
Intra-video Positive Pairs in Self-Supervised Learning for Ultrasound

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The videographic nature of ultrasound offers flexibility for defining the similarity
2 relationship between pairs of images for self-supervised learning (SSL). In this
3 study, we investigated the effect of utilizing proximal, distinct images from the
4 same ultrasound video as pairs for joint embedding SSL. Additionally, we intro-
5 duced a sample weighting scheme that increases the weight of closer image pairs
6 and demonstrated how it can be integrated into SSL objectives. Named *Intra-Video*
7 *Positive Pairs* (IVPP), the method surpassed previous ultrasound-specific con-
8 trastive learning methods' average test accuracy on COVID-19 classification with
9 the POCUS dataset by $\geq 1.3\%$. Investigations revealed that some combinations of
10 IVPP hyperparameters can lead to improved or worsened performance, depending
11 on the downstream task.

12 1 Introduction

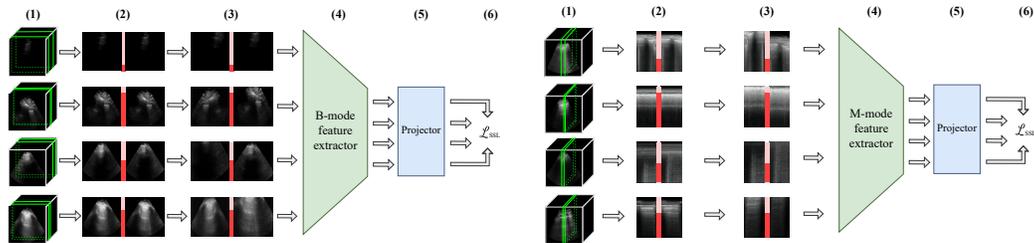
13 Deep learning has been extensively studied as a means to automate diagnostic tasks in medical
14 ultrasound (US) [1–5]. Barriers to the development of such systems include lack of open access
15 to large datasets, along with the cost and expertise required to label vast amounts of institutionally
16 acquired examinations. Additionally, many US examinations in acute care are not archived or
17 documented [6, 7]. Joint embedding self-supervised learning (SSL) has been explored as a means
18 to leverage unlabelled data for representation learning with US [8–11]. A common way to define
19 positive pairs for SSL is to apply two stochastic transformations to an image, producing two distorted
20 views with similar content. Being a video-based examination, US offers a unique opportunity to
21 compose alternative pairwise relationships. However, due to the dynamic nature of US, all frames in
22 a US video may not possess the same label for all downstream US interpretation tasks. Moreover, US
23 videos from the same examination or patient may bear a striking resemblance to each other.

24 In this study, we examined the effect of proximity and sample weighting of intra-video positive pairs
25 for common SSL methods applied to US. We evaluated on two tasks: one for B-mode US (i.e., US
26 video composed from several scan lines) and M-mode US (i.e., images depicting the evolution of one
27 US scan line over time). Our contributions are summarized as follows:

- 28 • A method for sampling US intra-video positive pairs for joint embedding SSL.
- 29 • A sample weighting scheme for joint embedding SSL methods that weighs positive pairs
30 according to the temporal or spatial distance between them in their video of origin
- 31 • A comprehensive assessment of intra-video positive pairs integrated with contrastive and
32 non-contrastive SSL methods, as measured by downstream performance in B-mode and
33 M-mode lung US classification tasks.

34 Figure 1 summarizes the methods proposed in this study. Although previous studies have formulated
35 contrastive learning methods with intra-video positive pairs for US [8, 10, 12, 13], the authors believe

36 there are no preceding studies that investigated the effect of sampling multiple images from the
 37 same US video in non-contrastive learning. More generally, we believe that this study is the first to
 38 integrate sample weights into non-contrastive objectives.



(a) For B-mode ultrasound, positive pairs are temporally separated images from the same video. (b) For M-mode ultrasound, positive pairs are spatially separated images from the same video.

Figure 1: An overview of the methods introduced in this study. Positive pairs of images separated by no more than a threshold are sampled from the same B-mode video (1). Sample weights inversely proportional to the separation between each image (red bars) are calculated for each pair (2). Random transformations are applied to each image (3). Images are sent to a neural network consisting of a feature extractor (4) and a projector (5) connected in series. The outputs are used to calculate the self-supervised objective \mathcal{L}_{SSL} (6).

39 2 Methods

40 **Datasets:** As done in previous studies on on US-specific joint embedding methods [8, 10, 12, 13], we
 41 evaluate on the public POCUS lung US dataset [14]. This dataset contains 140 publicly sourced US
 42 videos (2116 images) labelled for three classes: COVID-19 pneumonia, non-COVID-19 pneumonia,
 43 and normal lung. Pretraining is conducted on the public Butterfly dataset, which contains 22
 44 unlabelled lung ultrasound videos [15]. We also utilize a private dataset of 25 917 parenchymal lung
 45 US videos (5.9×10^6 images), hereafter referred to as *ParenchymalLUS*. Of these videos, 20 000
 46 had no labels. We evaluated on two binary classification tasks: A-lines versus B-line classification
 47 (i.e., AB), and lung sliding classification (i.e., LS). Details on *ParenchymalLUS* and descriptions of the
 48 downstream tasks can be found in Appendix A.

49 **Intra-video Positive Pairs (IVPP):** Clinically relevant patterns commonly surface and disappear
 50 within the same US video as the US probe and/or the patient move. Without further knowledge of
 51 the US examinations in an unlabelled dataset, we conjectured that it may be safest to only assume
 52 that positive pairs are intra-video images that are close to each other. Our method distinguishes itself
 53 from prior work by only considering proximal frames to be positive pairs and disregarding distant
 54 intra-video pairs, treating them as neither positive nor negative pairs.

55 For B-mode US videos, we define positive pairs as intra-video images x_1 and x_2 that are separated
 56 by no more than δ_t seconds (Figure 2a). To accomplish this, x_1 is randomly selected from the video’s
 57 images, and x_2 is randomly drawn from the set of images within δ_t seconds of x_1 . Videos with higher
 58 frame rates provide more positive pair candidates, potentially increasing the diversity of pairs due
 59 to naturally occurring noise. A similar sampling scheme is applied for M-mode US images. Like
 60 previous studies, we define M-mode images as vertical slices through the time axis of a B-mode video,
 61 taken at a specific x-coordinate [11, 16, 17]. Accordingly, M-mode images are columns of B-mode
 62 pixels for every frame, concatenated horizontally. We define positive pairs as M-mode images whose
 63 x-coordinates differ by no more than δ_x pixels (Figure 2b). To avoid resolution differences, B-mode
 64 videos are resized to the same dimensions prior to sampling M-mode images.

65 **Sample Weights:** The chance that intra-video images are semantically related increases as temporal
 66 or spatial separation decreases. To temper the effect of unrelated positive pairs, we applied sample
 67 weights to positive pairs in the SSL objective according to their temporal or spatial distance. Sample
 68 weights were incorporated into each SSL objective trialled in this study: SimCLR [18], Barlow
 69 Twins [19], and VICReg [20]. Appendix B details how sample weights are calculated and integrated
 70 into the SSL objectives.

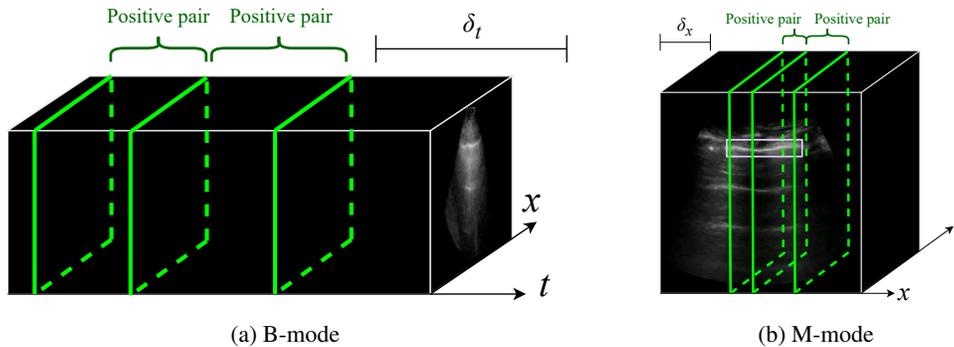


Figure 2: For B-mode ultrasound, positive pairs are frames in the same video that are within δ_t seconds of each other. For M-mode ultrasound, positive pairs are M-mode images originating from the same B-video that are no more than δ_x pixels apart. In the context of lung ultrasound, M-mode images should intersect the pleural line (outlined in mauve).

71 3 Experiments & Evaluation

72 **Fine-tuning Performance:** We evaluated our
 73 approach on each of the downstream tasks and
 74 report the mean cross-validation accuracy for
 75 COVID, along with the test set AUC for AB and
 76 LS. Pretraining and training protocols can be
 77 found in Appendix C. Multiple values of the
 78 threshold parameter were investigated, with and
 79 without sample weights. For the COVID and AB
 80 tasks, we examined $\delta_t \in \{0, 0.5, 1, 1.5\}$ sec-
 81 onds. For LS we explored $\delta_x \in \{0, 5, 10, 15\}$
 82 pixels. As shown in Figure 3, mean cross-
 83 validation accuracy of each fine-tuned model
 84 peaked at nonzero values of δ_t . Sample weights
 85 decreased performance when $\delta = 0.5$, but were
 86 helpful when $\delta = 1.0$ and $\delta = 1.5$. Figure 4
 87 gives fine-tuning performance for AB and LS. We
 88 observed no discernible trend for the effect of
 89 sample weights that was consistent for any task,
 90 pretraining method, δ_t , or δ_x . A striking finding
 91 across AB and LS was that SimCLR consistently outperformed Barlow Twins and VICReg, which
 92 are both non-contrastive methods. Evaluation via linear classification was also conducted, revealing
 93 similar results (see Appendix D).

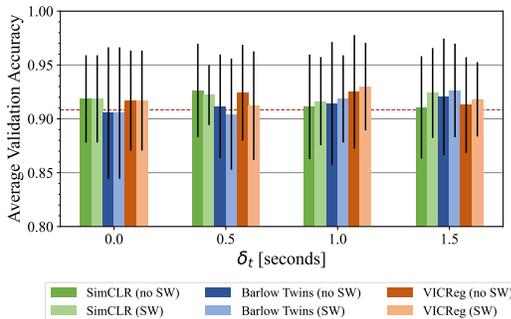


Figure 3: Average test accuracy across 5-fold cross validation on the POCUS dataset, pretrained using a variety of values for δ_t , with and without sample weights. The dashed line indicates initialization with ImageNet-pretrained weights.

94 **Comparison to Baselines:** We compared IVPP (with the best hyperparameter assignments) for
 95 each SSL objective against USCL [8] and UCL [10], which are preexisting US contrastive learning
 96 methods. As shown in Table 1, IVPP outperformed all baseline methods on POCUS, regardless of the
 97 pretraining objective. USCL and IVPP with SimCLR performed comparably on the AB task. On the
 98 LS task, which is more fine-grained and had a stronger class imbalance, IVPP with SimCLR achieved
 99 the greatest performance. Non-contrastive methods were unremarkable, achieving lower performance
 100 than networks initialized with ImageNet-pretrained weights.

101 **Label Efficiency:** We devised an experiment to compare the performance of models pretrained using
 102 different IVPP hyperparameters in low-label settings. The ParenchymalLUS training set was split
 103 by patient identifier into 20 subsets. Pretrained models were fine-tuned on each subset, resulting in
 104 a population of 20 test set metrics for each hyperparameter combination. As is visible in Figure 5,
 105 SimCLR obtained the best performance by a large margin. As detailed in Appendix D, the differences
 106 between some means were statistically significant.

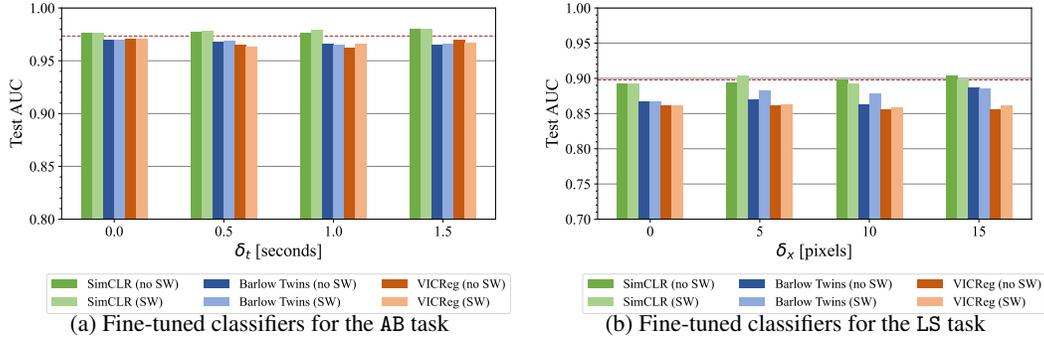


Figure 4: ParenchymaLUS test set AUC for models fine-tuned for the AB and LS binary classification tasks and pretrained with a variety of intra-video positive pair thresholds, with and without sample weights (SW). The dashed line indicates initialization with ImageNet-pretrained weights.

Dataset	POCUS		ParenchymaLUS	
	Task	COVID Mean (std) test accuracy	AB Test AUC	LS Test AUC
Random initialization		0.881 (0.050)	0.954	0.790
ImageNet initialization		0.908 (0.043)	0.973	0.898
USCL [8]		0.905 (0.044)	0.979	0.874
US UCL [10]		0.901 (0.054)	0.967	0.809
IVPP [SimCLR]		0.926 (0.043)	0.980	0.903
IVPP [Barlow Twins]		0.921 (0.054)	0.969	0.887
IVPP [VICReg]		0.930 (0.046)	0.971	0.862

Table 1: Performance of fine-tuned models pretrained using IVPP compared to US-specific contrastive learning methods, USCL and UCL, and to baseline Random and ImageNet initializations.

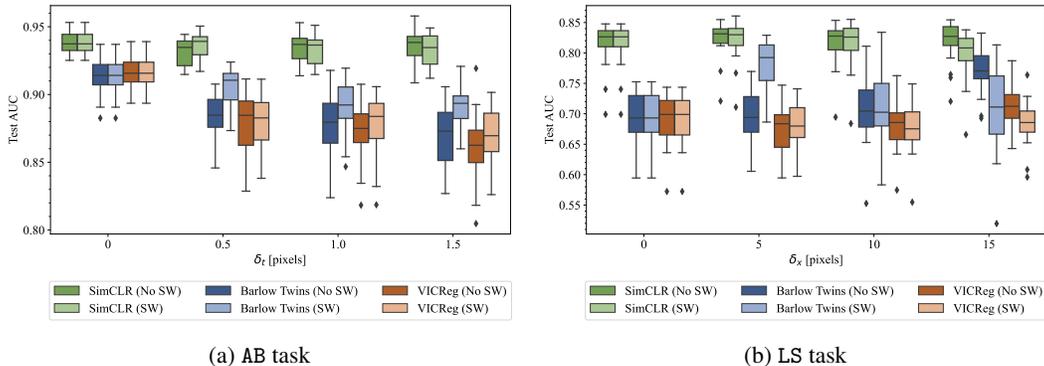


Figure 5: Distributions of test AUC for each pretraining method and assignment to δ , with and without sample weights. Each experiment is repeated 20 times on disjoint subsets of the training set, each containing all images from a group of patients.

107 4 Discussion

108 Overall, the results indicated that the optimal assignment for IVPP hyperparameters may be problem-
 109 specific. First, IVPP may improve performance on downstream ultrasound interpretation tasks;
 110 however, practitioners are advised to include a range of values of δ with and without sample weights
 111 in their hyperparameter search. Subsequent work could explore IVPP for other downstream tasks
 112 in US outside of the lung. Second, SimCLR outperformed non-contrastive methods across multiple
 113 tasks – contrary to our initial belief. Future work assessing non-contrastive methods for tasks in US
 114 examinations or alternative imaging modalities would shed light on the utility of non-contrastive
 115 methods outside the typical evaluation setting of photographic images.

References

- 116
- 117 [1] M. R. Whitson and P. H. Mayo, "Ultrasonography in the emergency department," *Critical Care*, vol. 20,
118 pp. 1–8, 2016.
- 119 [2] Y. H. Lau and K. C. See, "Point-of-care ultrasound for critically-ill patients: A mini-review of key
120 diagnostic features and protocols," *World Journal of Critical Care Medicine*, vol. 11, no. 2, p. 70, 2022.
- 121 [3] N. J. Soni, R. Arntfield, and P. Kory, *Point-of-Care Ultrasound*. Philadelphia: Elsevier, second ed., 2020.
- 122 [4] R. Sood, A. F. Rositch, D. Shakoor, E. Ambinder, K.-L. Pool, E. Pollack, D. J. Mollura, L. A. Mullen, and
123 S. C. Harvey, "Ultrasound for breast cancer detection globally: a systematic review and meta-analysis,"
124 *Journal of global oncology*, vol. 5, pp. 1–17, 2019.
- 125 [5] E. S. Yim and G. Corrado, "Ultrasound in sports medicine: relevance of emerging techniques to clinical
126 care of athletes," *Sports medicine*, vol. 42, pp. 665–680, 2012.
- 127 [6] M. K. Hall, J. Hall, C. P. Gross, N. J. Harish, R. Liu, S. Maroongroge, C. L. Moore, C. C. Raio, and R. A.
128 Taylor, "Use of point-of-care ultrasound in the emergency department: insights from the 2012 medicare
129 national payment data set," *Journal of Ultrasound in Medicine*, vol. 35, no. 11, pp. 2467–2474, 2016.
- 130 [7] R. Kessler, J. R. Stowell, J. A. Vogel, M. M. Liao, and J. L. Kendall, "Effect of interventional program on
131 the utilization of pacs in point-of-care ultrasound," *Journal of digital imaging*, vol. 29, pp. 701–705, 2016.
- 132 [8] Y. Chen, C. Zhang, L. Liu, C. Feng, C. Dong, Y. Luo, and X. Wan, "USCL: Pretraining Deep Ultrasound
133 Image Diagnosis Model Through Video Contrastive Representation Learning," in *Medical Image Com-
134 puting and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg,
135 France, September 27–October 1, 2021, Proceedings, Part VIII 24*, pp. 627–637, Springer, 2021.
- 136 [9] D. Anand, P. Annangi, and P. Sudhakar, "Benchmarking self-supervised representation learning from a
137 million cardiac ultrasound images," in *2022 44th Annual International Conference of the IEEE Engineering
138 in Medicine & Biology Society (EMBC)*, pp. 529–532, IEEE, 2022.
- 139 [10] S. Basu, S. Singla, M. Gupta, P. Rana, P. Gupta, and C. Arora, "Unsupervised contrastive learning of
140 image representations from ultrasound videos with hard negative mining," in *International Conference on
141 Medical Image Computing and Computer-Assisted Intervention*, pp. 423–433, Springer, 2022.
- 142 [11] B. VanBerlo, B. Li, A. Wong, J. Hoey, and R. Arntfield, "Exploring the utility of self-supervised pretraining
143 strategies for the detection of absent lung sliding in m-mode lung ultrasound," in *Proceedings of the
144 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3076–3085, 2023.
- 145 [12] Y. Chen, C. Zhang, C. H. Ding, and L. Liu, "Generating and weighting semantically consistent sample
146 pairs for ultrasound contrastive learning," *IEEE Transactions on Medical Imaging*, 2022.
- 147 [13] C. Zhang, Y. Chen, L. Liu, Q. Liu, and X. Zhou, "Hico: Hierarchical contrastive learning for ultrasound
148 video model pretraining," in *Proceedings of the Asian Conference on Computer Vision*, pp. 229–246, 2022.
- 149 [14] J. Born, G. Brändle, M. Cossio, M. Disdier, J. Goulet, J. Roulin, and N. Wiedemann, "POCOVID-Net:
150 Automatic Detection of COVID-19 From a New Lung Ultrasound Imaging Dataset (POCUS)," *arXiv
151 preprint arXiv:2004.12084*, 2020.
- 152 [15] Butterfly Network, "Covid-19 ultrasound gallery." [https://www.butterflynetwork.com/covid19/
153 covid-19-ultrasound-gallery](https://www.butterflynetwork.com/covid19/covid-19-ultrasound-gallery), 2020. Accessed: September 20, 2020.
- 154 [16] M. Jaščur, M. Bundzel, M. Malík, A. Dzian, N. Ferenčík, and F. Babič, "Detecting the absence of lung
155 sliding in lung ultrasounds using deep learning," *Applied Sciences*, vol. 11, no. 15, p. 6976, 2021.
- 156 [17] B. VanBerlo, D. Wu, B. Li, M. A. Rahman, G. Hogg, B. VanBerlo, J. Tschirhart, A. Ford, J. Ho, J. McCauley,
157 *et al.*, "Accurate assessment of the lung sliding artefact on lung ultrasonography using a deep learning
158 approach," *Computers in Biology and Medicine*, vol. 148, p. 105953, 2022.
- 159 [18] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual
160 representations," in *International conference on machine learning*, pp. 1597–1607, PMLR, 2020.
- 161 [19] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy
162 reduction," in *International Conference on Machine Learning*, pp. 12310–12320, 2021.
- 163 [20] A. Bardes, J. Ponce, and Y. LeCun, "VICReg: Variance-Invariance-Covariance Regularization for Self-
164 Supervised Learning," in *International Conference on Learning Representations*, 2022.
- 165 [21] R. Arntfield, D. Wu, J. Tschirhart, B. VanBerlo, A. Ford, J. Ho, J. McCauley, B. Wu, J. Deglint, R. Chaud-
166 hary, *et al.*, "Automation of Lung Ultrasound Interpretation via Deep Learning for the Classification of
167 Normal Versus Abnormal Lung Parenchyma: A Multicenter Study," *Diagnostics*, vol. 11, no. 11, p. 2049,
168 2021.
- 169 [22] B. VanBerlo, D. Smith, J. Tschirhart, B. VanBerlo, D. Wu, A. Ford, J. McCauley, B. Wu, R. Chaudhary,
170 C. Dave, *et al.*, "Enhancing annotation efficiency with machine learning: Automated partitioning of a lung
171 ultrasound dataset by view," *Diagnostics*, vol. 12, no. 10, p. 2351, 2022.

- 172 [23] B. VanBerlo, B. Li, J. Hoey, and A. Wong, "Self-supervised pretraining improves performance and
173 inference efficiency in multiple lung ultrasound interpretation tasks," *arXiv preprint arXiv:2309.02596*,
174 2023.
- 175 [24] D. A. Lichtenstein and Y. Menu, "A bedside ultrasound sign ruling out pneumothorax in the critically ill:
176 lung sliding," *Chest*, vol. 108, no. 5, pp. 1345–1348, 1995.
- 177 [25] D. A. Lichtenstein, G. Mezière, N. Lascols, P. Biderman, J.-P. Courret, A. Gepner, I. Goldstein, and
178 M. Tenoudji-Cohen, "Ultrasound diagnosis of occult pneumothorax," *Critical care medicine*, vol. 33, no. 6,
179 pp. 1231–1238, 2005.
- 180 [26] D. A. Lichtenstein, *Whole body ultrasonography in the critically ill*. Springer Science & Business Media,
181 2010.
- 182 [27] S. Azizi, B. Mustafa, F. Ryan, Z. Beaver, J. Freyberg, J. Deaton, A. Loh, A. Karthikesalingam, S. Kornblith,
183 T. Chen, *et al.*, "Big self-supervised models advance medical image classification," in *Proceedings of the*
184 *IEEE/CVF international conference on computer vision*, pp. 3478–3488, 2021.
- 185 [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the*
186 *IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- 187 [29] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan,
188 *et al.*, "Searching for MobileNetV3," in *Proceedings of the IEEE/CVF international conference on computer*
189 *vision*, pp. 1314–1324, 2019.
- 190 [30] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, K. Keutzer, and C.-J. Hsieh,
191 "Large batch optimization for deep learning: Training bert in 76 minutes," *arXiv preprint arXiv:1904.00962*,
192 2019.

193 **A ParenchymalLUS Dataset and Downstream Tasks**

194 ParenchymalLUS is a subset of a larger database of de-identified lung US videos that was partially
 195 labelled for previous work [17, 21]. Access to this database was permitted via ethical approval by
 196 [redacted].¹ The labelled portion of ParenchymalLUS was split by patient identifier into training,
 197 validation, and test sets. Its unlabelled portion consists of 20 000 videos from the unlabelled pool of
 198 videos in the database that were predicted to contain a parenchymal view of the lungs by a previously
 199 trained lung US view classifier [22]. Below are descriptions of the lung US binary classification tasks
 200 for which labels were available in ParenchymalLUS.

201 **A-line vs. B-line Classification (AB):** A-lines and B-lines are two cardinal artifacts in B-mode lung
 202 US that can provide quick information on the status of a patient’s lung tissue. A-lines are reverberation
 203 artifacts that are indicative of normal, clear lung parenchyma [3]. On lung US, they as horizontal
 204 lines deep to the pleural line. Conversely, B-lines are indicative of diseased lung tissue [3]. Generally,
 205 the two are mutually exclusive. We evaluate on the binary classification task of A-lines versus B-lines
 206 on lung US, as was done in previous work benchmarking joint embedding SSL methods for lung US
 207 tasks [23].

208 **Lung Sliding Classification: (LS)** Lung sliding is a dynamic artifact that, when observed on a
 209 parenchymal lung US view, rules out the possibility of a pneumothorax at the site of the probe [24].
 210 The absence of lung sliding is suggestive of pneumothorax, warranting further investigation. On
 211 B-mode US, lung sliding manifests as a shimmering of the pleural line [24]. The presence or absence
 212 of lung sliding is also appreciable on M-mode lung US images that intersect the pleural line [25, 26].
 213 We evaluate on the binary lung sliding classification task, where positive pairs are 3-second M-mode
 214 images originating from the same B-mode video that intersect the pleural line. Following prior studies,
 215 we estimate the horizontal bounds of the pleural line using a previously trained object detection
 216 model [17] and use the top half of qualifying M-mode images, in decreasing order of total pixel
 217 intensity [11].

218 Table 2 gives the composition of the ParenchymalLUS dataset, along with the number of examples
 219 labelled for the AB and LS tasks.

		UNLABELLED	LABELLED		
			Train	Validation	Test
Total	Patients	5204	1540	330	329
	Videos	20 000	4123	858	936
	Images	4 611 063	927 889	191 437	208 648
AB Labels	Videos	—	2100/998	441 / 197	512/213
	Images	—	484 287 / 216 505	99 132 / 40 608	116 648 / 42 122
LS Labels	Videos	—	3169/477	631/103	707/96
	Images	—	727 205/96 771	146 322/23 218	166 753/21 911

Table 2: Breakdown of ParenchymalLUS at the video and image level. x / y indicates the number of negative and positive labelled examples available for each task, respectively. Video labels apply to each image within the video. Note that some videos were not labelled for both tasks.

220 **B Sample Weights**

221 For a positive pair of B-mode images occurring at times t_1 and t_2 or M-mode images occurring at
 222 positions x_1 and x_2 , the sample weight is calculated as follows:

$$w = \frac{\delta_t - |t_2 - t_1| + 1}{\delta_t + 1} \qquad w = \frac{\delta_x - |x_2 - x_1| + 1}{\delta_x + 1} \qquad (1)$$

¹Omitted to protect anonymity during the review process.

223 Sample weights were incorporated into each SSL objective trialled in this study. Accordingly, we
 224 modified the objective functions for SimCLR, Barlow Twins, and VICReg in order to weigh the
 225 contribution to the loss differently based on sample weights. To the authors’ knowledge, this study
 226 is the first to propose sample weighting schemes for the aforementioned self-supervised learning
 227 methods.

228 The SimCLR objective can be easily modified by multiplying L_i , the per-example NT-Xent loss for
 229 the i^{th} positive pair, by sample weight w_i .

$$\mathcal{L}_{\text{SimCLR}} = \frac{1}{N} \sum_{i=1}^N w_i L_i \quad (2)$$

230 For VICReg [20], the invariance term is weighted with w_i for each positive pair in a batch. The
 231 invariance term is then calculated as follows:

$$s(Z_1, Z_2) = \frac{1}{N} \sum_{i=1}^N w_i \|Z_{1_i} - Z_{2_j}\|_2^2 \quad (3)$$

232 where Z_1 and Z_2 are batches of predicted embeddings for corresponding positive pairs; that is, Z_{1_i}
 233 and Z_{2_i} correspond to one positive pair. The entire VICReg objective can then be calculated as

$$\mathcal{L}_{\text{VICReg}}(Z_1, Z_2) = \underbrace{\lambda s(Z_1, Z_2)}_{\text{Invariance term}} + \underbrace{\mu(v(Z_1) + v(Z_2))}_{\text{Variance term}} + \underbrace{\nu(c(Z_1) + c(Z_2))}_{\text{Covariance term}} \quad (4)$$

234 where λ , μ , and ν are weights for each term. Since frames were sampled uniformly at random,
 235 $\mathbb{E}[w] \simeq 0.5$. Accordingly, we doubled λ when pretraining VICReg with sample weights.

236 For the Barlow Twins objective, weighting was applied to each positive pair in the invariance term by
 237 computing the weighted normalized cross correlation matrix $\mathcal{C}_W \in \mathbb{R}^{D \times D}$ between the weighted-
 238 mean-centered normalized batches of embeddings, Z_1 and Z_2 . For a batch of embeddings Z , the
 239 calculation for the weighted mean \bar{Z} and standard deviations $\sigma(Z)$ across the batch dimension was
 240 performed as follows:

$$\bar{Z} = \frac{\sum_{i=1}^N w_i Z_i}{\sum_{i=1}^N w_i} \quad \sigma(Z) = \sqrt{\frac{\sum_{i=1}^N w_i (Z_i - \bar{Z})^2}{\sum_{i=1}^N w_i}} \quad (5)$$

$$\mathcal{C} = \frac{1}{N} \left(\frac{Z_1 - \bar{Z}_1}{\sigma(Z_1)} \right)^T \left(\frac{Z_2 - \bar{Z}_2}{\sigma(Z_2)} \right) \quad \sigma(Z) = \sqrt{\frac{\sum_{i=1}^N w_i (Z_i - \bar{Z})^2}{\sum_{i=1}^N w_i}} \quad (6)$$

241 where w_i is the sample weight for the i^{th} positive pair in the batch. The redundancy reduction term
 242 should still be calculated using the normalized cross correlation matrix \mathcal{C} , since its purpose is to
 243 decorrelate the embedding dimensions. In the original Barlow Twins, the normalized cross correlation
 244 matrix is employed for both terms. The Barlow Twins objective then becomes

$$\mathcal{L}_{\text{BT}} = \underbrace{\sum_{d=1}^D (1 - \mathcal{C}_{W_{d,d}})^2}_{\text{Invariance term}} + \underbrace{\lambda \sum_{d=1}^D \sum_{\substack{e=1 \\ e \neq d}}^D \mathcal{C}_{d,e}^2}_{\text{Redundancy reduction term}} \quad (7)$$

245 C Pretraining and Training Protocols

246 Unless otherwise stated, all feature extractors are initialized with ImageNet-pretrained weights.
247 Similar studies concentrating on medical imaging have observed that this practice improves down-
248 stream performance when compared to random initialization [11, 27]. Moreover, we designate fully
249 supervised classifiers initialized with ImageNet-pretrained weights as a baseline against which to
250 compare models pretrained with SSL.

251 Evaluation on POCUS follows a similar protocol employed in prior works [8, 10]. Feature extractors
252 with the ResNet18 architecture [28] are pretrained on the Butterfly dataset. Prior to training on
253 the POCUS dataset, a 3-node fully connected layer with softmax activation was appended to the
254 pretrained feature extractor. Five-fold cross validation is conducted with POCUS by fine-tuning
255 the final three layers of the pretrained feature extractor. Unlike prior works, we adopt the average
256 across-folds validation accuracy, instead of taking the accuracy of the combined set of validation set
257 predictions across folds. Presenting the results in this manner revealed the high variance of model
258 performance across folds, which may be due to the benchmark dataset’s small video sample size.

259 All experiments with ParenchymalLUS utilize the MobileNetV3-Small architecture as the feature
260 extractor, which outputs a 576-dimensional representation vector [29]. Feature extractors are pre-
261 trained on the union of the unlabelled videos and labelled training set videos in ParenchymalLUS.
262 Performance is assessed via test set classification metrics. Prior to training on the downstream task,
263 a single-node fully connected layer with sigmoid activation was appended to the pretrained feature
264 extractor. We report the performance of linear classifiers trained on the frozen feature extractor’s
265 representations, along with classifiers that are fine-tuned end-to-end.

266 For each joint embedding method, the projectors were multilayer perceptrons with two 768-node
267 layers, outputting 768-dimensional embeddings. Pretraining is conducted for 500 epochs using the
268 LARS optimizer [30] with a batch size of 384 and a learning rate schedule with warmup and cosine
269 decay as in [20].

270 B-mode and M-mode images were resized to 224×224 and 224×112 pixels respectively using
271 bilinear interpolation. The reason that the width of M-modes was standardized to a smaller value was
272 because the height of B-mode images often far exceeded the number of frames in a 3-second segment
273 of B-mode video. Since feature extractors were initialized with pretrained weights, pixel intensities
274 were mean-centered and normalized using the mean and standard deviation of the ImageNet dataset.

275 All IVPP pretraining runs were subjected to stochastic data augmentations after intra-video positive
276 pairs were sampled. Each image was subjected to the following sequence of stochastic transformations
277 for data augmentation:

- 278 1. Randomly located crop of a fraction of the image’s area in the range $[0.4, 1.0]$, followed by
279 resizing to the original image dimensions. For B-mode images, the width/height aspect ratio
280 was confined to the range $[0.8, 1.25]$, while M-mode crops were confined to $[0.4, 0.6]$.
- 281 2. With probability 0.5, horizontal reflection.
- 282 3. With probability 0.5, random brightness change in the range $[-0.25, 0.25]$.
- 283 4. With probability 0.5, random contrast change in the range $[-0.25, 0.25]$.
- 284 5. With probability 0.25, random Gaussian blur with a kernel size of 5 pixels and a standard
285 deviation uniformly sampled from the range $[0.1, 2.0]$.

286 Training runs for the COVID and AB tasks with B-mode images also utilized this data augmentation
287 pipeline. For LS training runs with M-modes, the minimum allowable crop area was increased to
288 95% of the image’s area to ensure the pleural line was almost always visible.

289 Source code is available at [redacted]².

²Public GitHub repository URL redacted to preserve anonymity.

290 **D Additional Performance Details**

291 **Linear Classification Performance:** In addition to fine-tuning experiments, linear classifiers were
 292 trained using the feature vectors outputted by the pretrained models. Figure 6 summarizes the results
 293 obtained for the AB and LS tasks for each of the hyperparameter combinations studied.

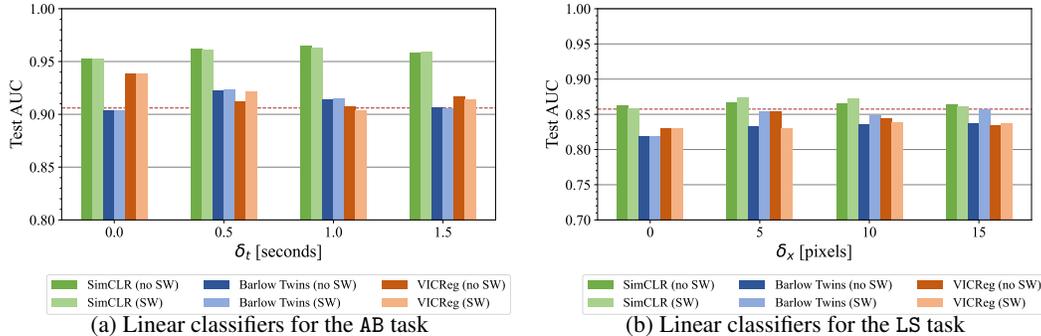


Figure 6: ParenchymalLUS test set AUC for linear classifiers trained on the AB and LS binary classification tasks and pretrained with a variety of intra-video positive pair thresholds, with and without sample weights (SW). The dashed line indicates initialization with ImageNet-pretrained weights.

294 **Label Efficiency Experiments:** Inspection of the central moments and boxplots from each distribu-
 295 tion (Figure 5) indicated that the normality and equal variance assumptions for ANOVA were not
 296 violated. For each pretraining method, a two-way repeated-measures analysis of variance (ANOVA)
 297 was performed to determine whether the mean test AUC scores across values of δ and sample weight
 298 usage were different. The independent variables were δ and the presence of sample weights, while
 299 the dependent variable was test AUC. Whenever the null hypothesis of the ANOVA was rejected,
 300 post-hoc paired t -tests were performed to compare the following:

- 301 • Pretraining with nonzero δ against standard positive pair selection ($\delta = 0$)
- 302 • For the same nonzero δ value, sample weights against no sample weights

303 For each group of post-hoc tests, the Bonferroni correction was applied to establish a family-wise
 304 error rate of $\alpha = 0.05$. To ensure that each training subset was independent, we split the dataset
 305 by anonymous patient identifier. This was a necessary step because intra-video images are highly
 306 correlated, along with videos from the same patient. As a result, the task became substantially
 307 more difficult than naively sampling 5% of training images because the volume *and* heterogeneity
 308 of training examples was reduced by training on a small fraction of examples from a small set of
 309 patients.

310 Table 3 gives the mean and standard deviation of each set of trials, for each hyperparameter combi-
 311 nation. For each task and each pretraining method, the ANOVA revealed significant interaction
 312 effects ($p \leq 0.05$). Accordingly, all intended post-hoc t -tests were performed to ascertain (1)
 313 which combinations of hyperparameters were different from the baseline setting of augmenting the
 314 same frame twice ($\delta = 0$) and (2) values of δ where the addition of sample weights changes the
 315 outcome. First, we note that SimCLR was the only pretraining method that consistently outperformed
 316 full supervision with ImageNet-pretrained weights. Barlow Twins and VICReg pretraining – both
 317 non-contrastive methods – resulted in worse performance.

318 For the AB task, no combination of intra-video positive pairs or sample weights resulted in statistically
 319 significant improvements compared to dual distortion of the same image ($\delta_t = 0$). For Barlow Twins
 320 and VICReg, several nonzero δ_t resulted in significantly worse mean test AUC. Sample weights
 321 consistently made a difference in Barlow Twins across δ_t values, but only improved mean test AUC
 322 for $\delta_t = 1$ and $\delta_t = 1.5$.

323 Different trends were observed for the LS task. SimCLR with $\delta_x = 5$ and no sample weights improved
 324 mean test AUC compared to the baseline where $\delta_x = 0$. No other combination of hyperparameters
 325 resulted in a significant improvement. For Barlow Twins, multiple IVPP hyperparameter combinations

326 resulted in improved mean test AUC over the baseline. No IVPP hyperparameter combinations
 327 significantly improved the performance of VICReg.

Pretrain Method	AB			LS				
	δ_t	SW	Mean (std) test AUC	δ_x	SW	Mean (std) test AUC		
SimCLR	0		0.938 (0.007)	0		0.812 (0.037)		
	0.5		0.931 (0.010) *	5		0.824 (0.030) *		
	0.5		0.936 (0.007) †	5		0.820 (0.033)		
	1		0.934 (0.011)	10		0.815 (0.035)		
	1		0.933 (0.011)	10		0.816 (0.037)		
	1.5		0.936 (0.013)	15		0.819 (0.034)		
	1.5		0.932 (0.012)	15		0.798 (0.039) *†		
Barlow Twins	0		0.914 (0.014)	0		0.693 (0.044)		
	0.5		0.914 (0.010) *	5		0.694 (0.040)		
	0.5		0.883 (0.017) *†	5		0.780 (0.040) *†		
	1		0.877 (0.022) *	10		0.705 (0.051)		
	1		0.891 (0.018) *†	10		0.706 (0.066)		
	1.5		0.870 (0.024) *	15		0.769 (0.037) *		
	1.5		0.892 (0.015) *†	15		0.707 (0.071) †		
VICReg	0		0.917 (0.011)	0		0.690 (0.042)		
	0.5		0.879 (0.024) *	5		0.675 (0.036)		
	0.5		0.879 (0.021) *	5		0.679 (0.038)		
	1		0.872 (0.023) *	10		0.680 (0.039)		
	1		0.876 (0.024) *	10		0.675 (0.040)		
	1.5		0.860 (0.026) *	15		0.710 (0.036)		
	1.5		0.870 (0.021) *†	15		0.685 (0.039) †		
None (ImageNet-pretrained)			0.896 (0.017)			0.783 (0.028)		
None (random initialization)			0.774 (0.051)			0.507 (0.022)		

* Significantly different ($p < 0.05$) than baseline for the pretraining method where $\delta = 0$

† Significantly different ($p < 0.05$) for particular δ when sample weights are applied, compared to no sample weight

Table 3: ParenchymalLUS test AUC for the the AB and LS tasks when trained using examples from 5% of the patients in the training set. Twenty trials were performed for each pretraining method, value of δ , with and without sample weights (SW). Mean and standard deviation of the test AUC across trials are reported for each condition.