X-rays and External Clinical Settings. In *Proceedings of the Conference on Health, Inference, and Learning (Association for Computing Machinery)*, pp. 125–132.
- Jin, B.T., Palleti, R., Shi, S., Ng, A.Y., Quinn, J.V., Rajpurkar, P., and Kim, D. (2022). Transfer learning enables prediction of myocardial injury from continuous single-lead electrocardiography. *J. Am. Med. Inf. Assoc.* 29, 1908–1918.
- Rajpurkar, P., and Lungren, M.P. (2023). The Current and Future State of AI Interpretation of Medical Images. *N. Engl. J. Med. Overseas. Ed.* 388, 1981–1990.
- Tiu, E., Talius, E., Patel, P., Langlotz, C.P., Ng, A.Y., and Rajpurkar, P. (2022). Expert-level detection of pathologies from unannotated chest X-ray images via self-supervised learning. *Nat. Biomed. Eng.* 6, 1399–1406.
- Agarwal, N., Moehring, A., Pranav, R., and Salz, T. (2023). Combining Human Expertise with Artificial Intelligence: Experimental Evidence from Radiology.
- Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Illcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., et al. (2019). CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1901.07031>.
- Smit, A., Jain, S., Rajpurkar, P., Pareek, A., Ng, A.Y., and Lungren, M.P. (2020). CheXbert: Combining Automatic Labelers and Expert Annotations for Accurate Radiology Report Labeling Using BERT. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2004.09167>.
- Pino, P., Parra, D., Besa, C., and Lagos, C. (2021). Clinically Correct Report Generation from Chest X-Rays Using Templates. In *Machine Learning in Medical Imaging (Springer)*, pp. 654–663.
- Miura, Y., Zhang, Y., Tsai, E.B., Langlotz, C.P., and Jurafsky, D. (2020). Improving Factual Completeness and Consistency of Image-to-Text Radiology Report Generation. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2010.10042>.
- Chen, Z., Song, Y., Chang, T.-H., and Wan, X. (2020). Generating Radiology Reports via Memory-driven Transformer. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2010.16056>.
- Endo, M., Krishnan, R., Krishna, V., Ng, A.Y., and Rajpurkar, P. (2021). Retrieval-Based Chest X-Ray Report Generation Using a Pre-trained Contrastive Language-Image Model. In *Machine Learning for Health (PMLR)*, pp. 209–219.
- Yan, A., He, Z., Lu, X., Du, J., Chang, E., Gentili, A., McAuley, J., and Hsu, C.-N. (2021). Weakly Supervised Contrastive Learning for Chest X-Ray Report Generation. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2109.12242>.
- Nicolson, A., Dowling, J., and Koopman, B. (2022). Improving Chest X-Ray Report Generation by Leveraging Warm-Starting. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2201.09405>.
- Zhou, H.-Y., Chen, X., Zhang, Y., Luo, R., Wang, L., and Yu, Y. (2022). Generalized radiograph representation learning via cross-supervision between images and free-text radiology reports. *Nat. Mach. Intell.* 4, 32–40.
- Ramesh, V., Chi, N.A., and Rajpurkar, P. (2022). Improving Radiology Report Generation Systems by Removing Hallucinated References to Non-existent Priors. In *Proceedings of Machine Learning for Health (PMLR)*, pp. 456–473.
- Jeong, J., Tian, K., Li, A., Hartung, S., Behzadi, F., Calle, J., Osayande, D., Pohlen, M., Adithan, S., and Rajpurkar, P. (2023). Multimodal Image-Text Matching Improves Retrieval-based Chest X-Ray Report Generation. In *Proceedings of Medical Imaging with Deep Learning (MIDL)*.
- Li, M., Lin, B., Chen, Z., Lin, H., Liang, X., and Chang, X. (2023). Dynamic Graph Enhanced Contrastive Learning for Chest X-Ray Report Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3334–3343.
- Yang, S., Wu, X., Ge, S., Zhou, S.K., and Xiao, L. (2022). Knowledge matters: Chest radiology report generation with general and specific knowledge. *Med. Image Anal.* 80, 102510.
- Nguyen, H., Nie, D., Badamdorj, T., Liu, Y., Zhu, Y., Truong, J., and Cheng, L. (2021). Automated Generation of Accurate & Fluent Medical X-ray Reports. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3552–3569.

22. Hossain, M.Z., Sohel, F., Shiratuddin, M.F., and Laga, H. (2018). A Comprehensive Survey of Deep Learning for Image Captioning. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1810.04020>.
23. Boag, W., Kané, H., Rawat, S., Wei, J., and Goehler, A. (2021). A Pilot Study in Surveying Clinical Judgments to Evaluate Radiology Report Generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*.
24. Jain, S., Agrawal, A., Saporta, A., Truong, S.Q.H., Duong, D.N., Bui, T., Chambon, P., Zhang, Y., Lungren, M.P., Ng, A.Y., et al. (2021). RadGraph: Extracting Clinical Entities and Relations from Radiology Reports. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2106.14463>.
25. Johnson, A.E.W., Pollard, T.J., Berkowitz, S.J., Greenbaum, N.R., Lungren, M.P., Deng, C.-Y., Mark, R.G., and Horng, S. (2019). MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Sci. Data* 6, 317–318.
26. Johnson, A.E.W., Pollard, T.J., Greenbaum, N.R., Lungren, M.P., Deng, C.-Y., Peng, Y., Lu, Z., Mark, R.G., Berkowitz, S.J., and Horng, S. (2019). MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1901.07042>.
27. Goldberger, A.L., Amaral, L.A., Glass, L., Hausdorff, J.M., Ivanov, P.C., Mark, R.G., Mietus, J.E., Moody, G.B., Peng, C.-K., and Stanley, H.E. (2000). PhysioBank, PhysioToolkit, and PhysioNet. *Circulation* 101, E215–E220. <https://doi.org/10.1161/01.cir.101.23.e215>.
28. Papineni, K., Roukos, S., Ward, T., and Zhu, W.J. (2002). BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*.
29. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., and Artzi, Y. (2019). BERTScore: Evaluating Text Generation with BERT. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1904.09675>.
30. Lavie, A., and Agarwal. (2007). Meteor: an automatic metric for MT evaluation with high levels of correlation with human judgments. In *Proceedings of the Second Workshop on Statistical Machine Translation*.
31. Vedantam, R., Zitnick, C.L., and Parikh, D. (2014). CIDEr: Consensus-based Image Description Evaluation. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1411.5726>.
32. Lin, C.-Y. (2004). ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, pp. 74–81.
33. Anderson, P., Fernando, B., Johnson, M., and Gould, S. (2016). SPICE: Semantic Propositional Image Caption Evaluation. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1607.08822>.
34. Monshi, M.M.A., Poon, J., and Chung, V. (2020). Deep learning in generating radiology reports: A survey. *Artif. Intell. Med.* 106, 101878.
35. Zhou, Y., Huang, L., Zhou, T., Fu, H., and Shao, L. (2021). Visual-Textual Attentive Semantic Consistency for Medical Report Generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3985–3994.
36. Wang, X., Zhang, Y., Guo, Z., and Li, J. (2018). ImageSem at ImageCLEF 2018 Caption Task: Image Retrieval and Transfer Learning.
37. Li, C.-Y., Liang, X., Hu, Z., and Xing, E.P. (2019). Knowledge-Driven Encode, Retrieve, Paraphrase for Medical Image Report Generation. *AAAI* 33, 6666–6673.
38. Yu, F., Endo, M., Krishnan, R., Pan, I., Tsai, A., Reis, E.P., Fonseca, E.K.U., Lee, H., Shakeri, Z., Ng, A., et al. (2023). Radiology Report Expert Evaluation (ReXVal) Dataset. <https://doi.org/10.13026/2fp8-qr71>.
39. Yu, K., and Rayan-Krishnan. (2023). rajpurkarlab/CXR-Report-Metric: v1.1.0. <https://doi.org/10.5281/zenodo.7936166>.
40. Goldberg-Stein, S., Frigini, L.A., Long, S., Metwalli, Z., Nguyen, X.V., Parker, M., and Abujudeh, H. (2017). ACR RADPEER Committee White Paper with 2016 Updates: Revised Scoring System, New Classifications, Self-Review, and Subspecialized Reports. *J. Am. Coll. Radiol.* 14, 1080–1086.
41. Lee, H.G., Sholle, E., Beecy, A., Al'Aref, S., and Peng, Y. (2021). Leveraging Deep Representations of Radiology Reports in Survival Analysis for Predicting Heart Failure Patient Mortality. *Proc. Conf.* 2021, 4533–4538.
42. Wadden, D., Wennberg, U., Luan, Y., and Hajishirzi, H. (2019). Entity, Relation, and Event Extraction with Contextualized Span Representations.

**Patterns, Volume 4**

## **Supplemental information**

### **Evaluating progress in automatic chest**

#### **X-ray radiology report generation**

**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**



# Supplemental information

**Stanford**

**Study Instructions**

**Goal:** Your goal is to judge the diagnostic accuracy of candidate report impressions based on a reference impression section.

**Study Setup:** We're looking to determine the accuracy of report impressions generated by hypothetical models that have high performance on popular metrics. You will be shown 50 studies with a reference impression section and various candidate report impression sections.


For each candidate, you will be asked how many clinically significant errors and clinically insignificant errors the report makes. An error can be one of the following:

- False prediction of finding
- Omission of finding
- Incorrect location/position of finding
- Incorrect severity of finding
- Mention of comparison that is **not** present in the reference impression
- Omission of comparison describing a change from a previous study

*Note: Some candidate reports will often be the same. In these cases, make sure to give the same scores to the candidates.*

When you are ready to begin the study, press the next page button below. You can always go back to this page to review the error guidelines.

→

Powered by Qualtrics 

**Figure S1. Radiologist evaluation survey instructions and interface on Qualtrics.**

## Stanford

Study #1

**Reference impression: Multiple chronic appearing left-sided rib fractures. No pneumothorax. Blunting of the costophrenic angle on the right likely represents pleural scarring and a small effusion, not significantly changed from \_\_\_.**

**Candidate 1: Blunting of the right costophrenic angle may be due to small pleural effusion .**

How many of the following errors does this report make:

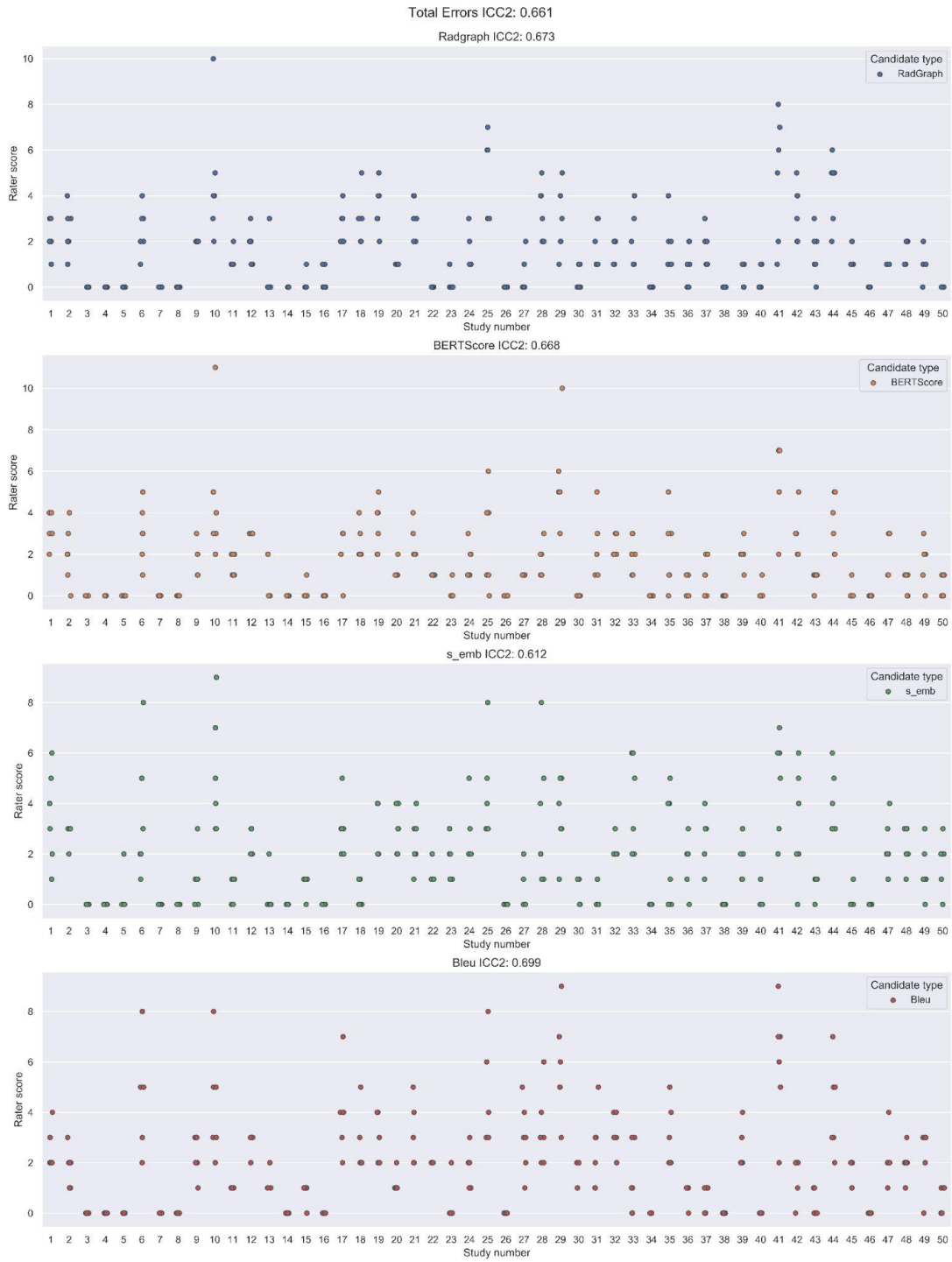
	Clinically significant	Clinically insignificant
False prediction of finding	<input type="text"/>	<input type="text"/>
Omission of finding	<input type="text"/>	<input type="text"/>
Incorrect location/position of finding	<input type="text"/>	<input type="text"/>
Incorrect severity of finding	<input type="text"/>	<input type="text"/>
Mention of comparison that is <b>not</b> present in the reference impression	<input type="text"/>	<input type="text"/>
Omission of comparison describing a change from a previous study	<input type="text"/>	<input type="text"/>

**Candidate 2: No focal consolidation. Minimal blunting of the left costophrenic angle may represent a trace pleural effusion.**

How many of the following errors does this report make:

	Clinically significant	Clinically insignificant
False prediction of finding	<input type="text"/>	<input type="text"/>
Omission of finding	<input type="text"/>	<input type="text"/>
Incorrect location/position of finding	<input type="text"/>	<input type="text"/>
Incorrect severity of finding	<input type="text"/>	<input type="text"/>
Mention of comparison that is <b>not</b> present in the reference impression	<input type="text"/>	<input type="text"/>
Omission of comparison describing a change	<input type="text"/>	<input type="text"/>

**Figure S2. Interface for evaluating a pair of a test report (denoted as “Reference Impression”) and a metric-oracle report (denoted as “Candidate 1”).** The survey asks radiologists to input the number of clinically significant and insignificant errors for six error categories.



**Figure S3. Dotplot of the radiologist total error scores on the 50 studies and corresponding intraclass correlation.** Candidate scores are split up by metric-oracle method. Each dot represents a single radiologist's score for a candidate report.

**Table S1. Coverage of pathologies, as determined by the CheXpert labels in MIMIC-CXR, for the 50 randomly sampled reports in the radiologist experiment.** The counts listed for “No Finding” refer to explicit labels for “No Finding”: namely 1.0 for positive mentions, 0.0 for negative mentions and -1.0 for uncertain mentions. There are edge cases where explicit labels for “No Finding” were suppressed, for instance if a pathology not included in the CheXpert label set was mentioned. If we also count studies where no other pathology had a positive label along with explicit positive labels for “No Finding” then we have 15 counts, instead of 14.

	Positive mentions	Negative mentions	Uncertain mentions
Atelectasis	8	0	2
Cardiomegaly	9	2	1
Consolidation	3	1	1
Edema	6	4	2
Enlarged Cardiomediastinum	1	0	0
Fracture	2	0	0
Lung Lesion	3	0	2
Lung Opacity	16	0	0
No Finding	14	0	0
Pleural Effusion	9	5	1
Pleural Other	3	0	0
Pneumonia	2	4	9
Pneumothorax	1	6	0
Support Devices	8	0	0

**Table S2. Per-radiologist Kendall rank correlation coefficient (tau-b) values quantifying metric-radiologist alignment.**

	Radiologist 1	Radiologist 2	Radiologist 3	Radiologist 4	Radiologist 5	Radiologist 6
BERTScore sig. and insig. errors	0.454 [95% CI 0.374 0.527]	0.441 [95% CI 0.362 0.517]	0.535 [95% CI 0.458 0.605]	0.442 [95% CI 0.348 0.521]	0.454 [95% CI 0.373 0.533]	0.511 [95% CI 0.424 0.590]
BERTScore sig. errors	0.410 [95% CI 0.322 0.490]	0.337 [95% CI 0.251 0.418]	0.540 [95% CI 0.469 0.600]	0.456 [95% CI 0.367 0.533]	0.419 [95% CI 0.342 0.497]	0.487 [95% CI 0.400 0.567]
RadGraph sig. and insig. errors	0.505 [95% CI 0.426 0.573]	0.499 [95% CI 0.429 0.566]	0.539 [95% CI 0.474 0.602]	0.491 [95% CI 0.415 0.554]	0.507 [95% CI 0.428 0.579]	0.451 [95% CI 0.364 0.526]
RadGraph sig. errors	0.474 [95% CI 0.398 0.556]	0.351 [95% CI 0.255 0.437]	0.540 [95% CI 0.476 0.596]	0.523 [95% CI 0.451 0.588]	0.499 [95% CI 0.424 0.566]	0.426 [95% CI 0.336 0.507]
BLEU sig. and insig. errors	0.398 [95% CI 0.307 0.483]	0.422 [95% CI 0.338 0.497]	0.475 [95% CI 0.392 0.550]	0.386 [95% CI 0.286 0.469]	0.412 [95% CI 0.319 0.492]	0.463 [95% CI 0.374 0.547]
BLEU sig. errors	0.345 [95% CI 0.248 0.434]	0.255 [95% CI 0.167 0.338]	0.472 [95% CI 0.389 0.542]	0.399 [95% CI 0.302 0.476]	0.357 [95% CI 0.264 0.444]	0.412 [95% CI 0.315 0.501]
CheXbert sig. and insig. errors	0.440 [95% CI 0.348 0.526]	0.424 [95% CI 0.332 0.509]	0.478 [95% CI 0.392 0.558]	0.503 [95% CI 0.420 0.581]	0.489 [95% CI 0.406 0.566]	0.451 [95% CI 0.364 0.530]
CheXbert sig. errors	0.392 [95% CI 0.300 0.475]	0.263 [95% CI 0.168 0.350]	0.451 [95% CI 0.362 0.529]	0.428 [95% CI 0.338 0.509]	0.411 [95% CI 0.330 0.490]	0.407 [95% CI 0.318 0.488]

**Table S3(a). Radiologist evaluation of metric-oracles in terms of total number of errors in six error categories, averaged over 6 radiologists and 50 studies.**

	Error 1	Error 2	Error 3	Error 4	Error 5	Error 6	Total
BLEU	0.807	0.550	0.113	0.133	0.140	0.097	1.840
CheXbert	0.597	0.443	0.227	0.197	0.150	0.093	1.707
BERTScore	0.477	0.523	0.183	0.153	0.077	0.113	1.527
RadGraph	0.427	0.573	0.147	0.160	0.077	0.110	1.493

**Table S3(b). Radiologist evaluation of metric-oracles in terms of number of clinically significant errors in six error categories, averaged over 6 radiologists and 50 studies.**

	Error 1	Error 2	Error 3	Error 4	Error 5	Error 6	Total
BLEU	0.607	0.353	0.087	0.093	0.107	0.077	1.323
CheXbert	0.430	0.263	0.193	0.176	0.107	0.077	1.247
BERTScore	0.363	0.310	0.147	0.117	0.053	0.083	1.073
RadGraph	0.300	0.343	0.133	0.143	0.053	0.070	1.043



**Table S4(a). Multiple hypothesis testing outputs in terms of *total number of clinically significant and insignificant errors in false prediction of finding*.** Significance of BLEU having a *more prominent failure mode* than BERTScore and RadGraph F1 in terms of *total number of clinically significant and insignificant errors in false prediction of finding*, as determined by the Benjamini-Hochberg Procedure with False Discovery Rate (FDR) of 1%.

	One-sided two-sample t test p-value	Benjamini-Hochberg Procedure critical value with False Discovery Rate (FDR) of 1%	Whether result is significant
BLEU > CheXbert	3.79e-3	2.50e-3	N
BLEU > BERTScore	9.50e-6	1.67e-3	Y
BLEU > RadGraph	1.07e-7	8.33e-4	Y
CheXbert > BLEU	9.96e-1	8.33e-3	N
CheXbert > BERTScore	7.65e-2	4.17e-3	N
CheXbert > RadGraph	6.39e-3	3.33e-3	N
BERTScore > BLEU	1.00e0	9.17e-3	N
BERTScore > CheXbert	9.51e-1	6.67e-3	N
BERTScore > RadGraph	2.24e-1	5.00e-3	N
RadGraph > BLEU	1.00e0	1.00e-2	N
RadGraph > CheXbert	9.94e-1	7.50e-3	N
RadGraph > BERTScore	7.76e-1	5.83e-3	N

**Table S4(b). Multiple hypothesis testing outputs in terms of *clinically significant errors in false prediction of finding*.** Significance of BLEU having a *more prominent failure mode* than BERTScore and RadGraph F1 in terms of *clinically significant errors in false prediction of finding*, as determined by the Benjamini-Hochberg Procedure with False Discovery Rate (FDR) of 1%.

	One-sided two-sample t test p-value	Benjamini-Hochberg Procedure critical value with False Discovery Rate (FDR) of 1%	Whether result is significant
BLEU > CheXbert	3.68e-3	2.50e-3	N
BLEU > BERTScore	1.48e-4	1.67e-3	Y
BLEU > RadGraph	6.44e-7	8.33e-4	Y
CheXbert > BLEU	9.96e-1	8.33e-3	N
CheXbert > BERTScore	1.34e-1	5.00e-3	N
CheXbert > RadGraph	9.77e-3	3.33e-3	N
BERTScore > BLEU	1.00e0	9.17e-3	N
BERTScore > CheXbert	8.66e-1	5.83e-3	N
BERTScore > RadGraph	1.33e-1	4.17e-3	N
RadGraph > BLEU	1.00e0	1.00e-2	N
RadGraph > CheXbert	9.90e-1	7.50e-3	N
RadGraph > BERTScore	8.67e-1	6.67e-3	N

**Table S4(c). Multiple hypothesis testing outputs in terms of *total number of clinically significant and insignificant errors in incorrect location/position of finding*.** Significance of BLEU having a *less prominent failure mode* than CheXbert vector similarity in terms of *total number of clinically significant and insignificant errors in incorrect location/position of finding*, as determined by the Benjamini-Hochberg Procedure with False Discovery Rate (FDR) of 1%.

	One-sided two-sample t test p-value	Benjamini-Hochberg Procedure critical value with False Discovery Rate (FDR) of 1%	Whether result is significant
BLEU < CheXbert	4.83e-4	8.33e-4	Y
BLEU < BERTScore	1.60e-2	2.50e-3	N
BLEU < RadGraph	1.51e-1	5.00e-3	N
CheXbert < BLEU	1.00e0	1.00e-2	N
CheXbert < BERTScore	8.90e-1	7.50e-3	N
CheXbert < RadGraph	9.89e-1	9.17e-3	N
BERTScore < BLEU	9.84e-1	8.33e-3	N
BERTScore < CheXbert	1.10e-1	3.33e-3	N
BERTScore < RadGraph	8.63e-1	6.67e-3	N
RadGraph < BLEU	8.49e-1	5.83e-3	N
RadGraph < CheXbert	1.14e-2	1.67e-3	N
RadGraph < BERTScore	1.37e-1	4.17e-3	N

**Table S4(d). Multiple hypothesis testing outputs in terms of *clinically significant errors in incorrect location/position of finding*.** Significance of BLEU having a *less prominent failure mode* than CheXbert vector similarity in terms of *clinically significant errors in incorrect location/position of finding*, as determined by the Benjamini-Hochberg Procedure with False Discovery Rate (FDR) of 1%.

	One-sided two-sample t test p-value	Benjamini-Hochberg Procedure critical value with False Discovery Rate (FDR) of 1%	Whether result is significant
BLEU < CheXbert	2.74e-4	8.33e-4	Y
BLEU < BERTScore	2.07e-2	1.67e-3	N
BLEU < RadGraph	5.75e-2	3.33e-3	N
CheXbert < BLEU	1.00e0	1.00e-2	N
CheXbert < BERTScore	9.25e-1	6.67e-3	N
CheXbert < RadGraph	9.67e-1	8.33e-3	N
BERTScore < BLEU	9.79e-1	9.17e-3	N
BERTScore < CheXbert	7.53e-2	4.17e-3	N
BERTScore < RadGraph	6.65e-1	5.83e-3	N
RadGraph < BLEU	9.42e-1	7.50e-3	N
RadGraph < CheXbert	3.32e-2	2.50e-3	N
RadGraph < BERTScore	3.35e-1	5.00e-3	N

**Table S5. The average, 95% confidence interval and range of metric scores of impression-generating models, including metric-oracle models, CXR-RePaiR and the random retrieval baseline model.**

	BLEU	BERTScore	CheXbert vector similarity	RadGraph F1	RadCliQ
BLEU metric-oracle	0.557 [95% CI 0.547 0.567] Range [0.009, 1.000]	0.661 [95% CI 0.652 0.670] Range [-0.266, 1.000]	0.689 [95% CI 0.678 0.699] Range [-0.088, 1.000]	0.476 [95% CI 0.464 0.489] Range [0.000, 1.000]	0.081 [95% CI 0.044 0.118] Range [-1.441, 2.567]
BERTScore metric-oracle	0.491 [95% CI 0.479 0.503] Range [0.000, 1.000]	0.721 [95% CI 0.714 0.729] Range [0.033, 1.000]	0.738 [95% CI 0.728 0.748] Range [-0.050, 1.000]	0.498 [95% CI 0.486 0.511] Range [0.000, 1.000]	-0.095 [95% CI -0.129 -0.062] Range [-1.441, 2.162]
CheXbert metric-oracle	0.381 [95% CI 0.367 0.395] Range [0.000, 1.000]	0.573 [95% CI 0.563 0.585] Range [-0.225, 1.000]	0.954 [95% CI 0.952 0.957] Range [-0.013, 1.000]	0.403 [95% CI 0.390 0.417] Range [0.000, 1.000]	0.052 [95% CI 0.017 0.086] Range [-1.441, 2.543]
RadGraph metric-oracle	0.366 [95% CI 0.356 0.377] Range [0.000, 1.000]	0.541 [95% CI 0.533 0.549] Range [-0.102, 1.000]	0.739 [95% CI 0.729 0.748] Range [-0.028, 1.000]	0.677 [95% CI 0.668 0.686] Range [0.000, 1.000]	-0.020 [95% CI -0.051 0.009] Range [-1.441, 2.500]
CXR-RePaiR	0.055 [95% CI 0.053 0.057] Range [0.000, 0.383]	0.193 [95% CI 0.188 0.198] Range [-0.402, 0.633]	0.379 [95% CI 0.370 0.387] Range [-0.146, 0.973]	0.090 [95% CI 0.086 0.095] Range [0.000, 1.000]	1.642 [95% CI 1.625 1.659] Range [-0.645, 2.904]
Random retrieval of impression	0.048 [95% CI 0.044 0.053] Range [0.000, 1.000]	0.222 [95% CI 0.216 0.227] Range [-0.326, 1.000]	0.269 [95% CI 0.259 0.279] Range [-0.226, 1.000]	0.050 [95% CI 0.045 0.055] Range [0.000, 1.000]	1.755 [95% CI 1.733 1.776] Range [-1.441, 3.111]

**Table S6. The average, 95% confidence interval and range of metric scores of findings-generating models, including M<sup>2</sup> Trans.**

	BLEU	BERTScore	CheXbert vector similarity	RadGraph F1	RadCliQ
M <sup>2</sup> Trans	0.220 [95% CI 0.214 0.224] Range [0.001, 0.643]	0.386 [95% CI 0.380 0.391] Range [-0.127, 0.738]	0.452 [95% CI 0.441 0.463] Range [-0.089, 0.987]	0.244 [95% CI 0.238 0.250] Range [0.000, 0.823]	1.059 [95% CI 1.037 1.083] Range [-0.701, 2.279]
Random retrieval of findings	0.123 [95% CI 0.119 0.127] Range [0.000, 0.778]	0.323 [95% CI 0.318 0.328] Range [-0.138, 0.845]	0.235 [95% CI 0.225 0.245] Range [-0.217, 0.920]	0.105 [95% CI 0.101 0.109] Range [0.000, 0.866]	1.553 [95% CI 1.533 1.573] Range [-0.938, 2.537]

**Table S7. The average, 95% confidence interval and range of metric scores of models that jointly generate findings and impression sections, including R2Gen, WCL and CvT2DistilGPT2.**

	BLEU	BERTScore	CheXbert vector similarity	RadGraph F1	RadCliQ
R2Gen	0.137 [95% CI 0.133 0.141] Range [0.000, 0.826]	0.271 [95% CI 0.266 0.276] Range [-0.318, 0.853]	0.286 [95% CI 0.276 0.295] Range [-0.204, 0.993]	0.134 [95% CI 0.130 0.138] Range [0.000, 0.883]	1.552 [95% CI 1.533 1.571] Range [-1.038, 3.063]
WCL	0.144 [95% CI 0.141 0.148] Range [0.000, 0.612]	0.275 [95% CI 0.271 0.279] Range [-0.332, 0.745]	0.309 [95% CI 0.300 0.318] Range [-0.214, 0.991]	0.143 [95% CI 0.139 0.147] Range [0.000, 0.694]	1.511 [95% CI 1.493 1.529] Range [-0.608, 3.060]
CvT2DistilGPT2	0.143 [95% CI 0.139 0.147] Range [0.000, 0.826]	0.280 [95% CI 0.275 0.285] Range [-0.351, 0.842]	0.335 [95% CI 0.326 0.344] Range [-0.224, 0.993]	0.154 [95% CI 0.149 0.159] Range [0.000, 0.883]	1.463 [95% CI 1.443 1.484] Range [-1.027, 2.909]
Random retrieval of jointly the findings and impression	0.100 [95% CI 0.097 0.103] Range [0.000, 0.498]	0.256 [95% CI 0.252 0.261] Range [-0.383, 0.809]	0.190 [95% CI 0.182 0.198] Range [-0.281, 0.940]	0.090 [95% CI 0.087 0.093] Range [0.000, 0.726]	1.726 [95% CI 1.710 1.741] Range [-0.724, 3.067]