039

040

050

051

052

053

054

055

056

057

058

059

060

068

069

070

071

072

073

074

075

BRADD: Balancing Representations with Anomaly Detection and Diffusion

Anonymous CVPR submission

Paper ID *****

Abstract

001 Self-supervised learning (SSL) has allowed for advancements in language processing and computer vision, as un-002 003 labelled data is available in large quantities. However, imbalances in training datasets can lead to strong biases 004 in the learned features of pre-trained models. Previous 005 results show that pre-training using imbalanced data can 006 also hurt downstream performance. We propose a data-007 008 centric approach: our method trains on the data, finds underrepresented samples, and uses diffusion to gener-009 ate novel data complementing the underrepresented im-010 ages. Our proposed method, BRADD (Balancing Represen-011 012 tations through Anomaly Detection and Diffusion), utilizes 013 distance-based outlier detection to identify regions of the 014 embedding space that are underrepresented in each training cycle. Experimental results on ImageNet-100-LT demon-015 strate that BRADD consistently outperforms both balanced 016 and imbalanced baselines, with significant improvements 017 018 on fine-grained classification tasks. Detailed ablation studies confirm that both out-of-distribution sample selection 019 and diffusion-based generation contribute substantially to 020 the effectiveness of our approach, offering a promising al-021 ternative to model-centric solutions for addressing imbal-022 023 ance in self-supervised learning.

024 1. Introduction

Self-supervised learning (SSL) methods have emerged 025 as powerful techniques for learning transferable features 026 027 across diverse tasks [2, 3, 17]. However, their performance is significantly affected by dataset imbalance [1, 028 12, 23]. For instance, SimCLR underperforms on long-029 tailed datasets due to insufficient negative samples [23], 030 while joint-embedding methods like VICReg assume uni-031 032 form clustering, hampering performance on imbalanced data [1]. 033

Existing solutions typically adopt model-centric approaches, such as ensemble learning [23], incorporating arbitrary feature priors [1], or modifying training dynamics [12]. These approaches, however, often require prior knowledge of dataset distributions or extensive hyperparameter tuning, contradicting the unsupervised nature of SSL.

We propose BRADD (Balancing Representations with 041 Anomaly Detection and Diffusion), a data-centric approach 042 to address imbalance in SSL. BRADD divides training into 043 cycles, identifies underrepresented samples via OOD detec-044 tion after each cycle, and augments them using diffusion 045 models for subsequent training. This approach (1) elimi-046 nates the need for prior dataset knowledge and (2) dynam-047 ically balances the latent space to avoid suboptimal local 048 minima. 049

Experiments on ImageNet-100-LT across multiple architectures (ResNet-50, ViT-S, ViT-B) and SSL methods (Sim-CLR, DINO, MoCo) demonstrate that BRADD consistently outperforms both balanced and imbalanced baselines on diverse downstream tasks. BRADD achieves significant improvements on fine-grained tasks (up to 11.7% on Oxford Flowers) and surpasses state-of-the-art methods on CIFAR-10 (73.3%) and CIFAR-100 (45.4%), showing that our datacentric approach offers a promising alternative to modelcentric solutions.

2. Related Work

Self-Supervised Learning (SSL) leverages unlabeled data061to learn meaningful representations for downstream tasks.062In computer vision, three main approaches are: (1) masked063prediction [7], (2) contrastive learning like MoCo [6] and064SimCLR [3], and (3) self-distillation methods such as065DINO [2]. These techniques have enabled large-scale training of models with emergent abilities [24].067

Training on Imbalanced Datasets often leads to inferior performance and bias toward majority classes [10]. While SSL methods are more robust to imbalance than supervised approaches [13], they still exhibit diminished performance on imbalanced data. Current solutions are primarily model-centric, attempting to learn arbitrary feature priors [1]. Data-centric approaches that directly complement imbalanced datasets remain underexplored.

Out-Of-Distribution (OOD) Detection is crucial for en-
hancing model robustness by maintaining high-quality076077

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

datasets [25]. Distance-based approaches [21] identify OOD samples by measuring their distance from indistribution samples in the embedding space. These nonparametric methods offer flexibility without distributional
assumptions, making them suitable for imbalanced, unlabeled data.

084 Diffusion-Based Image Generation models systematically add and then remove noise to generate high-quality images 085 [8]. Stable Diffusion [20], particularly Stable Diffusion 2 086 UnClip, leverages CLIP embeddings to generate semanti-087 088 cally similar images to the input context. Its ability to per-089 form image-to-image generation while preserving semantic content makes it well-suited for augmenting self-supervised 090 091 learning datasets.

3. Method: BRADD

We propose BRADD (Balancing Representations with
Anomaly Detection and Diffusion), a data-centric approach
to address imbalance in self-supervised learning. Unlike
model-centric approaches, BRADD identifies and augments
underrepresented concepts in the data distribution itself.

Algorithm 1 BRADD: Balancing Representations through Anomaly Detection and Diffusion

1: **Input:** model m, SSL algorithm a, diffusion model m_d , cycles N_C , epochs per cycle N_E , dataset D, samples per point N_{Aug} , OOD samples N_{OOD}

O

=

- 2: repeat
- 3: Train m using a for N_E epochs on D
- 4: Compute embeddings $E = \{m(x_i) | x_i \in D\}$

$$\{x_i | \operatorname{rank}(d_{kNN}(m(x_i), E)) \le N_{OOD}\}$$

6: for $x \in O$ do $7 \quad D = D \cup \{m_d(x, \epsilon_i) | j = 1, ..., N_{Aug}\}$

9: **until** N_C cycles completed

098BRADD alternates between self-supervised pre-training099and dataset augmentation phases (Algorithm 1). After100training for N_E epochs, we identify underrepresented data101points by computing k-nearest neighbor distances in the em-102bedding space. Rather than using a percentile threshold, we103select the top- N_{OOD} samples with highest k-NN distances,104providing precise control over augmentation.

For each identified OOD point, we generate N_{Aug} new samples using Stable Diffusion 2 UnCLIP [19], which preserves semantic content while introducing sufficient variation. We use k=5 for k-NN computation, $N_{OOD} = 500$, and $N_{Aug} = 5$, adding 2,500 new images per cycle across 5 cycles with 20 epochs each (100 total epochs).

4. Experiments

To evaluate the performance of our proposed method112BRADD (Balancing Representations through Automated113Detection and Diffusion), we conduct experiments on114ImageNet-100-LT. We test various backbone architectures,115SSL methods, and implementation details through compre-
hensive ablation studies, and compare against state-of-the-
art approaches.118

4.1. ImageNet-100 Experiments

ImageNet-100 is a subset of ImageNet with 100 randomly selected classes [22]. Following previous work [14], we introduce imbalance by creating a long-tailed distribution (ImageNet-100-LT) that follows a Pareto distribution with $\alpha = 6$. The dataset contains around 15 thousand images. **Models.** We evaluate our method across multiple architectures: ResNet-50 (25.6M parameters), ViT-Small (22.1M parameters), and ViT-Base (86.6M parameters).

Self-Supervised Learning Methods. Our primary experiments use SimCLR [3] with a temperature of 0.5, and we later compare with DINO [2] and MoCo [5] in our ablation studies.

Downstream Evaluation. We evaluate the learned features using both linear probing and K-nearest neighbor (KNN) classification across multiple datasets: CIFAR-10, CIFAR-100 [11], Stanford Cars [4], FGVC Aircraft [15], Oxford Flowers [16], and Oxford-IIIT Pets [18].

Baselines. We train two baselines: (1) a model trained on a balanced subset of ImageNet-100 with as much data uniformly removed as in ImageNet-100-LT and (2) a model trained on the imbalanced ImageNet-100-LT.

Proposed Method. Our proposed method BRADD starts with ImageNet-100-LT and trains for multiple cycles, where each cycle consists of training epochs followed by OOD detection and generation steps.

4.2. Experimental Results

Table 1 compares our method against balanced and imbalanced baselines using a ViT-B backbone with SimCLR. The results demonstrate that:

- While imbalance causes only slight performance degradation compared to the balanced setting on some datasets, our method consistently outperforms both baselines across all datasets.
- BRADD achieves substantial improvements in linear probing, with gains of up to 11.72% on Oxford Flowers and 7.53% on Oxford-IIIT Pets compared to the imbalanced baseline.
- Our method also shows consistent improvements in KNN classification, demonstrating the enhanced quality of the learned feature space.

Table 2 shows BRADD compared to state-of-the-art meth-ods using ResNet-50 with SimCLR (500 epochs):

194

195

196

197

198

205



Figure 1. The BRADD algorithm augments the most underrepresented datapoints in a dataset based on OOD detection.

Table 1. Comparison to balanced and imbalanced baselines (ViT-B, SimCLR, ImageNet-100-LT, 100 epochs)

	C-	10	C-	100	C	ars	Airc	raft	Flov	vers	Pe	ets
Setting	Lin.	KNN	Lin.	KNN	Lin.	KNN	Lin.	KNN	Lin.	KNN	Lin.	KNN
Balanced	59.90	57.13	33.63	28.34	4.29	3.31	9.14	4.83	31.25	26.04	23.64	13.04
Imbalanced	59.89	57.25	33.02	27.78	3.80	3.43	8.75	5.37	27.50	25.42	22.55	11.96
BRADD (ours)	64.00	60.73	38.42	31.56	6.13	4.54	11.55	6.72	39.22	30.89	30.08	14.69

- Our method achieves superior linear probing performance on CIFAR-10 (73.34%) and CIFAR-100 (45.38%) compared to previous methods.
- While TS [12] performs better on KNN classification for most datasets, BRADD shows competitive performance across all benchmarks.

168 4.3. Ablation Studies

To analyze the effectiveness of different components in ourapproach, we conducted extensive ablation studies:

SSL Method. Table 3a demonstrates that SimCLR consistently outperforms DINO and MoCo across all datasets
when using our method, with substantial margins particu-

- 174 larly on fine-grained classification tasks.
- Backbone Architecture. As shown in Table 3b, we compared ViT-S, ViT-B, and ResNet-50 backbones. ResNet-50 achieves the best linear probing performance on CIFAR-10

(71.17%), while ViT architectures perform better on finegrained datasets, with ViT-B showing the strongest perfor-

180 mance on Oxford-IIIT Pets (30.08%).

181 Sample Selection Strategy. In Table 3c, we compare uni182 form sampling versus our OOD-based selection. The re183 sults confirm that OOD-based selection provides consistent
184 performance gains, validating our hypothesis that targeting
185 underrepresented regions of the feature space is more effec186 tive than random augmentation.

187 Sample Generation Method. Table 3d compares re-

population (adding back removed samples) versus gener-
ation using Stable Diffusion. While both approaches im-
prove over the baseline, Stable Diffusion generation yields
better results, particularly on fine-grained datasets like Ox-
ford Flowers.188
189
190

Number of Cycles. As shown in Table 3e, we tested different numbers of cycles (2, 5, 10, 20) while keeping the total number of generated images fixed. Performance generally improves with more cycles, with the best results at 20 cycles for CIFAR-10 (68.68% linear probing accuracy) and 5 cycles for Oxford-IIIT Pets (30.08%).

Number of OOD Samples. Table 3f analyzes the effect of
selecting different numbers of samples per cycle for aug-
mentation. We find that 500 samples per cycle provides the
optimal balance, while selecting too many samples (2500)
degrades performance, suggesting that focusing on the most
out-of-distribution samples is important.199
200
201
202

5. Conclusion

Our comprehensive experiments demonstrate that BRADD 206 effectively mitigates the negative impact of dataset imbal-207 ance in self-supervised learning. By strategically detect-208 ing regions of the feature space that are underrepresented 209 and augmenting them with generated samples, our method 210 achieves consistent improvements across different architec-211 tures, SSL methods, and downstream tasks. The ablation 212 studies confirm that both OOD-based sample selection and 213

224

225

226

227

228

229

230

231

232

	C-	-10	C-	100	Ca	ars	Airc	eraft	Flov	wers	Pe	ets
Method	Lin.	KNN	Lin.	KNN	Lin.	KNN	Lin.	KNN	Lin.	KNN	Lin.	KNN
SDCLR [9]	68.72	65.16	38.71	35.53	7.84	5.21	11.24	10.11	11.87	31.14	40.22	21.98
TS [12]	71.26	66.76	43.90	35.64	10.91	5.58	12.95	10.85	32.05	3.36	47.01	23.20
BRIDGE (ours)	73.34	61.76	45.38	32.26	8.95	4.73	10.95	7.92	26.14	24.94	34.68	16.05

Table 2. **Comparison to SOTA** (ResNet50, SimCLR, ImageNet-100-LT, 500 epochs)

Table 3. Ablations of the key components of our method.

(a) Pretraining (ViT-B, ImageNet-100-LT, 100 epochs)										
C-10 Flowers Pets										
backbone	Lin.	KNN	Lin.	KNN	Lin.	KNN				
SimCLR	63.40	60.73	39.22	30.89	30.08	14.69				
DINO	47.97	37.36	20.92	12.46	10.30	6.00				
MoCo	49.13	42.86	15.03	17.63	14.36	7.50				

(c) Sample Selection (ViT-B, SimCLR, ImageNet-100-LT)

		C-10		Flov	vers	Pets	
SSL	Method	Lin.	KNN	Lin.	KNN	Lin.	KNN
SimCLR	uniform	62.01	58.81	31.37	29.56	23.85	13.25
	OOD	63.40	60.73	39.22	30.89	30.08	14.69
MoCoV3	uniform	51.38	41.08	21.79	16.16	12.90	7.14
	OOD	52.42	42.13	19.87	17.12	11.83	11.56

(e) Number of Cycles (ViT-S, SimCLR, ImageNet-100-LT, 100 Epochs, 10k created images total)

	C-10		Flov	wers	Pets		
#Cycles	Lin.	KNN	Lin.	KNN	Lin.	KNN	
2	64.08	61.15	39.87	31.63	26.02	14.83	
5	63.40	60.73	39.22	30.89	30.08	14.69	
10	64.46	60.63	31.38	31.28	27.91	14.23	
20	68.68	62.02	35.95	31.56	26.56	14.72	

diffusion-based generation contribute significantly to the ef-fectiveness of our approach.

216 References

- [1] M Assran, R Balestriero, Q Duval, F Bordes, I Misra, P Bojanowski, P Vincent, M Rabbat, and N Ballas. The hidden uniform cluster prior in self-supervised learning. URL https://arxiv. org/abs/2210.07277, 19:26, 2022a. 1
- [2] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou,Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerg-

(b) Backbone
(SimCLR, ImageNet-100-LT)

	C-10		Flov	vers	Pets		
backbone	Lin.	KNN	Lin.	KNN	Lin.	KNN	
ViT-S	64.48	62.13	44.44	34.78	27.91	15.94	
R50	71.17	59.06	26.80	21.07	26.29	15.43	
ViT-B	63.40	60.73	39.22	30.89	30.08	14.69	

(d) Sample Generation (ViT-B, SimCLR, ImageNet-100-LT)

		C-10		Flov	vers	Pets		
SSL	Method	Lin.	KNN	Lin.	KNN	Lin.	KNN	
SimCLR	Re-Popul. Stable-Diff.	63.36 63.40	59.68 60.73	30.07 39.22	29.24 30.89	26.02 30.08	14.12 14.69	
MoCoV3	Re-Popul. Stable-Diff.	48.76 52.42	38.68 42.13	24.56 19.87	15.56 17.12	19.23 11.83	15.50 11.56	

(f) Number of most OOD samples selected for aug. per cycle
(ViT-B, SimCLR, ImageNet-100-LT, 100 Epochs, 2500 created
images per cycle)

	C-10		Flov	vers	Pets		
#Samples	Lin.	KNN	Lin.	KNN	Lin.	KNN	
5	62.09	59.04	33.33	28.55	25.47	13.79	
500	63.40	60.73	39.22	30.89	30.08	14.69	
1250	58.09	56.10	28.10	24.21	27.10	11.50	
2500	57.03	55.75	24.18	23.15	23.57	11.42	

ing properties in self-supervised vision transformers. *CoRR*, abs/2104.14294, 2021. 1, 2

- [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 1, 2
- [4] Afshin Dehghan, Syed Zain Masood, Guang Shu, and Enrique. G. Ortiz. View independent vehicle make, model and color recognition using convolutional neural network, 2017.
 2

248

249

250

251

257 258

259

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

- 233 [5] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and 234 Ross B. Girshick. Momentum contrast for unsupervised vi-235 sual representation learning. CoRR, abs/1911.05722, 2019. 236 2
- 237 [6] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross 238 Girshick. Momentum contrast for unsupervised visual rep-239 resentation learning. In Proceedings of the IEEE/CVF con-240 ference on computer vision and pattern recognition, pages 241 9729-9738, 2020. 1
- 242 [7] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr 243 Dollár, and Ross Girshick. Masked autoencoders are scalable 244 vision learners. In Proceedings of the IEEE/CVF conference 245 on computer vision and pattern recognition, pages 16000-246 16009, 2022. 1
 - [8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. 2
 - [9] Ziyu Jiang, Tianlong Chen, Bobak Mortazavi, and Zhangyang Wang. Self-damaging contrastive learning, 2021. 4
- 252 [10] Justin M Johnson and Taghi M Khoshgoftaar. Survey on 253 deep learning with class imbalance. Journal of Big Data, 6 254 (1):1-54, 2019. 1
- 255 [11] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple 256 layers of features from tiny images. 2009. 2
- [12] Anna Kukleva, Moritz Böhle, Bernt Schiele, Hilde Kuehne, and Christian Rupprecht. Temperature schedules for selfsupervised contrastive methods on long-tail data, 2023. 1, 3, 260 4
- 261 [13] Hong Liu, Jeff Z HaoChen, Adrien Gaidon, and Tengyu Ma. 262 Self-supervised learning is more robust to dataset imbalance. 263 arXiv preprint arXiv:2110.05025, 2021. 1
- 264 [14] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, 265 Boqing Gong, and Stella X Yu. Large-scale long-tailed 266 recognition in an open world. In Proceedings of the 267 IEEE/CVF conference on computer vision and pattern 268 recognition, pages 2537-2546, 2019. 2
- 269 [15] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. 270 Fine-grained visual classification of aircraft. Technical re-271 port, 2013. 2
- 272 [16] Maria-Elena Nilsback and Andrew Zisserman. Automated 273 flower classification over a large number of classes. In 2008 274 Sixth Indian Conference on Computer Vision, Graphics Im-275 age Processing, pages 722-729, 2008. 2
- 276 [17] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy 277 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, 278 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. 279 Dinov2: Learning robust visual features without supervision. 280 arXiv preprint arXiv:2304.07193, 2023. 1
- 281 [18] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In IEEE Conference on Com-282 283 puter Vision and Pattern Recognition, 2012. 2
- [19] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 284 285 Patrick Esser, and Björn Ommer. High-resolution image syn-286 thesis with latent diffusion models, 2021. 2
- 287 [20] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 288 Patrick Esser, and Björn Ommer. High-resolution image syn-289 thesis with latent diffusion models, 2022. 2

- [21] Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-290 of-distribution detection with deep nearest neighbors. In In-291 ternational Conference on Machine Learning, pages 20827-292 20840. PMLR, 2022. 2 293
- [22] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Con-294 trastive multiview coding. In Computer Vision-ECCV 2020: 295 16th European Conference, Glasgow, UK, August 23–28, 296 2020, Proceedings, Part XI 16, pages 776-794. Springer, 297 2020. 2 298
- [23] Yonglong Tian, Olivier J Henaff, and Aäron Van den Oord. Divide and contrast: Self-supervised learning from uncurated data. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 10063-10074, 2021.
- [24] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022. 1
- [25] Steven Euijong Whang, Yuji Roh, Hwanjun Song, and Jae-Gil Lee. Data collection and quality challenges in deep learning: A data-centric ai perspective. The VLDB Journal, 32(4): 791-813, 2023. 2