BISSL: BILEVEL OPTIMIZATION FOR SELF SUPERVISED PRE-TRAINING AND FINE-TUNING

Anonymous authors

Paper under double-blind review

ABSTRACT

In this work, we present BiSSL, a first-of-its-kind training framework that introduces bilevel optimization to enhance the alignment between the pretext pretraining and downstream fine-tuning stages in self-supervised learning. BiSSL formulates the pretext and downstream task objectives as the lower- and upperlevel objectives in a bilevel optimization problem and serves as an intermediate training stage within the self-supervised learning pipeline. By more explicitly modeling the interdependence of these training stages, BiSSL facilitates enhanced information sharing between them, ultimately leading to a backbone parameter initialization that is better suited for the downstream task. We propose a training algorithm that alternates between optimizing the two objectives defined in BiSSL. Using a ResNet-18 backbone pre-trained with SimCLR on the STL10 dataset, we demonstrate that our proposed framework consistently achieves improved or competitive classification accuracies across various downstream image classification datasets compared to the conventional self-supervised learning pipeline. Qualitative analyses of the backbone features further suggest that BiSSL enhances the alignment of downstream features in the backbone prior to fine-tuning.

025 026 027

024

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

In the absence of sufficient labeled data, self-supervised learning (SSL) has emerged as a promising approach for training deep learning models. Rather than relying solely on labeled data, the SSL 031 framework aims to learn representations from unlabeled data which proves beneficial for subsequent use on various downstream tasks. These representations are learned by solving a pretext task, which 033 utilizes supervisory signals extracted from the unlabeled data itself. Extensive efforts has gone into 034 designing effective pretext tasks, achieving state-of-the-art or competitive performance in various fields such as computer vision (Chen et al., 2020b; Bardes et al., 2022; He et al., 2020; Grill et al., 2020; Caron et al., 2020; 2021; He et al., 2022; Oquab et al., 2024), audio signal processing (Schnei-037 der et al., 2019; Baevski et al., 2020; Hsu et al., 2021; Niizumi et al., 2021; Chung & Glass, 2018; 038 Chung et al., 2019; Yadav et al., 2024) and natural language processing (Devlin et al., 2019; Lewis et al., 2019; Brown et al., 2020; He et al., 2021; Touvron et al., 2023).

040 Making a self-supervised pre-trained backbone suitable for a downstream task typically involves 041 attaching additional layers that are compatible with that task, followed by fine-tuning the entire or 042 parts of the composite model in a supervised manner (Zhai et al., 2019; Dubois et al., 2022). When a 043 backbone is pre-trained on a distribution that differs from the distribution of the downstream data, the 044 representations learned during pre-training may not be initially well-aligned with the downstream task. During fine-tuning, this distribution misalignment could cause relevant semantic information, learned during the pre-training phase, to vanish from the representation space (Zaiem et al., 2024; 046 Chen et al., 2020a; Boschini et al., 2022). A potential strategy for alleviating the negative effects of 047 these distribution discrepancies would be to enhance the alignment between the pretext pre-training 048 and downstream fine-tuning stages. However, since the conventional SSL pipeline treats these stages as two disjoint processes, this poses a significant challenge in devising a strategy that enhances such alignment while not compromising on the benefits that SSL offers. 051

Meanwhile, bilevel optimization (BLO) has risen as a powerful tool for solving certain optimization
 problems within deep learning. It entails a main optimization problem constrained by the solution
 to a secondary optimization problem that depends on the parameters of the main objective. This hi-



Figure 1: The conventional self-supervised learning pipeline alongside the proposed pipeline involving BiSSL. The symbols θ and ϕ represent obtained backbone and task-specific attached head parameters, respectively. When they are transmitted to the respective training stages, they are used as initializations.

erarchical setup causes the solutions of both optimization problems to depend on each other, either directly or implicitly, which has proven advantageous in deep learning tasks that optimize multiple inter-dependent objectives simultaneously (Zhang et al., 2023a). Notable mentions of tasks within deep learning where BLO has proven useful are parameter pruning (Zhang et al., 2022b), invariant risk minimization (Arjovsky et al., 2019; Zhang et al., 2023b), meta-learning (Rajeswaran et al., 2019; Finn et al., 2017), adversarial robustness (Zhang et al., 2021), hyper-parameter optimization (Franceschi et al., 2018) and coreset selection (Borsos et al., 2020).

In this study, we propose BiSSL, a novel training framework that leverages BLO to enhance the 085 alignment between the pretext pre-training and downstream fine-tuning stages in SSL. Acting as an intermediate training stage within the SSL pipeline, BiSSL frames the pretext and downstream 087 task objectives as the lower- and upper-level objectives in a BLO problem - a challenging approach 088 that has not been explored until now. The objectives in BiSSL are connected by substituting the 089 lower-level backbone solution for the upper-level backbone parameters, while simultaneously en-090 forcing the lower-level backbone solution to resemble the upper-level backbone parameters. This 091 approach more explicitly captures the interdependence between pretext pre-training and downstream 092 fine-tuning, potentially leading to a lower-level backbone better aligned with the downstream task. Figure 1 compares the conventional SSL pipeline with our suggested pipeline involving BiSSL. Additionally, we propose a training algorithm for BiSSL and demonstrate that it consistently improves 094 or maintains comparable downstream performance across a range of image classification datasets. 095 For our experiments, we use SimCLR (Chen et al., 2020b) to pre-train a ResNet-18 backbone (He 096 et al., 2016) on the unlabeled partition of the STL10 dataset (Coates et al., 2011), a setup offering suitable model capacity and dataset complexity while being less resource-intensive than larger-scale 098 alternatives. The code implementation and pre-trained model weights are publicly available.¹

100 101

102 103

104

105

071 072

073

074

075 076

084

2 RELATED WORK

Bilevel Optimization in Self-Supervised Learning Bilevel optimization (BLO) refers to a constrained optimization problem, where the constraint itself is a solution to another optimization problem, which depends on the parameters of the "main" optimization problem. The general BLO

¹https://github.com/ICLR25-10484/ICLR25_10484_BiSSL

problem is formulated as

110 111 $\min_{\boldsymbol{\xi}} f(\boldsymbol{\xi}, \boldsymbol{\psi}^*(\boldsymbol{\xi})) \quad \text{s.t.} \quad \boldsymbol{\psi}^*(\boldsymbol{\xi}) \in \operatorname*{argmin}_{\boldsymbol{\psi}} g(\boldsymbol{\xi}, \boldsymbol{\psi}), \tag{1}$

where f and q are referred to as the upper-level and lower-level objectives, respectively. While 112 the lower objective g has knowledge of the parameters $\boldsymbol{\xi}$ from the upper-level objective, the upper-113 level objective f possesses full information of the lower objective q itself through its dependence 114 on the lower-level solution $\psi^*(\xi)$. Some works have incorporated bilevel optimization within self-115 supervised learning. Gupta et al. (2022) suggest formulating the contrastive self-supervised pretext 116 task as a bilevel optimization problem, dedicating the upper-level and lower-level objectives for up-117 dating the backbone and projection head parameters respectively. Other frameworks such as the 118 Local and Global (LoGo) (Zhang et al., 2022a) and Only Self-Supervised Learning (OSSL) Boonlia 119 et al. (2022) utilize auxiliary models, wherein the lower-level objective optimizes the parameters of 120 the auxiliary model, while the upper-level objective is dedicated to training the feature extraction 121 model. MetaMask (Li et al., 2022b) introduces a meta-learning based approach, where the upperlevel learns masks that filter out irrelevant information from inputs that are provided to a lower-level 122 self-supervised contrastive pretext task. Chen et al. (2023) introduces a pseudo-BLO setup where 123 the upper-level optimization still benefits from knowledge of the lower-level objective, but the pa-124 rameters of the lower-level objective are fixed during training. In Somayajula et al. (2023), a two-125 staged BLO problem is proposed to fine-tune self-supervised pre-trained large language models in 126 low-resource scenarios. Their approach focuses on solving downstream tasks while simultaneously 127 learning a task-dependent similarity structure. BLO-SAM (Zhang et al., 2024) is tailored towards 128 fine-tuning the segment anything model (SAM) (Kirillov et al., 2023) by interchangeably alternating 129 between learning (upper-level) prompt embeddings and fine-tuning the (lower-level) segmentation 130 model. The aforementioned frameworks integrate bilevel optimization into either the pre-training or 131 fine-tuning stage exclusively and are tailored towards specific pretext or downstream tasks. In con-132 trast, our proposed BiSSL employs a BLO problem that comprehensively incorporates *both* training stages of pretext pre-training and downstream fine-tuning, without being confined to any specific 133 type of pretext or downstream task. 134

Priming Pre-Trained Backbones Prior To Fine-Tuning Previous works have demonstrated that 136 downstream performance can be enhanced by introducing techniques that modify the backbone be-137 tween the pre-training and fine-tuning stages. Contrastive Initialization (COIN) (Pan et al., 2022) 138 introduces a supervised contrastive loss, to be utilised on backbones pre-trained with contrastive 139 SSL techniques. Noisy-Tune (Wu et al., 2022) perturbs the pre-trained backbone with tailored noise 140 before fine-tuning. Speaker-invariant clustering (Spin) (Chang et al., 2023) utilizes speaker dis-141 entanglement and vector quantization for improving speech representations for speech signal spe-142 cific downstream tasks. RIFLE (Li et al., 2020) conducts multiple fine-tuning sessions sequentially, 143 where the attached downstream specific layers are re-initialized in between every session. Unlike 144 BiSSL, these techniques either do not incorporate knowledge of both the pretext task and down-145 stream task objectives and their relationship or do so only implicitly.

146 147

148

135

3 PROPOSED METHOD

149 3.1 NOTATION

We denote the unlabeled pretext dataset $\mathcal{D}^P = \{\mathbf{z}_k\}_{k=1}^{C_P}$ and labeled downstream dataset $\mathcal{D}^D = \{\mathbf{x}_l, \mathbf{y}_l\}_{k=1}^{C_D}$, respectively, where $\mathbf{z}_k, \mathbf{x}_l \in \mathbb{R}^N$. Let $f_{\boldsymbol{\theta}} : \mathbb{R}^N \to \mathbb{R}^M$ denotes a feature extracting backbone with trainable parameters $\boldsymbol{\theta}$ and $p_{\boldsymbol{\phi}} : \mathbb{R}^M \to \mathbb{R}^P$ a task specific projection head with trainable parameters $\boldsymbol{\phi}$. Given pretext and downstream models $g_{\boldsymbol{\phi}_P} \circ f_{\boldsymbol{\theta}_P}$ and $h_{\boldsymbol{\phi}_D} \circ f_{\boldsymbol{\theta}_D}$ with $\boldsymbol{\theta}_P, \boldsymbol{\theta}_D \in \mathbb{R}^L$, we denote the pretext and downstream training objectives $\mathcal{L}^P(\boldsymbol{\theta}_P, \boldsymbol{\phi}_P; \mathcal{D}^P)$ and $\mathcal{L}^D(\boldsymbol{\theta}_D, \boldsymbol{\phi}_D; \mathcal{D}^D)$, respectively. To simplify notation, we omit the dataset specification from the training objectives, e.g. $\mathcal{L}^D(\boldsymbol{\theta}_D, \boldsymbol{\phi}_D) := \mathcal{L}^D(\boldsymbol{\theta}_D, \boldsymbol{\phi}_D; \mathcal{D}^D)$.

158 159

160

3.2 Optimization Problem Formulation

161 The conventional setup of self-supervised pre-training directly followed by supervised fine-tuning relies on using a single backbone model with parameters θ . In that instance, we minimize

s.

162 $\mathcal{L}^{P}(\theta, \phi_{P})$ to produce a backbone parameter configuration θ^{*} which is then used as an initial-163 ization when subsequently minimizing the downstream training objective $\mathcal{L}^{D}(\theta, \phi_{D})$. We deviate 164 from this by instead considering θ_{P} and θ_{D} as two separate parameter vectors that are strongly cor-165 related. In continuation, we suggest combining the two traditionally separate optimization problems 166 of pretext and downstream training into a joint optimization problem through bilevel optimization 167 called BiSSL. We formulate BiSSL as

$$\min_{\boldsymbol{\theta}_{D},\boldsymbol{\phi}_{D}} \mathcal{L}^{D}\left(\boldsymbol{\theta}_{P}^{*}\left(\boldsymbol{\theta}_{D}\right),\boldsymbol{\phi}_{D}\right) + \gamma \mathcal{L}^{D}\left(\boldsymbol{\theta}_{D},\boldsymbol{\phi}_{D}\right)$$
(2)

168 169 170

191

192

196 197

205 206

210

211

215

t.
$$\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D}) \in \operatorname*{argmin}_{\boldsymbol{\theta}_{P}} \min_{\boldsymbol{\phi}_{P}} \mathcal{L}^{P}(\boldsymbol{\theta}_{P}, \boldsymbol{\phi}_{P}) + \lambda r(\boldsymbol{\theta}_{D}, \boldsymbol{\theta}_{P})$$
 (3)

171 with $\gamma \in \mathbb{R}_+$ and r being some convex regularisation objective weighted by $\lambda \in \mathbb{R}_+$ enforcing 172 similarity between θ_D and θ_P . The upper-level training objective in equation 2 is tasked with minimizing the downstream task objective \mathcal{L}^{P} , while the lower-level objective in equation 3 aims to minimize the pretext task objective \mathcal{L}^{P} while also ensuring its backbone remains similar to the 173 174 175 upper-level backbone. As seen in the left term of equation 2, the backbone parameters $\theta_{P}^{*}(\theta_{D})$ 176 are transferred into the downstream training objective, mirroring how the backbone is transferred in 177 the conventional SSL pipeline. Although the second term of equation 2 is not strictly necessary, it 178 has empirically shown to improve stability and aid convergence of the upper-level solution during 179 training. Unlike the traditional SSL setup, the backbone solution of the pretext objective $\theta_{P}^{*}(\theta_{D})$ is now a function of the parameters of the downstream backbone θ_D , as the lower-level problem is dependent on the upper-level backbone parameters. 181

182 As the upper-level objective in equation 2 depends on the solution $\theta_P^*(\theta_D)$ of the lower-level ob-183 jective in equation 3, this enables the incorporation of information from the pretext objective when 184 solving the upper-level optimization problem. By including a regularization objective r that enforces 185 similarity between the lower-level and upper-level backbone parameters, this setup is hypothesized to guide the lower-level to achieve a configuration of model backbone parameters that is more ben-186 eficial for subsequent conventional fine-tuning on the downstream task. To more precisely under-187 stand how the pretext objective influences the downstream training procedure in this setup, we delve 188 deeper into the expression of the gradient of the upper-level training objective in equation 2 in the 189 following subsection. 190

3.3 UPPER-LEVEL DERIVATIVE

Given the upper-level objective $F(\theta_D, \phi_D) := \mathcal{L}^D(\theta_P^*(\theta_D), \phi_D) + \gamma \mathcal{L}^D(\theta_D, \phi_D)$ from equation 2, its derivative with respect to θ_D is given by

$$\frac{\mathrm{d}F}{\mathrm{d}\theta_D} = \underbrace{\frac{\mathrm{d}\theta_P^*(\theta_D)}{\mathrm{d}\theta_D}}_{\mathrm{IG}} \nabla_{\theta} \mathcal{L}^D(\theta, \phi_D)|_{\theta = \theta_P^*(\theta_D)} + \gamma \nabla_{\theta} \mathcal{L}^D(\theta, \phi_D)|_{\theta = \theta_D}. \tag{4}$$

Due to the dependence of the lower-level solution on the upper-level parameters, the first term of equation 4 includes the implicit gradient (IG) of the implicit function $\theta_P^*(\theta_D)$. To simplify notation, we let $\nabla_{\boldsymbol{\xi}} h(\boldsymbol{\xi})|_{\boldsymbol{\xi}=\boldsymbol{\psi}} := \nabla_{\boldsymbol{\xi}} h(\boldsymbol{\psi})$ when it is clear from context which variables are differentiated with respect to. Following an approach similar to Rajeswaran et al. (2019), with details on the derivations and underlying assumptions outlined in Section A.1 of Appendix A, the IG in equation 4 can be explicitly expressed as

$$\frac{\mathrm{d}\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})}{\mathrm{d}\boldsymbol{\theta}_{D}}^{T} = -\nabla_{\boldsymbol{\theta}_{D}\boldsymbol{\theta}_{P}}^{2} r(\boldsymbol{\theta}_{D}, \boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})) \left[\nabla_{\boldsymbol{\theta}}^{2} \left(\frac{1}{\lambda} \mathcal{L}^{P}(\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D}), \boldsymbol{\phi}_{P}) + r(\boldsymbol{\theta}_{D}, \boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})) \right) \right]^{-1}.$$
 (5)

A common convex regularization objective, which will also be the choice in the subsequent experiments of this work, is $r(\boldsymbol{\xi}, \boldsymbol{\psi}) = \frac{1}{2} \|\boldsymbol{\xi} - \boldsymbol{\psi}\|_2^2$. Using this regularization objective simplifies equation 5 down to

$$\frac{\mathrm{d}\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})}{\mathrm{d}\boldsymbol{\theta}_{D}}^{T} = \left[\frac{1}{\lambda}\nabla_{\boldsymbol{\theta}}^{2}\mathcal{L}^{P}(\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D}),\boldsymbol{\phi}_{P}) + I_{L}\right]^{-1},\tag{6}$$

where I_L is the $L \times L$ -dimensional identity matrix. Hence the upper-level derivative in equation 4 can be expressed as

$$\frac{\mathrm{d}F}{\mathrm{d}\boldsymbol{\theta}_D} = \left[\frac{1}{\lambda}\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}^P(\boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D), \boldsymbol{\phi}_P) + I_L\right]^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{L}^D(\boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D), \boldsymbol{\phi}_D) + \gamma \nabla_{\boldsymbol{\theta}} \mathcal{L}^D(\boldsymbol{\theta}_D, \boldsymbol{\phi}_D).$$
(7)

216 The inverse Hessian-vector product in the left term of equation 7 is computationally infeasible to 217 calculate directly, hence it is approximated using the conjugate gradient (CG) method (Nazareth, 218 2009; Shewchuk, 1994). While CG is established as a successful approach for approximating the 219 inverse Hessian-vector products in previous works (Pedregosa, 2016; Zhang et al., 2021; Rajeswaran 220 et al., 2019), it still introduces significant computational overhead due to its need for iterative evaluations of multiple Hessian vector products. Future work may explore alternative methods that offer 221 more efficient approximations without compromising downstream task performance. We employ a 222 layer-wise implementation of the CG method based on that of Rajeswaran et al. (2019) and refer to 223 their work for more details on applying CG in a deep learning setup with BLO. For a comprehensive 224 overview of other common methods used to approximate the upper-level derivative in BLO, we refer 225 to Zhang et al. (2023a). 226

With an explicit expression of the IG in equation 6, we can interpret the impact of the scaling factor 227 λ from equation 3 and equation 7: When λ is very large, the dependence of lower-level objective on 228 the upper-level parameters θ_D is also very large. This effectively drives the lower-level backbone 229 parameters toward the trivial solution $\theta_P^*(\theta_D) = \theta_D$. Meanwhile, the IG in equation 6 approxi-230 mately equals I_L , thereby diminishing the influence of the lower-level objective on the upper-level 231 gradient in equation 7. This roughly makes the task of the upper-level equivalent to conventional 232 fine-tuning. Conversely, if λ is very small, the lower-level objective in equation 3 effectively defaults 233 to conventional pretext task training. Additionally, the implicit gradient in equation 6 would consist 234 of numerically tiny entries, making the optimization of the first term in the upper-level objective in 235 equation 2 equivalent to probing of the downstream head on the frozen pretext backbone $\theta_P^*(\theta_D)$. 236

3.4 TRAINING ALGORITHM AND PIPELINE

Algorithm 1 BiSSL Training Algorithm

1: Input: Backbone and projection head parameter initializations θ , ϕ_P , ϕ_D . Training objectives $\mathcal{L}^{\bar{P}}, \mathcal{L}^{D}$. Weights $\lambda, \bar{\gamma} \in \mathbb{R}_{+}$. Optimizers opt_P, opt_D. Number of training stage alternations $T \in \mathbb{N}$ with upper and lower-level iterations $N_U, N_L \in \mathbb{N}$. Upper-level backbone adaption frequency $N_a \in \mathbb{N}$ and strength $\alpha \in [0, 1]$.

▷ Lower-level

2: Initialize
$$\theta_P \leftarrow \theta$$
 and $\theta_D \leftarrow \theta$.

3: for t = 1, ..., T do for $n = 1, ..., N_L$ do

4:

Compute $\mathbf{g}_{\boldsymbol{\phi}_P} = \nabla_{\boldsymbol{\phi}} \mathcal{L}^P(\boldsymbol{\theta}_P, \boldsymbol{\phi})|_{\boldsymbol{\phi} = \boldsymbol{\phi}_P}$. 5: Compute $\mathbf{g}_{\boldsymbol{\theta}_P} = \nabla_{\boldsymbol{\theta}} \mathcal{L}^P(\boldsymbol{\theta}, \boldsymbol{\phi}_P)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_P} + \lambda \nabla_{\boldsymbol{\theta}} r(\boldsymbol{\theta}_D, \boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_P}$. 6:

Update $\phi_P \leftarrow \operatorname{opt}_P(\phi_P, \mathbf{g}_{\phi_P})$ and $\theta_P \leftarrow \operatorname{opt}_P(\theta_P, \mathbf{g}_{\theta_P})$. 7:

if $t \mod N_a \equiv 0$ then 8:

 $\boldsymbol{\theta}_D \leftarrow (1-\alpha)\boldsymbol{\theta}_D + \alpha \boldsymbol{\theta}_P.$ 9:

10: for $n = 1, ..., N_U$ do ▷ Upper-level Compute $\mathbf{g}_{\phi_D} = \nabla_{\boldsymbol{\phi}} \mathcal{L}^D(\boldsymbol{\theta}_P, \boldsymbol{\phi})|_{\boldsymbol{\phi}=\phi_D} + \gamma \nabla_{\boldsymbol{\phi}} \mathcal{L}^D(\boldsymbol{\theta}_D, \boldsymbol{\phi})|_{\boldsymbol{\phi}=\phi_D}.$ Compute $\mathbf{v} = \nabla_{\boldsymbol{\theta}} \mathcal{L}^D(\boldsymbol{\theta}, \phi_D)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_P}.$ 11: 12: Approximate $\mathbf{v}_{IG} \approx \left[I_M + \frac{1}{\lambda} \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}^P(\boldsymbol{\theta}, \boldsymbol{\phi}_P)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_P}\right]^{-1} \mathbf{v}.$ Compute $\mathbf{g}_{\boldsymbol{\theta}_D} = \mathbf{v}_{IG} + \gamma \nabla_{\boldsymbol{\theta}} \mathcal{L}^D(\boldsymbol{\theta}, \boldsymbol{\phi}_D)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_D}.$ Update $\boldsymbol{\phi}_D \leftarrow \operatorname{opt}_D(\boldsymbol{\phi}_D, \mathbf{g}_{\boldsymbol{\phi}_D})$ and $\boldsymbol{\theta}_D \leftarrow \operatorname{opt}_D(\boldsymbol{\theta}_D, \mathbf{g}_{\boldsymbol{\theta}_D}).$ 13: ⊳ Use CG 14: 15:

16: **Return:** Backbone Parameters θ_P .

263 264

237

238 239

240

241

242

243

244 245

246

247

248

249

250

251

253

254 255

256

257

258

259

260

261 262

265 Algorithm 1 outlines the proposed training algorithm, which iteratively alternates between solv-266 ing the lower-level (equation 3) and upper-level (equation 2) optimization problems in BiSSL. The 267 lower-level training optimizes the pretext task objective, while additionally including the gradient of the regularization term r for the backbone parameter updates, complying with equation 3. For the 268 upper-level training, the gradient with respect to the backbone parameters as represented by the left 269 term on the right-hand side in equation 7, is approximated using the CG method. Additionally, the

pretext backbone parameters θ_P are weighted by α and added to the downstream backbone parameters θ_D every N_a alternations to further enforce similarity between them, which empirically has shown to aid convergence during training.

273 From Section A.1 in Appendix A, we get that $\theta_P^*(\theta_D)$ must fulfill the stationary condition 274 $\nabla_{\boldsymbol{\theta}} \left(\mathcal{L}^{P}(\boldsymbol{\theta}, \phi_{P}) + \lambda r(\boldsymbol{\theta}_{D}, \boldsymbol{\theta}) \right) |_{\boldsymbol{\theta} = \boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})} = \mathbf{0}$ to justify the explicit expression of the implicit gra-275 dient in equation 6. This means that executing Algorithm 1 using random initializations of θ and 276 ϕ_P will likely not suffice. The same applies to ϕ_D , as a random initialization of ϕ_D typically leads 277 to rapid initial changes of the backbone parameters θ_D during fine-tuning. This would then likely 278 violate the assumed stationary condition due to the dependence between θ_D and θ_P through the reg-279 ularization objective r. Figure 1 illustrates the suggested pipeline alongside the conventional SSL 280 pipeline. First conventional pretext pre-training is performed on the unlabeled dataset \mathcal{D}^P to obtain 281 initializations of θ and ϕ_P . Next, the downstream head is fitted on top of the frozen backbone θ us-282 ing the downstream dataset \mathcal{D}^D , which provides an initialization of the downstream head parameters ϕ_D . Then, BiSSL training is conducted as outlined in Algorithm 1, yielding an updated configura-283 tion of backbone parameters $\theta_{P}^{*}(\theta_{D})$. These updated backbone parameters are subsequently used 284 as an initialization for the final supervised fine-tuning on the downstream task. 285

286 287

4 EXPERIMENTS AND RESULTS

288 289 290

4.1 DATASETS

291 The STL10 dataset (Coates et al., 2011) is used throughout the experiments. It comprises two par-292 titions: 100.000 unlabeled images and 13.000 labeled images with 10 classes in total whereas 5000 293 and 8000 are assigned for training and testing, respectively. All images are natural images of resolution 96×96 , with the unlabeled partition drawn from a similar but broader distribution than the 295 labeled partition. This dataset strikes a balance between complexity and computational feasibility, 296 offering higher resolution and more diverse content than smaller datasets like CIFAR10 (Krizhevsky, 297 2012) while being less resource-intensive than larger-scale datasets such as ImageNet (Deng et al., 298 2009). For ease of reference, STL10U and STL10L will denote the unlabeled and labeled parti-299 tions, respectively. In all experiments, STL10U will be employed for self-supervised pre-training. For downstream fine-tuning and evaluation, we leverage a varied set of natural image classification 300 datasets that encompass a wide array of tasks, including general image classification, fine-grained 301 recognition across species and objects, scene understanding, and texture categorization. The datasets 302 include STL10L, Oxford 102 Flowers (Nilsback & Zisserman, 2008), StanfordCars (Yang et al., 303 2015), FGVC Aircraft (Maji et al., 2013), Describable Textures Dataset (DTD) (Cimpoi et al., 2014), 304 Oxford-IIIT Pets (Parkhi et al., 2012), FashionMNIST (Xiao et al., 2017), CIFAR10 (Krizhevsky, 305 2012), CIFAR100 (Krizhevsky, 2012), Caltech-101 (Li et al., 2022a), Food 101 (Bossard et al., 306 2014), SUN397 scene dataset (Xiao et al., 2010), Caltech-UCSD Birds-200-2011 (CUB200) (Wah 307 et al., 2011) and PASCAL VOC 2007 (Everingham et al.). All downstream datasets are split into 308 training, validation, and test partitions, with details on how these assignments are made provided in 309 Section B.1 of Appendix B.

- 310
- 311 4.2 IMPLEMENTATION DETAILS312
- 313 4.2.1 BASELINE SETUP 314

Pretext Task Training The SimCLR (Chen et al., 2020b) pretext task with temperature $\tau = 0.5$ is used for pre-training a ResNet-18 (He et al., 2016) backbone model. We selected this widely adopted architecture due to its proven ability to extract high-quality visual representations while maintaining relatively low computational requirements, striking an effective balance between performance and resource efficiency. On top of the backbone, a projection head is used, consisting of two fully connected layers with batch normalization (Ioffe & Szegedy, 2015) and ReLU (Agarap, 2018) followed by a single linear layer. Each layer consists of 256 neurons.

The image augmentation scheme follows the approach used in Bardes et al. (2022), with minor modifications: The image size is set to 96×96 instead of 224×224 , and the minimal ratio of the random crop is adjusted accordingly to 0.5 instead of 0.08. 324 The implementation of the LARS optimizer (You et al., 2017) from Bardes et al. (2022) is employed, 325 with a "trust" coefficient of 0.001, a weight decay of 10^{-6} and a momentum of 0.9. The learning rate 326 increases linearly during the first 10 epochs, reaching a peak base learning rate of 4.8, followed by 327 a cosine decay towards 0 with no restarts (Loshchilov & Hutter, 2017) for the remaining epochs. A 328 batch size of 1024 is used and, unless otherwise specified, pre-training is conducted for 600 epochs.

330 **Fine-Tuning on the Downstream Task** For downstream fine-tuning, a single linear layer is attached to the output of the pre-trained backbone. The training procedure utilizes the cross-entropy 331 loss, the SGD optimizer with a momentum of 0.9, and a cosine decaying learning rate scheduler 332 without restarts (Loshchilov & Hutter, 2017). Fine-tuning is conducted for 400 epochs with a batch 333 size of 256. An augmentation scheme similar to the fine-tuning augmentation scheme in Bardes 334 et al. (2022) is employed, where images are center cropped and resized to 96×96 pixels with a 335 minimal crop ratio of 0.5, followed by random horizontal flips. 336

A random grid search of 200 hyper-parameter configurations for the learning rates and weight decays 337 is conducted, where one model is fine-tuned for each configuration. Base learning rates and weight 338 decays are log-uniformly sampled over the ranges of 0.0001 to 1.0 and 0.00001 to 0.01, respectively. 339 Validation data accuracy is evaluated after each epoch. The hyper-parameter configuration yielding 340 the best balance between high validation accuracy and low validation loss is considered the optimal 341 hyper-parameter configuration.² The corresponding optimal hyper-parameters for each downstream 342 dataset are documented in Table 2 of Appendix B. 343

For subsequent evaluation on the test data, we train 10 models with different random seeds, each 344 using the considered optimal hyper-parameter configurations. During the training of each respective 345 model, the model parameters are stored after each epoch if the top-1 validation accuracy (or 11-346 point mAP for the VOC07 dataset) has increased compared to the previous highest top-1 validation 347 accuracy achieved during training. Top-1 and top-5 test data accuracies (or 11-point mAP for the 348 VOC07 dataset) are evaluated for each of the 10 models, from which the calculated means and 349 standard deviations of these accuracies are documented. 350

4.2.2 BISSL SETUP 352

353 In this section, we detail each stage of the proposed training pipeline for BiSSL, as outlined in the right part of Figure 1. 354

355

351

329

Pretext Warm-up The backbone θ and projection head ϕ_P are initialized by self-supervised pre-356 training using a setup almost identical to the baseline pretext task training setup in Section 4.2.1. 357 The only difference is that this training stage is conducted for 500 epochs instead of 600 epochs, 358 and that the peak base learning rate is set to 1.0 instead of 4.8. This adjustment is made because the 359 BiSSL training stage will conduct what is roughly equivalent to 100 pretext epochs, as detailed more 360 specifically in the composite configuration paragraph below. This ensures that the total number of 361 pretext pre-training steps is comparable to those conducted in the baseline setup. 362

363 **Downstream Head Warm-up** The training setup for the downstream head warm-up closely mir-364 rors the fine-tuning setup of Section 4.2.1. The main difference is that only the linear downstream 365 head is fitted on top of the now frozen backbone obtained from the pretext warm-up. Learning rates 366 and weight decays are initially selected based on those listed in Table 2, with adjustments made as needed when preliminary testing indicated a potential for improved convergence. These values 367 are provided in Table 3 in Appendix B. The authors recognize that more optimal hyper-parameter 368 configurations may exist and leave further exploration of this for future refinement. The downstream 369 head warm-up is conducted for 20 epochs with a constant learning rate. 370

371 **Lower-level of BiSSL** The training configuration for the lower-level primarily follows the setup 372 described for pretext pre-training in Section 4.2.1, with the modifications outlined here. As specified 373 in equation 3, the lower-level loss function is the sum of the pretext task objective \mathcal{L}^P (in our case, 374 the NT-Xent loss from SimCLR (Chen et al., 2020b)) and the regularization term $r(\theta_D, \theta_P) =$ 375

³⁷⁶ ²In certain scenarios during the experiments, the configuration that achieved the highest validation accuracy 377 also yielded a notably higher relative validation loss. To ensure better generalizability, an alternative configuration with a more favorable trade-off was selected in these cases.

Dataset	Top-1 Accuracy (*: 11-point mAP) Top-		p-5 Accuracy			
	BiSSL	Only FT	Avg Diff	BiSSL	Only FT	Avg Diff
STL10L	90.2 ± 0.1	90.3 ± 0.1	-0.1	99.7 ± 0.0	99.6 ± 0.0	+0.1
Flowers	74.8 ± 0.2	73.4 ± 0.4	+1.4	89.8 ± 0.3	90.0 ± 0.4	-0.2
Cars	73.0 ± 0.4	72.7 ± 0.5	+0.3	91.5 ± 0.3	91.4 ± 0.4	+0.1
Aircrafts	46.9 ± 0.5	46.1 ± 0.9	+0.8	78.9 ± 0.4	79.3 ± 0.6	-0.4
DTD	51.8 ± 0.5	49.3 ± 0.5	+2.5	79.9 ± 0.3	79.1 ± 0.4	+0.8
Pets	67.8 ± 0.2	65.0 ± 0.5	+2.8	92.3 ± 0.3	90.7 ± 0.3	+1.6
FMNIST	94.3 ± 0.2	94.1 ± 0.1	+0.2	100.0 ± 0.0	100.0 ± 0.0	0.0
CIFAR10	93.9 ± 0.1	93.8 ± 0.1	+0.1	99.9 ± 0.0	99.8 ± 0.0	+0.1
CIFAR100	73.0 ± 0.1	73.2 ± 0.2	-0.2	93.7 ± 0.1	92.8 ± 0.1	+0.9
Caltech-101	80.6 ± 0.7	78.1 ± 0.5	+2.5	95.5 ± 0.2	94.7 ± 0.2	+0.8
Food	72.0 ± 0.2	71.7 ± 0.2	+0.3	90.4 ± 0.1	90.4 ± 0.1	0.0
SUN397	41.1 ± 0.2	40.0 ± 0.3	+1.1	71.0 ± 0.2	69.9 ± 0.4	+1.1
CUB200	47.1 ± 0.4	45.7 ± 0.4	+1.4	72.1 ± 0.3	70.7 ± 0.6	+1.4
VOC07	$*60.4\pm0.1$	$*58.6\pm0.3$	+1.8	_	_	-

Table 1: Test classification accuracies. Accuracies significantly different from their counterparts are
 marked in bold font.

 $\frac{1}{2}||\boldsymbol{\theta}_D - \boldsymbol{\theta}_P||_2^2$. Based on early experiments, the regularization weight $\lambda = 0.001$ was selected, as it appeared to strike a well-balanced compromise between the convergence rates of both the lower- and upper-level objectives. The lower-level is trained for the equivalent of approximately 100 conventional pre-training epochs, with further details provided in the composite configuration paragraph. Each time the BiSSL training alternates back to the lower-level, the first 5 batches used for lower-level training are stored. These stored batches are utilized to approximate the Hessian of the lower-level objective when approximating the upper-level gradient. Further details are specified in Section B.3 of Appendix B and the paragraph below.

Composite Configuration Details of BiSSL As outlined in Algorithm 1, both lower- and upper-level backbone parameters θ_P and θ_D are initialized with the backbone parameters obtained during the pretext warm-up, and the training procedure alternates between solving the lower- and upper-level optimization problems. In this experimental setup, the lower-level performs $N_L = 20$ gradient steps before alternating to the upper-level, which then conducts $N_U = 8$ gradient steps. A total of T = 500 training stage alternations are executed. As the STL10U dataset with the current batch size of 1024 amounts to a total of 98 training batches without replacement, these T = 500 training stage alternations roughly equal 100 conventional pretext epochs. Section B.4 in Appendix B outlines further details on how data batches are handled during training. The upper-level backbone adaptation frequency and strength are set to $N_a = 100$ and $\alpha = 0.1$, respectively. Additionally, gradient normalization is employed on gradients exceeding ℓ_2 -norms of 10.

Fine-Tuning on the Downstream Task Subsequent downstream fine-tuning is conducted in a manner identical to that described in the 'Fine-Tuning on the Downstream Task' paragraph of section 4.2.1. Table 4 in Appendix B lists the considered optimal hyper-parameter configurations for each dataset.



Figure 2: Test classification accuracies on the Flowers dataset for separate models pre-trained for different durations, comparing the conventional and BiSSL training pipelines. BiSSL consistently achieves higher top-1 accuracy than the baseline after sufficient pre-training.

4.3 DOWNSTREAM TASK PERFORMANCE

454 The impact of using BiSSL compared to the conventional self-supervised training pipeline is bench-455 marked by evaluating classification accuracies on the various specified downstream datasets. Ta-456 ble 1 presents the means and standard deviations of top-1 and top-5 classification accuracies (or 457 the 11-point mAP on the VOC2007 dataset) on these downstream test datasets, comparing results 458 obtained from the conventional SSL pipeline with those achieved using the BiSSL pipeline. The results demonstrate that training with BiSSL significantly improves either top-1 or top-5 classification 459 accuracy, or 11-point mAP in the case of VOC07, on 10 out of 14 datasets, with no single result 460 showing a significant decline in performance compared to the baseline. 461

462

448

449

450 451 452

453

4.3.1 PERFORMANCE OVER VARYING PRE-TRAINING EPOCHS

464 To further assess the robustness of BiSSL, we conduct experiments varying the duration of self-465 supervised pre-training while keeping the remaining experimental setup unchanged. The Flowers 466 dataset, which demonstrated substantial benefits from BiSSL in terms of top-1 classification ac-467 curacies, was chosen for these experiments. We continue to compare models with BiSSL applied 468 after 100 fewer pre-training epochs than the baseline, ensuring the total pretext training duration 469 is approximately equal across methods, as described in Section 4.2.2. Figure 2 depicts the final fine-tuned test accuracy achieved by separate models pre-trained for varying durations. The results 470 indicate that BiSSL consistently outperforms the baseline once a sufficient duration of pre-training is 471 reached. This suggests that the BiSSL's benefits are not contingent on the amount of pre-training, but 472 rather that it provides a more efficient learning trajectory, stemming from the enhanced information 473 sharing it facilitates between the pretext and downstream tasks. 474

475 476

4.4 VISUAL INSPECTION OF LATENT FEATURES

477 To gain deeper insight into how BiSSL affects the representations learned compared to conventional 478 pretext pre-training, we perform a qualitative visual inspection of latent spaces. This involves com-479 paring features processed by backbones trained solely by pretext pre-training to those derived from 480 lower-level backbones obtained after conducting BiSSL, each trained as described in the "Pretext 481 Task Training" and "Lower-level of BiSSL" paragraphs in Section 4.2.1 and 4.2.2, respectively. By 482 comparing these features, we aim to assess whether BiSSL nudges the latent features toward be-483 ing more semantically meaningful for the downstream task. The t-Distributed Stochastic Neighbor Embedding (t-SNE) (Cieslak et al., 2020) technique is employed for dimensionality reduction. Fur-484 ther details regarding the experimental setup are outlined in Section C.1 of Appendix C. Figure 3 485 illustrates the results on the flowers dataset, indicating that BiSSL yields backbones with improved



Figure 3: Visualization of features from backbones trained using pretext pre-training exclusively and backbones derived from lower-level backbones obtained after applying BiSSL, respectively. Features are extracted from the test partition of the flowers dataset. Each color represents a different class. Details are outlined in Section C.1 of Appendix C

downstream feature alignment. Further plots on a selection of downstream datasets in Section C.1 reinforce this finding, also demonstrating that this trend persists even for datasets where BiSSL did not impose any classification accuracy improvements.

5 CONCLUSION

511 512

499 500

501

502

503 504 505

506

507

508 509 510

513 This study integrates pretext pre-training and downstream fine-tuning into a unified bilevel optimiza-514 tion problem, from which the BiSSL training framework is proposed. BiSSL explicitly models the 515 inheritance of backbone parameters from the pretext task, enhancing the transfer of relevant information between the pretext and downstream tasks. We propose a practical training algorithm and 516 pipeline that incorporates BiSSL as an intermediate stage between pretext pre-training and down-517 stream fine-tuning. Experiments across various image classification datasets demonstrate that BiSSL 518 consistently achieves improved or comparable downstream classification performance relative to the 519 conventional self-supervised learning pipeline. Additionally, our findings indicate that in instances 520 where BiSSL improves performance, this improvement remains consistent regardless of the pre-521 text pre-training duration. Further analysis suggests that BiSSL enhances the downstream semantic 522 richness of learned representations, as evidenced by qualitative inspections of latent spaces. BiSSL 523 marks a potential advancement towards enhancing the alignment between the pretext pre-training 524 and downstream fine-tuning stages, revealing a new direction for self-supervised learning algorithm 525 designs that leverage bilevel optimization.

526 527 528

5.1 FUTURE WORK

529 Formulating the self-supervised pipeline as a bilevel optimization problem offers various strategies 530 with trade-offs in computational complexity and theoretical justification. While this study presents 531 a promising approach for improving downstream performance, further investigation of alternative 532 formulations is needed to identify setups are are potentially more optimal. Although BiSSL is theo-533 retically applicable to any downstream task and model size, our experiments focused on small-scale 534 image classification due to resource constraints. Therefore, it remains uncertain whether BiSSL 535 can scale to larger setups and tasks. Additionally, a potential future advancement would integrating 536 more novel methods for solving BLO problems, which promise benefits in terms of reduced com-537 putational costs and improved solution convergence (Zhang et al., 2023a; Yang et al., 2021; Choe et al., 2023; Huang, 2024). Lastly, the current BiSSL framework relies on full access to pre-training 538 data and pretext tasks. Future research could investigate the use of only a subset of pre-training data and alternative pretext tasks to maintain BiSSL's benefits under these conditions.

540 REFERENCES

550

564 565

566

575

583

- Abien Fred Agarap. Deep learning using rectified linear units (relu). arXiv preprint
 arXiv:1803.08375, 2018.
- Martín Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
 ArXiv, abs/1907.02893, 2019.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A frame work for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022.
- Harshita Boonlia, Tanmoy Dam, Md Meftahul Ferdaus, Sreenatha G Anavatti, and Ankan Mullick.
 Improving self-supervised learning for out-of-distribution task via auxiliary classifier. In 2022
 IEEE International Conference on Image Processing (ICIP), pp. 3036–3040. IEEE, 2022.
- Zalán Borsos, Mojmir Mutny, and Andreas Krause. Coresets via bilevel optimization for continual
 learning and streaming. *Advances in neural information processing systems*, 33:14879–14890,
 2020.
- Matteo Boschini, Lorenzo Bonicelli, Angelo Porrello, Giovanni Bellitto, Matteo Pennisi, Simone Palazzo, Concetto Spampinato, and Simone Calderara. Transfer without forgetting. In *European Conference on Computer Vision*, pp. 692–709. Springer, 2022.
 - Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.
- 567 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 568 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel 569 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, 570 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, 571 Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, 572 and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, 573 R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, 574 volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jegou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In 2021 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9630–9640, 2021. doi: 10.1109/ICCV48922.2021.00951.
- Heng-Jui Chang, Alexander H. Liu, and James Glass. Self-supervised Fine-tuning for Improved
 Content Representations by Speaker-invariant Clustering. In *INTERSPEECH 2023*, pp. 2983–2987. ISCA, August 2023. doi: 10.21437/Interspeech.2023-847.
- Can (Sam) Chen, Jingbo Zhou, Fan Wang, Xue Liu, and Dejing Dou. Structure-aware protein self-supervised learning. *Bioinformatics*, 39(4):btad189, 04 2023. ISSN 1367-4811. doi: 10.1093/bioinformatics/btad189.
- Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, Ting Liu, and Xiangzhan Yu. Recall and
 learn: Fine-tuning deep pretrained language models with less forgetting. In Bonnie Webber,
 Trevor Cohn, Yulan He, and Yang Liu (eds.), *EMNLP (1)*, pp. 7870–7881. Association for Computational Linguistics, 2020a. ISBN 978-1-952148-60-6.

- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020b.
- Sang Keun Choe, Willie Neiswanger, Pengtao Xie, and Eric Xing. Betty: An automatic differentiation library for multilevel optimization. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=LV_MeMS38Q9.
- Yu-An Chung and James Glass. Speech2Vec: A Sequence-to-Sequence Framework for Learning
 Word Embeddings from Speech. In *Proc. Interspeech 2018*, pp. 811–815, 2018. doi: 10.21437/
 Interspeech.2018-2341.
- Yu-An Chung, Wei-Ning Hsu, Hao Tang, and James Glass. An Unsupervised Autoregressive Model
 for Speech Representation Learning. In *Proc. Interspeech 2019*, pp. 146–150, 2019. doi: 10.
 21437/Interspeech.2019-1473.
- Matthew C. Cieslak, Ann M. Castelfranco, Vittoria Roncalli, Petra H. Lenz, and Daniel K. Hartline.
 t-distributed stochastic neighbor embedding (t-sne): A tool for eco-physiological transcriptomic analysis. *Marine Genomics*, 51:100723, 2020. ISSN 1874-7787. doi: https://doi.org/10.1016/j. margen.2019.100723.
- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In
 Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2014.
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík (eds.), *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pp. 215–223, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier archical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition,
 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019.
- Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Dar rell. Decaf: A deep convolutional activation feature for generic visual recognition. In Eric P. Xing
 and Tony Jebara (eds.), *Proceedings of the 31st International Conference on Machine Learning*,
 volume 32 of *Proceedings of Machine Learning Research*, pp. 647–655, Bejing, China, 22–24
 Jun 2014. PMLR.
- A.L. Dontchev and R.T. Rockafellar. *Implicit Functions and Solution Mappings: A View from Variational Analysis*. Springer Series in Operations Research and Financial Engineering. Springer New York, 2014. ISBN 9781493910373.
 - Yann Dubois, Tatsunori Hashimoto, Stefano Ermon, and Percy Liang. Improving self-supervised learning by characterizing idealized representations. *ArXiv*, abs/2209.06235, 2022.
- M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html.
- 643
 644
 644
 645
 645
 Chen Fan, Parikshit Ram, and Sijia Liu. Sign-MAML: Efficient model-agnostic meta-learning by signSGD. In *Fifth Workshop on Meta-Learning at the Conference on Neural Information Processing Systems*, 2021.
- 646 647

637

638

624

625

626

627

Chelsea Finn, P. Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 2017.

677

686

687

689

690 691

- 648 Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazzi, and Massimiliano Pontil. Bilevel 649 programming for hyperparameter optimization and meta-learning. In International conference on 650 machine learning, pp. 1568–1577. PMLR, 2018. 651
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena 652 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, 653 et al. Bootstrap your own latent-a new approach to self-supervised learning. Advances in neural 654 information processing systems, 33:21271–21284, 2020. 655
- 656 Kartik Gupta, Thalaiyasingam Ajanthan, Anton van den Hengel, and Stephen Gould. Understanding 657 and improving the role of projection head in self-supervised learning. ArXiv, abs/2212.11491, 658 2022. 659
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-660 nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 661 770–778, 2016. 662
- 663 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for 664 unsupervised visual representation learning. In Proceedings of the IEEE/CVF conference on 665 computer vision and pattern recognition, pp. 9729-9738, 2020. 666
- 667 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In Proceedings of the IEEE/CVF conference on computer 668 vision and pattern recognition, pp. 16000-16009, 2022. 669
- 670 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert 671 with disentangled attention. In International Conference on Learning Representations, 2021. 672
- 673 Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, 674 and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked 675 prediction of hidden units. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 676 29:3451-3460, 2021.
- Feihu Huang. Optimal Hessian/Jacobian-free nonconvex-PL bilevel optimization. In Ruslan 678 Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and 679 Felix Berkenkamp (eds.), Proceedings of the 41st International Conference on Machine Learn-680 ing, volume 235 of Proceedings of Machine Learning Research, pp. 19598–19621. PMLR, 21–27 681 Jul 2024. 682
- 683 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In International conference on machine learning, pp. 448–456. 684 pmlr, 2015. 685
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. 688 Segment anything. 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 3992-4003, 2023.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. University of Toronto, 05 692 2012.
- 693 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdel rahman Mohamed, Omer 694 Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-695 training for natural language generation, translation, and comprehension. In Annual Meeting 696 of the Association for Computational Linguistics, 2019. 697
 - Fei-Fei Li, Marco Andreeto, Marc'Aurelio Ranzato, and Pietro Perona. Caltech 101, Apr 2022a.
- Jiangmeng Li, Wenwen Qiang, Yanan Zhang, Wenyi Mo, Changwen Zheng, Bing Su, and Hui 700 Xiong. Metamask: Revisiting dimensional confounder for self-supervised learning. Advances in 701 Neural Information Processing Systems, 35:38501–38515, 2022b.

702 703 704	Xingjian Li, Haoyi Xiong, Haozhe An, Cheng-Zhong Xu, and Dejing Dou. Rifle: Backpropagation in depth for deep transfer learning through re-initializing the fully-connected layer. In <i>International Conference on Machine Learning</i> , pp. 6010–6019. PMLR, 2020.
705 706 707	Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In Inter- national Conference on Learning Representations, 2017.
708 709	S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013.
710 711 712	J. L. Nazareth. Conjugate gradient method. <i>WIREs Computational Statistics</i> , 1(3):348–353, 2009. doi: https://doi.org/10.1002/wics.13.
713 714 715	Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Byol for audio: Self-supervised learning for general-purpose audio representation. In 2021 International Joint Conference on Neural Networks (IJCNN), pp. 1–8. IEEE, 2021.
716 717 718	Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In <i>Indian Conference on Computer Vision, Graphics and Image Processing</i> , Dec 2008.
719 720 721 722 723 724	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without super- vision. <i>Transactions on Machine Learning Research</i> , 2024. ISSN 2835-8856.
725 726 727	Haolin Pan, Yong Guo, Qinyi Deng, Hao-Fan Yang, Yiqun Chen, and Jian Chen. Improving fine- tuning of self-supervised models with contrastive initialization. <i>Neural networks : the official</i> <i>journal of the International Neural Network Society</i> , 159:198–207, 2022.
728 729 730	Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , 2012.
731 732 733 734	Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In Maria Florina Bal- can and Kilian Q. Weinberger (eds.), <i>Proceedings of The 33rd International Conference on Ma- chine Learning</i> , volume 48 of <i>Proceedings of Machine Learning Research</i> , pp. 737–746, New York, New York, USA, 20–22 Jun 2016. PMLR.
735 736 737	Aravind Rajeswaran, Chelsea Finn, Sham M. Kakade, and Sergey Levine. Meta-learning with implicit gradients. In <i>Neural Information Processing Systems</i> , 2019.
738 739 740	Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli. wav2vec: Unsupervised Pre-Training for Speech Recognition. In <i>Proc. Interspeech 2019</i> , pp. 3465–3469, 2019. doi: 10.21437/Interspeech.2019-1873.
741 742	Jonathan R Shewchuk. An introduction to the conjugate gradient method without the agonizing pain. Technical report, Carnegie Mellon University, 1994.
743 744 745	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014.
746 747 748 749 750	Sai Ashish Somayajula, Lifeng Jin, Linfeng Song, Haitao Mi, and Dong Yu. Bi-level finetuning with task-dependent similarity structure for low-resource training. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pp. 8569–8588, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.544.
751 752 753 754	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
755	C Wah S Branson P Welinder P Perona and S Belongie. The caltech-uosd birds-200-2011

755 C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

- Chuhan Wu, Fangzhao Wu, Tao Qi, Yongfeng Huang, and Xing Xie. Noisytune: A little noise can help you finetune pretrained language models better. In *Annual Meeting of the Association for Computational Linguistics*, 2022.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
 Large-scale scene recognition from abbey to zoo. In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 3485–3492, 2010. doi: 10.1109/CVPR.2010.
 5539970.
- Sarthak Yadav, Sergios Theodoridis, Lars Kai Hansen, and Zheng-Hua Tan. Masked autoencoders
 with multi-window local-global attention are better audio learners. In *The Twelfth International Conference on Learning Representations*, 2024.
- Junjie Yang, Kaiyi Ji, and Yingbin Liang. Provably faster algorithms for bilevel optimization. In
 A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural In- formation Processing Systems*, 2021.
- Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. A large-scale car dataset for fine grained categorization and verification. In 2015 IEEE Conference on Computer Vision and Pattern
 Recognition (CVPR), pp. 3973–3981, 2015. doi: 10.1109/CVPR.2015.7299023.
- Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv: Computer Vision and Pattern Recognition*, 2017.
- Salah Zaiem, Titouan Parcollet, and Slim Essid. Less forgetting for better generalization: Exploring continual-learning fine-tuning methods for speech self-supervised representations. *arXiv preprint arXiv:2407.00756*, 2024.
- Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4l: Self-supervised semisupervised learning. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 1476–1485, 2019.
- Li Zhang, Youwei Liang, Ruiyi Zhang, Amirhosein Javadi, and Pengtao Xie. BLO-SAM: Bi-level optimization based finetuning of the segment anything model for overfitting-preventing semantic segmentation. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 59289–59309. PMLR, 21–27 Jul 2024.
- Tong Zhang, Congpei Qiu, Wei Ke, Sabine Süsstrunk, and Mathieu Salzmann. Leverage your
 local and global representations: A new self-supervised learning strategy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16580–16589, 2022a.
- Yihua Zhang, Guanhua Zhang, Prashant Khanduri, Min-Fong Hong, Shiyu Chang, and Sijia Liu.
 Revisiting and advancing fast adversarial training through the lens of bi-level optimization. In *International Conference on Machine Learning*, 2021.
- Yihua Zhang, Yuguang Yao, Parikshit Ram, Pu Zhao, Tianlong Chen, Mingyi Hong, Yanzhi Wang, and Sijia Liu. Advancing model pruning via bi-level optimization. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 18309–18326. Curran Associates, Inc., 2022b.
- Yihua Zhang, Prashant Khanduri, Ioannis C. Tsaknakis, Yuguang Yao, Min-Fong Hong, and Sijia
 Liu. An introduction to bi-level optimization: Foundations and applications in signal processing
 and machine learning. *ArXiv*, abs/2308.00788, 2023a.
- Yihua Zhang, Pranay Sharma, Parikshit Ram, Mingyi Hong, Kush R. Varshney, and Sijia Liu. What
 is missing in IRM training and evaluation? challenges and solutions. In *The Eleventh Interna- tional Conference on Learning Representations*, 2023b.
- Nicolas Zucchet and João Sacramento. Beyond backpropagation: bilevel optimization through implicit differentiation and equilibrium propagation. *Neural Computation*, 34(12):2309–2346, 2022.

A THEORETICAL INSIGHTS AND FRAMEWORK COMPARISONS IN BISSL

A.1 DERIVATION OF THE IMPLICIT GRADIENT 813

Assume the setup of the BiSSL optimization problem described in equation 2 and equation 3. In the following derivations, we will assume that ϕ_P is fixed, allowing us to simplify the expressions involved. To streamline the notation further, we continue to use the convention $\nabla_{\boldsymbol{\xi}} h(\boldsymbol{\xi})|_{\boldsymbol{\xi}=\boldsymbol{\psi}} :=$ $\nabla_{\boldsymbol{\xi}} h(\boldsymbol{\psi})$, when it is clear from context which variables are differentiated with respect to. Under these circumstances, we then define the lower-level objective from equation 3 as

834 835 836

851 852

853

854

855 856

857

858 859

860

862

810

811

$$G(\boldsymbol{\theta}_D, \boldsymbol{\theta}_P) := \mathcal{L}^P(\boldsymbol{\theta}_P, \boldsymbol{\phi}_P) + \lambda r(\boldsymbol{\theta}_D, \boldsymbol{\theta}_P).$$
(8)

Recalling that r is a *convex* regularization objective, adequate scaling of λ effectively "convexi-821 fies" the lower-level objective G, a strategy also employed on the lower-level objective in previous 822 works (Rajeswaran et al., 2019; Zhang et al., 2022b; 2023a). This is advantageous because assuming 823 convexity of G ensures that for any $\theta_D \in \mathbb{R}^L$, there exists a corresponding $\hat{\theta}_P \in \mathbb{R}^L$ that satisfies 824 the stationary condition $\nabla_{\theta_P} G(\theta_D, \hat{\theta}_P) = 0$. In other words, we are assured that a minimizer of 825 $G(\boldsymbol{\theta}_D, \cdot)$ exists for all $\boldsymbol{\theta}_D \in \mathbb{R}^L$. Now, further assume that $\nabla_{\boldsymbol{\theta}_P} G(\boldsymbol{\theta}_D, \boldsymbol{\theta}_P)$ is continuously differen-826 tiable and that the Hessian matrix $\nabla^2_{\theta_P} G(\theta_D, \hat{\theta}_P)$ is invertible for all $\theta_D \in \mathbb{R}^L$. Under these condi-827 tions, the implicit function theorem (Dontchev & Rockafellar, 2014; Zucchet & Sacramento, 2022) 828 guarantees the existence of an implicit unique and *differentiable* function $\theta_P^* : \mathcal{N}(\theta_D) \to \mathbb{R}^L$, with 829 $\mathcal{N}(\boldsymbol{\theta}_D)$ being a neighborhood of $\boldsymbol{\theta}_D$, that satisfies $\boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D) = \hat{\boldsymbol{\theta}}_P$ and $\nabla_{\boldsymbol{\theta}_P} G(\tilde{\boldsymbol{\theta}}_D, \boldsymbol{\theta}_P^*(\tilde{\boldsymbol{\theta}}_D)) = \mathbf{0}$ for 830 all $\hat{\boldsymbol{\theta}}_D \in \mathcal{N}(\boldsymbol{\theta}_D)$. 831

As the lower-level solution $\theta_P^*(\theta_D)$ is indeed a differentiable function under these conditions, this justifies that the expression

$$\frac{\mathrm{d}}{\mathrm{d}\boldsymbol{\theta}_D} \nabla_{\boldsymbol{\theta}_P} [G(\boldsymbol{\theta}_D, \boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D))] = \mathbf{0}$$

is valid for all $\theta_D \in \mathbb{R}^L$. By applying the chain rule, the expression becomes

$$\nabla^2_{\boldsymbol{\theta}_D \boldsymbol{\theta}_P} G(\boldsymbol{\theta}_D, \boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D)) + \frac{\mathrm{d}\boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D)}{\mathrm{d}\boldsymbol{\theta}_D}^T \nabla^2_{\boldsymbol{\theta}_P} G(\boldsymbol{\theta}_D, \boldsymbol{\theta}_P^*(\boldsymbol{\theta}_D)) = \mathbf{0}.$$

Recalling that $\nabla^2_{\theta_P} G(\theta_D, \theta_P^*(\theta_D))$ is assumed to be invertible, the IG $\frac{\mathrm{d}\theta_P^*(\theta_D)}{\mathrm{d}\theta_D}^T$ can be isolated

$$\frac{\mathrm{d}\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})}{\mathrm{d}\boldsymbol{\theta}_{D}}^{T} = -\nabla_{\boldsymbol{\theta}_{D}\boldsymbol{\theta}_{P}}^{2} G(\boldsymbol{\theta}_{D},\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})) \left[\nabla_{\boldsymbol{\theta}_{P}}^{2} G(\boldsymbol{\theta}_{D},\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D}))\right]^{-1}$$

and by substituting the expression for G from equation 8, the expression becomes

$$\frac{\mathrm{d}\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})}{\mathrm{d}\boldsymbol{\theta}_{D}}^{T} = -\nabla_{\boldsymbol{\theta}_{D}\boldsymbol{\theta}_{P}}^{2}r(\boldsymbol{\theta}_{D},\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})) \left[\nabla_{\boldsymbol{\theta}_{P}}^{2}\left(\frac{1}{\lambda}\mathcal{L}^{P}(\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D}),\boldsymbol{\phi}_{P}) + r(\boldsymbol{\theta}_{D},\boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D}))\right)\right]^{-1}.$$
 (9)

To summarize, given the following assumptions:

- The lower-level pretext projection head parameters ϕ_P are fixed.
- G is convex such that $\nabla_{\boldsymbol{\theta}_{P}} G(\boldsymbol{\theta}_{D}, \boldsymbol{\theta}_{P}^{*}(\boldsymbol{\theta}_{D})) = \mathbf{0}$ is fulfilled for every $\boldsymbol{\theta}_{D} \in \mathbb{R}^{L}$.
- The Hessian matrix $\nabla^2_{\theta_P} G(\theta_D, \theta_P^*(\theta_D))$ exists and is invertible for all $\theta_D \in \mathbb{R}^L$.

Then, the IG $\frac{d\theta_P^*(\theta_D)}{d\theta_D}^T$ can be explicitly expressed by equation 9. The authors acknowledge that an explicit expression for the IG without fixing ϕ_P is achievable, though this is left for future exploration.

861 A.2 DISTINCTION FROM BILEVEL OPTIMIZATION IN META-LEARNING

⁸⁶³ While bilevel optimization (BLO) has been applied in meta-learning frameworks such as MAML (Finn et al., 2017), Sign-MAML (Fan et al., 2021) and iMAML (Rajeswaran et al., 2019),

864 BiSSL represents a distinct application and implementation of BLO, tailored for the challenges of 865 self-supervised learning (SSL). In the aforementioned works, BLO is primarily utilized to address 866 few-shot learning scenarios, focusing on efficiently adapting models to new tasks with minimal 867 labeled data. Conversely, BiSSL applies BLO to concurrently manage the more complex task of 868 self-supervised pretext pre-training on unlabeled with downstream fine-tuning on labeled data. Another key distinction is that in meta-learning, the upper- and lower-level objectives are closely related, with the upper-level objective formulated as a summation of the lower-level tasks. In contrast, 870 BiSSL involves fundamentally distinct objectives at each level, utilizing separate datasets and tasks 871 for pre-training and fine-tuning. This design allows BiSSL to better align the pre-trained model 872 with the requirements of a specific downstream task. Conversely, the BLO in meta-learning aims 873 to broadly generalize across a wide range of tasks, prioritizing adaptability rather than task-specific 874 optimization. Additionally, unlike BiSSL, the meta-learning frameworks discussed reinitialize the 875 lower-level backbone parameters with a copy of the upper-level parameters at every iteration. In 876 BiSSL, the closest comparable mechanism is the occasional update of the upper-level backbone us-877 ing an EMA update with the lower-level parameters (see Algorithm 1), though this occurs far less 878 frequently.

879 880

881

B EXPERIMENTAL DETAILS

882 B.1 DATASET PARTITIONS

The Caltech-101 (Li et al., 2022a) dataset does not come with a pre-defined train/test split, so the same convention as previous works is followed (Chen et al., 2020b; Donahue et al., 2014; Simonyan & Zisserman, 2014), where 30 random images per class are selected for the training partition, and the remaining images are assigned for the test partition. For the DTD (Cimpoi et al., 2014) and SUN397 (Xiao et al., 2010) datasets, which offer multiple proposed train/test partitions, the first splits are used, consistent with the approach in Chen et al. (2020b).

For downstream hyperparameter optimization, portions of the training partitions from each respective labeled dataset are designated as validation datasets. The FGVC Aircraft (Maji et al., 2013),
Oxford 102 Flowers (Nilsback & Zisserman, 2008), DTD, and Pascal VOC 2007 (Everingham et al.)
datasets already have designated validation partitions. For all the remaining labeled datasets, the validation data partitions are randomly sampled while ensuring that class proportions are maintained.
For the multi-attribute VOC07 dataset, sampling is performed with class balance concerning the first attribute present in each image. Roughly 20% of the training data is allocated for validation.

- 898 B.2 DOWNSTREAM TASK FINE-TUNING OF THE BASELINE SETUP
- In Table 2, the learning rates and weight decays used for each respective downstream dataset of the experiments described in Section 4.2.1 are outlined.
- 901 902 903

- B.3 DOWNSTREAM HEAD WARMUP AND UPPER-LEVEL OF BISSL
- Table 3 outlines the learning rates and weight decays used for the downstream head warm-up and upper-level of BiSSL of each respective downstream dataset, as described in the BiSSL experimental setup of Section 4.2.2.
- The first term of the upper-level gradient equation 7 is approximated using the Conjugate Gradi-908 ent (CG) method (Nazareth, 2009; Shewchuk, 1994). Our implementation follows a similar struc-909 ture to that used in Rajeswaran et al. (2019), employing $N_c = 5$ iterations and a dampening term 910 $\lambda_{\text{damp}} = 10$. Given matrix A and vector v, the CG method iteratively approximates $A^{-1}v$, which re-911 quires evaluation of multiple matrix-vector products Ad_1, \ldots, Ad_{N_c} . In practice, storing the matrix 912 \hat{A} (in our case, the Hessian $\nabla^2_{\theta_P} \mathcal{L}^P(\theta_P^*(\theta_D), \phi_P)$) in its full form is typically infeasible. Instead, 913 a function that efficiently computes the required matrix-vector products instead of explicitly storing 914 the matrix is typically utilized. For our setup, this function is detailed in Algorithm 2, showing 915 how the K stored lower-level batches (we use K = 5 as previously outlined in the "Lower-level of BiSSL' paragraph in Section 4.2.2) are used to calculate Hessian-vector products. This approach en-916 sures that the output of the CG algorithm is an approximation of the inverse Hessian-vector product 917 in the first term of Equation equation 7 as intended.

Table 2: Hyper-parameter configurations used for downstream fine-tuning after conventional pretext pre-training yielding the highest top-1 classification accuracies (11-point mAP for the VOC2007 dataset).

011			
922	Dataset	Learning Rate	Weight Decay
923	STI 101	0.0126	0.001
924	SILIUL	0.0130	0.001
925	Flowers	0.113	0.00226
926	Cars	0.035	0.00658
927	Aircrafts	0.0167	0.00996
928	DTD	0.0262	0.00332
929	Dete	0.0005	0.00470
930	Pets	0.0235	0.00472
931	FashionMNIST	0.0009	0.00829
932	CIFAR10	0.0067	0.00128
933	CIFAR100	0.005	0.00127
934	Caltech 101	0.0006	0.00002
935	Cancen-101	0.0030	0.00302
936	Food	0.015	0.00699
937	SUN397	0.0097	0.00121
938	CUB200	0.0722	0.00568
939	VOC2007	0.0108	0.00894
940			0.00001

Table 3: Hyper-parameters used for the Downstream Head Warm-up and Upper-level of BiSSL.

944	Dataset	Learning Rate	Weight Decay
945	CTI 101	0.015	0.01
946	SILIUL	0.015	0.01
947	Flowers	0.05	0.01
948	Cars	0.035	0.007
949	Aircrafts	0.015	0.01
950	DTD	0.015	0.0075
951	Pets	0.03	0.005
953	FashionMNIST	0.05	0.004
954	CIFAR10	0.03	0.006
955	CIFAR100	0.03	0.001
956	Caltech-101	0.03	0.007
958	Food	0.015	0.01
959	SUN397	0.03	0.002
960	CUB200	0.05	0.005
961	VOC2007	0.03	0.006
962	1002001	0.00	0.000

963

941 942

943

918

964 965

966

B.4 COMPOSITE CONFIGURATION OF BISSL

To avoid data being reshuffled between every training stage alternation, the respective batched lowerand upper-level training datasets are stored in separate stacks from which data is drawn. The stacks are only "reset" when the number of remaining batches is smaller than the number of gradient steps required before alternating to the other level. For example, the lower-level stack is reshuffled every fourth training stage alternation. If the downstream dataset does not provide enough data for making $N_U = 8$ batches with non-overlapping data points, the data is simply reshuffled every time

1:	Input: Input vector v. Mod	el parameters θ_P , ϕ_P . Training objective \mathcal{L}^P . Lower-level data
	batches $[\mathbf{z}_1, \ldots, \mathbf{z}_K]$. Regula	fization weight λ and dampening λ_{damp} .
2:	Initialize $\mathbf{v} \leftarrow 0$	▷ Initialize Hessian vector product v
3:	for $k = 1, \ldots, K$ do	·
4:	$\pi(\boldsymbol{\theta}_P) \leftarrow (\nabla_{\boldsymbol{\theta}} \mathcal{L}^P (\boldsymbol{\theta}, \boldsymbol{\phi}_P))$	$(\mathbf{z}_k) _{\boldsymbol{\theta} = \boldsymbol{\theta}_{\mathcal{P}}} ^T \mathbf{v}$
5:	$\mathbf{g} \leftarrow \nabla_{\boldsymbol{\theta}} \pi(\boldsymbol{\theta}) \big _{\boldsymbol{\theta} = \boldsymbol{\theta}_P}$	$\triangleright \text{ Memory efficient calculation of } \nabla^2_{\theta} \mathcal{L}^P(\theta, \phi_P; \mathbf{z}_k) _{\theta = \theta_P} \mathbf{x}$
6:	$\mathbf{y} \leftarrow \mathbf{y} + rac{1}{K}\mathbf{g}$	
7:	$\mathbf{y} \leftarrow \mathbf{v} + rac{1}{\lambda + \lambda_{ ext{down}}} \mathbf{y}$	
	2 - Champ	
8:	Return: $f_H(\mathbf{v}) := \mathbf{y}$	

the remaining number of data points is smaller than the upper-level batch size (256 images in these experiments).

989 **B.5** DOWNSTREAM FINE-TUNING AFTER BISSL

The learning rates and weight decays used for downstream fine-tuning after BiSSL for each respec-992 tive downstream dataset are outlined in Table 4. Section 4.2.2 outlines the experimental setup.

Table 4: Hyper-parameter configurations used for downstream fine-tuning after BiSSL leading to the highest top-1 classification accuracies (11-point mAP for the VOC2007 dataset).

997	Dataset	Learning Rate	Weight Decay
998	STI 101	0.0005	0.00043
999	STETUE	0.0000	0.00040
1000	Flowers	0.0009	0.00005
1001	Cars	0.0293	0.00851
1002	Aircrafts	0.011	0.00612
1003	DTD	0.0008	0.00764
1004	Pate	0.0002	0.00543
1005		0.0002	0.00343
1006	FashionMNIST	0.0022	0.00876
1007	CIFAR10	0.0003	0.00991
1008	CIFAR100	0.0012	0.00422
1009	Caltech-101	0.0008	0.00011
1010		0.0000	0.00011
1011	Food	0.006	0.0095
1012	SUN397	0.0011	0.00028
1013	CUB200	0.035	0.00868
1014	VOC2007	0.0002	0.00015
1015		1	
1016			

986

987

988

990 991

993 994

995

996

1017 1018 1019

ADDITIONAL RESULTS С

C.1 VISUAL INSPECTION OF LATENT FEATURES 1020

1021 Test data features of the downstream test data processed by backbones trained through conventional 1022 pretext pre-training are compared against those trained with BiSSL. This allows for an inspection of 1023 the learned representations prior to the final fine-tuning stage. 1024

During the evaluation, it is important to note that the batch normalization layers (Ioffe & Szegedy, 1025 2015) of the pre-trained backbones utilize the running mean and variance inferred during train-





