

# Learning to Exploit Temporal Structure for Biomedical Vision–Language Processing

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†

Microsoft Health Futures

## Abstract

Self-supervised learning in vision–language processing (VLP) exploits semantic alignment between imaging and text modalities. Prior work in biomedical VLP has mostly relied on the alignment of single image and report pairs even though clinical notes commonly refer to prior images. This does not only introduce poor alignment between the modalities but also a missed opportunity to exploit rich self-supervision through existing temporal content in the data. In this work, we explicitly account for prior images and reports when available during both training and fine-tuning. Our approach, named BioViL-T, uses a CNN–Transformer hybrid multi-image encoder trained jointly with a text model. It is designed to be versatile to arising challenges such as pose variations and missing input images across time. The resulting model excels on downstream tasks both in single- and multi-image setups, achieving state-of-the-art (SOTA) performance on (I) progression classification, (II) phrase grounding, and (III) report generation, whilst offering consistent improvements on disease classification and sentence-similarity tasks. We release a novel multi-modal temporal benchmark dataset, MS-CXR-T, to quantify the quality of vision–language representations in terms of temporal semantics. Our experimental results show the advantages of incorporating prior images and reports to make most use of the data.

## 1. Introduction

Self-supervision from image–text pairs has enabled the development of flexible general-purpose vision–language models both in the general domain [40, 53, 77] and for specialised domains such as biomedicine and radiology

\*These authors contributed equally.

†Corresponding author: ozan.oktay@microsoft.com

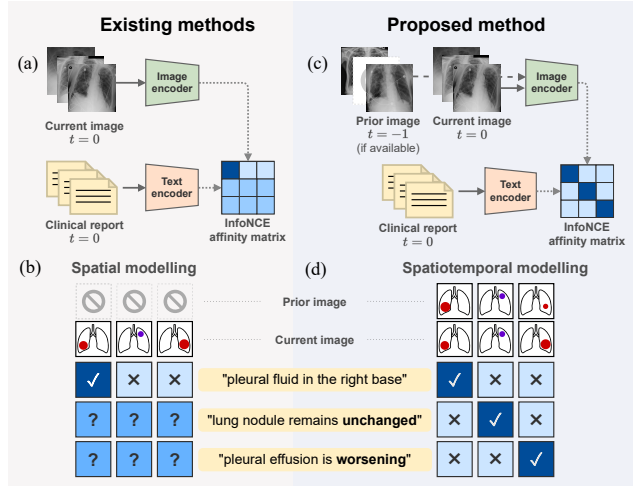


Figure 1. (a) Existing visual–language pre-training approaches [9, 32, 81] often use only a single image for contrastive learning (e.g., InfoNCE [49]). (b) In such settings, discarding the temporal connectivity of images limits the alignment of image–text pairs as shown with the affinity matrix, leading to suboptimal pre-training and missed opportunity to create additional model supervision for free. (c, d) Our approach exploits this domain knowledge by learning to incorporate a series of images and correlate them to reports, leading to pre-trained models that can generalise to a wider range of downstream tasks whilst achieving SOTA performance.

[9, 32, 81]. Vision–language processing (VLP) has shown that cross-modal supervision can provide a richer signal for training both image [19] and text [9] models. However, the success of VLP relies on paired samples sharing semantics, i.e., given an image and text pair, the text should describe the image with minimal extraneous detail [15, 16, 35].

In this regard, VLP in biomedicine and radiology poses a distinctive challenge, as reports routinely include comparisons to prior imaging studies [3, 47, 57]. Without knowl-

edge of this prior image<sup>1</sup>, temporal information in the text modality, e.g. “Pneumonia is improving”, could pertain to any image containing “Pneumonia”, producing ambiguity during contrastive training (Figure 1). Despite this, the existing VLP work to date considers alignment between only single images and reports [9, 32, 46, 81], going so far as to remove temporal content from reports in training data to prevent ‘hallucinations’ in downstream report generation [54]. However, temporal information can provide complementary self-supervision, solely by exploiting existing structure, and without requiring any additional data.

In this work, we neither ignore nor remove temporal information in the text modality, but explicitly account for it during pre-training. Rather than treating all image–report pairs in the dataset as independent, we exploit temporal correlations by making prior images available for comparison to a given report. To learn from this structure, we develop a temporal VLP pre-training framework named *BioViL-T*. A core component is its new multi-image encoder that can handle the absence of prior images and potential spatial misalignment between images across time. *BioViL-T* takes into account prior images where available, removing cross-modal ambiguity as illustrated in Fig. 1. Linking multiple images during pre-training proves beneficial to both image and text models: we report state-of-the-art (SOTA) performance on both temporal image classification and report generation. In the latter case, we show that prefixing the prior *report* substantially increases performance, again reflecting the value of prior information. We emphasise that the benefit is not restricted to temporal downstream tasks: our approach also achieves SOTA on non-temporal tasks of pneumonia detection [60] and phrase grounding [10], underscoring the value of a cleaner learning signal during VLP without needing to modify or add to the training dataset. Our contributions can be summarised as follows:

- We introduce a novel pre-training framework called *BioViL-T*. It leverages the temporal relationship of samples to self-supervise VLP models, making commonly used biomedical VLP models (e.g., [9, 32, 81]) more applicable to a wider range of downstream tasks without compromising performance on existing benchmarks.
- We develop a generic multi-image encoder that handles missing image inputs and incorporates longitudinal information without requiring explicit image registration.
- We achieve SOTA results in chest X-ray (CXR) report generation, temporal image classification, and phrase grounding downstream benchmarks by accounting for prior context in self-supervised training and fine-tuning.
- We release a new multimodal benchmark dataset, *MS-CXR-T*, curated by an expert radiologist. It enables

benchmarking of CXR VLP models in terms of temporal semantics extracted from image and text data.

## 2. Related work

**Vision–language processing** Self-supervised VLP can significantly reduce the need for manual labels required for the training of image encoders [19, 53]. The availability of large-scale paired image–text datasets has thus led to rapid development of general-purpose VLP models. Objectives include contrastive and discriminative image–text matching [40, 53, 69] including local variants [32, 76], auto-regressive (AR) captioning [4, 39, 77] and multi-modal masked modelling objectives [13, 40, 61].

**Biomedical vision–language processing** Paired medical image–report datasets were originally used for supervised learning via (typically) automated label extraction from clinical reports [33, 63, 70]. Using such datasets, advances in general-domain *self-supervised* VLP have been demonstrated to benefit biomedical imaging applications [9, 32, 81]. Work has incorporated ideas from general-domain VLP such as the original CLIP-style cross-modal contrastive objective [81], multi-modal masking with merged co-attention on image–text representations [46], and adaptations to the data of the domain. For example, a radiology report may have sparse image-specific details, prompting a local modification to the contrastive loss enabling alignment between text tokens and image patches [32]. Domain-specific pre-training of the text model is shown to benefit biomedical VLP [9], and preferential masking of medical terms during masked language modelling (MLM) was explored [75]. Here we use a local loss and domain-specific pre-training of the text model, but did not find a benefit to preferential masking. Similarly, cross-attention [22] is used rather than merged co-attention for image-guided MLM.

**Longitudinal modelling of medical images** While prior images are used in unimodal *supervised* longitudinal analysis of medical images [37, 58, 68, 74], temporal information has not directly been employed for self-supervision. The closest work exploits patient metadata to select positive or negative examples in unimodal contrastive learning [67, 79].

Existing models typically employ either late fusion of global image representations [58, 64, 68, 74], which can miss fine-grained localised changes [32], or explicit spatial correspondence of features, using fixed spatial grids [48] or object detection [37]. Registering image pairs is commonly used for change detection in other contexts [17, 52, 59], and has been applied to medical imaging [5, 23]. For CXRs however, registration entails the ill-posed problem of aligning 2D projections of 3D geometry, which inevitably results in residual misalignment. Our approach does not rely on bounding boxes or explicit graph construction as it uses

<sup>1</sup>In the MIMIC-CXR v2 dataset [36], around 40% of reports explicitly reference a previous image. See Appendix B for details.

self-attention of visual tokens across time to handle any spatial misalignment.

**Self-supervision across time** Self-supervision has found applications on densely-sampled time series data (e.g., video) to capture temporal information [30, 55, 78, 80]. Our problem setting involves sparsely and sporadically sampled data where temporal pretext tasks are less applicable [2]. Similarly, it requires text supervision to enable both static and temporal learning, when temporal structure is present.

### 3. BioViL-T training framework

Our approach comprises a multi-image encoder designed to extract spatio-temporal features from sequences of images (Section 3.1) and a text encoder incorporating optional cross-attention on image features. The models are trained jointly with image-guided MLM and cross-modal global and local contrastive objectives (Section 3.2). The resulting image and text models are later adapted for uni- or multi-modal downstream tasks as described in Section 3.3. Implementation details are presented in Appendices E and F.

For a given image and report pair  $(\mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{txt}}^{\text{curr}})$ , the report  $\mathbf{x}_{\text{txt}}^{\text{curr}}$  describes the current image content and changes in reference to prior images. Our proposed formulation focuses on a single prior image; however, it can be generalised to *multiple* prior images depending on the application. Hence, we construct datasets by including the prior image whenever it exists<sup>2</sup>:  $(\mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}}, \mathbf{x}_{\text{txt}}^{\text{curr}}) \in \mathcal{D}_m$  or  $(\mathbf{x}_{\text{img}}^{\text{curr}}, \emptyset, \mathbf{x}_{\text{txt}}^{\text{curr}}) \in \mathcal{D}_s$  with the resulting dataset being a union of single and multi-image examples:  $\mathcal{D} = \mathcal{D}_m \cup \mathcal{D}_s$ .

#### 3.1. Extracting spatio-temporal image features

Clinical findings are often observed across different image regions and co-occur simultaneously, which requires dense level visual reasoning across time to capture both static and temporal features. In contrast to late global fusion [64] and bounding-box based approaches [37], BioViL-T leverages local correspondences between image regions across time using transformer self-attention blocks [21]. Thus our method does not require an explicit image registration step between time points.

We propose a hybrid CNN–Transformer encoder model due to its data efficiency and spatial flexibility of cross-attention across time points:  $E_{\text{img}} : \mathbb{R}^{W \times H} \rightarrow \mathbb{R}^{W' \times H' \times D_{\text{img}}}$  (e.g., ResNet-50 [31]) and  $A_{\text{img}} : \mathbb{R}^{T \times L \times D_{\text{img}}} \rightarrow \mathbb{R}^{L \times D_{\text{img}}}$  (e.g., transformer [21]), where  $W$ ,  $H$ , and  $T$  correspond to spatiotemporal dimensions,  $L = W'H'$  is the number of visual tokens per image, and  $D_{\text{img}}$  is the embedding dimension. Here  $E_{\text{img}}$  serves as a stem network [51] to provide visual token features of individual images. The CNN’s inductive biases [24, 51] en-

sure data efficiency of our hybrid model, making it ideal for smaller scale biomedical datasets.  $E_{\text{img}}$  is initialised with BioViL weights [9]. The main purpose of  $A_{\text{img}}$  is to capture patch embedding interactions across time when a prior image  $\mathbf{x}_{\text{img}}^{\text{prior}}$  is available and to aggregate them into a fixed-length token representation. Input visual tokens,  $\mathbf{H}_0^{\text{curr}} = \mathbf{P}^{\text{curr}} := E_{\text{img}}(\mathbf{x}_{\text{img}}^{\text{curr}})$ ,  $\mathbf{H}_0^{\text{prior}} := E_{\text{img}}(\mathbf{x}_{\text{img}}^{\text{prior}})$  are augmented with spatio-temporal positional encodings and flattened across the spatial dimensions. They are then processed by  $K$  transformer encoder [66] layers  $A$  as follows:

$$\begin{bmatrix} \mathbf{H}_k^{\text{curr}} \\ \mathbf{H}_k^{\text{prior}} \end{bmatrix} = A_k \left( \begin{bmatrix} \mathbf{H}_{k-1}^{\text{curr}} + \mathbf{S} + \mathbf{1}_L \otimes \mathbf{t}^{\text{curr}} \\ \mathbf{H}_{k-1}^{\text{prior}} + \mathbf{S} + \mathbf{1}_L \otimes \mathbf{t}^{\text{prior}} \end{bmatrix} \right), \quad (1)$$

for  $k = 1, \dots, K$ , where  $\mathbf{S} \in \mathbb{R}^{L \times D_{\text{img}}}$  denotes 2D sinusoidal positional encodings [12] and  $\mathbf{T} = [\mathbf{t}^{\text{curr}}, \mathbf{t}^{\text{prior}}] \in \mathbb{R}^{2 \times D_{\text{img}}}$  is its temporal counterpart, which is learnt (Fig. 2) [4]. The layer-normalised (LN) [6] output of the final transformer encoder block  $\mathbf{P}^{\text{diff}} := \text{LN}(\mathbf{H}_K^{\text{curr}})$  is an ‘aggregated’ representation of patch-level progression information anchored on the current image. Figure 3 shows attention roll-out [1] applied to  $\mathbf{P}^{\text{diff}}$  after pre-training, showing how the prior image contributes to the fused representation. Figure A.3 further highlights the robustness to variations in pose underlining that registration is not necessary for this encoder.

**Static-temporal feature decomposition** When a prior image is available the final image representation  $\mathbf{V} := \mathbf{P}^{\text{curr}} \oplus \mathbf{P}^{\text{diff}} \in \mathbb{R}^{W' \times H' \times 2D_{\text{img}}}$  is formed by concatenating two sets of features (similar to [7]): those from the current image alone ( $\mathbf{P}^{\text{curr}}$ ) and the temporal features from current and prior images ( $\mathbf{P}^{\text{diff}}$ ). In this way, self-attention is mainly required to cope with pose variations and patch comparisons across time in extracting temporal content, removing the need for registration or explicit spatial feature alignment. When no prior scan is available ( $\mathbf{x} \in \mathcal{D}_s$ ),  $A_{\text{img}}$  is not used and  $\mathbf{P}^{\text{diff}}$  is replaced by a learnable token  $\mathbf{p}^{\text{miss}} \in \mathbb{R}^{D_{\text{img}}}$ , replicated across the spatial dimensions. Section 4.5 later demonstrates that  $A_{\text{img}}$  highlights the value of feature decomposition for tasks such as phrase grounding which require well-localised features [10].

Hereafter, downstream tasks that require solely single image features,  $\mathbf{P}^{\text{curr}}$ , are referred to as *static tasks*, and the ones that benefit from additional progression information,  $\mathbf{P}^{\text{diff}}$ , as *temporal tasks*, e.g., report decoding.

#### 3.2. Text-supervision for spatio-temporal learning

Let  $\mathbf{w} = (w_1, \dots, w_M)$  denote a vector of  $M$  tokens of a report  $\mathbf{x}_{\text{txt}}$  after tokenisation. We first obtain contextualised token features  $E_{\text{txt}}(\mathbf{w}) \in \mathbb{R}^{M \times D_{\text{txt}}}$  by passing a sequence of text tokens  $\mathbf{w} = (w_1, \dots, w_M)$  through a BERT encoder  $E_{\text{txt}}$  [20]. The input sequence is prepended with either a [CLS] or [MLM] token associated with a downstream training objective, conditioning the output features

<sup>2</sup>The prior *report* is not included during pre-training as it may further reference an earlier study, reintroducing temporal ambiguity.

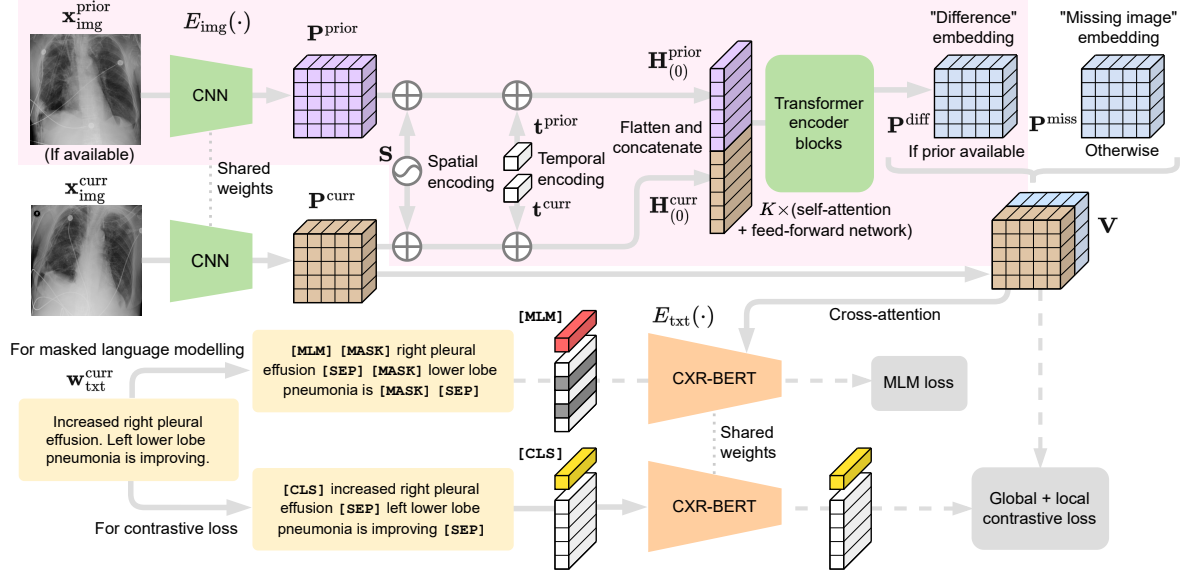


Figure 2. The proposed self-supervised VLP training framework BioViL-T: Image representations  $\mathbf{V}$  are extracted from single and multiple input scans (whenever available) using a hybrid CNN and transformer encoder [24, 51]. This design choice is to increase the data-efficiency and enable the fusion of temporal content without requiring image registration. They are later matched with their corresponding text representations obtained with CXR-BERT [9] using local [32] and global InfoNCE [49] training objectives. As an additional model supervision, multi-modal fused representations, obtained with cross-attention, are used for image-guided masked language modelling.

similar to [39, 42]. During training, we do two forward passes through  $E_{\text{txt}}$ : once with masking at 45% probability (for the MLM objective) and once without masking for contrastive learning, as shown in Figure 2. The text encoder is initialised with the weights of CXR-BERT<sup>3</sup> [9] canonical model, trained on domain-specific vocabulary and corpora.

Both text and image features are later projected into a joint latent space with  $\phi_{\text{txt}} : \mathbb{R}^{D_{\text{txt}}} \rightarrow \mathbb{R}^D$ , and similarly  $\mathbf{v}_{w,h}^{\text{proj}} := \phi_{\text{img}}(\mathbf{v}_{w,h})$  where  $\phi_{\text{img}} : \mathbb{R}^{D_{\text{img}}} \rightarrow \mathbb{R}^D$ , with  $\phi$  being a two-layer perceptron in our experiments.

**Contrastive objectives** Let  $\mathbf{r} := [E_{\text{txt}}(\mathbf{w})]_{[\text{CLS}]}$  denote the global representation of  $\mathbf{w}$ , with  $\mathbf{r}^{\text{proj}} := \phi_{\text{txt}}(\mathbf{r})$  its projected version. Given projected patch embeddings  $\mathbf{v}_{w,h}^{\text{proj}}$ , we can compute a global cosine similarity  $S_C(\bar{\mathbf{v}}^{\text{proj}}, \mathbf{r}^{\text{proj}})$  and a local similarity using weighted pairwise cosine similarities across text tokens and projected patch embeddings [32, 76]. These similarities are used in both global and local contrastive objectives with the InfoNCE loss [49, 53]. The local loss proves crucial both for static phrase-grounding and temporal image classification (see Table 7), highlighting the importance of localised self-supervision.

**Image-guided masked language modelling** Prior work [9, 46] has shown that biomedical visual-language learning benefits from an auxiliary task such as MLM since capturing the joint distribution of tokens can stabilise and improve

language understanding during joint learning. Given a batch  $\mathcal{B}$  of token vectors  $\mathbf{w}$ , it is often defined as the cross-entropy for predicting the randomly sampled masked tokens,  $m \in \{1, \dots, M\}$ ,  $\mathcal{L}_{\text{MLM}} = -\frac{1}{|\mathcal{B}|} \sum_{\mathbf{w} \in \mathcal{B}} \log p_{\theta}(\mathbf{w}_m | \mathbf{w}_{\setminus m})$ , where  $\theta$  are the weights of the text encoder  $E_{\text{txt}}$ .

In the absence of image information, however, certain masked findings and attributes are not readily predicted, e.g., “[MASK] is worsening”. As shown in the general domain [13], visual information can help disambiguate such masked predictions and provide additional cross-modal supervision. Thus, we use cross-attention [22, 66] to the image features  $\mathbf{v}_{w,h}^{\text{proj}}$  during this task. Specifically, for our image-guided MLM objective we model  $p_{\theta}(\mathbf{w}_m | \mathbf{w}_{\setminus m}, \mathbf{v}_{w,h}^{\text{proj}})$ .

### 3.3. Adaptations to downstream tasks

BioViL-T can be adapted to various downstream tasks. For phrase-grounding and zero-shot inference, we rely on  $S_C(\mathbf{r}^{\text{proj}}, \mathbf{v}_{w,h}^{\text{proj}})$  similar to [9, 32]. For multiple-text prompts, projected text embeddings are marginalised prior to  $\ell_2$ -normalisation [53]. To enable language decoding,  $\mathbf{v}_{w,h}^{\text{proj}}$  inputs are cross-attended by text queries  $\mathbf{w}$ , and causal-attention is utilised between text tokens [39, 66]. Differing from [9, 32, 81], we show that report generation tasks can greatly benefit from temporal joint latent space.

**Conditioning on prior reports** In contrast to existing work, we incorporate the prior report as a prompt to contextualise the report generation task:  $p_{\theta}(\mathbf{w}_{\text{txt}}^{\text{curr}} | \mathbf{w}_{\text{txt}}^{\text{prior}}, \mathbf{v}_{w,h}^{\text{proj}})$ ,

<sup>3</sup><https://huggingface.co/microsoft/BiomedVLP-CXR-BERT-general>



where  $\Phi$  are the multi-modal encoder-decoder network’s weights, and  $\mathbf{w}_{\text{txt}}^{\text{curr}}$ ,  $\mathbf{w}_{\text{txt}}^{\text{prior}}$  denote text tokens for current and prior reports respectively. This is analogous to fine-tuning GPT-3 [11] with prompts and instructions [71], but conditioning on both images and the previous report. A dedicated separation token [SEP] is added into the input sequence  $[\mathbf{w}_{\text{txt}}^{\text{prior}}, [\text{SEP}], \mathbf{w}_{\text{txt}}^{\text{curr}}]$ .

**Curation of imaging datasets** CXR datasets [36] often contain multiple image acquisitions  $\mathcal{Z} = \{\mathbf{x}_1^{\text{img}}, \dots, \mathbf{x}_Z^{\text{img}}\}$  in a single visit due to data quality issues such as a limited field-of-view or scanning the wrong body part (Figure A.4). Unlike [9, 32, 81], we conduct curation to choose higher quality images among the potential candidates instead of performing a random selection. For this step, a separate BioViL-T is trained on ‘clean’ studies with single acquisitions and later used in a zero-shot setting to detect out-of-distribution samples [26, 27] arising from the re-imaging process. The candidate  $\hat{z}$  is selected as follows:  $\hat{z} = \arg \max_{z \in \mathcal{Z}} S_C(\bar{\mathbf{v}}_z^{\text{proj}}, \mathbf{r}^{\text{proj}})$  s.t.  $|s_{\hat{z}} - s_{\mathcal{Z} \setminus \hat{z}}| > \delta$  for a margin  $\delta$ . This approach is applied to enhance the quality of the temporal classification dataset given its limited size.

## 4. Datasets & experiments

Here, we demonstrate BioViL-T’s data efficiency and adaptability to a wide range of applications, and show how the model achieves SOTA performance on various downstream tasks by learning from data instances linked across time, making effective use of domain priors and the available training data. Specifically, our model is evaluated on a diverse set of downstream tasks including zero- and few-shot static and temporal image classification, report generation, phrase-grounding [10], and sentence similarity.

**MS-CXR-T benchmark** We release a new multi-modal benchmark dataset<sup>4</sup>, *MS-CXR-T*, to evaluate chest X-ray VLP models on two distinct temporal tasks: image classification and sentence similarity. The former comprises multi-image and ground-truth label pairs ( $N = 1326$ ) across 5 findings, with classes corresponding to 3 states of disease progression for each finding: {Improving, Stable, Worsening}. The latter quantifies the temporal-semantic similarity of text embeddings extracted from pairs of sentences ( $N = 361$ ). The pairs can be either paraphrases or contradictions in terms of disease progression. The data for both tasks was manually annotated and reviewed by a board certified radiologist. Appendix C provides further details on its data distribution and annotation protocol.

**Datasets** For pre-training, we use the MIMIC-CXR v2 [28, 36] chest X-ray dataset, which contains longitudinal imaging studies with corresponding radiological reports,

<sup>4</sup>*MS-CXR-T* benchmark dataset can be accessed through PhysioNet: <https://aka.ms/ms-cxr-t>

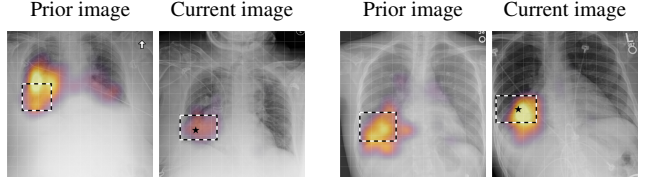


Figure 3. Attention rollout maps [1] from the reference patch (marked with  $\star$ ) to the current and prior images. The bounding boxes, annotated by a radiologist, show the extent of consolidation. Note that the reference patch attends to its anatomical neighbourhood in the prior image despite the misalignment between prior and current images. The grid ( $14 \times 14$ ) represents the patch tokens processed in the transformer encoder blocks.

see Fig. B.1 for the distribution of studies. We only use frontal view scans and discard samples where reports do not contain an IMPRESSION section. From this data, we gather 174.1k and 4.9k text-image pairs for training and validation respectively, with a majority of pairs including a prior image:  $|\mathcal{D}_m^{\text{train}}| = 118.8k$ ,  $|\mathcal{D}_s^{\text{train}}| = 55.3k$ . The text consists of the IMPRESSION section and, for MLM additionally the FINDINGS section if available. Note that *no manual labels* are used during pre-training and *no additional data* is used for the methods that leverage the link between current and prior images. For early stopping we track the validation loss, see Appendix E for implementation details.

Downstream evaluations are performed on a disjoint held-out test set shared across all tasks,  $|\mathcal{D}^{\text{test}}| = 2971$ . For report generation, we extend this test set with samples from healthy subjects ( $N = 815$ ) to match the prevalence of pathological studies used in prior work [14, 25, 45]. For fine-tuning on temporal image classification, we use labels from the Chest Imagenome dataset [72] as in [37] (statistics in Table F.2). In detail, we use the following benchmark datasets: (I) *MS-CXR* [10] for phrase grounding, (II) the RSNA Pneumonia dataset [60, 70] to test zero-shot and fine-tuned classification, (III) *MS-CXR-T* for temporal sentence similarity and temporal image classification.

**Comparison approaches** We compare our approach to other domain-specific SOTA pre-training frameworks [9, 32] specifically on phrase-grounding and zero-shot predictive performance. The non-temporal BioViL framework [9] is most similar to our approach and provides insight into non-temporal pre-training. We additionally compare to internal ablations such as removing the past report during report generation and masking prior images during phrase grounding. For SOTA performance comparison, various AR and nearest-neighbour (NN) based language decoding approaches are used as baselines: IFCC [45], R2Gen [14], CXR-RePaiR-2 [25], and CXR-RePaiR-Select [25].

For the temporal classification task, we compare against a baseline exploiting the BioViL image encoder [9], and an

Table 1. Results for report generation task: Predictions are evaluated in terms of lexical (BLEU-4, ROUGE) and factuality metrics (CHEXBERT, TEM). Approaches are grouped into two broad categories: nearest-neighbour (NN) and auto-regressive (AR). BioViL-T pre-training consistently yields improved decoding. Further, the consistent performance gains of using prior image and report demonstrate the importance of such domain priors. ‘PI / PR’ indicate usage of prior image and report, respectively.

	Method	Pre-training	PI / PR	BLEU-4	ROUGE	CHEXBERT	TEM
Z&F	CXR-RePaiR-2 [25]	BioViL	✗ / ✗	2.1	14.3	28.1	12.5
	Baseline (NN) [9]	BioViL	✗ / ✗	3.7	20.0	28.3	11.1
	Proposed (NN)	BioViL-T	✓ / ✗	4.5	20.5	29.0	13.0
AR	Baseline (AR) [9]	BioViL	✗ / ✗	7.5 ± 0.1	27.9 ± 0.1	29.3 ± 0.3	13.8 ± 0.1
	Proposed	BioViL-T	✓ / ✗	8.2 ± 0.1	28.7 ± 0.1	30.2 ± 0.7	16.0 ± 0.3
	Proposed	BioViL-T	✓ / ✓	<b>9.2 ± 0.3</b>	<b>29.6 ± 0.1</b>	<b>31.7 ± 1.0</b>	<b>17.5 ± 0.1</b>

Table 2. Temporal image classification results (repeated for 4 random seeds) on the *MS-CXR-T* benchmark for fully-supervised and zero-/few-shot (Z&F) learning settings, in terms of macro-accuracy across the three classes for each finding. Affine registration is performed for the baseline method (denoted with suffix ‘w/reg’), to partially address the pose variations across scans.

	Method (% of labels)	Pre-train	Consolidation	PI. effusion	Pneumonia	Pneumothorax	Edema
Z&F	BioViL-T prompt (0%)	Temporal	53.6 ± 1.9	59.7 ± 2.1	58.0 ± 3.9	34.9 ± 1.0	64.2 ± 1.5
	BioViL-T (10%)	Temporal	59.7 ± 2.4	62.4 ± 1.4	60.1 ± 2.1	35.3 ± 2.6	62.6 ± 1.7
Supervised	CNN + Transformer	ImageNet	44.0 ± 2.0	61.3 ± 1.6	45.1 ± 3.5	31.5 ± 3.1	65.5 ± 1.1
	CheXRelNet [37]	ImageNet	47	47	47	36	49
	BioViL [9]	Static	56.1 ± 1.5	62.3 ± 1.1	59.4 ± 1.0	41.7 ± 2.8	67.5 ± 0.8
	BioViL w/reg [9]	Static	56.0 ± 1.5	63.0 ± 0.9	60.2 ± 0.7	42.5 ± 2.7	67.5 ± 0.9
	BioViL-T w/out curation	Temporal	58.9 ± 1.7	65.5 ± 0.7	61.5 ± 2.2	44.4 ± 2.1	67.4 ± 0.8
	BioViL-T	Temporal	<b>61.1 ± 2.4</b>	<b>67.0 ± 0.8</b>	<b>61.9 ± 1.9</b>	42.6 ± 1.6	<b>68.5 ± 0.8</b>

approach that makes use of graph convolutions across regions of interest extracted from bounding boxes [37]. For BioViL, we perform affine image registration (with 4 DoF) for each pair of scans to cope with pose variations, and the encoded images are concatenated along the feature dimension and classified via a multilayer perceptron. For [37], we compare to the three-class setting. Lastly, we benchmark our final text model in isolation against domain specific SOTA models in a temporal sentence similarity task: CXR-BERT [9] and PubMedBert [29].

**Metrics** Due to class imbalance, we report macro-accuracy for temporal image classification. For phrase grounding, we use mean Intersection-Over-Union (mIoU) and Contrast-to-Noise-Ratio (CNR) [9]. The latter measures the discrepancies between cosine similarities inside and out of the bounding box region without requiring hard thresholds. To evaluate the quality of generated reports, we use both the standard lexical metrics, e.g., BLEU [50], ROUGE-L [41], and also domain-specific factuality metric: CheXbert<sup>5</sup> [62]. To directly probe the generation of change-related information, we introduce a new metric called temporal entity matching (TEM) to compute the match score of a fixed set of temporal entities (see Appendix D).

<sup>5</sup>The average of the weighted- $F_1$  score across 14 pathological observations labelled by CheXbert.

Table 3. Report generation results using the same train/test splits from [25], measured by lexical (BLEU-2) and factuality (CHEXBERT) metrics. Baseline results were also collected from [25]. Note the CHEXBERT score covers all 14 observations.

Method	Decoded sections	BLEU-2	CHEXBERT
R2gen [14]	Findings & Impression	21.20 ± 0.10	14.80 ± 0.30
IFCC [45]	Findings	21.70 ± 0.10	27.00 ± 0.40
CXR-RePaiR-Sel [25]	Impression	5.00 ± 0.10	27.40 ± 0.30
BioViL-T	Impression	15.86 ± 0.14	34.83 ± 0.73
BioViL-T	Findings & Impression	21.31 ± 0.19	<b>35.86 ± 0.35</b>

## 4.1. Temporal pre-training yields data efficiency

*Downstream tasks are enabled with minimal labels.*

The sections ‘NN’ and ‘Z&F’ on Tables 1 and 2 report zero- and few-shot performance on tasks benefitting from temporal information: temporal image classification and report generation. Here we measure the quality of the learnt joint latent space and the extent to which BioViL-T enables efficient use of raw data. For zero-shot classification we prompt the AR fine-tuned model with prefix: “[FINDING] is” and compare the next-token probability of words meaning ‘improving’, ‘stable’, and ‘worsening’ (Appendix F.4).

Without using any labelled data, Table 2 shows that the proposed AR-based approach already yields performance superior to prior fully-supervised work [37] on temporal image classification. With only 10% of labels, classification fine-tuning provides a further boost, indicating that BioViL-T produces a multi-image encoder readily adapted to temporal tasks. Similarly, in a zero-shot report-retrieval setting, the findings show that compared to temporally-agnostic pre-training, BioViL-T leveraging prior images improves across all metrics. Consistent with prior work [25], the retrieved reports already preserve factuality with high CheXbert scores, more-so than the other metrics which measure fine-grained specifics of phrasing. This demonstrates that the latent space captures the high-level semantics of the clinical features. Fine-grained phrasing however will be substantially improved by AR fine-tuning.

## 4.2. Achieving SOTA performance with BioViL-T

*A wide range of downstream tasks benefit substantially from temporally-aware pre-training.*

Through downstream adaptations and fine-tuning our model, we report SOTA performance on report generation and temporal image classification tasks. For the former, using both prior images *and* reports during fine-tuning substantially improves across metrics (Table 1). In particular, TEM metric results show that temporal context is key for accurately describing change in the generated report while avoiding hallucinations (see Table A.1 for examples). Comparing to published results on a comparable test split and

Table 4. Image classification results on RSNA Pneumonia Detection Benchmark [60] for train and test splits of 70% – 30% respectively.

Method	% of Labels	Supervision	Acc.	F1	AUROC
GLoRIA [32]	✗	Zero-shot	0.70	0.58	-
BioViL [9]	✗	Zero-shot	0.732	0.665	0.831
BioViL-T	✗	Zero-shot	<b>0.805</b>	<b>0.706</b>	<b>0.871</b>
BioViL [9]	1%	Few-shot	0.805	0.723	0.881
BioViL-T	1%	Few-shot	<b>0.814</b>	<b>0.730</b>	<b>0.890</b>

metrics (Sec. 4.1), we conclude that BioViL-T with fine-tuning achieves SOTA on report generation, producing reports that are lexically on par with prior work but substantially more factually accurate. Note that we do ‘vanilla’ AR fine-tuning to focus on the impact of the pre-trained encoders, so application-specific supervision [45] could be used in conjunction to further boost performance.

In temporal image classification (Tab. 2), BioViL-T pre-training outperforms the non-temporal baseline (BioViL) and improves on previously-reported results [37] by up to 20 percentage points (pp). Furthermore, baseline methods that rely on image registration (BioViL w/reg), underperform compared to the proposed approach. Further analysis reveals that errors tend to be in cases with disagreement between radiologists (Appendix A.2). We also note that pre-training is critical for a hybrid CNN-transformer model on this task, likely due to the small labelled dataset. Lastly, curation of temporal training data is observed to improve the classification results by .68 pp aggregated across the findings, see Appendix A.4 for details.

### 4.3. Static tasks benefit from temporal learning

*BioViL-T broadens the range of applicable downstream tasks whilst contributing to performance on static tasks.*

In this section, we demonstrate that performance improvements afforded by BioViL-T are not restricted to temporal tasks – *static* tasks also benefit. Table 4 reports results on zero- and few-shot pneumonia classification from single images [60], where BioViL-T establishes a new SOTA compared to prior work [9, 32].

We see a similar trend on the *MS-CXR* phrase grounding benchmark (Tab. 5). This task can be solved with single images, however we show that the inclusion of the prior image (where available) does not impair the performance of BioViL-T. Feature decomposition effectively preserves localised information from the current image.

### 4.4. Towards better sentence embedding quality

*Language models acquire increased temporal sensitivity.*

We hypothesise that text encoders learn temporal semantics through supervision from longitudinal image series. To verify this, RadNLI [45] and *MS-CXR-T* datasets are used in a zero-shot binary classification setting. Cosine similarity

Table 5. Results on *MS-CXR* benchmark [10] (5-runs with different seeds), “Multi-image” column indicates the input images used at test time.

Method	Multi-Image	Avg. CNR	Avg. mIoU
BioViL [9]	✗	1.07 ± 0.04	0.229 ± 0.005
+ Local loss [9, 32]	✗	1.21 ± 0.05	0.202 ± 0.010
BioViL-T	✗	<b>1.33 ± 0.04</b>	<b>0.243 ± 0.005</b>
BioViL-T	✓	<b>1.32 ± 0.04</b>	<b>0.240 ± 0.005</b>

Table 6. Results on *MS-CXR-T* sentence similarity benchmark.

Text Model	MS-CXR-T (361 pairs)		RadNLI (145 pairs)	
	Accuracy	ROC-AUC	Accuracy	ROC-AUC
PubMedBERT [29]	60.39	.542	81.38	.727
CXR-BERT-G [9]	62.60	.601	87.59	.902
CXR-BERT-S [9]	78.12	.837	89.66	.932
BioViL-T	<b>87.77 ± 0.5</b>	<b>.933 ± .003</b>	90.52 ± 1.0	<b>.947 ± .003</b>

of sentence pair embeddings [56] are treated as class-logits to label each pair either as paraphrase or contradiction. See Appendix F.6 for further details.

Our text model is benchmarked against SOTA domain-specific BERT models. Table 6 shows that the proposed framework greatly increases the sensitivity of sentence embeddings to temporal content whilst better capturing the static content (RadNLI). Note that CXR-BERT-Specialised [9] is learnt through single-images starting from the same canonical model, illustrating the substantial increase in temporal and static sensitivity due to BioViL-T pre-training.

### 4.5. Ablation experiments

In Table 7 we report extensive ablations across the multi-image encoder architecture, pre-training choices, and AR fine-tuning for report generation.

**Image encoder** Table 7 shows that decomposition of static and progression features is essential to ensure good performance on single-image tasks, such as phrase grounding. For temporal representations, on the other hand, positional encodings (T) are essential to disambiguate the order of scans, i.e., permutation variance across time.

**Model pre-training** The corresponding results are shown in the middle section of Table 7. The local contrastive loss proves crucial to ensure meaningful language supervision during pre-training, followed by the image-guided MLM objective. Lastly, use of the FINDINGS section results in only minor performance gains as the key findings are already captured in the IMPRESSION section.

**Report generation** The importance of prior image and report is demonstrated by the substantial drop in the “no prior image and report” ablation, confirming our hypothesis that temporal context is crucial for improving report quality. While both inputs are crucial for optimal performance,

Table 7. Ablation study on image encoder, pre-training settings, and report generation (one component at a time, and repeated for 4 random seeds). Note that for temporal classification, linear probing is applied to frozen image embeddings. In report generation, the baseline method is fine-tuned with both prior image and report.

	Ablation	Avg. CNR (mIoU)	Pl. Effusion Acc.
Encoder	Baseline	1.33 ± 0.02 (.248)	64.8 ± 0.6
	- Temporal pos. encoding	1.32 ± 0.02 (.242)	62.9 ± 1.0
	- Feature decomposition	1.11 ± 0.08 (.203)	64.0 ± 0.6
Pre-training	Baseline	1.33 ± 0.02 (.248)	64.8 ± 0.6
	- Use of findings section	1.32 ± 0.01 (.246)	63.8 ± 0.8
	- MLM loss	1.28 ± 0.02 (.238)	63.2 ± 0.7
	- Local contrastive loss	1.18 ± 0.02 (.236)	60.2 ± 0.6
	Ablation	ROUGE	TEM
Report gen.	Baseline	29.64 ± 0.08	17.54 ± 0.11
	- Prior image	29.35 ± 0.25	16.30 ± 0.40
	- Prior report	28.67 ± 0.12	16.00 ± 0.30
	- (Prior image and report)	27.78 ± 0.09	13.65 ± 0.48
	- Separation token	26.00 ± 0.40	15.50 ± 1.06

the prior report is more so because it summarises the image and provides a clearer signal. The prior image however cannot be dismissed entirely as it provides granular details which may not always be documented in a report. Finally, we found the separation token is crucial in differentiating between the predicted tokens for the current report and tokens from the prior report.

#### 4.6. Which tokens require a prior image in MLM?

We leverage the MLM objective in an inference setting to analyse the influence of prior images in predicting masked tokens. Inspired by the  $\Delta_{img}^{prior}$  loss of [8], we define  $\Delta_{img}^{prior}$  as the change in loss by conditioning the estimation with a prior image for a given token  $w$  as follows:

$$\Delta_{img}^{prior}(w) = l(w, \mathbf{x}_{img}^{curr}, \emptyset) - l(w, \mathbf{x}_{img}^{curr}, \mathbf{x}_{img}^{prior}) \quad (2)$$

where  $l(w, \mathbf{x}_{img}^{curr}, \mathbf{x}_{img}^{prior})$  is the cross-entropy of predicting the masked token  $w$  given visual features (MLM loss for a single token), averaged over sentences in which  $w$  appears.  $\Delta_{img}^{prior}$  is a measure of how much that token benefits from access to the prior image, as well as an assessment of the contribution of the prior image to the image representation. In Figure 4 we show the distribution of  $\Delta_{img}^{prior}$  as a function of token category (e.g., *Anatomy*, *Positional*; see F.5 for annotation details). For *Progression*-type terms in particular, the model heavily relies on the prior image for image-guided MLM. We further observe that this effect is specific to temporal tokens; as expected, those from other semantic categories do not consistently rely on the prior image.

## 5. Conclusion

In this paper, we introduced BioViL-T, a vision-language pre-training framework enabling alignment between text and multiple images. BioViL-T makes use of

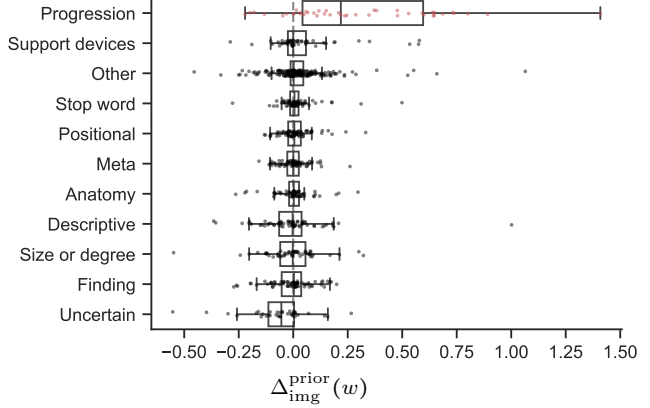


Figure 4. Mean token-level increase in image-guided MLM loss when prior image is discarded, grouped by token category. The prior image is excluded during inference to measure its impact on masked token predictions. *Progression* tokens are significantly better predicted when prior images are incorporated into image embeddings. The top five *Progression* tokens are ‘persist’, ‘improving’, ‘remains’, ‘unchanged’, and ‘residual’.

a novel multi-image encoder and explicitly decomposes static-temporal features to augment the current image representation with information from prior images. This enables the grounding of temporal references in the text. To our knowledge, this is the first method capable of leveraging the temporal content commonly present in biomedical text. It addresses an important limitation in existing VLP approaches, which simply discard such context. Also, incorporating such multi-modal temporal content provides strong learning signals to the model, resulting in richer representations and improved downstream performance.

We demonstrate the value of this paradigm through extensive experiments: BioViL-T excels on both static and temporal tasks, establishing new SOTA on report generation, temporal image classification, few/zero-shot pneumonia detection, and phrase grounding. Furthermore, we release a new multi-modal benchmark (*MS-CXR-T*) to measure the quality of image and text representations in terms of temporal semantics, enabling more diverse evaluation of biomedical VLP models. The corresponding model weights<sup>6</sup> and code<sup>7</sup> are publicly available.

Further exploration and evaluation are required on diverse datasets to characterise what kinds of tasks would benefit from a temporal modelling approach, and specifically from the proposed methodology.

**Acknowledgements:** We would like to thank Hannah Richardson, Hoifung Poon, Melanie Bernhardt, Melissa Bristow and Naoto Usuyama for their valuable feedback.

<sup>6</sup>Models can be found at: <https://aka.ms/biovil-t-model>

<sup>7</sup>Code can be found at: <https://aka.ms/biovil-t-code>



## References

- [1] Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, Online, July 2020. Association for Computational Linguistics. 3, 5, 13, 14
- [2] Monica N Agrawal, Hunter Lang, Michael Offin, Lior Gazit, and David Sontag. Leveraging time irreversibility with order-contrastive pre-training. In *International Conference on Artificial Intelligence and Statistics*, pages 2330–2353. PMLR, 2022. 3
- [3] Uwa O. Aideyan, Kevin Berbaum, and Wilbur L. Smith. Influence of prior radiologic information on the interpretation of radiographic examinations. *Academic Radiology*, 2(3):205–208, 1995. 1
- [4] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022. 2, 3
- [5] B.B. Avants, C.L. Epstein, M. Grossman, and J.C. Gee. Symmetric diffeomorphic image registration with cross-correlation: Evaluating automated labeling of elderly and neurodegenerative brain. *Medical Image Analysis*, 12(1):26–41, 2008. Special Issue on The Third International Workshop on Biomedical Image Registration – WBIR 2006. 2
- [6] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 3
- [7] Nadine Behrmann, Mohsen Fayyaz, Juergen Gall, and Mehdi Noroozi. Long short view feature decomposition via contrastive video representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9244–9253, 2021. 3
- [8] Yonatan Bitton, Gabriel Stanovsky, Michael Elhadad, and Roy Schwartz. Data efficient masked language modeling for vision and language. In *2021 Findings of the Association for Computational Linguistics, Findings of ACL: EMNLP 2021*, pages 3013–3028. Association for Computational Linguistics (ACL), 2021. 8
- [9] Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, Hoifung Poon, and Ozan Oktay. Making the most of text semantics to improve biomedical vision-language processing. In *European Conference on Computer Vision (ECCV)*, pages 1–21. Springer, 2022. 1, 2, 3, 4, 5, 6, 7, 17, 18, 19, 20
- [10] Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, Hoifung Poon, and Ozan Oktay. MS-CXR: Making the most of text semantics to improve biomedical vision-language processing (version 0.1). PhysioNet, 2022. 2, 3, 5, 7, 15
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 5, 19
- [12] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 3
- [13] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: UNiversal Image-TEXT Representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020. 2, 4
- [14] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, Nov. 2020. 5, 6
- [15] Ching-Yao Chuang, R Devon Hjelm, Xin Wang, Vibhav Vineet, Neel Joshi, Antonio Torralba, Stefanie Jegelka, and Yale Song. Robust contrastive learning against noisy views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16670–16681, 2022. 1
- [16] Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, Antonio Torralba, and Stefanie Jegelka. Debaised contrastive learning. *Advances in neural information processing systems*, 33:8765–8775, 2020. 1
- [17] R. C. Daudt, B. L. Saux, and Alexandre Boulch. Fully convolutional siamese networks for change detection. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 4063–4067, 2018. 2
- [18] Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310, 2016. 15
- [19] Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11162–11173, 2021. 1, 2
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- [21] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 3, 13
- [22] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuhang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, et al. An empirical study of training end-to-end vision-and-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18166–18176, 2022. 2, 4
- [23] Stanley Durrleman, Xavier Pennec, Alain Trounev, José Braga, Guido Gerig, and Nicholas Ayache. Toward a com-

- prehensive framework for the spatiotemporal statistical analysis of longitudinal shape data. *International Journal of Computer Vision*, 103:22–59, 2013. [2](#)
- [24] Stéphane d’Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *International Conference on Machine Learning*, pages 2286–2296. PMLR, 2021. [3](#), [4](#)
- [25] Mark Endo, Rayan Krishnan, Viswesh Krishna, Andrew Y Ng, and Pranav Rajpurkar. Retrieval-based chest x-ray report generation using a pre-trained contrastive language-image model. In *Machine Learning for Health*, pages 209–219. PMLR, 2021. [5](#), [6](#), [18](#)
- [26] Sepideh Esmailpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot out-of-distribution detection based on the pretrained model clip. In *Proceedings of the AAAI conference on artificial intelligence*, 2022. [5](#)
- [27] Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection. *Advances in Neural Information Processing Systems*, 34:7068–7081, 2021. [5](#)
- [28] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *Circulation*, 101(23):e215–e220, 2000. [5](#)
- [29] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pre-training for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021. [6](#), [7](#), [20](#)
- [30] Tengda Han, Weidi Xie, and Andrew Zisserman. Self-supervised co-training for video representation learning. *Advances in Neural Information Processing Systems*, 33:5679–5690, 2020. [3](#)
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. [3](#)
- [32] Shih-Cheng Huang, Liyue Shen, Matthew P Lungren, and Serena Yeung. GLoRIA: A multimodal global-local representation learning framework for label-efficient medical image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3942–3951, 2021. [1](#), [2](#), [4](#), [5](#), [7](#)
- [33] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, et al. CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. *Proceedings of the AAAI conference on artificial intelligence*, 33(01):590–597, 2019. [2](#)
- [34] Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Truong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew P Lungren, Andrew Y Ng, Curtis Langlotz, et al. Radgraph: Extracting clinical entities and relations from radiology reports. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. [17](#)
- [35] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021. [1](#), [13](#)
- [36] A. Johnson, T. Pollard, S.J. Berkowitz, R. Mark, and S. Horng. MIMIC-CXR database (version 2.0.0). PhysioNet, 2019. [2](#), [5](#)
- [37] Gaurang Karwande, Amarachi B Mbakwe, Joy T Wu, Leo A Celi, Mehdi Moradi, and Ismeni Lourentzou. Chexrelnet: An anatomy-aware model for tracking longitudinal relationships between chest x-rays. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022*, pages 581–591, 2022. [2](#), [3](#), [5](#), [6](#), [7](#)
- [38] Jiann-Shu Lee, Jing-Wein Wang, Hsing-Hsien Wu, and Ming-Zheng Yuan. A nonparametric-based rib suppression method for chest radiographs. *Computers & Mathematics with Applications*, 64(5):1390–1399, 2012. [20](#)
- [39] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–12900. PMLR, 17–23 Jul 2022. [2](#), [4](#)
- [40] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021. [1](#), [2](#)
- [41] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. [6](#)
- [42] Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. In *International Conference on Learning Representations*, 2018. [4](#)
- [43] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. [17](#), [19](#)
- [44] Bradley C Lowekamp, David T Chen, Luis Ibáñez, and Daniel Blezek. The design of simpleitk. *Frontiers in neuroinformatics*, 7:45, 2013. [20](#)
- [45] Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. Improving factual completeness and consistency of image-to-text radiology report generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5288–5304, Online, June 2021. Association for Computational Linguistics. [5](#), [6](#), [7](#), [18](#)

- [46] Jong Hak Moon, Hyungyung Lee, Woncheol Shin, Young-Hak Kim, and Edward Choi. Multi-modal understanding and generation for medical images and text via vision-language pre-training. *IEEE Journal of Biomedical and Health Informatics*, 2022. 2, 4
- [47] American College of Radiology (ACR). ACR practice guideline for communication of diagnostic imaging findings. *Practice guidelines & technical standards*, 2020. 1
- [48] Dong Yul Oh, Jihang Kim, and Kyong Joon Lee. Longitudinal change detection on chest X-rays using geometric correlation maps. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 748–756. Springer, 2019. 2
- [49] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 1, 4
- [50] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. 6
- [51] Namuk Park and Songkuk Kim. How do vision transformers work? In *International Conference on Learning Representations*, 2022. 3, 4
- [52] Daifeng Peng, Yongjun Zhang, and Haiyan Guan. End-to-end change detection for high resolution satellite images using improved UNet++. *Remote Sensing*, 11(11):1382, 2019. 2
- [53] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 1, 2, 4
- [54] Vignav Ramesh, Nathan Andrew Chi, and Pranav Rajpurkar. Improving radiology report generation systems by removing hallucinated references to non-existent priors. *arXiv preprint arXiv:2210.06340*, 2022. 2
- [55] Adria Recasens, Pauline Luc, Jean-Baptiste Alayrac, Luyu Wang, Florian Strub, Corentin Tallec, Mateusz Malinowski, Viorica Pătrăucean, Florent Althé, Michal Valko, et al. Broaden your views for self-supervised video learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1255–1265, 2021. 3
- [56] Nils Reimers, Iryna Gurevych, Nils Reimers, Iryna Gurevych, Nandan Thakur, Nils Reimers, Johannes Daxenberger, Iryna Gurevych, Nils Reimers, Iryna Gurevych, et al. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, pages 671–688. Association for Computational Linguistics, 2019. 7, 20
- [57] Liqa A Rousan, Eyhab Elobeid, Musaab Karrar, and Yousef Khader. Chest x-ray findings and temporal lung changes in patients with covid-19 pneumonia. *BMC Pulmonary Medicine*, 20(1):1–9, 2020. 1
- [58] Ruggiero Santeramo, Samuel Joseph Withey, and G. Montana. Longitudinal detection of radiological abnormalities with time-modulated LSTM. In *MICCAI 2018 Workshop on Deep Learning in Medical Imaging Analysis*, 2018. 2
- [59] Wenzhong Shi, Min Zhang, Rui Zhang, Shanxiong Chen, and Zhao Zhan. Change detection based on artificial intelligence: State-of-the-art and challenges. *Remote Sensing*, 12(10):1688, 2020. 2
- [60] George Shih, Carol C Wu, Safwan S Halabi, Marc D Kohli, Luciano M Prevedello, Tessa S Cook, Arjun Sharma, Judith K Amorosa, Veronica Arteaga, Maya Galperin-Aizenberg, et al. Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. *Radiology: Artificial Intelligence*, 1(1):e180041, 2019. 2, 5, 7
- [61] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15617–15629, 2022. 2
- [62] Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Ng, and Matthew Lungren. Combining automatic labelers and expert annotations for accurate radiology report labeling using BERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1500–1519, Online, Nov. 2020. Association for Computational Linguistics. 6
- [63] Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y Ng, and Matthew P Lungren. CheXbert: combining automatic labelers and expert annotations for accurate radiology report labeling using bert. *arXiv preprint arXiv:2004.09167*, 2020. 2, 16
- [64] Anuroop Sriram, Matthew Muckley, Koustuv Sinha, F. Shamout, Joelle Pineau, K. Geras, L. Azour, Y. Aphinyanaphongs, N. Yakubova, and William H. Moore. COVID-19 prognosis via self-supervised representation learning and multi-image prediction. *arXiv preprint arXiv:2101.04909*, 2021. 2, 3
- [65] Philippe Thévenaz and Michael Unser. Optimization of mutual information for multiresolution image registration. *IEEE transactions on image processing*, 9(12):2083–2099, 2000. 20
- [66] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3, 4
- [67] Yen Nhi Truong Vu, Richard Wang, Niranjana Balachandrar, Can Liu, Andrew Y Ng, and Pranav Rajpurkar. Medaug: Contrastive learning leveraging patient metadata improves representations for chest x-ray interpretation. In *Machine Learning for Healthcare Conference*, pages 755–769. PMLR, 2021. 2
- [68] Chuang Wang, Andreas Rimmer, Yu chi Hu, Neelam Tyagi, Jue Jiang, Ellen Yorke, Sadegh Riyahi, Gig S. Mageras, Joseph O. Deasy, and Pengpeng Zhang. Towards predicting the evolution of lung tumors during radiotherapy observed

- on a longitudinal MR imaging study via a deep learning algorithm. *Medical Physics*, 2019. 2
- [69] Wenhui Wang, Hangbo Bao, Li Dong, and Furu Wei. Vlm0: Unified vision-language pre-training with mixture-of-modality-experts. *arXiv preprint arXiv:2111.02358*, 2021. 2
- [70] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. ChestX-Ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2097–2106. IEEE Computer Society, 2017. 2, 5
- [71] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2021. 5
- [72] Joy Wu, Nkechinyere Agu, Ismini Lourentzou, Arjun Sharma, Joseph Paguio, Jasper Seth Yao, Edward Christopher Dee, William Mitchell, Satyananda Kashyap, Andrea Giovannini, Leo Anthony Celi, Tanveer Syeda-Mahmood, and Mehdi Moradi. Chest imagenome dataset (version 1.0.0). PhysioNet, 2021. 5, 15, 19
- [73] Kan Wu, Houwen Peng, Minghao Chen, Jianlong Fu, and Hongyang Chao. Rethinking and improving relative position encoding for vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10033–10041, 2021. 19
- [74] Yiwu Xu, Ahmed Hosny, Roman Zeleznik, Chintan Parmar, Thibaud P. Coroller, Idalid Ivy Franco, Raymond H. Mak, and Hugo J.W.L. Aerts. Deep learning predicts lung cancer treatment response from serial medical imaging. *Clinical Cancer Research*, 25:3266 – 3275, 2019. 2
- [75] Bin Yan and Mingtao Pei. Clinical-BERT: Vision-language pre-training for radiograph diagnosis and reports generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(3):2982–2990, 2022. 2
- [76] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *arXiv preprint arXiv:2111.07783*, 2021. 2, 4
- [77] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. 1, 2
- [78] Sukmin Yun, Jaehyung Kim, Dongyoon Han, Hwanjun Song, Jung-Woo Ha, and Jinwoo Shin. Time is matter: Temporal self-supervision for video transformers. In *International Conference on Machine Learning*, pages 25804–25816. PMLR, 2022. 3
- [79] Dwen Zeng, John N Kheir, Peng Zeng, and Yiyu Shi. Contrastive learning with temporal correlated medical images: A case study using lung segmentation in chest x-rays. In *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*, pages 1–7. IEEE, 2021. 2
- [80] Zhaoyang Zeng, Daniel McDuff, Yale Song, et al. Contrastive learning of global and local video representations. *Advances in Neural Information Processing Systems*, 34:7025–7040, 2021. 3
- [81] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020. 1, 2, 4, 5
- [82] Yuhao Zhang, Yuhui Zhang, Peng Qi, Christopher D Manning, and Curtis P Langlotz. Biomedical and clinical english model packages for the stanza python nlp library. *Journal of the American Medical Informatics Association*, 28(9):1892–1899, 2021. 16, 17, 18



## A. Additional Results and Analyses

### A.1. Qualitative analysis of generated reports

Table A.1 shows example reports generated with BioViL-T and BioViL models, which are compared to the reference radiologist’s reports. In comparison with BioViL which only models the current image, BioViL-T shows the benefit from incorporating prior study information and is able to provide factually more accurate reports especially in terms of describing temporal progression of the findings. This is showcased in the first two examples in the table: In the first row, BioViL-T is able to comment on not only the presence of the pleural effusion but also its improvement while BioViL fails to mention the change. In the second example, BioViL-T is able to correctly identify that there is no relevant change by comparing with the previous study, while BioViL wrongly hallucinates the tube in the current image as a new placement. BioViL-T can also avoid hallucination of the temporal information when there is no prior study. For instance, in the third example, BioViL-T correctly acknowledges that there is no prior image and generates the report based on information from the single current image, while BioViL hallucinates a non-existent prior study and wrongly generates temporal descriptions in the report.

### A.2. Further analysis on temporal classification

A subset of the *MS-CXR-T* benchmark dataset is re-annotated by an expert radiologist by blinding them to the existing ground-truth labels and displaying only pairs of images obtained from each subject. With the new set of labels, the analysis focuses on measuring the correlation between inter-rater agreement and image model’s prediction errors. Figure A.1 shows the dependency between the two where the x-axis corresponds to the cross entropy loss between the *MS-CXR-T* benchmark labels and model predictions. We observe lower model performance in cases with smaller inter-rater reliability for the three classes in the dataset, indicating that the model’s prediction errors occur more often for the cases where experts may disagree with each other.

### A.3. Self-attention visualisation

In Figure A.2, we show examples of self-attention rollout [1] maps for pleural effusion and consolidation, including radiologist-annotated bounding boxes surrounding the corresponding pathology in each prior and current image.

To model the attention flow through the transformer encoder block, we first average each attention weight matrix across all heads, subsequently we multiply the matrices between every two layers. For every block we add the identity matrix in order to model the residual connections. Last, we only keep the top 10 % of attention weights per block to reduce noise in the final rollout map. In contrast to [21], we

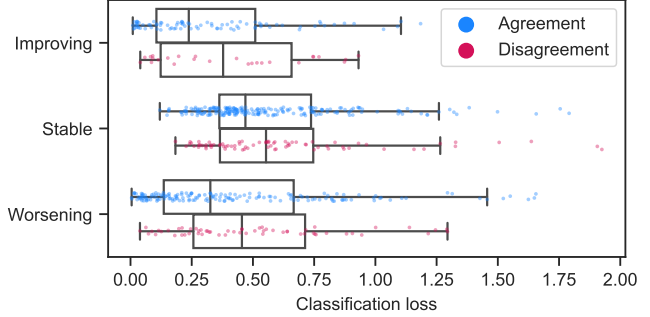


Figure A.1. Cross entropy between model predictions and *MS-CXR-T* temporal classification labels. ‘Disagreement’ indicates cases for which annotations differed amongst radiologists. Model performance is higher for cases with low ambiguity (‘Agreement’).

do not visualize the rollout map with respect to a [CLS] token. Instead, we choose a reference image patch from the center of the radiologist-annotated bounding boxes, marked with  $\star$  in Figure A.2.

We find that the rollout maps in Figure A.2 are in good agreement with radiologist-annotated bounding boxes, i.e., the reference patch attends to other patches within the bounding boxes in the prior and current image. In addition, we find that BioViL-T is robust to pose variations, e.g., in Figure A.2 (a) we show that despite the vertical shift between prior and current image, the reference patch attends to the correct image patches in the prior image.

To further assess the robustness of BioViL-T against pose variations between prior and current images, we performed multiple rotations to the prior image within a pair and computed rollout maps from the same reference patch in the current image. Figure A.3 shows that BioViL-T consistently attends to the corresponding anatomical region independently of the spatial transformation applied, demonstrating that registration is not needed.

### A.4. Data curation of imaging datasets

Large datasets often contain instances that are mislabelled or out of distribution [35]. We used BioViL-T to perform pairwise ranking of instances in MIMIC-CXR (Section 3.3,  $\delta = 0.2$ ) and selected representative examples found in the dataset. Our method is able to select the most appropriate image for a range of different image-acquisition or image-processing issues (Figure A.4).

We found that many lateral acquisitions in the dataset were unexpectedly labelled as frontal (Figure A.4a). Some images contained only noise (Figure A.4b), non-human samples (Figures A.4d and A.4e) or incorrect anatomy (Figure A.4g). Often, acquisitions with an incomplete field of view (FOV) (i.e., the lungs are not completely visible) were repeated (Figure A.4c). Lastly, post-processed images were

Reference (by radiologist)	BioViL	BioViL-T
No evidence of acute cardiopulmonary process. <span style="border: 1px solid blue; padding: 2px;">Decreased</span> right pleural effusion.	Small right pleural effusion.	Small right pleural effusion, <span style="border: 1px solid blue; padding: 2px;">decreased</span> since ...
<span style="border: 1px solid blue; padding: 2px;">No relevant change</span> as compared to the previous image. No evidence of pneumonia or other parenchymal pathology. Normal size of the heart. No pleural effusions.	<span style="border: 1px solid brown; padding: 2px;">In comparison</span> with the study of ..., there has been <span style="border: 1px solid brown; padding: 2px;">placement</span> of a right ij port - a - cath that extends to the lower svc. No evidence of post procedure pneumothorax. The cardiac silhouette is within normal limits and there is no vascular congestion or pleural effusion. There is some asymmetry of opacification at the bases, more prominent on the left. In the appropriate clinical setting, this could possibly represent a developing consolidation.	As compared to the previous radiograph, <span style="border: 1px solid blue; padding: 2px;">no relevant change</span> is seen. The lung volumes are normal. Normal size of the cardiac silhouette. Normal hilar and mediastinal structures. No pneumonia, no pulmonary edema, no pleural effusions.
<span style="border: 1px solid blue; padding: 2px;">No previous images</span> . The cardiac silhouette is within normal limits and there is no evidence of vascular congestion, pleural effusion, or acute focal pneumonia.	<span style="border: 1px solid brown; padding: 2px;">In comparison</span> with the study of ..., there is <span style="border: 1px solid brown; padding: 2px;">little change</span> and no evidence of acute cardiopulmonary disease. No pneumonia, vascular congestion, or pleural effusion.	<span style="border: 1px solid blue; padding: 2px;">No previous images</span> . The cardiac silhouette is within normal limits and there is no vascular congestion, pleural effusion, or acute focal pneumonia.

Table A.1. Comparison between reports generated by radiologists, BioViL using only a single current image and BioViL-T using both the current and previous study. BioViL-T with access to longitudinal information can generate more accurate reports with more precise details on the progression of findings (as in the first and second example) while avoiding hallucination (in the third example). Blue box highlights the correct temporal information and brown box highlights incorrect temporal information including hallucination.

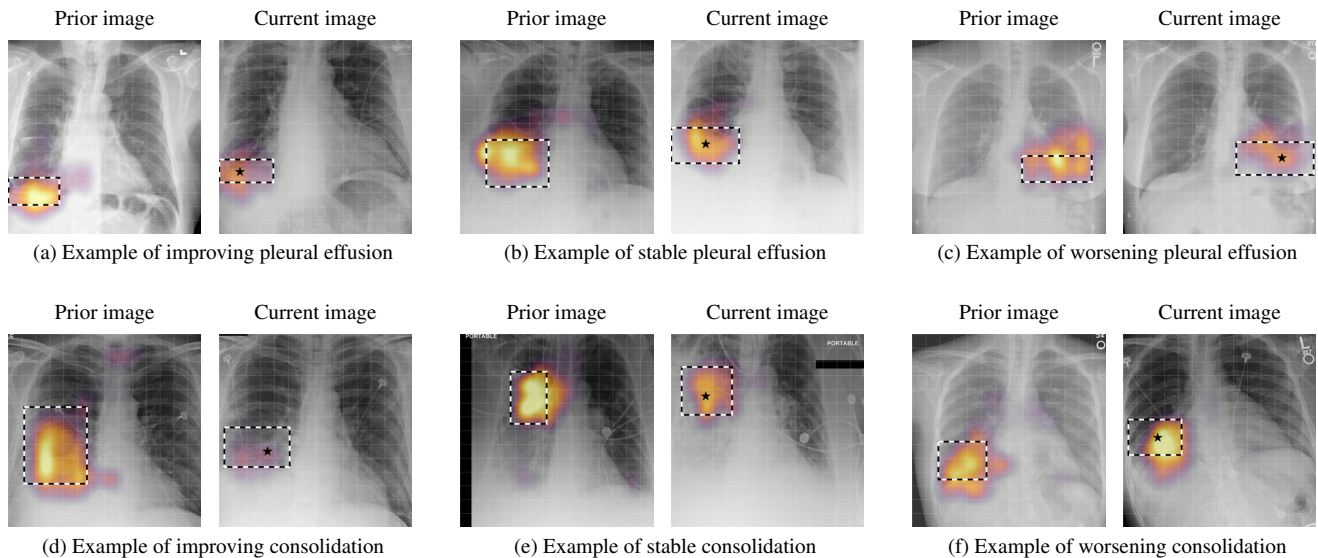


Figure A.2. Self-attention rollout maps [1] from the reference patch (marked with  $\star$ ) to the current and prior images, overlaid on example cases of (a) improving, (b) stable and (c) worsening pleural effusion (top row) and consolidation (bottom row). The bounding boxes, annotated by a radiologist, show the area corresponding to the pathology. The centre patch in the bounding box for the current image was selected as reference. The grid ( $14 \times 14$ ) represents the visual tokens processed in the transformer encoder blocks.

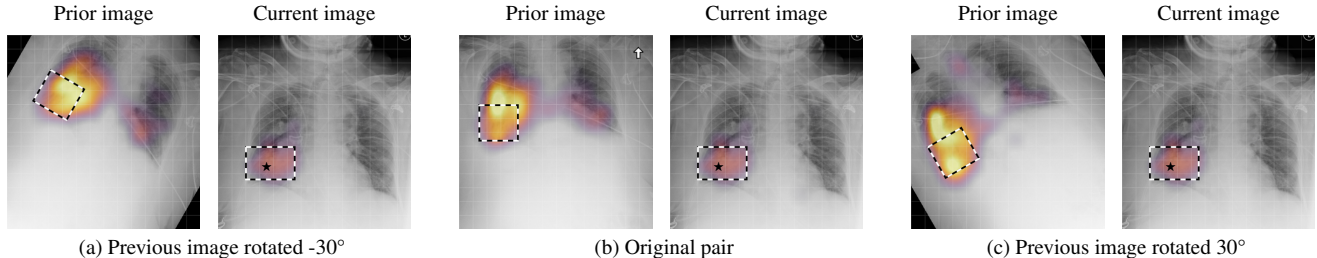


Figure A.3. Comparison of roll-out maps computed after applying in-plane spatial rotations to the prior image. The reference visual token (★) attends to the corresponding anatomical region annotated by an expert independent of the underlying spatial transformation.

detected by the algorithm such as contrast-enhanced scans (Figure A.4i) that are not often used for diagnostic purposes in clinical practice.

### A.5. Phrase-grounding on external data

We have additionally conducted a robustness analysis on an out-of-distribution dataset. For this purpose, a small set of expert labels (N=137 bounding-box-caption pairs) were collected on Open-Indiana CXR dataset [18] for phrase grounding on the same set of abnormalities as *MS-CXR* benchmark [10]. The dataset differs in terms of text token distribution, demographics, and disease prevalence. The experiment was performed with the same methods and setup described in Section 4.3. The results show that the performance gains due to temporal pre-training is observed to be consistent on external datasets.

Table A.2. Multi-modal phrase-grounding results obtained on a subset of Open-Indiana CXR dataset [18] image-text pairs. “Multi-image” column indicates the input images used at test time. The results are reported in terms of micro-averages owing to the limited number of samples in some classes.

Method	Pre-Train	Multi-Image	Avg. CNR	Avg. mIoU
BioViL [9]	Static	✗	$1.19 \pm 0.04$	$0.259 \pm 0.003$
BioViL-T	Temporal	✗	<b><math>1.53 \pm 0.05</math></b>	<b><math>0.289 \pm 0.006</math></b>

## B. Temporal aspects of the MIMIC-CXR v.2 dataset

Subjects in the MIMIC-CXR dataset often have multiple associated studies that happened at different times. A study, sometimes referred to as an ‘exam’ or ‘procedure’, refers to “one or more images taken on a single visit to a medical facility”<sup>8</sup>. To assess pathology progression, radiologists compare images (also referred to as ‘scans’ or ‘series’) from different studies. In the MIMIC-CXR dataset, each study (with one or more images) is accompanied by the report written by the radiologist. Figure B.1 represents the distribution of studies per subject within MIMIC-CXR and

the corresponding cumulative distribution function, showing that 67 % of the subjects have at least two different associated studies (and therefore at least two images acquired at different stages of the disease).

Another way to quantify temporal information in MIMIC-CXR is through the progression labels provided by the Chest ImaGenome dataset [72]. These progression labels are extracted from the reports and thus identify the cases when the radiologist explicitly describes changes. We found that in MIMIC, around 40 % of the reports are associated with a progression label from any of the available findings defined by ImaGenome.

## C. *MS-CXR-T* benchmark

### C.1. Temporal image classification

The *MS-CXR-T* temporal image classification contains progression labels for five findings (Consolidation, Edema, Pleural Effusion, Pneumonia and Pneumothorax) across three progression classes (Improving, Stable, and Worsening). This benchmark builds on the publicly available Chest ImaGenome gold and Chest ImaGenome silver datasets [72] which provide progression labels automatically derived from radiology reports. We collected a set of studies that are part of the ImaGenome silver dataset, excluding any studies that had been previously verified as part of the ImaGenome gold dataset. Additionally, we excluded studies where there are multiple progression labels for a single pathology (e.g. left pleural effusion has increased, right pleural effusion remains stable). We conducted a review process of the selected candidates, asking a board certified radiologist to either accept or reject the label. To inform their review of the labels, the radiologist was given access to the radiology report for the current image, and the sentence from which the auto generated label had been extracted.

After collecting our curated labels and labels from the ImaGenome gold dataset, we matched the report-based labels to specific image pairs, performing a second data curation step to create the image dataset. To ensure the diagnostic quality of all images in the dataset, if a study had multiple frontal scans we performed a quality control step

<sup>8</sup> Adapted from <https://ncithesaurus.nci.nih.gov/>

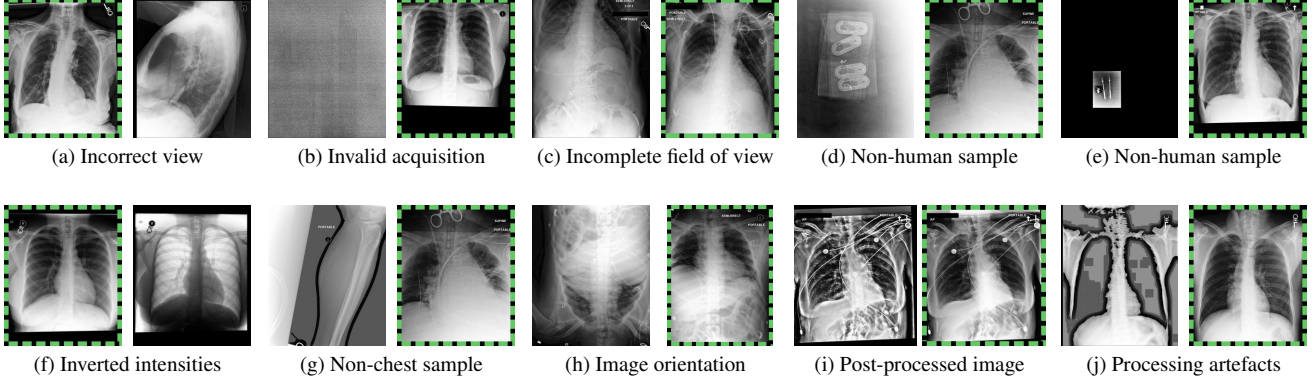


Figure A.4. Pairwise ranking of images performed by the proposed data curation method (see Section 3.3) on images from the MIMIC-CXR v2 dataset. Images highlighted with dashed green rectangles are automatically selected by our method and used for training to improve model’s downstream performance. The rejected image samples may not be appropriate for training due to image acquisition or image processing issues as shown in each subfigure above.

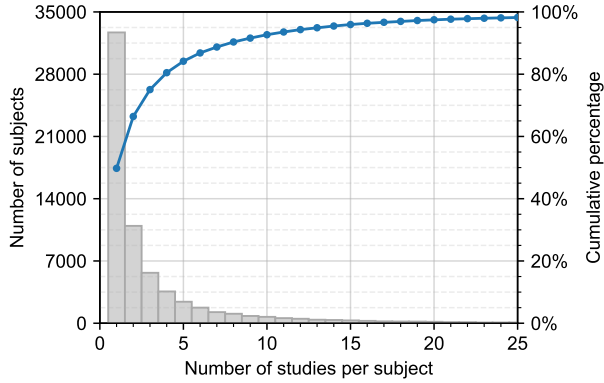


Figure B.1. Number of studies per subject in the MIMIC-CXR dataset. A study, sometimes referred to as an ‘exam’ or ‘procedure’, refers to “one or more images taken on a single visit to a medical facility” (adapted from <https://ncithesaurus.nci.nih.gov/>). Note that 67 % of subjects have at least two studies that happened at different times.

asking a radiologist to select the best image for each study. Fig. F.1 shows examples from the benchmark across different pathologies and progression labels.

The class distribution for the image classification task in *MS-CXR-T* is shown in Tab. C.1. As seen in the table, the class distribution of the dataset skews towards the *stable* and *worsening* classes. This could be explained as patients are more likely to get a chest X-ray scan when their condition is stable or deteriorating as opposed to when there is an improvement in patient condition.

## C.2. Temporal sentence similarity

In this section, we describe the process of creating the *MS-CXR-T* temporal sentence similarity benchmark, which consists of pairs of paraphrase or contradiction sentences in

Table C.1. *MS-CXR-T* temporal image classification benchmark: Showing the distribution of multi-image studies across different clinical findings, distribution of classes {Improving, Stable, Worsening} per finding, and number of subjects.

Findings	# of annotation pairs	Class distribution	# of subjects
Consolidation	201	14% / 42% / 44%	187
Edema	266	31% / 26% / 43%	241
Pleural effusion	411	19% / 49% / 32%	370
Pneumonia	237	8% / 25% / 67%	218
Pneumothorax	211	15% / 55% / 30%	148
Total	1326	18% / 40% / 42%	800

Table C.2. *MS-CXR-T* temporal sentence similarity benchmark: Number of paraphrase and contradiction examples in the full dataset and across the RadGraph and Swaps subsets.

Subset	# of paraphrase pairs	# of contradiction pairs	Total
Radgraph	42	75	117
Swaps	99	145	244
Total	141	220	361

terms of disease progression. We create this dataset using two different methods, *RadGraph* where paraphrase and contradiction sentence pairs are discovered by analysing graph representations of sentences and *Swaps* where paraphrases and contradictions are created by swapping out temporal keywords in the sentence.

To create this dataset, we first collected a set of sentences from the MIMIC dataset, using the Stanza constituency parser [82] to extract individual sentences from reports. Using the CheXbert labeller [63], we filtered this set to sentences that described one of seven pathologies - Atelectasis, Consolidation, Edema, Lung Opacity, Pleural Effusion, Pneumonia or Pneumothorax. We then filtered to sentences which contained at least one mention of a temporal key-



Table C.3. Examples of paraphrase and contradiction sentence pairs from the *MS-CXR-T* temporal sentence similarity benchmark. The examples are selected from the RadGraph and Swaps subsets (see Appendix C.2).

	Label	Sentence 1	Sentence 2
Swaps	Paraphrase	“Unchanged small-to-moderate right pleural effusion.”	“Stable small-to-moderate right pleural effusion.”
	Contradiction	“Interval worsening of the right-sided pneumothorax.”	“Interval resolution of the right-sided pneumothorax.”
RadGraph	Paraphrase	“There has also been a slight increase in left basal consolidation.”	“There is slight interval progression of left basal consolidation.”
	Contradiction	“Right mid and lower lung consolidations are unchanged.”	“There has been worsening of the consolidation involving the right mid and lower lung fields.”

word. Using this sentence pool, paraphrase and contradiction pairs were constructed in two ways. (I) We paired sentences from the sentence pool by matching on RadGraph [34] entities, relaxing the matching constraint only for temporal keywords and possible mentions of pathologies. (II) We swapped out temporal keywords in a sentence to create sentence pairs, choosing swap candidates from the top 5 masked token predictions from CXR-BERT-Specialized [9] provided they were temporal keywords. After creating candidate sentence pairs, we manually filtered out sentence pairs with ambiguous differences in terms of disease progression. A board certified radiologist then annotated each sentence pair as either paraphrase or contradiction. Sentences were filtered out in the annotation process if (I) they were not clear paraphrases or contradictions (II) the sentences differed in meaning and this difference was not related to any temporal information (III) they were not grammatically correct. The distribution of sentence pairs across the paraphrase and contradiction classes are described in Table C.2, see Table C.3 for examples from the benchmark.

## D. Temporal entity matching

To quantify how well the generated report describes progression-related information, we propose a new metric, namely temporal entity matching (TEM) score.

### D.1. Metric Formulation

We first extract entities (tagged as “observation” or “observation\_modifier”) from the text by running the named entity recognition model in the Stanza library [82]. Within the extracted entities, we manually curated a list of temporal entities that indicate progression (Appendix D.2). The list is reviewed by an expert radiologist. Given extracted temporal entities  $E$  in  $N$  pairs of reference and generated reports, we calculate global precision ( $p_E$ ) and global recall ( $r_E$ ), which are later used to compute the TEM score. It is defined as the harmonic mean of precision and recall (also

known as the F1 score).

$$p_E = \frac{\sum_{i=1}^N |E_{gen}^i \cap E_{ref}^i|}{\sum_{i=1}^N |E_{gen}^i|} \quad (3)$$

$$r_E = \frac{\sum_{i=1}^N |E_{gen}^i \cap E_{ref}^i|}{\sum_{i=1}^N |E_{ref}^i|} \quad (4)$$

### D.2. List of temporal keywords

The list of temporal keywords used to compute the TEM score are as follows: {bigger, change, cleared, constant, decrease, decreased, decreasing, elevated, elevation, enlarged, enlargement, enlarging, expanded, greater, growing, improved, improvement, improving, increase, increased, increasing, larger, new, persistence, persistent, persisting, progression, progressive, reduced, removal, resolution, resolved, resolving, smaller, stability, stable, stably, unchanged, unfolded, worse, worsen, worsened, worsening, unaltered}.

## E. Architecture and implementation details

### E.1. Hyper-parameters

The models are trained in a distributed setting across 8 GPU cards. For pre-training, we use a batch size of 240 (30 \* 8 GPUs) and the AdamW optimizer [43]. We use a linear learning rate scheduler with a warm-up proportion of 0.03 and base learning rate of  $2 \times 10^{-5}$ . We train for a maximum of 50 epochs and use validation set loss for checkpoint selection. The overall loss is a sum of components with weighting factors: global contrastive (1.0), local contrastive (0.5), and image-guided MLM (1.0) respectively, see Sec. 3.1 for further details on their formulation.

Following [9] we use sentence permutation as text-based data augmentation. Similarly, spelling errors in the reports are corrected prior to tokenisation of the text data<sup>9</sup>. For image augmentations, note that we apply the same augmentation to current and prior images to prevent severe misalignment. We resize the shorter edge to 512 and centre-crop to (448, 448). We apply random affine transformations

<sup>9</sup>[https://github.com/farrell236/mimic-cxr/blob/master/txt/section\\_parser.py](https://github.com/farrell236/mimic-cxr/blob/master/txt/section_parser.py)

(rotation up to  $30^\circ$  and shear up to  $15^\circ$ ) and colour jitter (brightness and contrast).

## E.2. Training infrastructure

We train with distributed data processing (DDP) on eight NVIDIA Tesla V100s with 32GB of memory each. To handle inconsistently-present prior images with DDP, we define a custom batch sampler. This sampler is a mixture of two samplers, in proportion to their dataset coverage: a sampler which produces batches with *only* multi-image examples –  $(\mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}}, \mathbf{x}_{\text{txt}}^{\text{curr}}) \in \mathcal{D}_m$  and one with only single-image examples –  $(\mathbf{x}_{\text{img}}^{\text{curr}}, \emptyset, \mathbf{x}_{\text{txt}}^{\text{curr}}) \in \mathcal{D}_s$ . Each GPU then processes a batch which is entirely single or multi-image, avoiding branching logic within the forward pass and enabling an efficient single pass through the CNN to process all input images (current or prior) by concatenating them along the batch dimension.

We confirmed that although the custom sampler theoretically impacts the order in which the dataset is traversed, it has a negligible effect on training metrics relative to fully random sampling. Since we train on eight GPUs and collect negatives across all GPUs during contrastive training, each update involves on average a representative mixture of both single-image and multi-image samples.

Finally, following [9] we use the DICOM images from MIMIC-CXR to avoid JPEG compression artefacts.

## F. Adaptation and experimentation details

### F.1. Fine-tuning BioViL-T for report generation

During fine-tuning of BioViL-T for report generation, we minimise the cross entropy loss to maximise the log likelihood of the report in an autoregressive manner given the input images. The model is initialised from the pretrained weights of the image encoder and the text encoder. Similar to the cross-modal masked language modelling task, we additionally train a linear projection layer to map the projected patch embeddings to the same hidden dimension of the text encoder, and we train cross-attention layers in each transformer block. The difference from the masked language modelling task is that we change the bidirectional self-attention to unidirectional causal attention that can only access the past tokens. If trained with prior report, we pass the prior report as prefix to condition the generation of the current report (the current and prior report are separated by [SEP]), and we only back-propagate the gradients from the loss on the tokens in the current report.

For all experiments, we train the model for 100 epochs and we chose the best checkpoint according to metrics on the validation set. We performed grid search for learning rate in  $[10^{-5}, 2 \times 10^{-5}, 5 \times 10^{-5}]$  and found  $2 \times 10^{-5}$  to be optimal. We ran each experiment with 3 random seeds and report mean and standard deviation.

Table F.1. Results for report generation task: Predictions are evaluated on NEM. The approaches are grouped into two broad categories: NN (Nearest Neighbour) and AR (Auto-Regressive). BioViL-T pre-training consistently yields superior decoding performance. Further, the use of prior image and report consistently yield performance gains demonstrating the importance of such domain priors.

	Method	Pre-training	Prior Img/Report	NEM
NN	CXR-RePaiR-2 [25]	BioViL	✗ / ✗	13.36
	Baseline (NN) [9]	BioViL	✗ / ✗	16.25
	Proposed (NN)	BioViL-T	✓ / ✗	17.55
AR	Baseline (AR) [9]	BioViL	✗ / ✗	24.27 ± 0.22
	Proposed	BioViL-T	✓ / ✗	25.50 ± 0.04
	Proposed	BioViL-T	✓ / ✓	<b>26.95 ± 0.17</b>

In addition to the metrics we reported in the main text, we also evaluate the generated reports by named entity metric (NEM). This metric was defined in [45] to measure the accuracy of reporting clinically relevant entities in the generated reports (Similar to how TEM is computed to measure the match of temporal entities in our study). Following [45], we extract entities (tagged as “observation” or “observation\_modifier”) from the text by running the named entity recognition model in the Stanza library [82]. The results are presented in Tab. F.1.

### F.2. Nearest-neighbour-based report retrieval

The joint latent space learnt by BioViL-T can also be used to directly perform report retrieval without requiring task-specific model fine-tuning. Given the test image, we retrieve its semantically closest report from the training set in the joint latent space. Specifically, we encode each test image with the image model in BioViL-T and collect its projected image embeddings, and similarly we encode all the reports in the training data with their projected text embeddings. For each test study, we compute cosine similarity between the test image embedding and all the text embeddings from the training set in the joint latent space, and we retrieve the closest text embedding and use its corresponding report as the prediction. To evaluate the retrieval performance, we use the same decoding metrics on the retrieved reports and report results in the top section of Table 1. In a separate set of experiments, we also tried performing nearest neighbour search only within the image embedding space by retrieving the report associated with the closet image embedding, but this yielded sub-optimal performance compared with using the joint latent space.

### F.3. Fine-tuning for temporal image classification

In this section, we describe the training dataset and fine-tuning procedure for the fully supervised and few-shot settings of the temporal image classification task. For this task, we finetune BioViL-T on a subset of the Chest ImaGenome

silver dataset [72] to predict progression labels for 5 different pathologies. To create our training dataset, we filter out image pairs from this dataset where there are multiple directions of progression of a single pathology in the image-pair. We additionally perform an automatic data curation step to choose higher quality image pairs when possible, as described in 3.3. Table F.2 shows the number of training samples and label distribution for the training dataset.

Table F.2. Statistics of the training dataset used for downstream fine-tuning on temporal image classification.

Findings	# labelled pairs	Class distribution	# of subjects
Consolidation	7012	15% / 42% / 43%	3308
Edema	14170	28% / 33% / 39%	4813
Pleural effusion	26320	16% / 53% / 31%	6838
Pneumonia	8471	12% / 29% / 59%	4197
Pneumothorax	3795	21% / 57% / 22%	1161

For the fully supervised setting, we add a multilayer classification head to the BioViL-T image encoder and fine-tune the model independently for each pathology. We use weighted cross entropy loss with a batch size of 128 and the AdamW optimizer [43]. During parameter optimisation, positional encodings and missing-image embeddings are exempt from weight decay penalty as in [73]. We train for 30 epochs, with a linear learning rate schedule, a warmup proportion of 0.03 and a base learning rate of  $1 \times 10^{-5}$ . For data augmentation, we first resize the shorter edge of the image to 512 and centre crop to (448, 448). We apply random horizontal flips, random cropping, random affine transformations (rotation up to  $30^\circ$ , shear up to  $15^\circ$ ), colour transforms (brightness and contrast) and Gaussian noise.

For the few-shot setting we tune only a single-layer linear head on the BioViL-T image encoder and freeze the rest of the encoder. We initialise the weight matrix of the linear head with values from encoded text prompts [9] for each of the three progression classes, and the bias matrix is initialised with zeros. To train, we again use weighted cross entropy loss, with a batch size of 32 and the AdamW optimizer. We use a learning rate of  $1 \times 10^{-3}$  and train for 40 epochs. For data augmentation, we resize the shorter edge of the image to 448 and center crop to (448, 488). We apply random horizontal flips, random affine transformations (rotation up to  $45^\circ$  and shear up to  $25^\circ$ ), colour transforms (brightness and contrast). As in the pre-training step, we always synchronise image data augmentations to apply the identical transforms to the current and prior images.

#### F.4. Auto-regressive prompting for zero-shot temporal image classification

Following the GPT-3 style language prompting [11], we prompt the fine-tuned AR language decoding model with the template: “[FINDING] is” and infer the next token

Table F.3. Prompting the AR language decoding model for zero-shot image classification. The list above shows the mapping from decoded tokens to progression classes.

Target class	Tokens
Improving	better, cleared, decreased, decreasing, improved, improving, reduced, resolved, resolving, smaller
Stable	constant, stable, unchanged
Worsening	bigger, developing, enlarged, enlarging, greater, growing, increased, increasing, larger, new, progressing, progressive, worse, worsened, worsening

to perform temporal classification for each of the five findings. The mapping from the predicted next token to the three progression classes is characterised by a short list of tokens provided in Table F.3. After computing the posterior for each token in the list, the obtained values are normalised across the three classes, and the class with the highest score is selected as the prediction. The corresponding results are reported in Table 2.

#### F.5. Further analysis of image-guided MLM

In Section 4.6 we used a simplified notation for the computation of  $\Delta_{\text{img}}^{\text{prior}}(m)$  for ease of exposition – here we provide further detail. Recall that  $\mathbf{w} = (w_1, \dots, w_M)$  is a sequence of tokens and  $\mathbf{w}_{\setminus m}$  is that sequence with token  $m$  masked. Let  $p_\theta(\mathbf{w}_m | \mathbf{w}_{\setminus m}, \mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}})$  be the text model’s predicted probability of token  $m$  given  $\mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}}$ , and  $\mathbf{w}_{\setminus m}$  ( $\theta$  are the weights of the model). Then,  $l(w, p_\theta(\mathbf{w}_m | \mathbf{w}_{\setminus m}, \mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}}))$  is the cross-entropy loss of predicting token  $m$  given those inputs.

It is possible for different sentences in a report to refer to the same image finding. Since we mask single tokens at a time, to prevent information leakage from other sentences we consider each sentence in a report independently. Suppose report  $\mathbf{x}_{\text{txt}}^{\text{curr}}$  consists of  $S$  sentences, so we have  $\mathbf{x}_{\text{txt}}^{\text{curr}} = [\mathbf{w}^1, [\text{SEP}], \dots, [\text{SEP}], \mathbf{w}^S]$ , where  $\mathbf{w}^s$  is the tokens of sentence  $s$  and [SEP] separates sentences.

For a given sample  $(\mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}}, \mathbf{x}_{\text{txt}}^{\text{curr}}) \in \mathcal{D}_m$  in the test set indexed by  $i$ , we define

$$\delta_i(m) = \sum_{s \in S} [l(m, p_\theta(\mathbf{w}_m^s | \mathbf{w}_{\setminus m}^s, \mathbf{x}_{\text{img}}^{\text{curr}}, \emptyset)) - l(m, p_\theta(\mathbf{w}_m^s | \mathbf{w}_{\setminus m}^s, \mathbf{x}_{\text{img}}^{\text{curr}}, \mathbf{x}_{\text{img}}^{\text{prior}}))]$$

This is the MLM loss for predicting  $m$  given each sentence in the report with and without the prior image. Note that if  $m$  does not appear in a given sentence, its contribution to the sum is zero. The overall  $\Delta_{\text{img}}^{\text{prior}}(m)$  is computed across

Category	Description	Examples
Progression	Pertaining to change or progression	<i>bigger, cleared, new</i>
Support devices	Tubes, lines and implants	<i>nasogastric, pacemaker, cannula</i>
‘Other’	No clear category	<i>can, relevant, overall</i>
Stop word	‘Insignificant’ words	<i>the, no, of</i>
Positional	Localisation (not anatomical)	<i>right, lower, bilateral</i>
Meta	Pertaining to the report itself or practice of radiology	<i>evidence, radiograph, study</i>
Anatomy	Anatomical locations	<i>pulmonary, chest, mediastinal</i>
Descriptive	Qualitative appearance of a finding	<i>layering, focal, patchy</i>
Size or degree	Quantifying extent or severity	<i>extensive, moderate, severe</i>
Finding	Radiographic finding or pathology	<i>edema, penumonia, pneumothorax</i>
Uncertain	Expression of certainty or doubt	<i>may, possible, concerning</i>

Table F.4. Semantic categories used in Figure 4.

all samples:

$$\Delta_{\text{img}}^{\text{prior}} = \frac{1}{N_m} \left( \sum_{i \in \mathcal{D}_m^{\text{test}}} \delta_i(m) \right) \quad (5)$$

where  $N_m$  is the number of *sentences* in reports in  $\mathcal{D}_m^{\text{test}}$  in which token  $m$  appears. This estimate is subject to high variance when  $N_m$  is small. Hence, for Figure 4 we filter to tokens  $m$  with  $N_m \geq 10$ . We collected 931 tokens with  $N_m \geq 10$  from the validation set for manual annotation by a board-certified radiologist. The categories, shown in Figure 4 and described in Table F.4 are specific to the radiology domain.

## F.6. Sentence similarity experiment

The text models are evaluated in isolation to observe if their encoding is sensitive to key clinical observations. To achieve this, we assess the quality of sentence representations obtained from our text model by examining how well the contradiction and paraphrase pairs can be separated in the embedding space. Unlike the traditional NLI task where a model needs to be fine-tuned, here the models are probed in a zero-shot setting and the BERT output token embeddings are utilised. To do so, we encode the sentences from RadNLI and *MS-CXR-T* sentence similarity datasets with the [CLS] token from CXR-BERT-Specialised [9] and BioViL-T. For PubMedBERT [29] and CXR-BERT-General [9] which did not directly optimise the [CLS] token during pretraining, we follow [56] to average the token output embeddings to represent each sentence.

Cosine similarity is computed between the representations of each sentence pair in the dataset [56] and is used as logits for the binary classification between paraphrase and contradiction. Note that for RadNLI, we use the subset of ‘entailment’ and ‘contradiction’ pairs and discard the ‘neutral’ pairs to unify the task across the two datasets. Given the similarities for each sentence pair, we report ROC-AUC

and binary-accuracy. For the latter, a threshold value for each method is derived by setting aside a validation set. For this, we perform ten-fold cross validation and tune the threshold with step size of 0.005 on the validation set.

## F.7. Image registration algorithm

In Section 4.2, image registration is applied to pairs of images as a preprocessing step to enable a fair comparison for the baseline approaches (e.g., BioViL [9]). We performed bidirectional multi-scale registration between image pairs optimising an affine transformation (4 degrees of freedom), using mutual information (MI) [65] with 128 bins as the similarity criterion. In more detail, the spatial transformation is characterised by four parameters: two for translation, one for isotropic scaling, and one for rotation. The optimisation is repeated five times with different random seeds for initialisation, and the run with the highest MI is selected to determine the final spatial alignment. To better identify the correspondences between the scans, bilateral filtering is applied to each image before registration to remove detailed texture whilst preserving edge information [38]. Our implementation is based on the SimpleITK library [44].



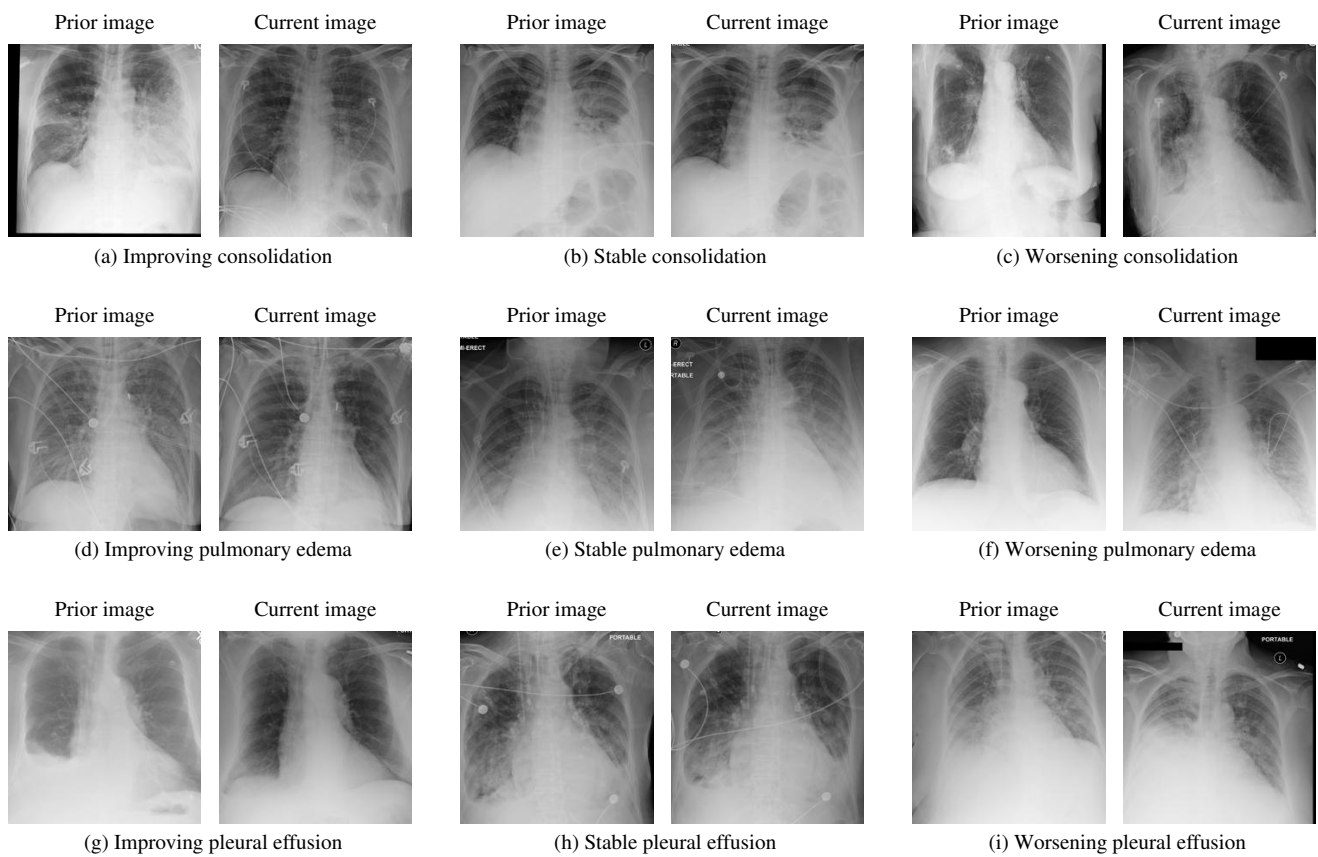


Figure F.1. Examples of image pairs in our *MS-CXR-T* benchmark.