
Anomaly Detection In The Wild: Can SSL Handle Strong Distribution Imbalances?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Accurate anomaly detection is critical in vision-based infrastructure inspection,
2 where it helps prevent costly failures and enhances safety. Self-Supervised Learn-
3 ing (SSL) offers a promising approach by learning robust representations from
4 unlabeled data. However, its application in anomaly detection remains underex-
5 plored. This paper addresses this gap by providing a comprehensive evaluation of
6 SSL methods for real-world anomaly detection, focusing on sewer infrastructure.
7 Using the Sewer-ML dataset, we evaluate lightweight models such as ViT-Tiny
8 and ResNet-18 across SSL frameworks, including BYOL, Barlow Twins, SimCLR,
9 DINO, and MAE, under varying class imbalance levels. Through 250 experiments,
10 we rigorously assess the performance of these SSL methods to ensure a robust
11 and comprehensive evaluation. Our findings highlight the superiority of joint-
12 embedding methods like SimCLR and Barlow Twins over reconstruction-based
13 approaches such as MAE, which struggle to maintain performance under class
14 imbalance. Furthermore, we find that the SSL model choice is more critical than
15 the backbone architecture. Additionally, we emphasize the need for better label-
16 free assessments of SSL representations, as current methods like RankMe fail to
17 adequately evaluate representation quality, making cross-validation without labels
18 infeasible. Despite the remaining performance gap between SSL and supervised
19 models, these findings highlight the potential of SSL to enhance anomaly detection,
20 paving the way for further research in this underexplored area of SSL applications.

21 1 Introduction

22 Self-Supervised Learning (SSL) is a machine learning paradigm where models are trained on
23 unlabeled data by creating surrogate labels through pretext tasks that exploit inherent structures or
24 patterns within the data. As a result, this approach enables learning meaningful representations that
25 can be applied to various downstream tasks without the need for explicit manual labeling [4]. Because
26 of this, SSL is particularly advantageous for semi-supervised anomaly detection problems where
27 obtaining labeled data is costly, labor-intensive, impossible, or undesirable [1]. Despite these benefits,
28 anomaly detection is frequently underrepresented in SSL research, with recent large-scale ablation
29 studies often neglecting its inclusion in benchmarking [10, 2]. In fact, common benchmarks such
30 as ImageNet [7] and CIFAR [16], are object-centric and do not accurately reflect the complexity of
31 real-world environments, where images are more diverse and less structured [11].

32 Closer to the anomaly detection scenario, recent works have started to stress-test SSL on more
33 realistic scenarios with uncurated data. Albeit still employing a classification task evaluation, it
34 has been shown that SSL may be sensitive to the quality of the data and in particular to balance
35 representations of the features to be learned [18, 3]. Clearly, such warning seems to go against the

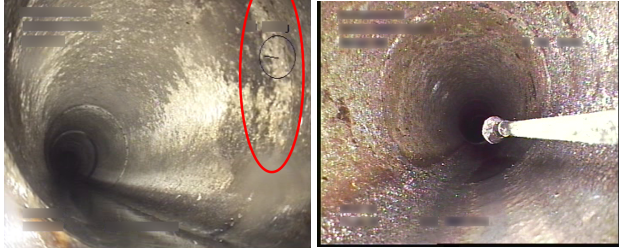


Figure 1: Samples from the Sewer-ML dataset [13]. The red circle highlights a defect associated with lateral reinstatement cuts, where improper cutting or misalignment can cause issues such as blockages, leaks, or structural weakness in the sewer system (for further details on defect types and image samples used for training and validation, refer to Appendix E).

36 ability of SSL to solve anomaly detection as such scenario assumes by nature that the model can
 37 learn the anomaly features from very few samples.

38 This gap in evaluation coupled with those recent sensitivity studies of SSL beg the following question:

39 *Can SSL successfully learn representations on highly imbalance datasets and capture salient features*
 40 *to solve anomaly detection tasks?*

41 To scientifically approach that question, we propose the use of the Sewer-ML dataset [13], a sewerage
 42 infrastructure dataset that contains 1.3 million images captured from video inspections, featuring
 43 17 different defect classes. Figure 1 presents examples of defect and non-defect images taken from
 44 Sewer-ML. To determine the robustness of self-supervised learning in handling class imbalances,
 45 we conducted 250 binary classification experiments where we systematically vary the proportion of
 46 defect samples in the train and validation datasets as 1%, 2%, 5%, and 15%.

47 We summarize our key findings below:

- 48 • **Superiority of joint-embedding methods over reconstruction methods:** we observe that
 49 MAE has difficulties to maintain performance with varying class imbalance as opposed to
 50 methods like SimCLR or BarlowTwins that compete with supervised baselines.
- 51 • **Impact of SSL methodology over backbone architecture:** the choice of backbone archi-
 52 tecture (Resnet or Transformer) is not critical compared to the choice of SSL method.
- 53 • **Need for better label-free assessment of SSL representations:** our findings indicate
 54 that methods such as RankMe fail to assess the richness of SSL representations making
 55 cross-validation without labels currently impossible on such task and dataset.

56 The code for this research will be open-sourced and is available at [Anomaly-Detection-In-The](#)
 57 [Wild_code.zip](#).

58 2 Controlled evaluation of self-supervised anomaly detection

59 We will first describe our methodology in detail and then provide results and discussions at the end of
 60 this section.

61 2.1 Methodology

62 Prior to conducting the SSL ablation study, an initial hyperparameter search was performed to optimize
 63 the data augmentation pipeline for the dataset. The focus was placed on tuning image resolution and
 64 augmentation settings to achieve a balance between model performance and computational efficiency.
 65 Detailed results and configurations are provided in Appendix A.

66 **Methods and models.** Our study conducts an ablation analysis on anomaly detection using self-
 67 supervised learning methods, with a particular emphasis on their robustness to distribution imbalances.
 68 We primarily focus on joint-embedding architectures—specifically Barlow Twins, SimCLR, BYOL,
 69 and DINO [19, 6, 12, 5]—which aim to learn an embedding space by aligning representations of
 70 different augmented views of the same input while avoiding collapse. To provide a comparative
 71 perspective, we also evaluate Masked Autoencoders (MAE) [14], a self-supervised approach that
 72 reconstructs missing parts of the input data. For our backbone architectures, we use lightweight
 73 models such as ViT-Tiny [8] and ResNet-18 [15]. Please refer to Appendix B to inquire about specific
 74 training hyperparameters used for each method.

75 **Controlled dataset imbalance.** To determine the robustness of self-supervised learning in handling
76 class imbalances, we conduct binary classification experiments with the Sewer-ML dataset [13], where
77 we systematically vary the proportion of defect samples in the dataset. Specifically, we assess the
78 performance of these methods under different levels of class imbalance—defect sample proportions
79 of 1%, 2%, 5%, and 15%. As a binary setting, Sewer-ML has an approximately balanced distribution
80 by design—45% and 47% defect proportion on the train and validation sets, respectively. We seek to
81 assess how well self-supervised learning approaches can learn from imbalanced distributions and to
82 evaluate their effectiveness on datasets with varying levels of class imbalance. We apply each of the
83 specified defect proportions to both the training and testing datasets, allowing us to evaluate model
84 performance across all combinations of class imbalance scenarios.

85 **Evaluation metrics.** In addition to reporting performance based solely on the F1 score, we also
86 evaluate our models using the metrics proposed by the benchmark. This includes a weighted F2
87 metric ($F2_{CIW}$) for identifying defects and a conventional F1 score ($F1_{Normal}$) for non-defect
88 classifications [13]. The F2 metric’s weights are assigned to each defect category according to their
89 economic significance. Additionally, they encourage the use of an F2 score to give greater emphasis
90 to recall over precision, recognizing that overlooking a defect incurs a higher economic cost than the
91 occurrence of a false positive.

92 2.2 Results

93 We present in Table 1 the performance of various SSL models compared to their supervised coun-
94 terparts. Training was conducted on imbalanced data, while validation adhered to a balanced set,
95 aligning with benchmark standards. Although previous studies have noted that CNN-based architec-
96 tures like ResNet-18 often outperform vision transformers (ViTs) [10], our results do not indicate a
97 significant difference in performance between these backbone architectures.

98 **Robustness of joint-embedding methods.** When examining individual SSL methods, SimCLR
99 achieves the best overall results especially when paired with ResNet-18 and defect proportion is
100 higher than 5%. When working with stronger distribution imbalances, particularly with 1%, BYOL
101 and Resnet-18 have the best results. This might be due to the training dynamics behind these methods,
102 in situations with moderate imbalance (like 5%), the contrastive approach of SimCLR can still
103 adequately separate the minority class from the majority, leveraging the discriminative power of the
104 contrastive loss. This can be counter-productive when working with extremely imbalance levels.
105 BYOL avoids this issue by not requiring explicit negative sampling, allowing it to maintain more
106 consistent performance even when the class imbalance becomes more extreme. Thus, while SimCLR
107 thrives with moderate imbalance, BYOL proves more resilient in handling extreme class disparities.
108 Another interesting insight is Barlow Twins’ competitive performance under extreme imbalance. Its
109 use of redundancy reduction loss, which both maximizes similarity between augmented views and
110 decorrelates the learned representations, may help avoid the pitfalls of overfitting to the dominant
111 class. DINO, however, shows a unique trend, performing poorly in 1% and 2% settings with ResNet-
112 18, but not with ViT-Tiny. This could be attributed to DINO’s focus on global feature learning through
113 knowledge distillation, which aligns better with ViT’s global attention approach.

114 **Failure of reconstruction-based methods.** Finally, the Masked Autoencoder (MAE) delivers the
115 weakest performance across all scenarios, especially at higher imbalance levels (e.g., 45%). MAE’s
116 reliance on reconstructing multiple classes might introduce noise that hampers its ability to generalize
117 under severe imbalance. In highly imbalanced data, MAE could face challenges in differentiating
118 between common and rare classes, possibly due to the reconstruction bias that favors the majority
119 class. This might also be tied to the model’s neural capacity, which can struggle to produce robust
120 representations when not sufficiently overparameterized.

121 Full tables with training and validation results are in Appendix C.2. Detailed analysis of each model’s
122 performance on imbalanced test distributions is provided in Appendix C.1. SSL monitoring metrics
123 including RankMe [9] and the mean and standard deviation of features and embeddings for DINO,
124 BYOL, and SimCLR—are discussed in Appendix D.

125 **Performance Trends on Imbalanced Validation Sets.** When evaluating the models on increasingly
126 imbalanced validation sets (Figure 2), there is a clear and consistent decline in performance across
127 most methods, as indicated by the F1 scores. As the validation set imbalance becomes more severe
128 (e.g., at 1% and 2% imbalance), no method is able to sustain strong performance, particularly in

Table 1: Validation metrics for each method and architecture across different imbalance levels. The percentages represent the imbalance levels applied during training, while the validation data proportion remained unchanged.

		RESNET-18					ViT-TINY				
		1%	2%	5%	15%	45%	1%	2%	5%	15%	45%
F1 \uparrow	<i>Supervised</i>	0.831	0.847	0.865	0.882	0.898	0.768	0.785	0.831	0.858	0.878
	<i>BYOL</i>	0.764	0.768	0.774	0.783	0.792	0.676	0.738	0.752	0.757	0.730
	<i>Barlow Twins</i>	0.768	0.775	0.789	0.796	0.803	0.682	0.738	0.747	0.769	0.772
	<i>DINO</i>	0.049	0.032	0.673	0.719	0.727	0.691	0.683	0.553	0.704	0.719
	<i>MAE</i>	-	-	-	-	-	0.513	0.636	0.570	0.575	0.024
	<i>SimCLR</i>	0.740	0.778	0.796	0.807	0.814	0.762	0.768	0.771	0.778	0.786
F2 $_{CIW}$ \uparrow	<i>Supervised</i>	0.838	0.849	0.880	0.904	0.879	0.837	0.812	0.831	0.870	0.881
	<i>BYOL</i>	0.837	0.844	0.819	0.825	0.785	0.680	0.768	0.760	0.781	0.725
	<i>Barlow Twins</i>	0.817	0.805	0.817	0.822	0.780	0.704	0.749	0.786	0.784	0.758
	<i>DINO</i>	0.566	0.402	0.712	0.777	0.752	0.733	0.736	0.690	0.724	0.695
	<i>MAE</i>	-	-	-	-	-	0.629	0.640	0.704	0.702	0.009
	<i>SimCLR</i>	0.775	0.813	0.826	0.822	0.813	0.766	0.795	0.798	0.801	0.788
F1 $_{Normal}$ \uparrow	<i>Supervised</i>	0.850	0.861	0.879	0.894	0.908	0.774	0.794	0.844	0.869	0.890
	<i>BYOL</i>	0.758	0.759	0.783	0.789	0.818	0.719	0.757	0.778	0.771	0.781
	<i>Barlow Twins</i>	0.782	0.796	0.799	0.807	0.828	0.716	0.760	0.766	0.785	0.806
	<i>DINO</i>	0.689	0.686	0.700	0.730	0.763	0.655	0.688	0.419	0.727	0.759
	<i>MAE</i>	-	-	-	-	-	0.471	0.350	0.553	0.389	0.677
	<i>SimCLR</i>	0.752	0.793	0.806	0.818	0.836	0.782	0.781	0.786	0.785	0.810

129 terms of accurately identifying the minority class. On the other hand, the F2 score remains relatively
 130 stable, likely due to the use of a weighted loss function based on defect proportions.

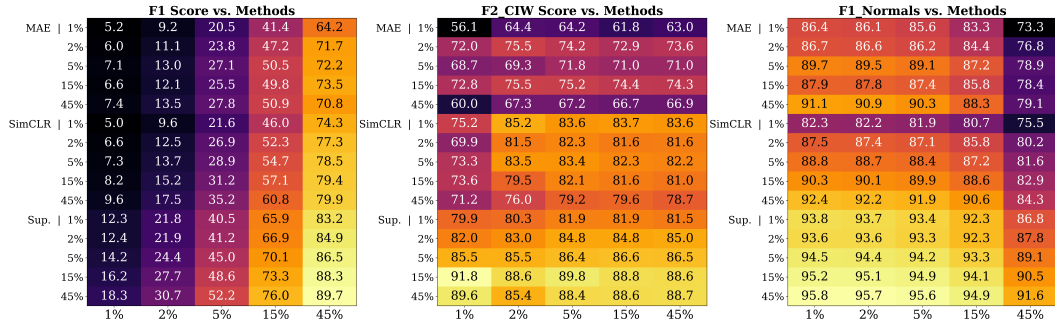


Figure 2: ResNet-18 validation performance heatmaps across imbalance levels. The x-axis represents the imbalance levels in the validation set, while the y-axis indicates the method and the imbalance level used during training.

131 3 Conclusions and discussion

132 Our study indicates that self-supervised learning (SSL) is effective for anomaly detection and remains
 133 robust even when facing significant distribution imbalances. We also find that the choice of backbone
 134 architecture is not the most critical factor in model performance, as neither ViT-Tiny nor ResNet-18
 135 consistently outperforms the other across all cases. In contrast, the selection of the SSL methodology
 136 significantly impacts performance, with substantial variations observed among different SSL model
 137 families. Therefore, for practitioners, choosing the appropriate SSL family is more crucial than
 138 selecting a specific backbone architecture. Furthermore, there is a pressing need to accurately measure
 139 the quality of representations produced by SSL models. Our findings indicate that the RankMe metric
 140 is ineffective for this purpose; it aims to assess the richness of representations but fails to correlate
 141 with actual performance. As shown in Appendix D, there is no correlation between performance and
 142 RankMe metrics, underscoring the necessity for better methods of evaluating representations quality.

References

- 143
- 144 [1] Leman Akoglu and Jaemin Yoo. Self-supervision for tackling unsupervised anomaly detection:
145 Pitfalls and opportunities. In *2023 IEEE International Conference on Big Data (BigData)*,
146 pages 1047–1051, Los Alamitos, CA, USA, dec 2023. IEEE Computer Society.
- 147 [2] Haider Al-Tahan, Quentin Garrido, Randall Balestriero, Diane Bouchacourt, Caner Hazirbas,
148 and Mark Ibrahim. UniBench: Visual Reasoning Requires Rethinking Vision-Language Beyond
149 Scaling. (arXiv:2408.04810), August 2024. arXiv:2408.04810 [cs].
- 150 [3] Mido Assran, Randall Balestriero, Quentin Duval, Florian Bordes, Ishan Misra, Piotr Bo-
151 janowski, Pascal Vincent, Michael Rabbat, and Nicolas Ballas. The Hidden Uniform Cluster
152 Prior in Self-Supervised Learning. In *The Eleventh International Conference on Learning
153 Representations, 2023*.
- 154 [4] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Gold-
155 stein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild,
156 Andrew Gordon Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed
157 Pirsiavash, Yann LeCun, and Micah Goldblum. A Cookbook of Self-Supervised Learning.
158 (arXiv:2304.12210), June 2023.
- 159 [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jegou, Julien Mairal, Piotr Bojanowski,
160 and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. In *2021
161 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9630–9640, 2021.
- 162 [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework
163 for Contrastive Learning of Visual Representations. In Hal Daumé III and Aarti Singh, edi-
164 tors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of
165 *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020.
- 166 [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-
167 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern
168 Recognition*, pages 248–255, 2009.
- 169 [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
170 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
171 Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for
172 Image Recognition at Scale. In *Proceedings of the 2021 International Conference on Learning
173 Representations (ICLR)*, 2021.
- 174 [9] Quentin Garrido, Randall Balestriero, Laurent Najman, and Yann LeCun. RankMe: Assessing
175 the Downstream Performance of Pretrained Self-Supervised Representations by Their Rank. In
176 *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org,
177 2023.
- 178 [10] Micah Goldblum, Hossein Souri, Renkun Ni, Manli Shu, Viraj Uday Prabhu, Gowthami
179 Somepalli, Prithvijit Chattopadhyay, Mark Ibrahim, Adrien Bardes, Judy Hoffman, Rama
180 Chellappa, Andrew Gordon Wilson, and Tom Goldstein. Battle of the Backbones: A Large-
181 Scale Comparison of Pretrained Models across Computer Vision Tasks. In *Thirty-seventh
182 Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2023*.
- 183 [11] Priya Goyal, Mathilde Caron, Benjamin Lefaudeaux, Min Xu, Pengchao Wang, Vivek Pai,
184 Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, and Piotr Bojanowski. Self-
185 supervised Pretraining of Visual Features in the Wild. (arXiv:2103.01988), March 2021.
186 arXiv:2103.01988 [cs].
- 187 [12] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena
188 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi
189 Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap Your Own
190 Latent: A New Approach to Self-Supervised Learning. *Proceedings of Advances in Neural
191 Information Processing Systems (NeurIPS)*, 2020.

- 192 [13] Joakim Bruslund Haurum and Thomas B. Moeslund. Sewer-ML: A Multi-Label Sewer De-
193 fect Classification Dataset and Benchmark. In *Proceedings of the IEEE/CVF Conference on*
194 *Computer Vision and Pattern Recognition (CVPR)*, pages 13456–13467, June 2021.
- 195 [14] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
196 Autoencoders Are Scalable Vision Learners. In *Proceedings of the IEEE/CVF Conference on*
197 *Computer Vision and Pattern Recognition (CVPR)*, pages 16000–16009, June 2022.
- 198 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for
199 Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
200 *Recognition (CVPR)*, pages 770–778, 2016.
- 201 [16] Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical report,
202 University of Toronto, 2009.
- 203 [17] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
204 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
205 Köpf, Edward Z. Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
206 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style,
207 High-Performance Deep Learning Library, 2019.
- 208 [18] Huy V. Vo, Vasil Khalidov, Timothée Darcet, Théo Moutakanni, Nikita Smetanin, Marc
209 Szafraniec, Hugo Touvron, camille couprie, Maxime Oquab, Armand Joulin, Herve Jegou,
210 Patrick Labatut, and Piotr Bojanowski. Automatic Data Curation for Self-Supervised Learning:
211 A Clustering-Based Approach. *Transactions on Machine Learning Research*, 2024.
- 212 [19] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stephane Deny. Barlow Twins: Self-
213 Supervised Learning via Redundancy Reduction. In Marina Meila and Tong Zhang, editors,
214 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Pro-*
215 *ceedings of Machine Learning Research*, pages 12310–12320. PMLR, 18–24 Jul 2021.

216 **A Images resolution and data augmentations ablation**

217 For the resolution, we trained both architectures in a supervised fashion, using no data augmentations
 218 other than resizing. As seen in Table 2 the resolution ablation study reveals a clear trend where
 219 increasing the image size leads to better model performance, as measured by validation F1 scores.
 220 However, this improvement comes at the cost of increased training time. Notably, the best trade-off
 221 between performance and runtime was achieved at a resolution of 224x224.

Table 2: Best validation F1 scores for different image sizes and architectures.

	64	128	224	384	512
ViT-Tiny	0.866	0.885	0.894	0.897	0.90
ResNet-18	0.877	0.892	0.901	0.907	0.909
Time/Epoch	28m	31m	37m	52m	69m

222 For the data augmentations, we explored how changes in color jitter and random cropping impacted
 223 model performance. To maintain computational efficiency, we used 25% of the dataset, running
 224 experiments with the ResNet-18 architecture. We adjusted the ColorJitter parameters using a
 225 variable t_val that ranged from 0.1 to 0.8. Specifically, brightness, contrast, and saturation were
 226 each set to t_val , while hue was set to half of that value ($hue = t_val / 2$). Furthermore, the
 227 minimum scale for RandomResizedCrop was varied from 0.08 to 0.71. As illustrated in Figure 3,
 228 the best performance, measured by F1 Score, was obtained when t_val was around 0.1 and 0.45 and
 229 when min_scale was around 0.395 and 0.605. Due to the small performance difference, and looking
 230 to introduce stronger augmentations when using SSL methodologies, we decided to use min_scale
 231 $= 0.395$ and $t_val = 0.275$ for further ablations.

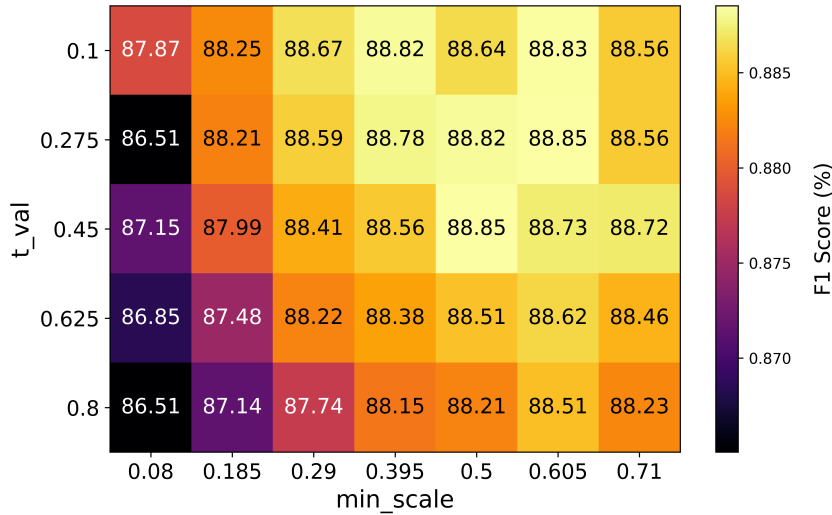


Figure 3: Val F1 Score by min_scale and t_val . This heatmap shows the performance variation in terms of F1 score, demonstrating the interaction between these two hyperparameters.

232 **B SSL configurations**

233 Training was conducted using PyTorch 2.0.2 [17] on a SLURM cluster equipped with NVIDIA L4
 234 GPUs. In total, we performed 250 experiments, comprising 74 training runs and 184 validation
 235 tests. The experiments were conducted over 45 epochs varying the optimizer and its corresponding
 236 parameters depending on the method at hand. A linear warmup was applied for the first 10 epochs,
 237 followed by a cosine scheduler with no restarts. The base and final decay rates (τ) were 0.996 and
 238 0.999, respectively, with a minimum learning rate of 1×10^{-6} .

239 During pre-training, a linear probe was trained, and the reported results reflect its final performance.
240 For all experiments, the linear probe was trained using a learning rate of 0.001, with a momentum of
241 0.9 and no weight decay applied. However, for MAE, a reduced learning rate of 0.0001 was necessary
242 to prevent numerical issues caused by 16-bit mixed precision. Similarly, DINO was trained using
243 32-bit precision due to instability on imbalanced cases.

244 Moreover, due to the imbalanced nature of the problem, we use a binary cross entropy loss with
245 positive weights. The positive weights are assigned using the inverse of the defect proportion, a
246 method commonly referred to as inverse class frequency weighting. Finally, every SSL method was
247 trained on two augmented views of the same image and no multicrop was used.

248 To ensure reproducibility, the remainder of the section contains the specific configurations adopted
249 for each method.

250 B.1 BYOL

```
251 optimization:  
252   optimizer: lars  
253   batch_size: 256  
254   lr: 0.2  
255   momentum: 0.9  
256   weight_decay: 1.5e-06  
257   exclude_bias_and_norm: true  
258   warmup_start_lr: 0.1  
259 ssl_settings:  
260   proj_hidden_dim: 2048  
261   pred_hidden_dim: 2048  
262   proj_output_dim: 256  
263   normalize_projector: false
```

264 B.2 DINO

```
265 optimization:  
266   optimizer: adamw  
267   batch_size: 128  
268   lr: 0.0005  
269   momentum: 0.9  
270   weight_decay: 0.0001  
271   warmup_start_lr: 3.0e-05  
272 ssl_settings:  
273   proj_hidden_dim: 2048  
274   proj_output_dim: 256  
275   num_prototypes: 8192  
276   clip_grad: true  
277   freeze_prototyper: 1  
278   use_bn_in_head: false  
279   norm_prototyper: false  
280   student_temperature: 0.1  
281   warmup_teacher_temperature: 0.04  
282   teacher_temperature: 0.07  
283   warmup_temperature_epochs: 15
```

284 B.3 MAE

```
285 optimization:  
286   optimizer: adamw  
287   batch_size: 256  
288   lr: 0.001  
289   momentum: 0.9  
290   weight_decay: 0.05  
291   warmup_start_lr: 0.0005  
292 ssl_settings:
```



```
293 mask_ratio: 0.75
294 norm_pix_loss: true
295 decoder_embed_dim: 192
296 decoder_depth: 12
297 decoder_num_heads: 12
```

298 **B.4 SimCLR**

```
299 optimization:
300   optimizer: lars
301   batch_size: 256
302   lr: 0.3
303   momentum: 0.9
304   weight_decay: 1.0e-06
305   warmup_start_lr: 0.15
306 ssl_settings:
307   proj_hidden_dim: 2048
308   proj_output_dim: 256
309   temperature: 0.2
```

310 **B.5 Barlow Twins**

```
311 optimization:
312   optimizer: lars
313   batch_size: 256
314   lr: 0.2
315   momentum: 0.9
316   weight_decay: 1.5e-06
317   warmup_start_lr: 0.1
318 ssl_settings:
319   proj_hidden_dim: 2048
320   proj_output_dim: 256
321   lambda: 0.0051
322   scale_loss: 0.024
```

323 **B.6 Supervised**

```
324 optimization:
325   optimizer: adamw
326   batch_size: 256
327   lr: 0.0005
328   momentum: 0.9
329   weight_decay: 0.0001
330   warmup_start_lr: 3.0e-05
```

331 **C Complementary tables**

332 When evaluating SSL methods trained and tested on various class imbalances, several important
333 observations emerged (refer to Figure 4). DINO, when trained on low defect proportions (1% and
334 2%), struggled significantly in identifying defects, despite performing well at detecting non-defects,
335 suggesting a limitation in handling extreme class imbalances. Additionally, the supervised baseline
336 consistently outperformed all SSL methods across every class imbalance, with the performance gap
337 widening as the class imbalance in the validation set increased. This highlights the need for further
338 exploration of SSL methods in real-world anomaly detection scenarios. Interestingly, Barlow Twins
339 demonstrated the strongest resilience in highly imbalanced settings, particularly when both training
340 and evaluation involved significant class imbalances. Specifically, Barlow Twins variants trained
341 on 1% and 2% defect proportions achieved the best performance among SSL methods under these
342 conditions, underscoring its robustness in scenarios with severe class imbalances.

343 When evaluating precision across varying validation imbalances, a degradation trend is evident as class
 344 imbalance increases (refer to Figure 5). Notably, BYOL, when trained on a nearly balanced setting
 345 (45% defect proportion), shows a larger performance gap compared to its 15% counterpart, indicating
 346 the potential influence of training balance. Meanwhile, the recall score remains consistently strong
 347 across all validation imbalance proportions, likely due to the weighted loss choice, demonstrating
 348 solid recall performance even under significant class imbalances.

349 C.1 Imbalanced validation setting

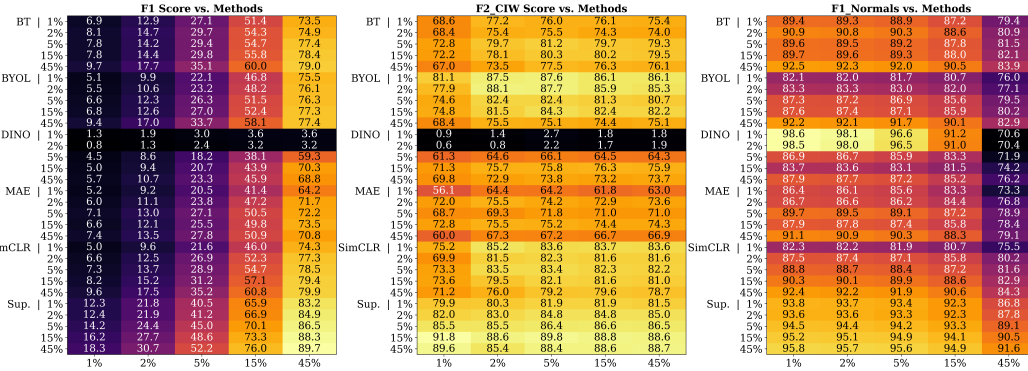


Figure 4: ResNet-18’s $F1$, $F2_{CIW}$, and $F1_{Normal}$ validation score heatmaps across imbalance levels. The x-axis represents the imbalance levels in the validation set, while the y-axis indicates the method and the imbalance level used during training.

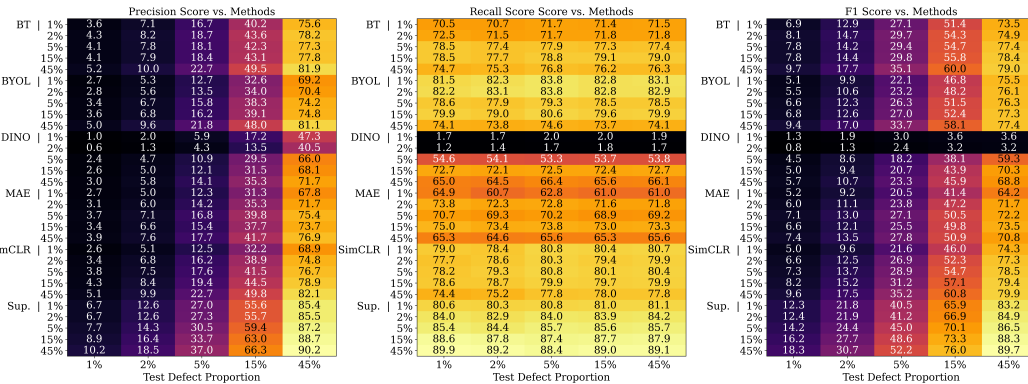


Figure 5: ResNet-18’s precision, recall, and $f1$ validation score heatmaps across imbalance levels. The x-axis represents the imbalance levels in the validation set, while the y-axis indicates the method and the imbalance level used during training.

350 C.2 Balanced validation setting

351 Overall, the training performance across the models show a clear linear degradation as the proportion
 352 of defects in the dataset decreases. This decline is primarily reflected in the precision scores, which
 353 drop significantly at lower defect proportions, such as 1% and 2%. The decreasing precision indicates
 354 that as fewer defects are present in the training set, the models fail to generalize properly, as the
 355 limited number of defect examples hinders their ability to accurately distinguish between defective
 356 and non-defective instances, leading to an increase in false positives. However, this trend is less
 357 evident in the validation set, where the precision remains relatively stable. This is due to the balanced
 358 nature of the validation set. On the other hand, recall remains fairly consistent across both the training
 359 and validation sets, even as the defect proportion decreases. This behavior can be explained by
 360 the use of a weighted loss function, which adjusts weights based on defect proportion. As a result, the
 361 model becomes biased towards positive predictions, preserving recall but impacting precision.

Table 3: Comprehensive performance metrics for ResNet-18 across both training and validation phases. The percentages represent the imbalance levels applied during training, while the validation default distribution remained unchanged.

		TRAIN					VAL				
		1%	2%	5%	15%	45%	1%	2%	5%	15%	45%
PRECISION \uparrow	<i>Supervised</i>	0.055	0.114	0.282	0.614	0.918	0.834	0.840	0.861	0.877	0.893
	<i>BYOL</i>	0.027	0.053	0.143	0.372	0.782	0.715	0.715	0.747	0.750	0.801
	<i>Barlow Twins</i>	0.032	0.067	0.167	0.410	0.805	0.749	0.769	0.764	0.772	0.811
	<i>DINO</i>	0.024	0.043	0.107	0.309	0.719	0.571	0.422	0.664	0.693	0.738
	<i>SimCLR</i>	0.028	0.066	0.168	0.433	0.808	0.715	0.762	0.770	0.786	0.817
RECALL \uparrow	<i>Supervised</i>	0.893	0.895	0.910	0.928	0.927	0.827	0.854	0.869	0.887	0.903
	<i>BYOL</i>	0.805	0.822	0.807	0.811	0.777	0.820	0.829	0.802	0.820	0.784
	<i>Barlow Twins</i>	0.771	0.782	0.791	0.802	0.777	0.787	0.781	0.815	0.821	0.796
	<i>DINO</i>	0.663	0.670	0.669	0.733	0.711	0.026	0.017	0.683	0.747	0.716
	<i>SimCLR</i>	0.748	0.794	0.808	0.808	0.793	0.766	0.795	0.824	0.829	0.812
F1 \uparrow	<i>Supervised</i>	0.104	0.202	0.431	0.739	0.922	0.831	0.847	0.865	0.882	0.898
	<i>BYOL</i>	0.052	0.099	0.243	0.510	0.779	0.764	0.768	0.774	0.783	0.792
	<i>Barlow Twins</i>	0.061	0.123	0.276	0.543	0.791	0.768	0.775	0.789	0.796	0.803
	<i>DINO</i>	0.046	0.080	0.185	0.434	0.715	0.049	0.032	0.673	0.719	0.727
	<i>SimCLR</i>	0.054	0.123	0.279	0.564	0.800	0.740	0.778	0.796	0.807	0.814
F2 _{CIW} \uparrow	<i>Supervised</i>	0.868	0.877	0.886	0.904	0.900	0.838	0.849	0.880	0.904	0.879
	<i>BYOL</i>	0.836	0.841	0.825	0.831	0.786	0.837	0.844	0.819	0.825	0.785
	<i>Barlow Twins</i>	0.828	0.822	0.828	0.825	0.788	0.817	0.805	0.817	0.822	0.780
	<i>DINO</i>	0.757	0.719	0.733	0.761	0.699	0.566	0.402	0.712	0.777	0.752
	<i>SimCLR</i>	0.792	0.821	0.834	0.830	0.807	0.775	0.813	0.826	0.822	0.813
F1 _{Normal} \uparrow	<i>Supervised</i>	0.916	0.922	0.933	0.939	0.936	0.850	0.861	0.879	0.894	0.908
	<i>BYOL</i>	0.826	0.821	0.849	0.846	0.821	0.758	0.759	0.783	0.789	0.818
	<i>Barlow Twins</i>	0.865	0.873	0.879	0.870	0.834	0.782	0.796	0.799	0.807	0.828
	<i>DINO</i>	0.832	0.812	0.820	0.808	0.769	0.689	0.686	0.700	0.730	0.763
	<i>SimCLR</i>	0.849	0.869	0.878	0.880	0.839	0.752	0.793	0.806	0.818	0.836

Table 4: Comprehensive performance metrics for ViT-Tiny across both training and validation phases. The percentages represent the imbalance levels applied during training, while the validation data proportion remained unchanged.

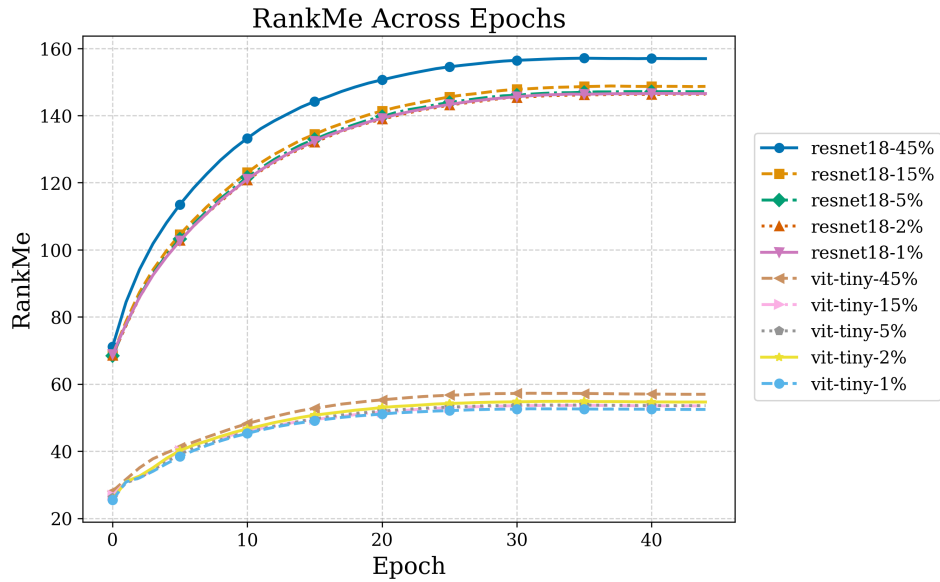
		TRAIN					VAL				
		1%	2%	5%	15%	45%	1%	2%	5%	15%	45%
PRECISION↑	<i>Supervised</i>	0.030	0.065	0.209	0.520	0.879	0.735	0.757	0.816	0.842	0.873
	<i>BYOL</i>	0.024	0.055	0.152	0.356	0.755	0.687	0.723	0.751	0.737	0.770
	<i>Barlow Twins</i>	0.024	0.058	0.136	0.375	0.780	0.681	0.729	0.733	0.753	0.792
	<i>DINO</i>	0.018	0.041	0.052	0.307	0.714	0.625	0.651	0.474	0.692	0.734
	<i>MAE</i>	0.010	0.022	0.063	0.154	0.223	0.469	0.503	0.530	0.479	0.232
	<i>SimCLR</i>	0.033	0.063	0.155	0.373	0.776	0.752	0.746	0.753	0.747	0.788
RECALL↑	<i>Supervised</i>	0.807	0.808	0.855	0.880	0.890	0.804	0.816	0.847	0.875	0.882
	<i>BYOL</i>	0.657	0.746	0.749	0.770	0.690	0.666	0.753	0.752	0.777	0.694
	<i>Barlow Twins</i>	0.677	0.731	0.754	0.779	0.744	0.682	0.748	0.762	0.787	0.752
	<i>DINO</i>	0.679	0.710	0.674	0.715	0.695	0.773	0.719	0.662	0.717	0.703
	<i>MAE</i>	0.496	0.500	0.588	0.705	0.013	0.567	0.864	0.617	0.719	0.013
	<i>SimCLR</i>	0.761	0.783	0.788	0.805	0.776	0.772	0.792	0.790	0.812	0.784
F1↑	<i>Supervised</i>	0.059	0.120	0.335	0.654	0.884	0.768	0.785	0.831	0.858	0.878
	<i>BYOL</i>	0.046	0.102	0.252	0.487	0.721	0.676	0.738	0.752	0.757	0.730
	<i>Barlow Twins</i>	0.046	0.107	0.231	0.506	0.762	0.682	0.738	0.747	0.769	0.772
	<i>DINO</i>	0.035	0.077	0.096	0.430	0.704	0.691	0.683	0.553	0.704	0.719
	<i>MAE</i>	0.020	0.043	0.114	0.253	0.025	0.513	0.636	0.570	0.575	0.024
	<i>SimCLR</i>	0.062	0.116	0.259	0.509	0.776	0.762	0.768	0.771	0.778	0.786
F2 _{C1W} ↑	<i>Supervised</i>	0.824	0.825	0.848	0.869	0.868	0.837	0.812	0.831	0.870	0.881
	<i>BYOL</i>	0.711	0.773	0.776	0.789	0.726	0.680	0.768	0.760	0.781	0.725
	<i>Barlow Twins</i>	0.708	0.759	0.787	0.804	0.760	0.704	0.749	0.786	0.784	0.758
	<i>DINO</i>	0.722	0.744	0.715	0.729	0.680	0.733	0.736	0.690	0.724	0.695
	<i>MAE</i>	0.624	0.394	0.711	0.683	0.014	0.629	0.640	0.704	0.702	0.009
	<i>SimCLR</i>	0.792	0.809	0.807	0.815	0.785	0.766	0.795	0.798	0.801	0.788
F1 _{Normal} ↑	<i>Supervised</i>	0.849	0.862	0.903	0.912	0.904	0.774	0.794	0.844	0.869	0.890
	<i>BYOL</i>	0.839	0.847	0.869	0.840	0.789	0.719	0.757	0.778	0.771	0.781
	<i>Barlow Twins</i>	0.834	0.858	0.850	0.851	0.813	0.716	0.760	0.766	0.785	0.806
	<i>DINO</i>	0.785	0.791	0.507	0.810	0.764	0.655	0.688	0.419	0.727	0.759
	<i>MAE</i>	0.615	0.709	0.692	0.461	0.695	0.471	0.350	0.553	0.389	0.677
	<i>SimCLR</i>	0.869	0.862	0.867	0.847	0.816	0.782	0.781	0.786	0.785	0.810

362 D Self-supervised monitoring metrics

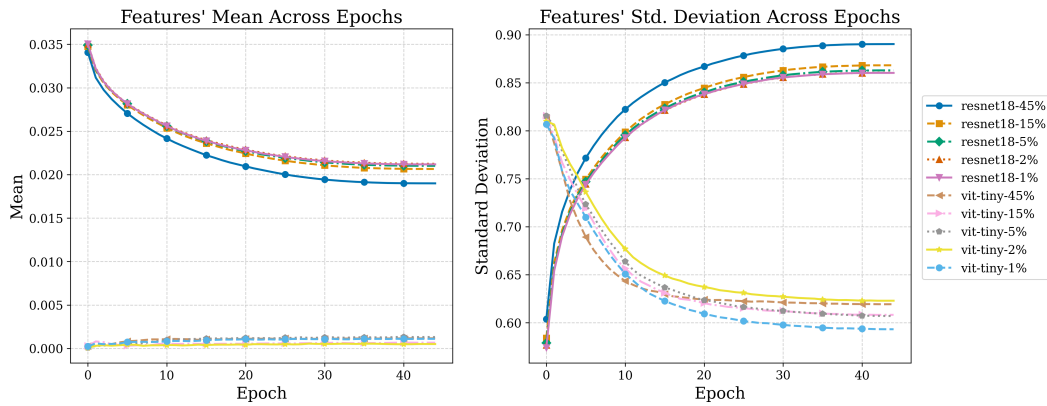
363 When evaluating SimCLR, BYOL, and DINO, significant differences emerge in the behavior of
364 ResNet-18 versus ViT-Tiny architectures. ResNet-18 consistently achieved higher RankMe values
365 and demonstrated steady improvement in feature standard deviation, indicating more stable and robust
366 feature representations. In contrast, ViT-Tiny models showed strong initial performance that often
367 degraded over time, particularly under SimCLR and BYOL, suggesting challenges in maintaining
368 feature consistency (see Figures 6, 7). However, under DINO, ViT-Tiny’s performance was more
369 stable, aligning closely with ResNet-18’s in terms of RankMe values (see Figure 8).

370 Overall, these findings suggest that ResNet architectures offer superior stability and consistency
371 across different self-supervised learning algorithms, whereas ViT-Tiny models may require tailored
372 optimization techniques to sustain their initial performance levels. This underscores the importance of
373 continuously monitoring self-supervised metrics to ensure robust feature learning and representation
374 stability.

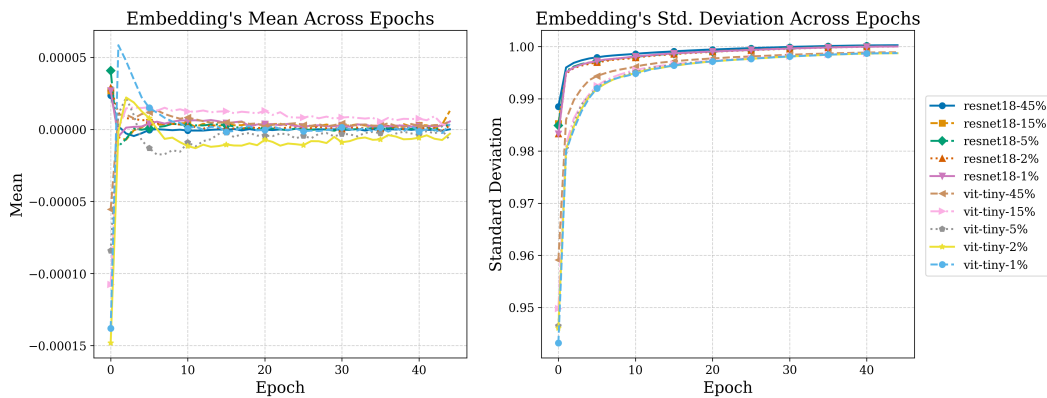
375 Furthermore, scatter plots comparing RankMe with train and validation F1 scores—grouped by SSL
376 method, model type, and defect proportion—show no correlation (Figures 9, 10) indicating that
377 RankMe is not a reliable predictor of downstream performance. This challenges RankMe’s intended
378 role and highlights a gap in its effectiveness within the SSL domain. Nonetheless, distinct RankMe
379 differences between ResNet-18 and ViT-Tiny architectures persist, reinforcing the architectural
380 distinctions previously observed.



(a) RankMe metric across epochs.

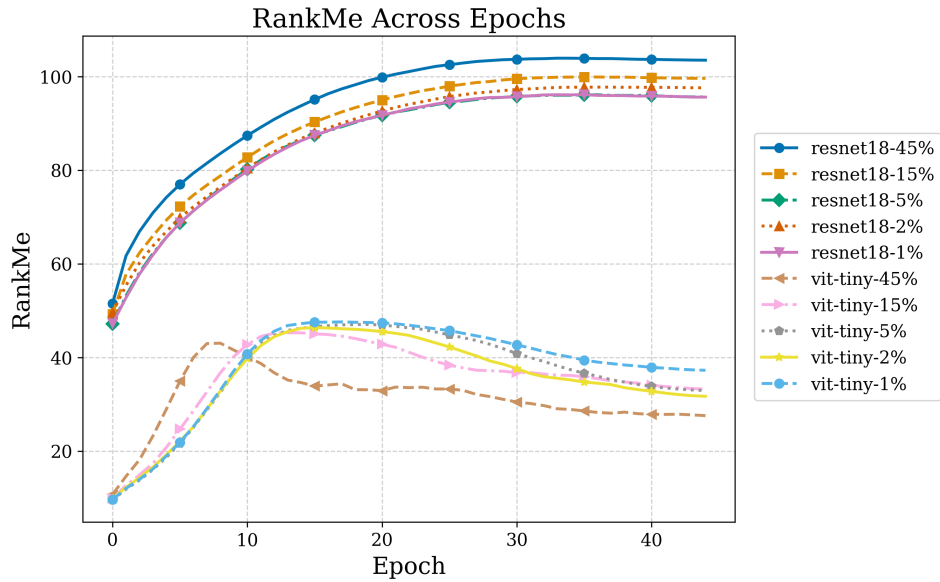


(b) Features' mean and standard deviation across epochs.

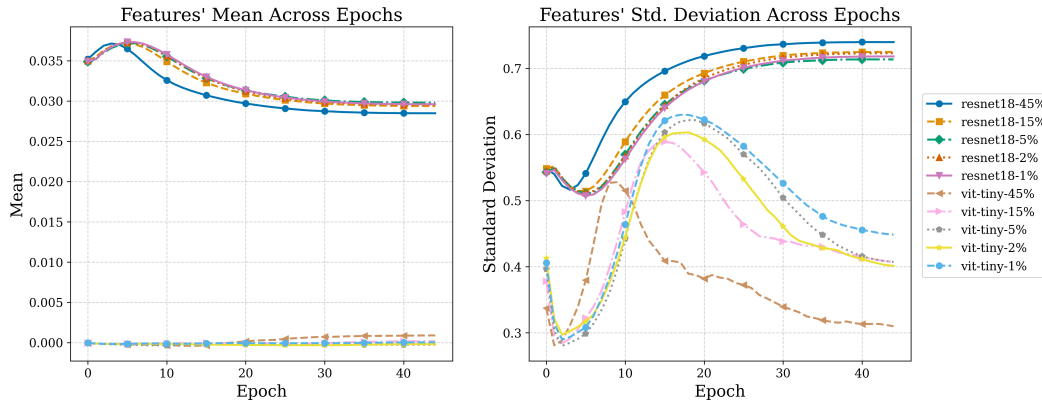


(c) Embedding's mean and standard deviation across epochs.

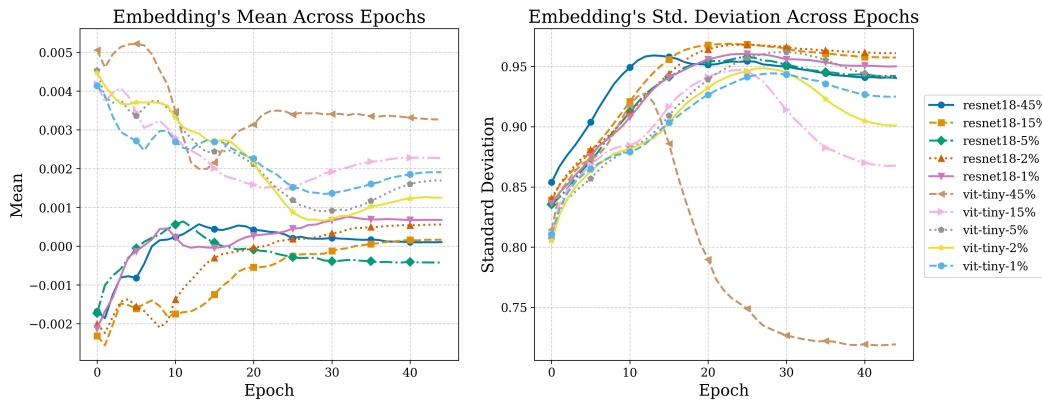
Figure 6: Self-supervised monitoring metrics for SimCLR.



(a) RankMe metric across epochs.

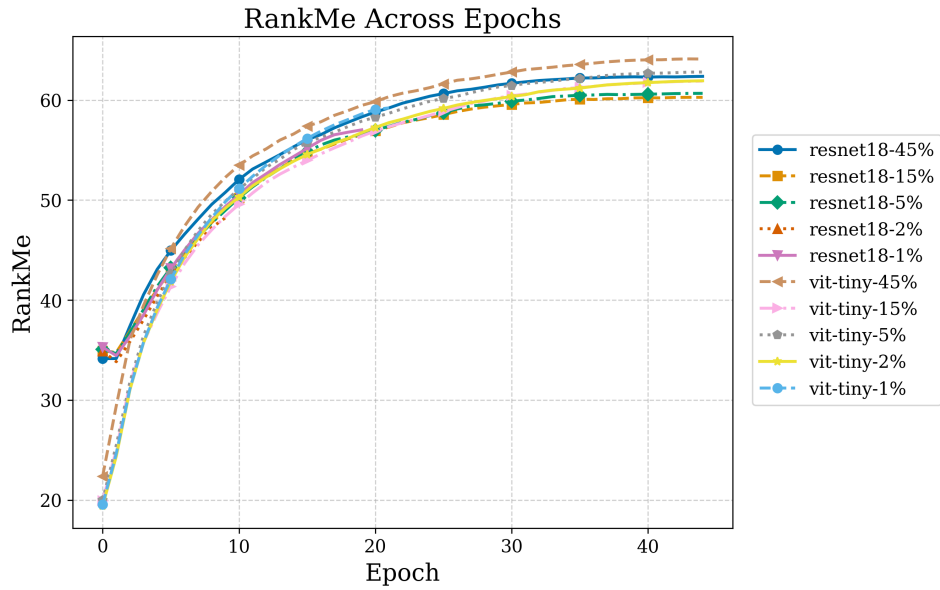


(b) Features' mean and standard deviation across epochs.

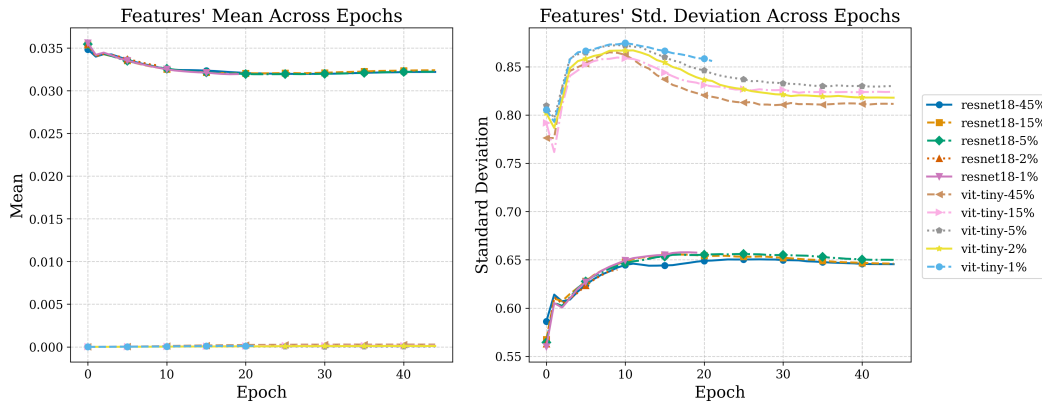


(c) Embedding's mean and standard deviation across epochs.

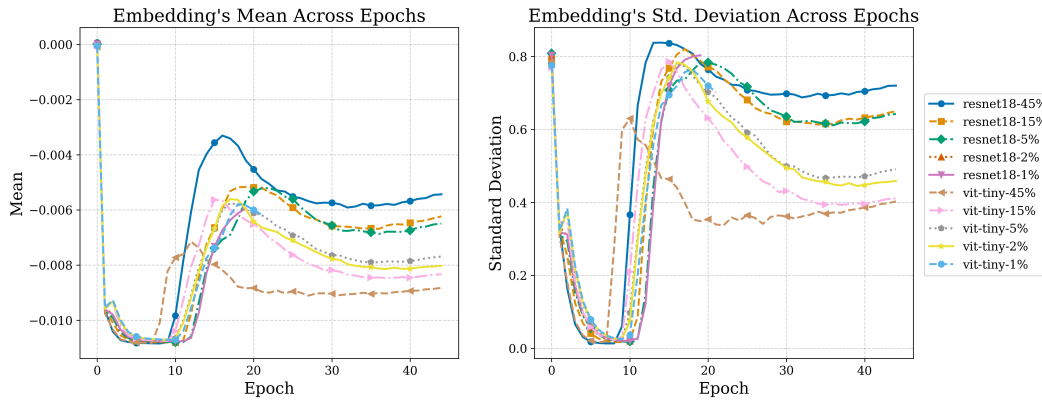
Figure 7: Self-supervised monitoring metrics for BYOL.



(a) RankMe metric across epochs.



(b) Features' mean and standard deviation across epochs.



(c) Embedding's mean and standard deviation across epochs.

Figure 8: Self-supervised monitoring metrics for DINO.

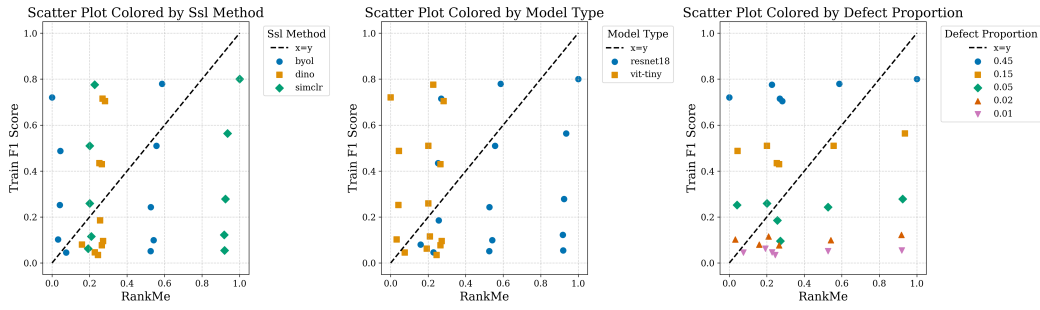


Figure 9: Train F1 scores plotted against re-scaled RankMe metrics. RankMe metrics were re-scaled to enable easier comparisons with F1 scores.

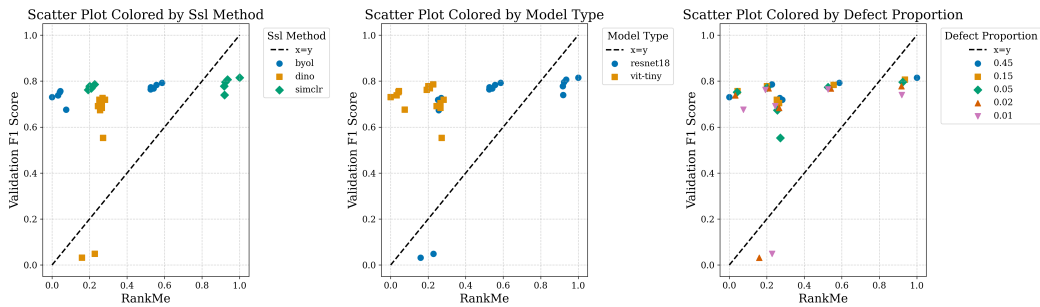


Figure 10: Validation F1 scores plotted against re-scaled RankMe metrics. RankMe metrics were re-scaled to enable easier comparisons with F1 scores.

385 **E Sewer-ML samples**











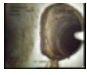

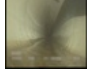
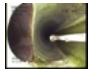





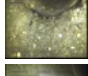
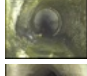

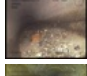
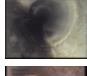
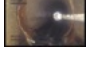




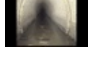




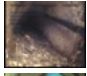










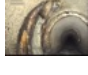



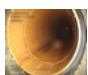





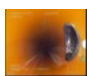









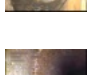
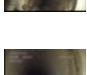


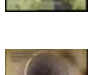


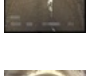


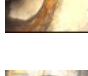














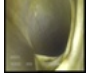

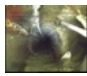














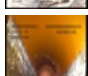



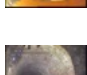


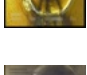


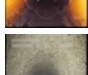


Defect Type	Train				Validation			
RB								
OS								
FS								
OB								
OK								
PH								
PB								
OP								
RO								
IN								
PF								
FO								
BE								
IS								
DE								

Table 5: Sample images by defect type from train and validation sets.