# AN X-RAY IS WORTH 15 FEATURES: SPARSE AU-TOENCODERS FOR INTERPRETABLE RADIOLOGY RE-PORT GENERATION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Radiological services are experiencing unprecedented demand, leading to increased interest in automating radiology report generation. Existing Vision-Language Models (VLMs) suffer from hallucinations, lack interpretability, and require expensive fine-tuning. Sparse Autoencoders (SAEs) have been shown to provide a principled approach to reverse-engineer a model's internal activations into discrete, verifiable components. Thus, we introduce SAE-Rad, the first instance of using mechanistic interpretability techniques explicitly for a downstream multi-modal reasoning task. SAE-Rad uses a novel SAE architecture to decompose latent representations from a pre-trained vision transformer into human-interpretable features. These features are then labelled using an off-theshelf language model and compiled into a full report for each image, eliminating the need for fine-tuning large models for this task. On the MIMIC-CXR dataset, SAE-Rad achieves competitive radiology-specific metrics compared to state-ofthe-art models while using significantly fewer computational resources for training. Qualitative analysis reveals that SAE-Rad learns meaningful visual concepts and generates reports aligning closely with expert interpretations. Our results suggest that SAEs can enhance multimodal reasoning in healthcare, providing a more interpretable alternative to existing VLMs.

032

005 006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

### 1 INTRODUCTION

Radiological services are essential to modern clinical practice, with demand rising rapidly. In the UK, the NHS performs over 43 million radiological procedures annually (Lewis et al., 2021), costing over £2 billion, and demand for scans more than doubled between 2012 and 2019 (NHS England & NHS Improvement, 2019). Much of the cost covers agency, bank, and overtime staff, but a national imaging strategy deemed it unsustainable. (NHS England & NHS Improvement, 2019).
Consequently, there's growing interest in (semi)-automating tasks like radiology report generation, augmentation, and summarization to assist clinicians (Zhu et al., 2024; Chen et al., 2024; Pérez-García et al., 2024), spurred by advances in multimodal text-vision modelling techniques.

Recent architectures combining vision encoders with pretrained Large Language Models (LLMs) 041 into multimodal Vision-Language Models (VLMs) excel in visual and language tasks (Liu et al., 042 2024b; 2023a; Li et al., 2024; Lin et al., 2023; Liu et al., 2023b). VLMs have been applied to 043 healthcare tasks, including radiology report generation (Hyland et al., 2023; Bannur et al., 2024; 044 Chen et al., 2024; Stock et al., 2024; Yang et al., 2024), often by mapping image representations into 045 the LLM's token embedding space (Pérez-García et al., 2024). Despite improvements from scaling 046 VLMs (Yang et al., 2024), hallucinations and disagreements with domain experts remain common 047 (Yildirim et al., 2024; Jeblick et al., 2023; Lee et al., 2023; Tanno et al., 2024a). Hallucinations 048 are unavoidable in LLMs (Kalai & Vempala, 2024; Xu et al., 2024), and this limits VLMs for radiology reporting, but other important considerations remain. For current state-of-the-art systems, it is necessary to finetune a multi-billion parameter LLM to perform visual instruction tuning (Liu 051 et al., 2024b), which is computationally intensive. Additionally, reports generated by VLMs may not be faithful to the underlying computations of the image encoder – we seek to design a framework 052 that is verifiably faithful to the image model by reverse engineering it's computations. This may yield more interpretable results, engendering more trust in automated radiology reporting models.



Figure 1: SAE-Rad overview. Panel A: We learn a set of sparsely activating features by training 072 a Sparse Autoencoder (SAE) on class tokens produced by a radiology-image encoder. Panel B: We retrieve the corresponding reference reports for highest activating images for a feature, from 074 which we can produce text descriptions of each feature. Panel C: We pass a new image through the radiology-image encoder and SAE encoder to retrieve the highest activating features. Text de-076 scriptions of these features are subsequently used by a pretrained large language model (LLM) to 077 generate a detailed radiology report.

054

056

061

063

064

067

068

071

073

075

080 To this end, we introduce SAE-Rad (Fig. 1), a framework which leverages sparse autoencoders 081 (SAEs) to learn human-interpretable features from image latents extracted from a pre-trained radiology image encoder (Pérez-García et al., 2024). We propose to automatically label these features 083 using pre-trained LLMs (Bricken et al., 2023) and then compile them into descriptive radiology reports. To the best of our knowledge, this is the first instance of using SAEs explicitly for a down-084 stream task requiring multi-modal reasoning. Our contributions are as follows: 1) Novel SAE: We 085 propose a hybrid architecture which builds on gated SAEs to achieve accurate reconstructions with comparable sparsity to state-of-the-art techniques; 2) Multimodal reasoning pipeline: We develop a 087 framework based on our SAE which leverages automated interpretability to label SAE features us-088 ing paired text data and ultimately produce radiology reports given an image; 3) Interpretability in *healthcare*: Our case study on healthcare data demonstrates the utility of mechanistic interpretabil-090 ity for downstream tasks in a relatively smaller and significantly more homogeneous dataset than 091 most general-domain text or natural-image datasets.

092 093

094 095

#### 2 **RELATED WORK**

**Multimodal reasoning** Multimodal reasoning methods like ScienceQA (Lu et al., 2022) intro-096 duced multimodal chain-of-thought (CoT) by zero-shot prompting models to generate both ratio-097 nales and answers. Multimodal-CoT (MM-CoT) (Zhang et al., 2023) extended this with a two-stage 098 framework that separates rationale generation and answer inference using two models of the same architecture. Duty-Distinct CoT (DDCoT) (Zheng et al., 2023) further factorizes rationales by decom-100 posing the initial question into sub-questions answered by a vision-language model (VLM). Other 101 divide-and-conquer approaches decompose questions into sub-questions but often require training 102 task-specific visual question generation (VQG) models and additional scoring models (Selvaraju 103 et al., 2020; Uehara et al., 2022; Wang et al., 2022). IdealGPT (You et al., 2023) iteratively decom-104 poses queries and uses a VLM to answer sub-questions, repeating the process if confidence is low. 105 Unlike these methods that rely on decomposing questions or generating rationales through additional models, our approach directly extracts and interprets features from pre-trained image encoders. This 106 enables faithful and transparent multimodal reasoning without the need for extensive fine-tuning or 107 supplementary VQG models.

108 **Radiological VLMs** Several works have finetuned or otherwise trained specialized foundation 109 models for radiological applications including Med-flamingo (Moor et al., 2023), Med-PaLM M (Tu 110 et al., 2024), LLava-Med (Li et al., 2024), Med-Gemini (Yang et al., 2024), Rad-DINO (Pérez-111 García et al., 2024), MAIRA-1 (Hyland et al., 2023), R2gengpt (Wang et al., 2023b), and Radiology-112 GPT (Liu et al., 2023c). Regarding radiology report synthesis, many models produce both the 'findings' and 'impression' sections of the reports (Chen et al., 2020; Jin et al., 2024; Yan et al., 113 2023), whilst others only produce the 'impression' section (Bannur et al., 2023), or the 'findings' 114 section (Tu et al., 2024; Miura et al., 2020; Delbrouck et al., 2022; Tanida et al., 2023; Nicolson 115 et al., 2023). As noted by others (Hyland et al., 2023; Yu et al., 2023; Jeong et al., 2024), studies 116 examining all three settings found that the choice of section(s) to report significantly affects the 117 performance metrics, making comparison between results difficult. As such, we focus on the most 118 common setting of producing the 'findings' section.

119

120 Mechanistic interpretability Bricken et al. (2023) demonstrated that SAEs could recover 121 monosemantic features by training on the residual stream of small transformers. The gated SAE 122 was a Pareto improvement over the baseline SAE in terms of sparsity (L0) and the loss recovered 123 (Rajamanoharan et al., 2024). There was concern that SAEs would not scale to frontier transformers 124 until recent work by Templeton et al. (2024), which trained SAEs on Claude 3 Sonnet and discovered 125 many monosemantic features. Similarly, Gao et al. (2024) trained SAEs on GPT-4 and introduced the Top-K training variant, while Lieberum et al. (2024) released Gemma Scope, a comprehensive 126 suite of SAEs for Gemma 2 models. Contemporaneously, SAEs were trained on the class tokens of 127 a CLIP vision transformer (Fry, 2024), InceptionV1 (Gorton, 2024), the conditioning embeddings 128 of diffusion models (Daujotas, 2024), and the vision transformer of a pathology foundation model 129 (Le et al., 2024). Other methods use an overcomplete basis and LASSO regressions on CLIP em-130 beddings (Bhalla et al., 2024). However, in all cases the discovered language/visual features were 131 not used to perform a downstream multimodal reasoning task.

132 133 134

135

136

137 138

139

# 3 BACKGROUND

In this section we give a brief overview of mechanistic interpretability, SAEs, and gated SAEs, before introducing our SAE-Rad framework.

### 3.1 MECHANISTIC INTERPRETABILITY AND SPARSE AUTOENCODERS (SAES)

140 Mechanistic interpretability Mechanistic interpretability research aims to identify, understand, 141 and verify the algorithms that an ML model implements by reverse engineering a model's compu-142 tations into human-interpretable components (Olah et al., 2020; Rajamanoharan et al., 2024). Clas-143 sical approaches attempted to achieve this by analysing the firing patterns of individual neurons, 144 which were interpreted as possible 'concept representations'. However, this was broadly ineffective 145 as neurons can be *polysemantic*, meaning that a single neuron may fire on many unrelated concepts (Rajamanoharan et al., 2024; Bolukbasi et al., 2021; Elhage et al., 2022a). Polysemantic neurons 146 are believed to arise during training due to the composition of both the linear representation and 147 superposition hypotheses. 148

- 149
- Linear representation and superposition hypotheses Motivated by a number of findings which 150 suggest that concept representations are linear (Gurnee et al., 2023; Olah et al., 2020; Park et al., 151 2023), the linear representation hypothesis states that neural networks represent concepts as direc-152 tions in *activation space* (Nanda et al., 2024). This hypothesis can be thought of as being composed 153 of two properties: 1) Linearity: That is, features are represented as directions; and 2) Decom-154 posability: We can understand neural network outputs as a composition of multiple independently 155 understandable features (Elhage et al., 2022b). The superposition hypothesis states that for an in-156 termediate representation of dimension n, neural networks will encode  $M \gg n$  concepts as linear 157 directions (Rajamanoharan et al., 2024; Elhage et al., 2022b). These directions form an overcom-158 plete basis of the activation space and must therefore necessarily overlap with each other. However, 159 a single input will only activate a sparse subset of these concepts, leading to minimal interference between the (non-orthogonal) concept directions (Gurnee et al., 2023; Rajamanoharan et al., 2024). 160 Recent work (Bricken et al., 2023) has proposed using SAEs to take features out of superposition 161 and learn monosemantic interpretable representations.

Sparse autoencoders (SAEs) SAEs attempt to 'undo' superposition by learning the sparse over-complete basis (Mallat & Zhang, 1993; Rajamanoharan et al., 2024) (or dictionary) of the activation space induced by superposition. SAEs attempt to learn both the concept directions and a sparse vector of coefficients for the inputs, that reflect how much each concept is activated for each input (Cunningham et al., 2023; Bricken et al., 2023). We will henceforth refer to such sparse vectors of coefficients as 'feature activations'.

We begin by defining the 'baseline SAE' described by Bricken et al. (2023). Let n be the dimension of the input and output (typically the input is the residual stream of a transformer, and the output is its reconstruction (Elhage et al., 2021)), and m be the SAE hidden layer dimension. Let s be the size of the dataset. Then given encoder weights and biases  $W^{\text{enc}} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b}^{\text{enc}} \in \mathbb{R}^m$ , and decoder weights and biases  $W^{\text{dec}} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{b}^{\text{dec}} \in \mathbb{R}^n$ , the encoding and decoding operations for a dataset  $X \in \mathbb{R}^{s,n}$  are

$$\mathbf{h}(\mathbf{x}) := \operatorname{ReLU}(W^{\operatorname{enc}}(\mathbf{x} - \mathbf{b}^{\operatorname{dec}}) + \mathbf{b}^{\operatorname{enc}})$$
(1)

$$\hat{\mathbf{x}}(\mathbf{h}(\mathbf{x})) := W^{\text{dec}}\mathbf{h}(\mathbf{x}) + \mathbf{b}^{\text{dec}}.$$
(2)

The loss function is then

174

175

176

181

182

183

184 185 186

201 202

203 204

212 213 214

$$\mathcal{L}(\mathbf{x}) := \frac{1}{|X|} \sum_{\mathbf{x} \in X} \left[ \|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{h}(\mathbf{x}))\|_2^2 + \lambda \|\mathbf{h}(\mathbf{x})\|_1 \right],\tag{3}$$

where  $\lambda$  is an L1 sparsity coefficient. The first term is a reconstruction error measured by the squared distance between the input and its reconstruction, and the second is an L1 regularization loss to induce sparsity. The column-wise norm of the decoder  $W^{dec}$  is constrained to unit norm during training (Bricken et al., 2023).

187 As can be seen in Eq. (3), SAEs jointly optimize two opposing objectives: 1) Reconstruction fidelity 188 and 2) L1 regularization as a proxy for sparsity (as measured by L0). This means the SAE is free 189 to trade-off some reconstruction fidelity in order to perform better on the sparsity penalty. One consequence of this is shrinkage (Wright & Sharkey, 2024). That is, for a fixed decoder, the sparsity 190 penalty pushes the feature activations h(x) towards zero whilst the squared distance loss encourages 191 h(x) to be large enough in order to produce high quality reconstructions. Thus, the standard SAE 192 will systematically underestimate the optimal magnitude of feature activations (and simply rescaling 193 these does not necessarily overcome this bias) (Rajamanoharan et al., 2024; Wright & Sharkey, 194 2024). 195

Rajamanoharan et al. (2024) propose a gated SAE which separates the encoding procedure into two
tasks: 1) Detecting which features should activate for a given input (this requires an L1 penalty if the
features are to be sparse); and 2) Estimating the magnitude of the feature activations (this does **not**require an L1 loss; indeed, including this loss here introduces a shrinkage bias). The architecture of
the gated encoder is

$$\mathbf{h}(\mathbf{x}) := \underbrace{\mathbb{I}[W^{\text{gate}}(\mathbf{x} - \mathbf{b}^{\text{dec}}) + \mathbf{b}_{\text{gate}}}_{\mathbf{h}_{\text{gate}}(\mathbf{x})} \odot \underbrace{\text{ReLU}(W^{\text{mag}}(\mathbf{x} - \mathbf{b}^{\text{dec}}) + \mathbf{b}_{\text{mag}})}_{\mathbf{h}_{\text{mag}}(\mathbf{x})}, \tag{4}$$

where  $\mathbb{I}[\bullet > 0]$  is an element-wise Heaviside step function and  $\odot$  is element-wise multiplication. The  $\mathbf{h}_{gate}$  sub-function learns which features should activate for a given input and  $\mathbf{h}_{mag}$  estimates the magnitude of activations for these features. Here,  $\pi_{gate}$  is referred to as the  $\mathbf{h}_{gate}$  sub-function's 'preactivations'. To minimize the number of additional parameters required,  $W^{mag}$  shares the same feature directions as  $W^{gate}$ , and is defined as  $W^{mag}_{i,j} := \exp(\mathbf{r}^{mag}_i) \cdot W^{gate}_{i,j}$ , where  $\mathbf{r}^{mag} \in \mathbb{R}^m$  is a vectorvalued scaling parameter. Letting  $RA(\cdot) := ReLU(\pi_{gate}(\cdot))$  denote the rectified pre-activations of the gating sub-function, the loss function is defined as

$$\mathcal{L}(\mathbf{x}) := \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{h}(\mathbf{x}))\|_{2}^{2}}_{\mathcal{L}_{reconstruct}} + \underbrace{\lambda \|\mathbf{R}\mathbf{A}(\mathbf{x})\|_{1}}_{\mathcal{L}_{sparsity}} + \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}_{frozen}(\mathbf{R}\mathbf{A}(\mathbf{x}))\|_{2}^{2}}_{\mathcal{L}_{aux}},$$
(5)

where  $\hat{\mathbf{x}}_{\text{frozen}}$  is a fixed copy of the decoder so that gradients from the auxiliary loss  $\mathcal{L}_{\text{aux}}$  do not back-propagate to the decoder weights or bias terms. The auxiliary term  $\mathcal{L}_{\text{aux}}$  ensures that  $\mathbf{h}_{\text{gate}}$  216 correctly identifies features necessary for reconstruction, as its (positive) pre-activations must be 217 able to reproduce the input. The sparsity term  $\mathcal{L}_{\text{sparsity}}$  applies an L1 penalty to the rectified pre-218 activations (and thus sparsity is only imposed on the gating sub-function), and the reconstruction 219 term serves the same function as in Eq. (3).

#### 4 SAE-RAD

220 221

222

224

225

226 227

228 229

232 233

238

239

243

244

245 246

248 249

In this section we introduce SAE-Rad. First, we describe the autoencoder architecture, which is based on the gated SAE described in Section 3.2. Then, we describe our end-to-end radiology report generation pipeline.

#### 4.1 SAE ARCHITECTURE

It was recently demonstrated that a lower overall SAE loss is achievable without constraining the L2 230 norm of the decoder weights, or centering the input based on the decoder bias (Conerly et al., 2024). 231 Concretely, they use the same decoder as in Eq. (2) and define the encoder as

$$\mathbf{h}(\mathbf{x}) := \operatorname{ReLU}(W^{\operatorname{enc}}\mathbf{x} + \mathbf{b}^{\operatorname{enc}}).$$
(6)

Note that the input x is no longer centered by subtracting the decoder bias  $b^{dec}$  as in Eq. (1). The 234 sparsity penalty in the loss also now includes the L2 norm of the columns of the decoder  $W^{dec}$ . 235 We present a novel proof of the equivalence between the equation below and the loss in Eq. (3) in 236 Appendix K: 237

$$\mathcal{L}(\mathbf{x}) := \|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{h}(\mathbf{x}))\|_{2}^{2} + \lambda \sum_{i} \mathbf{h}_{i}(\mathbf{x}) \cdot \|W_{\cdot,i}^{\text{dec}}\|_{2}.$$
(7)

The feature activation for a feature *i* is then  $\mathbf{h}_i(\mathbf{x}) \cdot \|W_{\cdot,i}^{\text{dec}}\|_2$ . The 'concept directions' are the unit-normalized decoder vectors  $\frac{W_{\cdot,i}^{\text{dec}}}{\|W_{\cdot,i}^{\text{dec}}\|_2}$ . 240 241

242

The SAE-Rad sparse autoencoder is a hybrid architecture which combines a gated encoder layer with unconstrained decoder norms. Its encoder is defined as

$$\mathbf{h}(\mathbf{x}) := \mathbb{I}[W^{\text{gate}}\mathbf{x} + \mathbf{b}^{\text{gate}} > 0] \odot \text{ReLU}(W^{\text{mag}}\mathbf{x} + \mathbf{b}^{\text{mag}}), \tag{8}$$

247 and the training objective given by

$$\mathcal{L}(\mathbf{x}) := \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{h}(\mathbf{x}))\|_{2}^{2}}_{\mathcal{L}_{\text{reconstruct}}} + \underbrace{\lambda \sum_{i} RA_{i}(\mathbf{x}) \cdot \|W_{\cdot,i}^{\text{dec}}\|_{2}}_{\mathcal{L}_{\text{sparsity}}} + \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}(RA(\mathbf{x}))\|_{2}^{2}}_{\mathcal{L}_{\text{aux}}}.$$
(9)

250 251

There are four architectural differences between the gated SAE and SAE-Rad: 1) The L2 norm 253 of the decoder weights are not constrained to unit norm; 2) The L2 norm of the decoder weights 254 are included in the sparsity loss term; 3) The input is not centered by subtracting the decoder bias 255 term; 4) We do not leverage stop gradient operations/decoder copying – this means that we allow gradients to back-propagate to the decoder parameters from the auxiliary loss. SAEs are trained 256 to balance reconstruction fidelity and sparsity, placing them on a Pareto frontier of these objectives 257 (Rajamanoharan et al., 2024). We therefore evaluate SAEs based on the L0 norm and the mean-258 squared error loss as defined in Eq. (3). We demonstrate that this novel SAE architecture outperforms 259 the architecture described by Conerly et al. (2024) on our dataset, achieving both a lower L0 and a 260 lower mean-squared error. As such, SAE-Rad is a more Pareto-optimal architecture in the present 261 context. Details of this comparison can be found in Appendix B.1. 262

263 264

### 4.2 SAE-RAD – AUTOMATED RADIOLOGY REPORTING PIPELINE

265 In this section we describe our pipeline to automate the task of radiology report generation. 266 Concretely, a radiographic image x is passed through a pre-trained and frozen vision encoder 267  $f_{img}(\cdot) : \mathbf{x} \mapsto \mathbf{z}$  to produce an image latent  $\mathbf{z}$ . We leverage the hybrid SAE architecture described in Section 4.1 to learn feature directions  $\frac{W_{\cdot,i}^{\text{dec}}}{\|W_{\cdot,i}^{\text{dec}}\|_2}$  and their associated activations  $\mathbf{h}(\mathbf{z}) \cdot \|W^{\text{dec}}\|_2$  from 268 269 the latents.

270 In order to generate a text-based report, we produced plain-English descriptions of the learned sparse 271 dictionary (i.e., a description of what each feature direction represents). To do this we performed 272 automated interpretability (Bricken et al., 2023) by using a pre-trained and frozen LLM to analyse 273 the ground-truth radiology reports of the highest activating images for each feature. Let  $\mathbf{X}_{ ext{highest}}^{(i)}$ 274 represent the set of images with the largest feature activations for feature *i*. For each image  $\mathbf{x} \in$ 275  $\mathbf{X}_{\text{highest}}^{(i)}$  there is an associated ground-truth radiology report  $r(\mathbf{x})$ . We collect these reports into 276 a set  $\mathcal{R}^{(i)} := \{r(\mathbf{x}) | \mathbf{x} \in \mathbf{X}_{highest}^{(i)}\}$ . We then utilized a pre-trained and frozen language model 277  $f_{\text{descriptor}}: \mathcal{R} \mapsto d$  to generate a description  $d^{(i)}$  for a feature *i* by analyzing the set  $\mathcal{R}^{(i)}$  as  $d^{(i)} =$ 278 279  $f_{\text{descriptor}}(\mathcal{R}^{(i)})$ . This process yielded a set of feature descriptions  $\{d^{(i)}\}_{i=1}^{M}$ , where M is the total number of features learned by the SAE. 280

For a new image x we identified the set of active features I(x) based on a threshold  $\tau$  as 282

$$I(\mathbf{x}) := \{ i | \mathbf{h}_i(f_{\text{img}}(\mathbf{x})) \cdot \| W^{\text{dec}} \|_2 > \tau \}.$$
(10)

The automated radiology report  $R(\mathbf{x})$  is then generated by a pre-trained and frozen LLM from the descriptions of the active features

$$R(\mathbf{x}) = f_{\text{generator}}(\{d^{(i)} | i \in I(\mathbf{x})\}).$$
(11)

#### 5 EXPERIMENTS

291 Our overarching hypotheses are that: 1) SAE features capture meaningful visual concepts even in ho-292 mogeneous datasets (such as is the case for chest radiographs); 2) The visual concepts captured by an 293 SAE can be appropriately described by a pre-trained LLM by use of automated interpretability techniques with paired text data; 3) Natural language descriptions of visual features in the latent space of 295 an SAE can be composed into high-quality radiology reports without explicit use (or training/finetuning) of a VLM for multimodal reasoning. We assessed these hypotheses with our automated 296 radiology reporting experiment (Section 5.1). We then conducted a number of ablation studies 297 to investigate the effects of model size, different sparsity constraints, and the inclusion of auxiliary 298 information (Section 5.2). Next, we performed a case study for image-based feature localization to 299 assess whether the features learned by our SAE relate appropriately to their corresponding location 300 in the images (Section 5.3). Finally, we conducted a **reader study** with a specialist radiologist to 301 assess the quality of our generated reports (Section 5.4). 302

#### 303 5.1AUTOMATED RADIOLOGY REPORTING 304

305 **Dataset description** We trained and evaluated all models on the MIMIC-CXR dataset (Johnson 306 et al., 2019), a public dataset of 227,835 radiographic studies for a total of 377,110 chest radio-307 graphs and associated written text reports. We linked all images to their DICOM metadata files to retrieve scan orientations. We only considered images for which metadata files exist, and retained 308 only antero-posterior(AP)/postero-anterior(PA) scans for training (these are 'head-on' scans, as op-309 posed to lateral ones). These are the default views for the generation of diagnostic reports given the 310 increased clarity and ability to visualize the relevant anatomy comprehensively (Hyland et al., 2023). 311 We then extracted the 'findings' section from each text report. The 'findings' are a natural language 312 description of all relevant negative and positive features for a given radiograph. Occasionally, the 313 'findings' section is placed into another section with the sub-heading of 'impression'. Datapoints 314 without either a findings or impression section were discarded. We used the recommended train/test 315 split for MIMIC-CXR, resulting in a total of 239,931 training and 3,403 test images.

316

281

283 284

285

290

317 **Evaluation metrics** We evaluated generated radiology reports using both general Natural Lan-318 guage Generation (NLG) metrics (BLEU-4 (Papineni et al., 2002), ROUGE-L (Lin, 2004), ME-319 TEOR (Banerjee & Lavie, 2005)) and radiology-specific metrics (RGER score (Delbrouck et al., 320 2022), CheXpert F1 score (Irvin et al., 2019)). While lexical metrics assess n-gram overlap and 321 word order, clinical metrics like RGER and CheXpert F1 attempt to evaluate factual completeness and consistency by analyzing entity-relationship graphs and predicting common chest X-ray 322 pathologies, respectively. NLG metrics can be inadequate for assessing radiology reports as they 323 don't account for clinical significance (Bannur et al., 2024), while radiology-specific metrics often rely on specialized models (Yu et al., 2023) or pre-specified findings classes (Smit et al., 2020; Bannur et al., 2024). To address these limitations, the RadFact framework Bannur et al. (2024) uses
LLMs to assess sentence-level factuality through bi-directional entailment verification with reference reports— offering a robust evaluation method without relying on pre-specified error types or specialized models. Additional details are provided in Appendix B.2.

Experimental setup The SAE-Rad framework was trained on class tokens produced by the Rad-DINO vision transformer, using an expansion factor of 64 resulting in a latent dimension of 49,152.
 The model was optimized using Adam with a learning rate peaking at 5e-5 and a sparsity penalty of 8e-3, trained for 200,000 steps with a batch size of 2048. Claude 3.5 Sonnet was used for automated feature interpretation and report generation, while RadFact evaluation employed Llama3-70B-Instruct. Additional experimental setup details are given in Appendix B.3.

Table 1: **Report generation performance on the official MIMIC-CXR test split.** BL4 = BLEU-4, RG-L = ROUGE-L, MTR = Meteor. Ma-5 (Macro-F1-5), Ma-14 (Macro-F1-14), Mi-5 (Micro-F1-5), and Mi-14 (Micro-F1-14) represent the clinical CheXbert labeler scores. Bolding represents best performance in the current study or between the upper bound models.

| Model         | RadFact $\uparrow$ | NLG Metrics ↑ |      |      | Clinical Metrics ↑ |      |       |      |       |  |
|---------------|--------------------|---------------|------|------|--------------------|------|-------|------|-------|--|
| Widder        | F1                 | BL4           | RG-L | MTR  | RGER               | Ma-5 | Ma-14 | Mi-5 | Mi-14 |  |
| Current study |                    |               |      |      |                    |      |       |      |       |  |
| Baseline      | 30.0               | 3.1           | 18.6 | 23.5 | 14.9               | 38.3 | 25.5  | 45.0 | 42.8  |  |
| CheXagent     | 36.9               | 3.7           | 21.5 | 21.1 | 18.0               | 31.5 | 22.5  | 38.6 | 38.1  |  |
| SAE-Rad (×64) | 37.2               | 1.9           | 17.1 | 29.1 | 18.2               | 47.2 | 34.3  | 54.4 | 53.2  |  |
| Upper bound n | odels              |               |      |      |                    |      |       |      |       |  |
| MAIRA-1       | 47.8               | 14.2          | 28.9 | 33.3 | 29.6               | 47.7 | 38.6  | 56.0 | 55.7  |  |
| MAIRA-2       | 50.4               | 23.1          | 38.4 | 41.7 | 39.6               | 50.4 | 41.6  | 59.1 | 58.1  |  |

Feature 1 dimension 714 Dextroscoliosis of the thoracic spine



Feature 3 DIMENSION 150 RIGHT-SIDED PLEURAL EFFUSION



Feature 2 dimension 152 BILATERAL LUNG OPACITIES & SIGNS OF PULMONARY



Feature 4 dimension 89 PACEMAKER DEVICE & ASSOCIATED LEADS WITHIN CHEST CAVITY



Figure 2: **SAE-Rad identifies clinically relevant and interpretable features within radiological images.** We illustrate a number of pathological and instrumentation features relevant for producing radiology reports. We add annotations (green arrows) to emphasize the presence of each feature.

Quantitative evaluation We compared SAE-Rad to the current state-of-the-art radiology reporting systems. CheXagent (Chen et al., 2024) is an instruction-tuned foundation model for CXRs
trained on 1.1M scans for question-answering and text-generation tasks. MAIRA-1 &-2 (Hyland
et al., 2023; Bannur et al., 2024) are VLMS based on the LLaVA 1.5 architecture (Liu et al.,
2024b;a). MAIRA-2 is trained on 510,848 CXRs from four datasets and sets the current state-of-theart for report generation. The MAIRA systems are not publicly available for result replication, and



Figure 3: **SAE-Rad accurately captures features reported by human radiologists and more**. Above, we showcase a side-by-side comparison between a ground-truth radiology report and one generated by SAE-Rad. The model successfully identifies key clinically relevant features. SAE-Rad also identifies additional details, such as a right-sided dialysis catheter, without hallucination (we annotate this feature with green arrows for emphasis). SAE-Rad can also miss features when compared to the reference report.

400

401 thus we quote their evaluation values directly as our upper-bound. CheXagent is publicly available, 402 and we therefore performed independent replications for this model for a direct comparison. The 403 'baseline' approach is a naïve method of report generation that uses the report of the closest image in 404 the MIMIC train split. Further details of the 'baseline' approach can be found in Appendix B.6. As 405 Table 1 demonstrates, SAE-Rad underperforms on generic NLG metrics such as BLEU-4. This is expected as we do not try to optimize for any specific 'writing style' by fine-tuning an LLM on the 406 407 reference reports from MIMIC-CXR. Conversely, SAE-Rad demonstrates strong performance on radiology-specific metrics which are clinically relevant, outperforming CheXagent by up to 52% in 408 the CheXpert F1 score (macro-averaged F1-14), and achieving 92.1% and 89.9% of the performance 409 of MAIRA-1 and MAIRA-2 on these scores, respectively. We additionally compare SAE-Rad to 14 410 alternative approaches for automated report generation in Appendix J, broadly outperforming them 411 across all clinical metrics. 412

413 **Qualitative investigation** Figure 2 illustrates randomly selected monosemantic visual features 414 from SAE-Rad. As can be seen, the SAE learns human-interpretable visual concepts despite the ho-415 mogeneity and relatively small size of the dataset. These include dextroscoliosis of the spine (Fig. 2; 416 feature 1), bilateral opacifications (Fig. 2; feature 2), unilateral pleural effusions (Fig. 2; feature 3), 417 and the presence of instrumentation - in this case a pacemaker (Fig. 2; feature 4). In Fig. 3, we 418 illustrate an example 'findings' section for a CXR with a number of pathological findings; SAE-419 Rad is capable of detecting multiple relevant pathologies for a given image. Like other radiology report generation systems, SAE-Rad can miss findings. However, it can also occasionally describe 420 a relevant finding which is otherwise missing from the reference report – an example relating to the 421 presence of a dialysis catheter is shown in Fig. 3. 422

- 423
- 424 5.2 ABLATION STUDIES

We conducted a set of additional experiments which characterize the effects of: 1) Varying the SAE
expansion factor; 2) Investigating less sparse ('dense') SAEs by reducing the L1 penalty coefficient;
and 3) Evaluating the relative benefits of including auxiliary information including the 'indication'
section of a report as well as previous reports, where available. The experimental setup for these
ablations is described in Appendix H.

As can be seen in Table 2, an expansion factor of  $\times 64$  produced a higher RadFact F1 score compared with both smaller ( $\times 32$ ) and larger ( $\times 128$ ) expansion factors. In addition, denser SAEs with a larger L0 norm underperformed sparser models. This suggests that concepts useful for
radiology report generation are likely to exist in balance between being too coarse or too finegrained. If the feature are too coarse, this may cause 'feature absorption'— an asymmetric form
of feature splitting that can negatively impact the interpretability of an SAE feature (Chanin et al.,
2024). Conversely if the features are too fine-grained, the features may be more difficult to accurately describe given insufficient amount of descriptive detail in the ground-truth reference reports.
Table 2 demonstrates that the addition of aux-

438 iliary information such as the indication, which 439 describes why the patient required the scan in 440 the first instance, can boost the RadFact F1 441 score, with a large boost to recall. However, 442 in our experiments this caused a small degra-443 dation to the precision sub-metric. This find-444 ing runs against prior work on the effect of in-445 cluding this section (Yu et al., 2023; Tu et al., 446 2024; Nguyen et al., 2023) and warrants additional investigation in future. We find that 447 adding both previous indications and prior stud-448 ies has a net positive effect on the quality of 449 generated reports. Additionally, we report NLG 450 and radiology-specific metrics for this ablation 451 experiment in Table 6 of Appendix H, These 452 additional analyses support the findings in Table 2. 453

Table 2: RadFact performance metrics for different SAE-Rad configurations. /w inds = with indication(s), /w inds + prev.reps = with indications and previous text reports.

| SAE-Rad Configuration            | Precision $\uparrow$ | Recall ↑ | F1 Score ↑ |
|----------------------------------|----------------------|----------|------------|
| ×128                             | 34.83                | 29.91    | 32.18      |
| ×64                              | 35.95                | 31.95    | 33.83      |
| ×32                              | 31.22                | 27.89    | 29.46      |
| ×128_dense                       | 32.57                | 27.06    | 29.56      |
| ×64_dense                        | 31.15                | 28.46    | 29.74      |
| ×32_dense                        | 32.02                | 28.61    | 30.22      |
| $\times 64$ /wo inds             | 38.78                | 28.22    | 32.67      |
| $\times 64$ /w inds              | 38.45                | 32.42    | 35.18      |
| $\times 64$ /w inds + prev. reps | 37.32                | 39.83    | 38.45      |

5.3 IMAGE-BASED FEATURE LOCALIZATION

456 To test whether SAE features relate appropriately to their corresponding image location we trained 457 a diffusion model conditioned on Rad-DINO class tokens (Pérez-García et al., 2024). We then pro-458 duced *counterfactual* images, where we ask, what would this image have looked like if a particular 459 feature were present or absent? If the SAE successfully captures visual features, the counterfac-460 tual images should reflect the targeted feature's presence or absence. To produce such images, we 461 passed a class token through the SAE, intervened on encoder activations, and reconstructed a coun-462 terfactual token via the decoder, which conditioned the diffusion model to project interventions into imaging space. We tested whether: 1) interventions alter the reconstructed class token accordingly, 463 2) changes affect only the targeted feature, and 3) features can be "added" or "removed" by manipu-464 lating the same activation. Fig. 4 shows the results for two features (cardiomegaly and pacemaker), 465 demonstrating that our interpretations accurately reflect their impact on model behaviour. This also 466 enables us to conduct *unsupervised segmentation*, where we semantically segment specific visual 467 features by adding them to an image and then taking the difference between the pre- and post-edit 468 images. Further details are given in Appendix D.

469 470

471

454

455

### 5.4 READER STUDY

In a reader study with a specialist radiologist evaluating the quality of the automated radiology reports, 328 sentences from 60 reports (SAE-Rad, CheXagent, and a baseline) were analyzed. SAE-Rad had up to 7% fewer edits than other models and demonstrated significantly fewer errors with clinical impact, particularly in the "significant" category, where SAE-Rad had almost half the rate compared to others. This highlights SAE-Rad's potential for radiology report generation in a real clinical scenario. The full study can be found in Appendix E.

478 479

480

## 6 DISCUSSION, LIMITATIONS, AND CONCLUSION

In this work, we introduced SAE-Rad, a novel framework that leverages sparse autoencoders to automate radiology report generation. Our approach directly decomposes image class tokens from a pre-trained radiology image encoder into human-interpretable features, which are then compiled into comprehensive radiology reports. The experimental results demonstrate that SAE-Rad achieves competitive performance on radiology-specific metrics, outperforming existing models like CheXagent (Chen et al., 2024) and approaching the performance of state-of-the-art systems such as



Figure 4: **SAE-Rad enables targeted counterfactual image generation and unsupervised segmentation with disentangled class tokens.** Row 1 examines a pacemaker, and Row 2 investigates cardiomegaly. Column 1 shows original MIMIC-CXR images, Column 2 shows model reconstructions, and Columns 3 and 4 depict counterfactuals by adding and removing features. The final column demonstrates unsupervised segmentation by comparing counterfactual and original images. Details are in Appendix D.

501 502

504

505

506

507

508

MAIRA-2 whilst being trained on a significantly smaller dataset and with a much lower training compute budget; see Appendix F. By reverse-engineering the computations of the image encoder, SAE-Rad provides a framework that is verifiably faithful to the underlying model, enhancing transparency and trustworthiness, which are critical considerations in the healthcare setting.

Qualitative analyses confirm that SAE-Rad successfully captures meaningful visual concepts in cluding the presence or absence of pathological features. These interpretable features contribute to generating detailed and accurate radiology reports, as evidenced by strong performance in the clinical evaluation metrics. Our ablation studies indicate that the choice of expansion factor and inclusion of auxiliary information, such as previous reports and indications, can significantly impact the quality of the generated reports, which is broadly in line with the previous literature on this topic (Bannur et al., 2024).

Our approach has limitations. First, SAE-Rad relies on pre-trained (frozen) models for both the 524 image encoder and the LLM in the interpretability pipeline, potentially introducing inherent biases. 525 Spurious correlations and other biases in the image encoder's latent space may prevent our method 526 from fully disentangling concepts. This may lead to *feature absorption* (Chanin et al., 2024), where 527 distinct concepts are inappropriately merged into a single SAE dimension. However, due to the 528 pipeline's modular nature, these biases can be mitigated by replacing either model without retrain-529 ing the SAE if the LLM is swapped. Additionally, SAE-Rad underperforms on general language 530 metrics like BLEU-4, suggesting that while the generated reports are clinically accurate, they may 531 lack the fluency and stylistic nuances of human-generated reports, particularly those in the MIMIC-CXR dataset. Improving these metrics through style-aware radiology report generation (Yan et al., 532 2023) is a natural avenue for future work. Another exciting direction is to prevent errors from LLM 533 hallucinations by having SAE features labeled by human inspection (Appendix I) and using a deter-534 ministic regex for report synthesis, thereby obviating the need for LLMs in our pipeline altogether. 535

Overall, this feasibility study presents a novel and effective approach to radiology report genera tion by leveraging mechanistic interpretability techniques to extract and utilize human-interpretable
 features from medical images. By providing a verifiably faithful representation of the underlying
 computations, SAE-Rad contributes to the development of more interpretable AI systems in health care.

# 540 REFERENCES

553

576

Anthropic. Claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet, 2024.
 Accessed: 2024-09-25.

- Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In <u>Proceedings of the acl workshop on intrinsic and extrinsic</u> evaluation measures for machine translation and/or summarization, pp. 65–72, 2005.
- Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fernando Pérez-García, Maximilian Ilse, Daniel C. Castro, Benedikt Boecking, Harshita Sharma, Kenza Bouzid, Anja Thieme, Anton Schwaighofer, Maria Wetscherek, Matthew P. Lungren, Aditya Nori, Javier Alvarez-Valle, and Ozan Oktay. Learning to exploit temporal structure for biomedical vision-language processing. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15016–15027, 2023.
- Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Anton Schwaighofer, Sam Bond-Taylor, Maximilian Ilse, Fernando Pérez-García, Valentina Salvatelli, Harshita Sharma, Felix Meissen, et al. Maira-2: Grounded radiology report generation. <u>arXiv preprint arXiv:2406.04449</u>, 2024.
- Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. The reversal curse: LLMs trained on "A is B" fail to learn "B is A". arXiv preprint arXiv:2309.12288, 2023.
- Usha Bhalla, Alex Oesterling, Suraj Srinivas, Flavio P Calmon, and Himabindu Lakkaraju. Interpreting clip with sparse linear concept embeddings (splice). arXiv preprint arXiv:2402.10376, 2024.
- Tolga Bolukbasi, Adam Pearce, Ann Yuan, Andy Coenen, Emily Reif, Fernanda Viégas, and Martin
   Wattenberg. An interpretability illusion for bert. arXiv preprint arXiv:2104.07143, 2021.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
   Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,
   Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina
   Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and
   Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary
   learning. <u>Transformer Circuits Thread</u>, 2023. URL https://transformer-circuits.pub/
   2023/monosemantic-features/index.html.
- 573 David Chanin, James Wilken-Smith, Tomáš Dulka, Hardik Bhatnagar, and Joseph Bloom. A is
   574 for absorption: Studying feature splitting and absorption in sparse autoencoders. <u>arXiv preprint</u> arXiv:2409.14507, 2024.
- Juan Manuel Zambrano Chaves, Shih-Cheng Huang, Yanbo Xu, Hanwen Xu, Naoto Usuyama, Sheng Zhang, Fei Wang, Yujia Xie, Mahmoud Khademi, Ziyi Yang, et al. Training small multimodal models to bridge biomedical competency gap: A case study in radiology imaging. <u>arXiv</u> preprint arXiv:2403.08002, 2024.
- Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via
   memory-driven transformer. arXiv preprint arXiv:2010.16056, 2020.
- Zhihong Chen, Maya Varma, Jean-Benoit Delbrouck, Magdalini Paschali, Louis Blankemeier, Dave Van Veen, Jeya Maria Jose Valanarasu, Alaa Youssef, Joseph Paul Cohen, Eduardo Pontes Reis, et al. Chexagent: Towards a foundation model for chest x-ray interpretation. <u>arXiv preprint</u> arXiv:2401.12208, 2024.
- Noel CF Codella, Ying Jin, Shrey Jain, Yu Gu, Ho Hin Lee, Asma Ben Abacha, Alberto Santamaria Pang, Will Guyman, Naiteek Sangani, Sheng Zhang, et al. Medimageinsight: An open-source embedding model for general domain medical imaging. <u>arXiv preprint arXiv:2410.06542</u>, 2024.
- Tom Conerly, Adly Templeton, Trenton Bricken, Jonathan Marcus, and Tom Henighan. Circuits Updates April 2024: Update on how we train SAEs, April 2024. URL https://transformer-circuits.pub/2024/april-update/index.html#training-saes. Accessed: 2024-09-23.

612

619

625

631

- 594 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-595 coders find highly interpretable features in language models. arXiv preprint arXiv:2309.08600, 596 2023. 597
- Gytis Daujotas. Case study: Interpreting, manipulating, and controlling 598 clip, 2024. URL https://www.lesswrong.com/posts/iYFuZo9BMvr6GgMs5/ case-study-interpreting-manipulating-and-controlling-clip. Accessed: 2024 -600 09-25. 601
- 602 Jean-Benoit Delbrouck, Pierre Chambon, Christian Bluethgen, Emily Tsai, Omar Almusa, and Cur-603 tis P Langlotz. Improving the factual correctness of radiology report generation with semantic 604 rewards. arXiv preprint arXiv:2210.12186, 2022.
- 605 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, et al. 606 A mathematical framework for transformer circuits. https://transformer-circuits.pub/ 2021/framework/index.html, 2021. Accessed: 2024-09-30. 608
- 609 Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, et al. 610 Softmax linear units: Designing for interpretability. https://transformer-circuits.pub/ 611 2022/solu/index.html, 2022a. Accessed: 2024-09-19.
- Nelson Elhage, Catherine Olsson, Tom Henighan, Neel Nanda, Scott Johnston, et al. A mathematical 613 framework for transformer circuits: Toy model. https://transformer-circuits.pub/2022/ 614 toy\_model/index.html#motivation, 2022b. Accessed: 2024-09-19. 615
- 616 Hugo Fry. Towards multimodal interpretability: Learning sparse interpretable features in vi-617 sion transformers, 2024. URL https://www.lesswrong.com/posts/bCtbuWraqYTDtuARg/ 618 towards-multimodal-interpretability-learning-sparse-2. Accessed: 2024-09-25.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya 620 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. arXiv preprint 621 arXiv:2406.04093, 2024. 622
- 623 Liv Gorton. The missing curve detectors of inceptionv1: Applying sparse autoencoders to incep-624 tionv1 early vision. arXiv preprint arXiv:2406.03662, 2024.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bert-626 simas. Finding neurons in a haystack: Case studies with sparse probing. arXiv preprint 627 arXiv:2305.01610, 2023. 628
- 629 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in 630 neural information processing systems, 33:6840-6851, 2020.
- Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang. Kiut: Knowledge-injected u-transformer 632 for radiology report generation. In Proceedings of the IEEE/CVF Conference on Computer Vision 633 and Pattern Recognition, pp. 19809-19818, 2023. 634
- 635 Huggingface. Diffusers: State-of-the-art diffusion models in pytorch. https://github.com/ 636 huggingface/diffusers, 2023. Version 0.19.0.
- Stephanie L Hyland, Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Mercy Ranjit, Anton 638 Schwaighofer, Fernando Pérez-García, Valentina Salvatelli, Shaury Srivastav, Anja Thieme, et al. 639 Maira-1: A specialised large multimodal model for radiology report generation. arXiv preprint 640 arXiv:2311.13668, 2023. 641
- 642 Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik 643 Marklund, Behzad Haghgoo, Robyn L. Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, 644 Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. CheXpert: A large chest radiograph dataset with 645 uncertainty labels and expert comparison. In Proceedings of the AAAI Conference on Artificial 646 Intelligence (AAAI 2019), volume 33, pp. 590–597. AAAI Press, July 2019. doi: 10.1609/aaai. 647 v33i01.3301590.

- 648 Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven QH Truong, Du Nguyen Duong, Tan Bui, 649 Pierre Chambon, Yuhao Zhang, Matthew P Lungren, Andrew Y Ng, et al. Radgraph: Extracting 650 clinical entities and relations from radiology reports. arXiv preprint arXiv:2106.14463, 2021. 651 Katharina Jeblick, Balthasar Schachtner, Jakob Dexl, Andreas Mittermeier, Anna Theresa Stüber, 652 Johanna Topalis, Tobias Weber, Philipp Wesp, Bastian Oliver Sabel, Jens Ricke, et al. Chat-653 gpt makes medicine easy to swallow: an exploratory case study on simplified radiology reports. 654 European radiology, pp. 1–9, 2023. 655 656 Jaehwan Jeong, Katherine Tian, Andrew Li, Sina Hartung, Subathra Adithan, Fardad Behzadi, Juan 657 Calle, David Osayande, Michael Pohlen, and Pranav Rajpurkar. Multimodal image-text matching improves retrieval-based chest x-ray report generation. In Medical Imaging with Deep Learning, 658 pp. 978-990. PMLR, 2024. 659 660 Haibo Jin, Haoxuan Che, Yi Lin, and Hao Chen. Promptmrg: Diagnosis-driven prompts for medical 661 report generation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, 662
- pp. 2607-2615, 2024. 663
- Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, 664 Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a 665 large publicly available database of labeled chest radiographs. arXiv preprint arXiv:1901.07042, 666 2019. 667
- 668 Adam Tauman Kalai and Santosh S Vempala. Calibrated language models must hallucinate. In 669 Proceedings of the 56th Annual ACM Symposium on Theory of Computing, pp. 160-171, 2024. 670
- Diederik P Kingma. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 671 2014. 672
- 673 Nhat Le, Ciyue Shen, Chintan Shah, Blake Martin, Daniel Shenker, Harshith Padigela, Jennifer 674 Hipp, Sean Grullon, John Abel, Harsha Vardhan Pokkalla, et al. Interpretability analysis on a 675 pathology foundation model reveals biologically relevant embeddings across modalities. arXiv 676 preprint arXiv:2407.10785, 2024. 677
- Peter Lee, Sebastien Bubeck, and Joseph Petro. Benefits, limits, and risks of gpt-4 as an ai chatbot 678 for medicine. New England Journal of Medicine, 388(13):1233-1239, 2023. 679
- 680 T Lewis, M Wood, and M Myers. Pathology girft programme national specialty report, 2021. 681

687

688

689

690

- 682 Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision as-683 sistant for biomedicine in one day. Advances in Neural Information Processing Systems, 36, 684 2024. 685
  - Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. arXiv preprint arXiv:2408.05147, 2024.
  - Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. Video-Ilava: Learning united visual representation by alignment before projection. arXiv preprint arXiv:2311.10122, 2023.
- 692 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In Text summarization 693 branches out, pp. 74-81, 2004. 694
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction 695 tuning. arXiv preprint arXiv:2310.03744, 2023a. 696
- 697 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual in-698 struction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 699 Recognition, pp. 26296–26306, 2024a. 700
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances 701 in neural information processing systems, 36, 2024b.

| 702<br>703<br>704               | Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, et al. Llava-plus: Learning to use tools for creating multimodal agents. <u>arXiv</u> preprint arXiv:2311.05437, 2023b.  |
|---------------------------------|---|
| 705<br>706<br>707<br>708        | Z Liu, A Zhong, Y Li, et al. Radiology-gpt: a large language model for radiology. arxiv 2306.08666<br>[preprint] https://arxiv. org/abs/2306.08666. published june 14, 2023. <u>Accessed November</u> , 14, 2023c.  |
| 709<br>710<br>711<br>712        | Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,<br>Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for<br>science question answering. <u>Advances in Neural Information Processing Systems</u> , 35:2507–2521,<br>2022.   |
| 713<br>714<br>715               | Stéphane G Mallat and Zhifeng Zhang. Matching pursuits with time-frequency dictionaries. <u>IEEE</u><br><u>Transactions on signal processing</u> , 41(12):3397–3415, 1993.  |
| 716<br>717                      | Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426, 2018.   |
| 718<br>719<br>720<br>721        | Jake Mendel. Sae feature geometry is outside the superposition hypoth-<br>esis. https://www.alignmentforum.org/posts/MFBTjb2qf3ziWmzz6/<br>sae-feature-geometry-is-outside-the-superposition-hypothesis, 2024. Accessed:<br>2024-11-18.   |
| 722<br>723<br>724               | Yasuhide Miura, Yuhao Zhang, Emily Bao Tsai, Curtis P Langlotz, and Dan Jurafsky. Improving factual completeness and consistency of image-to-text radiology report generation. <u>arXiv preprint</u> <u>arXiv:2010.10042</u> , 2020.  |
| 725<br>726<br>727<br>728        | Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril<br>Zakka, Eduardo Pontes Reis, and Pranav Rajpurkar. Med-flamingo: a multimodal medical few-<br>shot learner. In <u>Machine Learning for Health (ML4H)</u> , pp. 353–367. PMLR, 2023.   |
| 729<br>730<br>731               | Neel Nanda, Tom Henighan, Scott Johnston, et al. Scaling monosemanticity in language models.<br>https://transformer-circuits.pub/2024/scaling-monosemanticity/, 2024. Accessed:<br>2024-09-19.  |
| 732<br>733                      | Dang Nguyen, Chacha Chen, He He, and Chenhao Tan. Pragmatic radiology report generation. In <u>Machine Learning for Health (ML4H)</u> , pp. 385–402. PMLR, 2023.  |
| 735<br>736                      | NHS England and NHS Improvement. Transforming imaging services in England: a national strat-<br>egy for imaging networks. Technical report, NHS, 2019.  |
| 737<br>738                      | Aaron Nicolson, Jason Dowling, and Bevan Koopman. Improving chest x-ray report generation by<br>leveraging warm starting. <u>Artificial intelligence in medicine</u> , 144:102633, 2023.  |
| 739<br>740<br>741               | Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.<br>Zoom in: An introduction to circuits. <u>Distill</u> , 5(3):e00024–001, 2020.   |
| 742<br>743<br>744               | Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. <u>arXiv preprint arXiv:2304.07193</u> , 2023.   |
| 745<br>746<br>747<br>748        | Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In <u>Proceedings of the 40th annual meeting of the Association</u> for Computational Linguistics, pp. 311–318, 2002.  |
| 749<br>750                      | Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. <u>arXiv preprint arXiv:2311.03658</u> , 2023.   |
| 751<br>752<br>753<br>754<br>755 | Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In <u>Advances in Neural Information Processing Systems 32</u> , pp. 8024–8035. Curran Associates, Inc., 2019. |

756 Fernando Pérez-García, Harshita Sharma, Sam Bond-Taylor, Kenza Bouzid, Valentina Salvatelli, 757 Maximilian Ilse, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Matthew P Lungren, 758 et al. Rad-dino: Exploring scalable medical image encoders beyond text supervision. arXiv 759 preprint arXiv:2401.10815, 2024. 760 Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János 761 Kramár, Rohin Shah, and Neel Nanda. Improving dictionary learning with gated sparse autoen-762 coders. arXiv preprint arXiv:2404.16014, 2024. 763 764 Ramprasaath R Selvaraju, Purva Tendulkar, Devi Parikh, Eric Horvitz, Marco Tulio Ribeiro, Be-765 smira Nushi, and Ece Kamar. Squinting at vqa models: Introspecting vqa models with sub-766 questions. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 767 Recognition, pp. 10003–10011, 2020. 768 Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y Ng, and Matthew P Lungren. 769 Chexbert: combining automatic labelers and expert annotations for accurate radiology report la-770 beling using bert. arXiv preprint arXiv:2004.09167, 2020. 771 772 Raphael Stock, Stefan Denner, Yannick Kirchhoff, Constantin Ulrich, Maximilian Rouven Rokuss, 773 Saikat Roy, Nico Disch, and Klaus Maier-Hein. From generalist to specialist: Incorporating 774 domain-knowledge into flamingo for chest x-ray report generation. In Medical Imaging with Deep Learning, 2024. 775 776 Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. Interactive and explainable 777 region-guided radiology report generation. In Proceedings of the IEEE/CVF Conference on 778 Computer Vision and Pattern Recognition, pp. 7433-7442, 2023. 779 Ryutaro Tanno, David Barrett, Andrew Sellergren, Sumedh Ghaisas, Sumanth Dathathri, Abigail 780 See, Johannes Welbl, Karan Singhal, Shekoofeh Azizi, Tao Tu, et al. Consensus, dissensus 781 and synergy between clinicians and specialist foundation models in radiology report generation, 782 2024a. 783 784 Ryutaro Tanno, David GT Barrett, Andrew Sellergren, Sumedh Ghaisas, Sumanth Dathathri, Abi-785 gail See, Johannes Welbl, Charles Lau, Tao Tu, Shekoofeh Azizi, et al. Collaboration between 786 clinicians and vision-language models in radiology report generation. Nature Medicine, pp. 1-10, 787 2024b. 788 Sonnet Templeton, Tom Conerly, Jonathan Marcus, Tom Henighan, Joshua Batson, Chris Olah, and 789 Adam Jermyn. Scaling monosemanticity: Extracting interpretable features from claude 3, 2024. 790 URL https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html. 791 Accessed: September 24, 2024. 792 793 Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaekermann, Mohamed Amin, Pi-Chuan Chang, 794 Andrew Carroll, Charles Lau, Ryutaro Tanno, Ira Ktena, et al. Towards generalist biomedical ai. NEJM AI, 1(3):AIoa2300138, 2024. 796 Kohei Uehara, Nan Duan, and Tatsuya Harada. Learning to ask informative sub-questions for visual 797 question answering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 798 Recognition, pp. 4681-4690, 2022. 799 800 Ruonan Wang, Yuxi Qian, Fangxiang Feng, Xiaojie Wang, and Huixing Jiang. Co-vqa: Answering 801 by interactive sub question sequence. arXiv preprint arXiv:2204.00879, 2022. 802 Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. Metransformer: Radiology report gen-803 eration by transformer with multiple learnable expert tokens. In Proceedings of the IEEE/CVF 804 Conference on Computer Vision and Pattern Recognition, pp. 11558-11567, 2023a. 805 806 Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. R2gengpt: Radiology report generation 807 with frozen llms. Meta-Radiology, 1(3):100033, 2023b. 808 Zilong Wang, Xufang Luo, Xinyang Jiang, Dongsheng Li, and Lili Qiu. Llm-radjudge: Achieving 809 radiologist-level evaluation for x-ray report generation. arXiv preprint arXiv:2404.00998, 2024.

| 810<br>811<br>812                      | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.  |
|--|---|
| 813<br>814<br>815                      | Benjamin Wright and Lee Sharkey. Addressing feature suppression in saes. In <u>AI Alignment Forum</u> , pp. 16, 2024.   |
| 816<br>817<br>818                      | Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards generalist foun-<br>dation model for radiology. <u>arXiv preprint arXiv:2308.02463</u> , 2023.   |
| 819<br>820                             | Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models. <u>arXiv preprint arXiv:2401.11817</u> , 2024.  |
| 821<br>822<br>823<br>824               | Benjamin Yan, Ruochen Liu, David E Kuo, Subathra Adithan, Eduardo Pontes Reis, Stephen Kwak, Vasantha Kumar Venugopal, Chloe P O'Connell, Agustina Saenz, Pranav Rajpurkar, et al. Style-aware radiology report generation with radgraph and few-shot prompting. <u>arXiv preprint</u> <u>arXiv:2310.17811</u> , 2023.  |
| 825<br>826<br>827<br>828               | Lin Yang, Shawn Xu, Andrew Sellergren, Timo Kohlberger, Yuchen Zhou, Ira Ktena, Atilla Ki-<br>raly, Faruk Ahmed, Farhad Hormozdiari, Tiam Jaroensri, et al. Advancing multimodal medical<br>capabilities of gemini. <u>arXiv preprint arXiv:2405.03162</u> , 2024.  |
| 829<br>830<br>831                      | Shuxin Yang, Xian Wu, Shen Ge, Zhuozhao Zheng, S Kevin Zhou, and Li Xiao. Radiology report generation with a learned knowledge base and multi-modal alignment. <u>Medical Image Analysis</u> , 86:102798, 2023.   |
| 832<br>833<br>834<br>835<br>836<br>837 | Nur Yildirim, Hannah Richardson, Maria Teodora Wetscherek, Junaid Bajwa, Joseph Jacob,<br>Mark Ames Pinnock, Stephen Harris, Daniel Coelho De Castro, Shruthi Bannur, Stephanie<br>Hyland, et al. Multimodal healthcare ai: identifying and designing clinically relevant vision-<br>language applications for radiology. In <u>Proceedings of the CHI Conference on Human Factors in</u><br><u>Computing Systems</u> , pp. 1–22, 2024. |
| 838<br>839<br>840                      | Haoxuan You, Rui Sun, Zhecan Wang, Long Chen, Gengyu Wang, Hammad A Ayyubi, Kai-Wei<br>Chang, and Shih-Fu Chang. Idealgpt: Iteratively decomposing vision and language reasoning via<br>large language models. <u>arXiv preprint arXiv:2305.14985</u> , 2023.   |
| 841<br>842<br>843<br>844               | Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser<br>Ururahy Nunes Fonseca, Henrique Min Ho Lee, Zahra Shakeri Hossein Abad, Andrew Y Ng,<br>et al. Evaluating progress in automatic chest x-ray radiology report generation. <u>Patterns</u> , 4(9),<br>2023.   |
| 845<br>846<br>847                      | Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. <u>arXiv preprint arXiv:2302.00923</u> , 2023.   |
| 848<br>849<br>850                      | Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibei Yang. Ddcot: Duty-distinct chain-of-<br>thought prompting for multimodal reasoning in language models. <u>Advances in Neural Information</u><br><u>Processing Systems</u> , 36:5168–5191, 2023.  |
| 851<br>852<br>853<br>854               | Qingqing Zhu, Xiuying Chen, Qiao Jin, Benjamin Hou, Tejas Sudharshan Mathai, Pritam Mukher-<br>jee, Xin Gao, Ronald M Summers, and Zhiyong Lu. Leveraging professional radiologists' exper-<br>tise to enhance llms' evaluation for radiology reports. <u>arXiv preprint arXiv:2401.16578</u> , 2024.   |
| 855<br>856<br>857                      |   |
| 859<br>860                             |   |
| 861<br>862<br>863                      |   |

# Appendix

# Table of Contents

| 870 |   |    |
|-----|---|----|
| 871 | Appendix A - Reproducibility statement  |    |
| 872 | Appendix B - Automated radiology report experiment                                |    |
| 873 | B.1 SAE architecture comparison   |    |
| 874 | B.2 Evaluation metrics  |    |
| 875 | B.3 SAE-Rad additional experimental setup details                                 |    |
| 877 | B.4 Prompt used to annotate SAE features  |    |
| 878 | B.5 Prompt used to compose finding sections via SAE features                      |    |
| 879 | B.6 Baseline experimental setup   |    |
| 880 | Appendix C - Additional qualitative evaluations                                   |    |
| 881 | C.1 More radiology report generations   |    |
| 882 | C.2 Image examples of monosemantic features                                       |    |
| 883 | Appendix D - Image-based feature localization implementation                      |    |
| 885 | Appendix E - Reader study   |    |
| 886 | E.1 Results   |    |
| 887 | Appendix F - Compute efficiency comparisons                                       |    |
| 888 | Appendix G - Limits of automated quality assessment for radiology reports $\dots$ |    |
| 890 | Appendix H - Ablation experimental setup  |    |
| 891 | Appendix I - Interactive exploration of SAE features                              |    |
| 892 | Appendix J - Additional SAE-Rad comparisons                                       |    |
| 893 | Appendix K - Equivalence of loss functions between SAE architectures              | 43 |
| 894 | K 1 Motivation  | 43 |
| 806 | K 2 Explanation of the Conerly et al. loss  |    |
| 897 | K 3 Equivalence with the original loss  |    |
| 001 |   |    |

# 918 A REPRODUCIBILITY STATEMENT

We provide details on the experimental setup, hyperparameters, and data preprocessing steps to ensure the reproducibility of our results. We cannot share the MIMIC-CXR dataset, however it can be accessed through application here: https://physionet.org/content/mimic-cxr/2.1. 0/. Our codebase is available at: https://anonymous.4open.science/r/sae\_rad-5B77/.

#### 

## **B** AUTOMATED RADIOLOGY REPORT EXPERIMENT

## B.1 SAE ARCHITECTURE COMPARISON

In this section we compare the performance of our novel SAE architecture with the state-of-the-art architecture proposed by Conerly et al. (2024), which itself introduced a novel loss that outperformed prior works. The following hyperparameters were used to train both SAEs:

- Expansion factor of  $\times 64$ .
- Batch size of 2048.
- Learning rate of  $5 \times 10^{-5}$ .
- Linear warm-up of learning rate for the first 1% of training.
- Linear warm-down of learning rate for the last 20% of training.
- L1 coefficient warmup for the first 5% of training.
- Adam optimizer with no weight decay.
  - Trained for 200,000 optimization steps.

The L1 coefficient was increased for the SAE-Rad architecture in comparison to the Conerly et al. (2024) SAE in order to compensate for the additional auxiliary loss term used to train the SAE-Rad architecture. Table 3 displays the resulting metrics comparing the two SAEs - our proposed SAE attains both a lower L0 and a higher reconstruction accuracy.

|                           | L1 coefficient      | L0↓  | Explained Variance (%) $\uparrow$ |
|---------------------------|---------------------|------|-----------------------------------|
| Conerly et al. (2024) SAE | $5.3 	imes 10^{-3}$ | 15.7 | 81.3                              |
| SAE-Rad                   | $8.0 	imes 10^{-3}$ | 13.6 | 84.3                              |

Table 3: Comparison of SAE types.

### **B.2** EVALUATION METRICS

**NLG and classical radiology-specific metrics** We evaluated generated radiology reports using both general NLG metrics and radiology-specific metrics. For lexical evaluation, we report BLEU-4 (Papineni et al., 2002) for 4-gram overlap based on n-gram precision, ROUGE-L (Lin, 2004) for longest common subsequence matching, and METEOR (Banerjee & Lavie, 2005), which performs unigram matching using surface forms, stems, and meanings, computing scores based on precision, recall, and fragmentation assessing word order. Whilst widely reported, lexical metrics do not cap-ture factual completeness or consistency (Miura et al., 2020; Bannur et al., 2024) and we therefore also include classical clinical metrics. The RGER score (Delbrouck et al., 2022), based on the RadGraph model (Jain et al., 2021), evaluates entity-relationship graphs extracted from reports by matching entities and verifying relationships. Additionally, we report the CheXpert F1 score (Irvin et al., 2019), utilizing the CheXbert model (Smit et al., 2020) to predict 14 common pathologies in chest X-rays and calculating the harmonic mean of precision and recall between generated and reference texts. Following recommendations (Miura et al., 2020; Tu et al., 2024; Hyland et al., 2023), we provide micro- and macro-averaged scores over five key observations-atelectasis, car-diomegaly, consolidation, edema, and pleural effusion-as well as the F1-14 score encompassing all observations.

972 **RadFact** — A robust evaluation framework for radiology Natural language generation (NLG) 973 metrics are insufficient to appropriately assess radiology report generations as they do not ac-974 count for the clinical significance of each sentence, and instead treat all words equally (Bannur 975 et al., 2024). Conversely, radiology-specific metrics are often based on specialised models such as 976 CheXbert (Smit et al., 2020; Irvin et al., 2019) or RadGraph (Yu et al., 2023; Jain et al., 2021; Delbrouck et al., 2022), which are themselves limited in that rely on specialised models (Yu et al., 2023), 977 pre-specified findings classes (Smit et al., 2020; Bannur et al., 2024), and/or error types (Chaves 978 et al., 2024; Wang et al., 2024). To this end, Bannur et al. (2024) proposed the RadFact framework 979 which leverages LLMs to assess the factuality of *each sentence* by use of bi-directional entailment 980 verification with the ground-truth reference report. RadFact does not rely on pre-specified error 981 types or radiology-specialized models and is therefore significantly more robust for use to assess the 982 quality of generated radiology reports.

983 984 985

### B.3 SAE-RAD ADDITIONAL EXPERIMENTAL SETUP DETAILS

986 In our instantiation of the SAE-Rad framework (described in Sections 4.1 and 4.2), all bias terms 987  $\mathbf{b}^{enc}$ ,  $\mathbf{b}^{gate}$ , and  $\mathbf{b}^{mag}$  were initialized to zeros. The elements of the shared encoder weights matrix  $W^{\text{gate}}$  were initialized such that all rows point in random directions. The decoder weights  $W^{\text{dec}}$ 988 were initialized to  $W^{\text{gate}\top}$ . For a dataset of size S, we trained our SAE on the class tokens  $\mathbf{Z} \in$ 989  $\mathbb{R}^{S \times 768}$  produced by the Rad-DINO vision transformer (Pérez-García et al., 2024), a fine-tuned 990 Dino V2 (Oquab et al., 2023) on a large dataset of chest X-rays. We use the model weights released 991 at https://huggingface.co/microsoft/rad-dino. The dataset was shuffled and scaled by a 992 constant such that  $\mathbb{E}_{z \in Z}[\|z\|_2] = \sqrt{\dim(\overline{z})}$ , where  $\dim(z) = 768$ . We used an expansion factor of 993 64 and thus the latent dimension of our SAE is 49, 152. We used the Adam optimization algorithm 994 (Kingma, 2014) with  $\beta_1 = 0.9, \beta_2 = 0.999$ , and no weight decay. Our learning rate was increased 995 linearly over the first 1% of training to  $5 \times 10^{-5}$  and then decayed linearly to zero over the final 20% of training. Our sparsity penalty  $\lambda$  was linearly increased over the first 5% of training to  $8 \times 10^{-3}$ . 997 We trained our model for 200,000 steps. We used a batch size of 2048. We did not use ghost 998 gradients or resampling strategies as none of our features were dead at the end of training — in 999 other words, all learned concept directions were activated at least once when passing the dataset through the SAE. We trained all models on a single NVIDIA A6000 GPU. We used Claude 3.5 1001 Sonnet (Anthropic, 2024) for our automated interpretability pipeline; for a given SAE feature, we retrieved the 10 highest activating images for that feature and passed their ground-truth reference report to Claude with the task of extracting the most consistent information across the reports. The 1003 output description was used to label the feature. The full prompts are described in Appendix B.4. 1004

To generate a report from SAE features, we accrued all activating features and their descriptions, and passed them to Claude for concatenation into a full 'findings' paragraph. The full prompts are describe in Appendix B.5. We evaluate RadFact using Llama3-70B-Instruct.

1008 1009

1010

B.4 PROMPT USED TO ANNOTATE SAE FEATURES

In this section we describe the prompt used to instruct Claude 3.5 Sonnet (Anthropic, 2024) to annotate SAE features. Overall, we instruct Claude to identify commonalities between radiological reports corresponding to the 10 maximally activating examples for a given SAE feature. The prompt encourages the model to use chain-of-thought (Wei et al., 2022) reasoning in-order to ensure the quality of the feature annotations.

| 1016 <sup>1</sup> | You are an expert radiologist specializing in chest radiographs. We're |
|-------------------|--|
| 1017              | studying neurons in an image neural network, where each neuron         |
| 4040              | detects specific features in chest X-rays. I've identified the         |
| 1018              | radiology images that most strongly activate a particular neuron and   |
| 1019              | will provide you with their associated text radiology reports. Your    |
| 1020              | task is to analyze these reports and determine the common feature      |
| 1021              | that this neuron is detecting.   |
| 1000 2            | To arrive at the most accurate and precise explanation of what this    |
| 1022              | neuron is detecting you must engage in explicit chain of thought       |
| 1023              | neuron is detecting, you must engage in expirere endin of endagite     |
|                   | reasoning. Begin by thoroughly examining all provided radiology        |
| 1024              | reports, noting any patterns or commonalities. Pay close attention to  |
| 1025              | recurring terminology, described anatomical structures, and            |
|                   | consistent pathological findings. Consider how these elements might    |

1026 interrelate to form a singular, distinctive feature that the neuron 1027 could be identifying. Evaluate the context of chest radiographs and 1028 consider which aspects would be most significant or unique within this imaging modality. 1029 As you progress through your analysis, verbalize your thought process. 1030 3 Explain each step of your reasoning, from initial observations to 1031 intermediate conclusions, and finally to your overall assessment. 1032 This chain of thought approach will help ensure a comprehensive and 1033 well-reasoned final explanation. 1034 4 After this detailed analytical process, formulate a single, specific explanation of what the neuron is detecting. Your explanation should 1035 be as precise and fine-grained as possible, avoiding vague or general 1036 statements. Focus on specific features or combinations of features, 1037 using 'and' to connect multiple elements if necessary. Avoid using 1038 or' to list multiple possibilities. Refrain from explaining the pathology itself (e.g., avoid statements like "This feature 1039 represents X, which is characterized by..."). Base your explanation 1040 solely on the information provided in the reports, without additional 1041 medical knowledge that might not be captured by the neuron. 1042 It is crucial that you present your final explanation in the following 5 1043 format: 1044 6 \*This feature represents [your specific, detailed description of what the neuron is detecting]. 1045 The asterisk is absolutely essential. Your explanation must begin 7 1046 immediately after the asterisk, without any additional text, 1047 numbering, or preamble. The presence of this asterisk is critical for 1048 the proper processing of your response. Below are the radiology reports, listed in order of how strongly they 1049 8 activate the neuron. Use these to inform your analysis and final 1050 explanation: 1051 9 Report number 1: ... 1052 <sub>10</sub> Report number 2: ... 1053 11 1054 12 Report number 10: . 1055 Below is an example output from Claude in response to the above prompt. 1056 1057 Thank you for providing these detailed radiology reports. I'll analyze 1 1058 them to determine the common feature that this neuron is detecting. Let me walk you through my thought process: 1059 1060 2 1. First, I'll review each report, noting recurring elements: 3 1061 4 - Report 1: Mentions cardiomegaly and atherosclerotic calcifications 1062 in the thoracic aorta. 1063 5 Report 2: Notes mild cardiomegaly and a moderately tortuous, mildly 1064 calcified aorta. - Report 3: Describes an enlarged cardiomediastinal silhouette and a 1065 6 calcified, unfolded aorta. 1066 - Report 4: Mentions moderate enlargement of the cardiac silhouette. 1067 - Report 5: No specific mention of heart or aorta abnormalities. 1068 g - Report 6: Notes moderate cardiomegaly. 1069 10 - Report 7: Mentions moderately enlarged cardiac silhouette. Report 8: Describes moderate cardiomegaly and moderate tortuosity of 1070 11 **1071** 12 the thoracic aorta. - Report 9: No specific mention of heart or aorta abnormalities. 1072 <u>1</u>3 - Report 10: Notes an enlarged cardiac silhouette and a tortuous 1073 descending thoracic aorta. 1074 14 1075 15 2. Analyzing these observations, I notice two recurring themes: a) Enlarged cardiac silhouette or cardiomegaly 1076<sup>16</sup> 1077<sup>17</sup> b) Abnormalities of the thoracic aorta (tortuosity, calcification, unfolding) 1078 <sub>18</sub> 3. Considering the context of chest radiographs, both of these features 1079 19 would be significant and distinctive within this imaging modality.

1080 <sub>20</sub> 1081 <sub>21</sub> 4. The frequency of these observations is notable: **1082** 22 - Cardiomegaly or enlarged cardiac silhouette is mentioned in 7 out of 1083 10 reports. **1084** <sup>23</sup> - Aortic abnormalities are mentioned in 5 out of 10 reports. 1085 .24 25 5. Given that cardiomegaly is mentioned more frequently, it might be the 1086 primary feature the neuron is detecting. However, the aortic 1087 abnormalities are often mentioned in conjunction with cardiomegaly. 1088 26 6. It's possible that the neuron is detecting a combination of these 1089 27 features, as they often co-occur and are related to cardiovascular 1090 changes. 1091 <sub>28</sub> 1092 29 7. The neuron might be sensitive to the overall appearance of an enlarged cardiac silhouette combined with changes in the adjacent thoracic 1093 aorta, which together create a distinctive pattern on chest 1094 1095 <sub>30</sub> radiographs. 1096 <sub>31</sub> Based on this analysis, I believe the most precise and specific 1097 explanation for what this neuron is detecting is: 1098 32 1099 33 \*This feature represents an enlarged cardiac silhouette (cardiomegaly) in conjunction with thoracic aortic abnormalities, particularly 1100 tortuosity and calcification. 1101 1102 Finally, we use regex to parse the description using the asterisk as a separator. 1103 1104 B.5 PROMPT USED TO COMPOSE FINDING SECTIONS VIA SAE FEATURES 1105 1106 In this section we describe the prompt used to instruct Claude to compose the finding sections of 1107 a radiology report based on: 1) a patient's past radiology reports, if they exist, 2) the features that 1108 SAE-Rad has identified are present within the image and 3) the indication for the radiograph. Each 1109 feature is associated to an *importance score* in accordance to the magnitude of the SAE dimension 1110 of provenance — the higher the magnitude, the higher the importance score. The scores were nor-1111 malised to the interval between 0 and 1. Our prompt encourages Claude to focus on features with 1112 the highest scores as these were likely to be more significant. See Appendix C.1 for more example 1113 outputs following this prompt, alongside their corresponding radiographs. 1114 You are an expert radiologist specializing in chest radiographs. Your 1 1115 task is to write the findings section for a radiology report based on 1116 a chest X-ray image. To assist you, I may provide up to three of the 1117 patient's past radiology reports, if available. These might contain useful information related to the features of the current scan. I 1118 will also give you the indication (reason) for the current X-ray. 1119 Additionally, you'll receive text descriptions of features present in 1120 the current X-ray image, along with importance scores for each 1121 feature. Your primary focus should be on producing the findings 1122 section for the latest scan, given the features about that scan. Focus on features with higher importance scores, as these are more 1123 prominent in the image and should be emphasized. Assess the current 1124 features, and then judge whether it would be appropriate to relate 1125 them to information in previous scans, if provided. Do not explicitly 1126 mention dates and times from previous reports. Discuss the features 1127 present in the X-ray, along with their implications and any 1128 deductions you can make. Your response should constitute the ' findings' section of the radiology report, providing a comprehensive 1129 analysis of the current X-ray. All of the information is provided 1130 below: 1131 2 <patient\_history> 1132 3 <past\_report> Report number 1. This report was written 0 years, 2 days, 16 hours and 28 1133 4 minutes before the current chest x-ray

```
1134
      INDICATION: Left-sided pleuritic chest pain ...
    5
1135
       COMPARISON: Chest radiograph ___ and chest CT __
    6
1136 7
       IMPRESSION: Ill-defined patchy opacities in lung bases which may
           represent ...
1137
1138 <sup>8</sup>
       </past_report>
1139 <sup>9</sup><sub>10</sub>
       </patient_history>
1140 11
       <current_chest_x_ray>
1141 12
       <feature 1>
1142 13
       Feature number 1. Relative importance score 1.0:
1143 <sup>14</sup>
       This feature represents the absence of pneumothorax, characterized by
1144 <sub>15</sub>
           normal lung appearance at the pleural margins.
       </feature 1>
1145 <sub>16</sub>
1146 17
       <feature n>
       </current_chest_x_ray>
1147 18
       Using the information provided, compose the findings section of the
1148 19
           radiology report. Be aware that some of the described features may be
1149
            inaccurate or only loosely related to the actual characteristics
1150
           present in the X-ray. When faced with conflicting information, rely
1151
           on the importance scores or a majority consensus to determine which
           features are most likely correct. In your report, refrain from simply
1152
            listing the features. Avoid using the word 'feature' entirely in
1153
           your report. Keep the radiology report brief and to the point. The
1154
           reason for the current x-ray examination is provided below:
1155 <sub>20</sub>
1156 21
       <indication>
       Status post CABG.
1157 22
       </indication>
1158 <sup>23</sup>
1159 <sup>24</sup><sub>25</sub>
       Now write the findings section. This should be a single contiguous
1160
           paragraph with the findings of the X-ray radiology report. No more
1161
           than 5 to 6 sentences. Be concise and avoid simply listing the
           features. Do not respond with any additional text other than the
1162
           findings. Do not add any concluding statements at the end, only
1163
           include findings.
1164
```

1167

#### B.6 BASELINE EXPERIMENTAL SETUP

In this section, we provide details on the experimental details of our baseline report generation method. In summary, the baseline approach returns the 'closest' report in the train split. Explicitly, we encode a given radiograph using Rad-DINO into a single embedding vector  $\mathbf{z} \in \mathbb{R}^{768}$  by selecting the class token. Then, we identify the radiograph y with the closest Rad-DINO class token from the MIMIC-CXR train split X as measured by the Frobenius norm  $y = \arg \min_{p \in X} ||\mathbf{z}(x) - \mathbf{z}(p)||_2$ . We then return the original findings section from y as our predicted, baseline report.

1174

1176

1178

### 1175 C ADDITIONAL QUALITATIVE EVALUATIONS

1177 C.1 Additional radiology report generations

In this section we provide additional qualitative examples of SAE-Rad generated findings sections, highlighting findings jointly described between radiologist generated reports, missing findings by SAE-Rad, and correct findings highlighted by SAE-Rad not present in the original report. Descriptions and analyses of findings (dis)similarity were produced and reviewed by a clinician.

1183

1103

1184 1185

1186



1242

1243





# **A: Ground-Truth Report**

The lungs are hyperinflated but clear of consolidation. Linear opacity in the right mid to upper lung is compatible with scarring as well as changes of the posterior right ribs which are chronic. Blunting of the right lateral posterior costophrenic angle is chronic, potentially due to scarring or trace effusion. Blunting of the left posterior costophrenic angle suggests small pleural effusion. Cardiomediastinal silhouette is within normal limits. Prominent retrocardiac opacity on the right is compatible with a neo esophagus. No acute osseous abnormalities.

# **B: SAE-Rad Report**

The chest radiograph demonstrates linear opacities in the right lung, consistent with scarring or subsegmental atelectasis. There are focal areas of consolidation or atelectasis noted, which may represent residual changes from previous infections or aspiration events. The right posterior pleural sulcus appears blunted, possibly due to a small chronic pleural effusion or scarring, as noted in prior studies. The cardiac silhouette is within normal limits, and there is no evidence of acute pulmonary edema or pneumothorax. Emphysematous changes are present, reflecting the patient's known history of emphysema. Post-surgical changes from previous esophagectomy and gastric pull-through are noted and unchanged from prior examinations.

Figure 6: **Comparison between SAE-Rad Report and the Ground Truth (2)**. Hyperinflation in the lungs, right lung linear opacities, blunting of the right posterior pleural sulcus and the normal limits of the cardiac sihloutee are noted in both reports. SAE-Rad detects the presence of focal areas of consolidation or atelectasis, the patients known history of emphysema, and the gastric pull-through that remains unchanged from prior exams. SAE-Rad misses the changes of the posterior right ribs and the lack of acute osseous abnormalities.



Figure 7: **Comparison between SAE-Rad Report and the Ground Truth (3)**. Ground-truth and SAE-Rad reports identify the normal catdiomediastinal contours, clear lung fields, as well as the lack of pleural effusion or pneumothorax. SAE-Rad misses that the patient underwent CABG and that their heart, mediastinum and elevated right hemidiaphraghm remain unchanged. SAE-Rad correctly comments on the lack of focal consolidation and cardiopulmonary abnormalities, and notes acute findings with respect to the patients past diagnosis of melanoma. Additionally, SAE-Rad notes the presence of mediastinal clips and a slight reduction in right lung volume.









#### C.2 IMAGE EXAMPLES OF MONOSEMANTIC FEATURES

In this section we showcase highest activating images for a number of features, as well as the corresponding feature explanations. We highlight the variety of features captured by SAE-Rad, from instrumentation features in Fig. 12, Fig. 13, Fig. 14, as well as visual features such as radiograph inversion in Fig. 16, pathology-related features in Fig. 17, and small details such as piercings in Fig. 18. 



Figure 12: Maximally activating images for a feature corresponding to deep brain stimulators which are typically used to treat Parkinson's disease.



Figure 13: Maximally activating images for a feature corresponding to orthopaedic rods and screws.







Figure 15: Maximally activating images for a feature corresponding to female radiographs with no pathology detected.







Figure 17: Maximally activating images for a feature corresponding to bowel obstruction.

 Under review as a conference paper at ICLR 2025



Figure 18: Maximally activating images for a feature corresponding to piercings.

# 1890 D IMAGE-BASED FEATURE LOCALIZATION IMPLEMENTATION

1892 **Model Training** In order to visualise intervention effects, we trained a conditional denoising diffusion probabilistic model (DDPM) Ho et al. (2020) to generate chest radiographs given Rad-DINO 1894 tokens. Our DDPM architecture is a three block, 2D UNet, where each block consists of two resid-1895 ual sub-units with SiLU activation and group normalisation. The downward path consists of a first block of 2D convolutions, followed by spatial self-attention for the second and third blocks. The upward path consists of two spatial self-attention blocks followed by upsampling for the last block. 1897 The DDPM is conditioned by linearly projecting the CLS token  $\mathbf{z} \in \mathbb{R}^{768}$  to a vector  $\mathbf{v} \in \mathbb{R}^{1024}$ , 1898 broadcasting the result to a 2D tensor  $\mathbf{j} \in \mathbb{R}^{128 \times 128}$ , and concatenating  $\mathbf{j}$  channel-wise with a sam-1899 ple  $n \in \mathbb{R}^{128 \times 128}$  as an input to the UNet. The DDPM was trained for 330 epochs using a 90/10 1900 train/val split of the MIMIC-CXR dataset, with a batch size of 256 and an AdamW optimizer with 1901 a learning rate of 0.001. We follow standard diffusion model training procedures by randomly sam-1902 pling timesteps  $\{t \in \mathbb{Z} | 0 \le t \le 1000\}$ , noising an image  $\mathbf{x}_0 \subset \mathbf{X}$  for t steps to obtain noised image 1903  $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon$ , where  $\alpha_t$  is a scalar controlling the variance of  $\mathbf{x}_t$  as a function of t, and 1904  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . The objective function can be defined over the mean squared error between diffusion 1905 model  $\epsilon_{\theta}$  outputs conditioned on class tokens z:

$$\mathcal{L} = \arg\min_{\theta} \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim U(0, 1000), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[ \| \epsilon_{\theta}^{(t)} (\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t}, \mathbf{z}) - \epsilon \|_2^2 \right]$$
(12)

Training was implemented using PyTorch (Paszke et al., 2019) and the diffusers (Huggingface, 2023) package on a single LambdaLabs 8xNVIDIA H100 cluster.

1911 1912 Intervention Experiments We define a qualitative experiment to assess the visual interpretability 1913 of SAE-Rad features by evaluating whether activating or deactivating particular SAE features results 1914 in the expected visual changes. These interventions can be defined using a constant reassignment 1915 operator

$$do(i,\beta;\mathbf{h}(\mathbf{x})):\mathbf{h}(\mathbf{x})\mapsto \tilde{\mathbf{h}}(\mathbf{x})_{i,\beta}$$
(13)

which sets a particular feature in the SAE activation  $\mathbf{h}(\mathbf{x})$  at index *i* to a desired constant value  $\beta$ . Following Eq. (9), a given image token z may be reconstructed into an intervention token  $\tilde{\mathbf{z}}_{i,\beta}$  by decoding it's intervened-upon activation  $\tilde{\mathbf{h}}(\mathbf{x})_{i,\beta}$ 

1920 1921

1922

1916

1907 1908

$$\tilde{\mathbf{z}}_{i,\beta}(\mathbf{h}(\mathbf{z})) := W^{\text{dec}} \mathsf{do}(i,\beta;\mathbf{h}(\mathbf{z})) + \mathbf{b}^{\text{dec}}.$$
(14)

Given a feature of interest *i*, intervention quantity  $\beta$ , and token **z**, we can compute a single "counterfactual" image  $\tilde{I}_{i,\beta}$  by reverse denoising  $\tilde{\mathbf{z}}_{i,\beta} + \delta$  using  $\hat{\epsilon}_{\theta}$ , where  $\delta = \hat{\mathbf{z}} - \mathbf{z}$  accounts for the reconstruction error inherent to the SAE between the original token **z** and reconstructed  $\hat{\mathbf{z}}$ .

1926 Using this framework, we assessed the counterfactual cyclic consistency for a given feature by 1927 performing two sequential interventions: we generated  $\tilde{I}_{i,\beta}$  by *activating* feature *i* by  $\beta$  for a token 1928 z, and subsequently reverted the effects of the intervention by *deactivating*  $\beta$  at *i* in  $\tilde{z}_{i,\beta}$  to generate 1929  $\tilde{I}_{i,0}$ . Intuitively,  $|\tilde{I}_{i,0} - \tilde{I}_{i,\beta}|$  should represent the visual concepts of interest. In our experiments, 1930 we empirically set  $\beta = 15$ , and select features 311 (pacemaker feature) and 162 (cardiomegaly) 1931 our features of interest. Features were empirically identified by reviewing highest activating feature 1932 subsets for our final SAEs. Results are showcased in Fig. 4.

1933

# <sup>1934</sup> E READER STUDY

1935

In addition to qualitative analyses performed in Section 5.1, we performed a reader study to assess radiologist opinions of generated reports compared to currently reported methods. We followed the reader study design from MAIRA-2 Bannur et al. (2024) - in summary, we presented a radiologist with a radiograph, the indication for the exam, and the generated report, and subsequently tasked the radiologist with assessing reports sentence by sentence, where the radiologist could edit, delete or maintain the sentences. Where changes were made, we additionally tasked the radiologist to categorise errors as one or multiple of "omission", "misclassification", "overspecific", "incorrect location", and "other". We omitted the category "incorrect progression", as we do not include past radiographs for consideration as part of our reader assessment. Additionally, unlike the MAIRA-2 reader study, errors could be classified as multiple error categories, which allowed for increased granularity in labelling. Finally, errors were categorised according to their clinical implications errors could be "significant", "minor", or "none". The custom user interface (UI) used for the study is showcased in Fig. 19. For full details on error taxonomy, we refer the reader to Bannur et al. (2024).

| ><br>1/30 |             |                              |                                   |                   |            |
|-----------|-------------|------------------------------|-----------------------------------|-------------------|------------|
| Previous  |             |                              | Next                              |                   |            |
|           |             | Sentence 1                   |                                   |                   |            |
|           | 100 C 100 C | Edit sentence                |                                   |                   |            |
|           | S. 40       | There is no focal consolidat | ation, effusion, or pneumothorax. |                   |            |
|           |             | Error type:                  |                                   |                   |            |
|           |             | Omission Misci               | lassifica 🗌 Overspecific          | Incorrect Incorre | ct 🗌 Other |
|           |             | tion                         |                                   | Location Progres  | ssion      |
|           |             | Error Comments               |                                   |                   |            |
|           |             |                              |                                   |                   |            |
|           |             |                              |                                   |                   |            |
|           |             | Clinical implication:        |                                   |                   |            |
|           | X AND       | Significant                  | Minor                             | None              |            |
|           |             | Severity Comments            |                                   |                   |            |
|           | 8 P         |                              |                                   |                   |            |
| NO DESS   |             |                              |                                   |                   |            |
|           | eli E       |                              |                                   |                   |            |

Figure 19: Custom UI for radiograph evaluation.

Overall, a radiologist (ST3) reviewed twenty unique radiographs with findings section extracted from the ground-truth closest report baseline described in Appendix B.6, as well as findings generated by CheXagent and SAE-Rad; in total, sixty individual reports were assessed. The radiologist was blinded to the model that findings originated from during the study.

#### 1975 E.1 RESULTS

1957

1959

1961

1968 1969

1970

1971

1972

1973 1974

Overall, a total of 328 sentences were analysed, with an average of 5.5 sentences per report (CheX-agent: 5.0, SAE-Rad: 5.9, Baseline: 5.5). Of 328 sentences, 240 (73%) sentences were edited or deleted, with 88 (27%) sentences maintained as is. Each report, on average, required 4.0 number of edits. Split by model, CheXagent had 74 (74%) sentences with edits, the baseline totalled 85 (77%) sentences with edits, and SAE-Rad totalled 81 (69%) sentences with edits.

Of sentences that required edits, 128 (68%) sentences overall had no clinical impact, 72 (30%) had a minor impact, and 30 (12%) had a significant impact. By model, CheXagent had 38 (51.4%) sentences with edits with no clinical impact, 23 (31.0%) sentences with minor impact and 13 (17.5%) sentences with significant clinical impact, the baseline had 44 (51.7%) sentences with edits with no clinical impact, 30 (35.3%) sentences with minor impact and 11 (13%) sentences with significant clinical impact, and SAE-Rad had 56 (69.1%) sentences with edits with no clinical impact, 19 (23.5%) sentences with minor impact and 6 (7.4%) sentences with significant clinical impact.

1988 Sentences with edits were manually reviewed to assess their hallucinatory potential based off of the 1989 radiologist's feedback. We define hallucination as the presence of a clinical fact that is not verifi-1990 able or incorrect. This obviates edits which, for example, arise due to referencing previous studies, or edits resulting from minor grammatical differences or re-wording of clinical facts. Overall, we 1992 find that 124 (53.4%) of edits were classified as hallucinations, whilst 108 (46.5%) of edits were 1993 not classified as hallucinations. Of these edits, SAE-Rad produced the fewest hallucinations, with the baseline generating the most. Specifically, SAE-Rad produced 50 (64.1%) non-hallucinatory sentences, and 28 (35.9%) sentences that were classified as hallucinations. CheXagent produced 30 (41.1%) non-hallucinatory sentences, and 43 (58.9%) sentences that were classified as halluci-1996 nations. The baseline method produced 28 (34.5%) non-hallucinatory sentences, and 53 (65.4%) 1997 sentences that were classified as hallucinations.

In total, 160 (46.6%) errors were classified partly as "other", representing the majority class. 52 (15.1%) errors were omissions, 59 (17.2%) errors were mis-classifications, 57 (16.6%) errors were overspecific, and 15 (4.4%) errors were incorrect location. Error type and severity distribution are visualised in Fig. 20. Notably, whilst many "other" type errors were observed, the majority of these errors carried no clinical risk. In contrast, the "omission" and "misclassification" categories convey much higher risk, with a comparatively higher proportion of minor and significant severity categories. Overall, for sentences with edits, each sentence was labelled with an average  $1.40 \pm$ 0.68 error categories (SAE-Rad:  $1.42 \pm 0.70$ , CheXagent:  $1.35 \pm 0.69$ , baseline:  $1.42 \pm 0.66$  error categories per edited sentence).



Figure 20: **Distribution of error types for different models**. Notably, SAE-Rad achieves the lowest number of "significant" errors in comparison to alternate methods.

Overall, SAE-Rad requires 7% fewer edits than the baseline, and 5% fewer edits than CheXagent.
 Additionally, significant clinical impact error rates are almost half for reports generated using our method compared to other models, whilst maintaining comparable minor and no impact error types.
 Similarly to results reported by Bannur et al. (2024), the majority of errors have no clinical impact.
 Overall, this showcases the potential our proposed method for radiology report generation in a real clinical scenario.

## F COMPUTE EFFICIENCY COMPARISONS

In this section we provide and discuss the compute estimates for training and inference of SAE-Rad in comparison to both MAIRA 2 and CheXagent. In Table 4 we provide quantitative estimates of the 2035 compute efficiency - we estimate that SAE-Rad uses approximately 1000x less train compute, 100x 2036 less parameters and 4x less data. We note that while in this paper we combine the SAE text features 2037 into a report using Claude, our method enables the generation of reports using a much smaller LLM, 2038 or alternately without an LLM — for example, by concatenating the SAE feature text explanations 2039 into a report using only regex. This reduces the inference FLOPs and cost to a negligible quantity. 2040 In Table 4, we have estimated the lower bound inference compute requirements for SAE-Rad using 2041 this assumption.

|   | SAE-Rad   | cheXagent | MAIRA 2     |
|---|-----------|-----------|-------------|
| Number of unique training images          | 239,931   | 1,100,000 | 510,848     |
| Active training parameters                | 76M       | 8B        | 7B          |
| Estimated training Tera-FLOPs             | 60,000    | NA        | 100,000,000 |
| Estimated training cost                   | \$3.25    | NA        | \$4000      |
| Estimated inference Tera-FLOPs per report | 0.1       | 8         | 7           |
| Estimated inference cost per report       | \$0.00005 | \$0.001   | \$0.001     |

2049 2050

2022

2023 2024

2025

2026

2027

2028

2029

2030 2031

2032 2033

Table 4: Comparison of compute resources for different models.

#### 2052 LIMITS OF AUTOMATED QUALITY ASSESSMENT FOR RADIOLOGY REPORTS G

2053

2054 As described in Appendix B.2, NLG metrics such as BLEU-4 (Papineni et al., 2002), ROUGE-L 2055 (Lin, 2004), and METEOR (Banerjee & Lavie, 2005) do not account for the clinical relevance of the 2056 sentences composing a radiology report. Radiology-specific metrics such as the RGER score (Del-2057 brouck et al., 2022) and the CheXpert F1 score (Smit et al., 2020) were consequently developed, 2058 seeking to improve upon these limitations by using radiology-specialized models as backbones to compute clinically relevant commonalities between generated and reference reports. However, 2060 CheXpert and RGER rely on pre-specified findings classes and error types (Bannur et al., 2024), which limits their ability to assess the quality of generated reports when elements outside these 2061 specified categories are present. RadFact (Bannur et al., 2024) was recently proposed to address 2062 this limitation; by relying on the logical inference capabilities of large language models (LLMs), 2063 RadFact evaluates the correctness and completeness of generated reports outside of pre-specified 2064 classes. 2065

- RadFact introduces two key metrics: RadFact logical precision and RadFact logical recall, hence-2066 forth denoted as RadFact precision and recall. RadFact precision calculates the proportion of gener-2067 ated sentences that are entailed by the ground-truth report, assessing the truthfulness of the model's 2068 outputs by penalizing hallucinations. RadFact recall measures the proportion of ground-truth sen-2069 tences that are entailed by the generated report, evaluating the completeness of the generated report 2070 by penalizing omissions. 2071
- However, RadFact has several limitations. Firstly, it relies on an LLM to extract discrete statements 2072 from reports, which may introduce inaccuracies. The LLM must also perform bi-directional en-2073 tailment verification—a challenge for models like Llama3-70B-Instruct (RadFact's backbone) given 2074 the complex medical reasoning in the reports. Furthermore, LLMs face the reversal curse (Berglund 2075 et al., 2023): they excel at forward relationships (e.g., 'France's capital is Paris') but struggle with reverse ones ('Paris is the capital of what country?'), compounding the difficulty of bi-directional 2077 entailment verification. A particularly important limitation of RadFact is that it relies on the as-2078 sumption that the ground-truth report used as a basis for evaluation is a gold-standard label which 2079 exhaustively captures clinically relevant detail.
- 2080 Therefore, in cases where a predicted report correctly includes findings present in the image but 2081 omitted from the ground-truth report, RadFact precision will drop. This means that the metric may 2082 not fully reflect the model's ability to generate clinically accurate and comprehensive reports, when 2083 the ground-truth reports are incomplete. 2084
- We illustrate potential pitfalls of RadFact in Fig. 21 below. In this case, the ground-truth report fails 2085 to report several clinically relevant details observed in the SAE-Rad report, including, for example, 2086 the normal size of the cardiac silhouette, and the increased density in the right upper lung field. We 2087 provide additional examples of this phenomenon in Appendix C.1.
- We note that the level of detail for radiographic reports is likely impacted by multiple factors, such 2089 as report indication, patient history, urgency of the request, and care setting of the radiograph. Given 2090 we perform evaluation on MIMIC-CXR —a dataset of radiographs captured at an emergency depart-2091 ment (Johnson et al., 2019)— SAE-Rad is likely to report clinically correct and relevant statements 2092 which are not explicitly mentioned by radiologists. Therefore, reported metrics must be interpreted 2093 with care, and may be underestimates of actual clinical factuality and relevance of generated reports. 2094
- 2095 2096

2097

#### Η ABLATION EXPERIMENTAL SETUP

In this section we present the experimental setup for our ablation studies. We evaluated six SAEs on 2098 RadFact across three expansion factors and two sparsity levels. We evaluated the SAEs using a set 2099 of 60 randomly selected radiographs from the MIMIC-CXR test split. RadFact was evaluated using 2100 the "Llama3-70B-Instruct" model. Additionally, we report NLG and clinical metrics to assess the 2101 difference in quality of reports across models. In this study, we used an earlier prompting template 2102 to generate feature explanations and subsequent reports. Overall, we note that this change resulted in 2103 slightly worse quantitative metrics than those that were obtained using the final prompting strategy. 2104

All SAEs were trained for 200,000 steps with a batch size of 2048. Models were trained using an 2105 Adam optimizer with no weight decay and a learning rate of  $5 \times 10^{-5}$ , as well as a linear warm-up



detail than the ground-truth. SAE-Rad captures nearly all features noted in the ground-truth such as: fibronodular changes, a diagnosis of sarcoidosis and pleural effusions. However SAE-Rad also notes the lack of pneuomothorax, pulmonary edema and the increased density in the right upper lung field. Despite the clinical relevance of these elements this report likely scored poorly on RadFact due to the incompleteness of the ground-truth.

of the learning rate for the first 1% of training, linear warm-down of learning rate for the last 20% of training, and an L1 coefficient warm-up for the first 5% of training.

We summarise experimental differences between runs in Table 5, as well as resultant L0 and ex-plained variance metrics. The L0 and explained variance metrics in Table 5 provide insights into the quality of our SAEs. Notably, the "dense" variants exhibit higher L0 values, indicating reduced sparsity, but achieve greater explained variance. The different variants therefore represent different points on the Pareto frontier between sparsity and reconstruction quality.

Table 5: Varying hyperparameters and resulting metrics for the six SAEs.

|            | Expansion factor | L1 coefficient        | $\text{L0}\downarrow$ | Explained variance (%) $\uparrow$ |
|------------|------------------|-----------------------|-----------------------|-----------------------------------|
| ×128       | 128              | $6 \times 10^{-3}$    | 12.0                  | 90.9                              |
| ×64        | 64               | 8 × 10 <sup>-3</sup>  | 13.6                  | 84.3                              |
| ×32        | 32               | 9 × 10 <sup>-3</sup>  | 15.1                  | 79.2                              |
| ×128_dense | 128              | $3.75 \times 10^{-3}$ | 26.0                  | 93.2                              |
| ×64_dense  | 64               | 5 × $10^{-3}$         | 28.4                  | 88.5                              |
| ×32_dense  | 32               | $5.63 \times 10^{-3}$ | 31.8                  | 84.7                              |

Results for clinical and NLG metrics for the ablation study are presented in Table 6. Overall, we find that the ×64 SAE (i.e SAE-Rad) is the maximally performant model for 6/9 of the metrics, with an additional 2/9 metrics being the second best of all models.

#### INTERACTIVE EXPLORATION OF SAE FEATURES Ι

To investigate the geometry of SAE features (Mendel, 2024), we embed the feature directions from the SAE decoder into a 2-dimensional space using a UMAP (McInnes et al., 2018) (for which we set the number of neighbors to 50 and the minimum distance to 0.05). This analysis allows us to Table 6: Report generation performance with clinical metrics for ablation study radiographs.
BL4 = BLEU-4, RG-L = ROUGE-L, MTR = Meteor. Ma-5 (Macro-F1-5), Ma-14 (Macro-F1-14), Mi-5 (Micro-F1-5), and Mi-14 (Micro-F1-14) represent the clinical CheXbert labeler scores.
Bolding represents best performance on the development set, and italics represent second best model for a given metric.

| Model      | RadFact $\uparrow$ NLG Metrics $\uparrow$ |     |      | Clinical Metrics $\uparrow$ |      |      |       |      |       |
|------------|---|-----|------|-----------------------------|------|------|-------|------|-------|
| inouer     | F1  | BL4 | RG-L | MTR                         | RGER | Ma-5 | Ma-14 | Mi-5 | Mi-14 |
| ×32 dense  | 30.22                                     | 1.5 | 18.6 | 23.3                        | 19.5 | 51.8 | 33.2  | 55.0 | 52.6  |
| ×64 dense  | 29.74                                     | 1.9 | 19.5 | 23.9                        | 19.5 | 57.8 | 31.9  | 58.7 | 55.5  |
| ×128 dense | 29.56                                     | 1.7 | 16.1 | 22.8                        | 17.3 | 46.9 | 33.9  | 57.5 | 53.7  |
| ×32        | 29.46                                     | 2.2 | 16.2 | 23.7                        | 16.6 | 44.2 | 25.3  | 53.5 | 48.8  |
| ×64        | 33.83                                     | 2.4 | 17.1 | 24.4                        | 20.7 | 54.7 | 33.6  | 58.9 | 56.8  |
| ×128       | 32.18                                     | 1.8 | 16.4 | 24.1                        | 18.9 | 45.2 | 25.7  | 55.8 | 49.8  |

visualize whether features which have semantically similar descriptions are clustered in the latent representations of the SAE. Text information is not used to train the Rad-Dino vision encoder, SAEs, or the UMAP projections of the features. Therefore, assessing the clustering patterns based on text descriptions allows us to independently validate the quality of learned representations. We develop interactive tooling to evaluate the results of this analysis. Clicking on each point will show the text description of the feature and load highest activating images. The tool can be accessed using the following link: https://scatter-plot-app.vercel.app/. For convenience, we also illustrate a number of well-clustered feature categories in Fig. 22. 

It can be seen from Fig. 23 that the clustering enables the detection of possibly incorrect feature descriptions. The figure provides an example relating to a feature that captures the orientation of radiographs, rather than the currently described pathology. This approach can be used to ensure hallucination-free feature descriptions.



Figure 22: **UMAP analysis for several feature categories.** Subfigures show examples of UMAP embeddings for (a) Breathing Tube, (b) Clear Lungs, (c) Pacemaker, and (d) Post-Surgical cases.



Figure 23: **Illustration of an orientation cluster.** Here we show an example of how an automated description can be manually visualized and corrected. In this the case the feature actually represents sideways orientation radiographs. To explore this cluster please visit the interactive tool here: https://scatter-plot-app.vercel.app/

## J ADDITIONAL SAE-RAD COMPARISONS

We present additional comparisons to our method by collating reported clinical efficacy metrics andNLG metrics in Table 7. We find that SAE-Rad broadly outperforms existing methods across all clinical metrics.

Table 7: **Report generation performance on the official MIMIC-CXR test split on additional baselines.** BL4 = BLEU-4, RG-L = ROUGE-L. Ma-5 (Macro-F1-5), Ma-14 (Macro-F1-14), Mi-5 (Micro-F1-5), and Mi-14 (Micro-F1-14) represent the clinical CheXbert labeler scores as reported in the original studies. Bolding represents best performance in the current study or between the upper bound models.

| Model                                      |      | Metrics $\uparrow$ | Clinical Metrics ↑ |       |      |       |  |
|--|------|--------------------|--------------------|-------|------|-------|--|
| Model                                      | BL4  | RG-L               | Ma-5               | Ma-14 | Mi-5 | Mi-14 |  |
| GPT-4V Wu et al. (2023)                    | 1.9  | 13.2               | 19.6               | 20.4  | 25.8 | 35.5  |  |
| LLaVa-Med Wu et al. (2023)                 | 1.0  | 13.3               | 16.6               | 15.5  | 22.0 | 27.2  |  |
| CvT2Dist Codella et al. (2024)             | 12.7 | 28.6               | -                  | 30.7  | -    | 44.2  |  |
| LLaVa Codella et al. (2024)                | 1.3  | 13.8               | 17.5               | 15.4  | 23.4 | 22.9  |  |
| GPT-40 finetune Codella et al. (2024)      | 17.8 | 32.1               | 43.8               | 33.0  | 52.7 | 48.9  |  |
| GPT-40 mini finetune Codella et al. (2024) | 16.2 | 32.2               | 42.0               | 30.8  | 51.8 | 47.6  |  |
| R2GenGPT Wang et al. (2023b)               | 13.4 | 16.0               | -                  | -     | -    | 38.9  |  |
| METransformer Wang et al. (2023a)          | 12.4 | 29.1               | -                  | -     | -    | 31.1  |  |
| R2Gen Wang et al. (2023a)                  | 10.3 | 27.7               | -                  | -     | -    | 27.6  |  |
| R2GenCMN Wang et al. (2023a)               | 17.0 | 19.1               | -                  | -     | -    | 27.8  |  |
| MSAT Yang et al. (2023)                    | 11.1 | -                  | -                  | -     | -    | 33.9  |  |
| KIUT Huang et al. (2023)                   | 11.3 | 28.5               | -                  | -     | -    | 32.1  |  |
| RGRG Tanida et al. (2023)                  | 12.6 | 26.4               | -                  | -     | 54.7 | 44.7  |  |
| Flamingo-CXR Tanno et al. (2024b)          | 29.7 | 10.1               | -                  |       | 58.0 | 51.9  |  |
| SAE-Rad (x64)                              | 1.9  | 17.1               | 47.2               | 34.3  | 54.4 | 53.2  |  |
|  | >    |                    |                    |       |      |       |  |

2290 2291

2292

2298

2299 2300

2305

2306

2307

2311

2313

2316

2317

### K EQUIVALENCE OF LOSS FUNCTIONS BETWEEN SAE ARCHITECTURES

Here, we show that the original SAE loss function in Eq. 3 is isomorphic to the loss introduced by (Conerly et al., 2024) (displayed in Eq. 7). We show that they are related by a fundamental group of symmetries in the design of SAEs. In particular, we construct a group of symmetries that preserve the reconstructed activations of the SAE and leave the (Conerly et al., 2024) loss invariant. The authors modified the sparsity-inducing part of the loss, which is specified below:

$$\mathcal{L}_{ ext{sparsity}} = \sum_{i} \mathbf{h}_{i}(\mathbf{x}) \cdot \| W^{ ext{dec}}_{\cdot,i} \|_{2}$$

where  $W_{i,i}^{\text{dec}}$  is the decoder weight matrix and  $\mathbf{h}_i(\mathbf{x})$  is the activation value of the *i*<sup>th</sup> SAE feature.

For simplicity, we consider SAEs without any biases (though it should be noted that the explanation provided generalises to SAEs with biases). With this in mind, an SAE has the following architecture:

$$\begin{aligned} \mathbf{h}(\mathbf{x}_{in}) &= \text{ReLU}(W^{\text{enc}}\mathbf{x}_{in})\\ \mathbf{x}_{\text{out}} &= W^{\text{dec}}\mathbf{h}(\mathbf{x}_{in}) \end{aligned}$$

where  $\mathbf{x}_{in} \in \mathbb{R}^n$  is the input activation vector,  $W^{enc} \in \mathbb{R}^{m \times n}$  is the encoder matrix,  $W^{dec} \in \mathbb{R}^{n \times m}$ is the decoder matrix,  $\mathbf{h}(\mathbf{x}_{in})$  is the hidden activation of the SAE, and  $\mathbf{x}_{out} \in \mathbb{R}^n$  is the reconstructed activation. Note that there are no normalisation constraints on the decoder matrix.

2312 K.1 MOTIVATION

Let us begin by highlighting the fact that multiplication by positive constants commutes with the ReLU activation function:

$$\operatorname{ReLU}(\lambda x) = \lambda \operatorname{ReLU}(x), \forall x \in \mathbb{R}, \lambda > 0$$

2318 Motivated by this expression, we can define a second SAE with weights given by  $\tilde{W}^{enc}$ ,  $\tilde{W}^{dec}$  as 2319 follows:

$$\tilde{W}^{\rm enc} = {\rm diag}(\boldsymbol{\lambda}) W^{\rm enc}$$

$$ilde{W}^{ ext{dec}} = W^{ ext{dec}} ext{diag}(rac{1}{oldsymbol{\lambda}})$$

where  $\lambda \in \mathbb{R}^m$  is now a vector. This SAE is identical to the original SAE when viewed as a function  $\mathbf{x}_{in} \to \mathbf{x}_{out}$ . The only difference is that the hidden activations  $\tilde{\mathbf{h}}(\mathbf{x}_{in}) = \text{diag}(\lambda)\mathbf{h}(\mathbf{x}_{in})$  have been scaled by  $\text{diag}(\lambda)$ . Since this transformation does not change the output  $\mathbf{x}_{out}$ , it will not change the MSE reconstruction loss. This transformation will however change the  $l_1$  sparsity loss. The  $l_1$  loss changes by:

$$\tilde{l_1} = |\tilde{\mathbf{h}}(\mathbf{x}_{in})|_1 = |\text{diag}(\boldsymbol{\lambda})\mathbf{h}(\mathbf{x}_{in})|_2$$

It follows that if we do not impose the constraint of normalising the decoder weights, the SAE will use this symmetry in the limit  $\lambda \to 0$  to reduce  $l_1$  to 0. This explains the necessity of normalising the decoder weights.

2333 2334

2335

2329

#### K.2 EXPLANATION OF THE CONERLY ET AL. LOSS

Are there alternatives to normalising the decoder bias? What if instead we find a way to penalise the SAE for the transformation  $\lambda \to 0$  instead? We know that as  $\lambda \to 0$ ,  $||W_{\cdot,i}^{dec}||_2 \to \infty$ . Could we penalise the network by weighting the terms in the sum for  $l_1$  by the corresponding term  $||W_{\cdot,i}^{dec}||_2$ ? We know that the  $l_1$  loss is given by:

2340 2341

2342

2346 2347 2348

2354

2360

We also know that the components  $\mathbf{h}_i(\mathbf{x}_{in})$  transform in the opposite (covariant) way to decoder vectors  $W_{\cdot,i}^{\text{dec}}$ . In particular, the product  $\mathbf{h}_i(\mathbf{x}_{in})||W_{\cdot,i}^{\text{dec}}||_2$  remains invariant under our transformation. We therefore define an invariant loss function as follows:

$$\mathcal{L}_{ ext{sparsity}} = \sum_i \mathbf{h}_i(\mathbf{x}_{ ext{in}}) \cdot \| W^{ ext{dec}}_{\cdot,i} \|_2$$

 $l_1 = \sum_i |\mathbf{h}_i(\mathbf{x}_{\mathsf{in}})|$ 

We have therefore penalised the SAE in such a way that there is no insentive to transform by any value of  $\lambda$ ; both the MSE loss and sparsity loss are now completely invariant. This is precisely the loss function introduced by Conerly et al. (2024). If we use this loss function, we have found a group of symmetries that preserve the MSE and sparsity losses. This group is isomorphic to the direct product of the group ( $\mathbb{R}^+$ , ×) with itself *m*-times.

### 2355 K.3 EQUIVALENCE TO THE ORIGINAL LOSS

We now construct an equivalence to the old  $l_1$  sparsity loss function. Suppose we have an SAE trained with the loss in Eq. 7. We are now at will to transform with whatever value of  $\lambda$  we like since they define a group of symmetries preserving the new loss. Let us pick  $\lambda$  as follows:

$$\boldsymbol{\lambda}_i = ||W_{\cdot,i}^{\text{dec}}||_2$$

By using this transformation, we define an equivalent SAE in which the decoder weights are now normalised. In addition the resulting (invariant) sparsity loss function has the form of the original loss function, the  $l_1$  loss, since the decoder weights are now unit normalised.

This analysis shows that the Conerly et al. (2024) loss function is equivalent to the original SAE loss function but now has a natural invariance preserved under our symmetry group.

- 2367
- 2368
- 2370
- 2371
- 2372
- 2373
- 2374
- 2375