# MODEL DEVELOPMENTAL SAFETY: A RETENTION-CENTRIC METHOD AND APPLICATIONS IN VISION-LANGUAGE MODELS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

Paper under double-blind review

#### ABSTRACT

In the real world, a learning-enabled system usually undergoes multiple cycles of model development to enhance the system's ability to handle difficult or emerging tasks, which involve collecting new data, training a new model and validating the model. This continual model development process raises a significant issue that the model development for acquiring new or improving existing capabilities may inadvertently lose capabilities of the old model, also known as catastrophic forgetting. Existing continual learning studies focus on mitigating catastrophic forgetting by trading off performance on previous tasks and new tasks to ensure good average performance. However, they are inadequate for many applications especially in safety-critical domains, as failure to strictly preserve the good performance of the old model on some tasks not only introduces risks and uncertainties but also imposes substantial expenses in the re-improving and re-validation of existing properties. To address this issue, we introduce **model developmental** safety as a guarantee of a learning system such that in the model development process the new model should strictly retain the existing protected capabilities of the old model while improving its performance on target tasks. To ensure the model developmental safety, we present a retention-centric framework by formulating the model developmental safety as data-dependent constraints. Under this framework, we study how to develop a pretrained vision-language model, specifically the CLIP model, for acquiring new capabilities or improving existing capabilities of image classification. We propose an efficient constrained optimization algorithm with theoretical guarantee and use its insights to finetune a CLIP model with task-dependent heads for promoting the model developmental safety. Our experiments on improving vision perception capabilities on autonomous driving and scene recognition datasets demonstrate the efficacy of the proposed approach.

## 7 1 INTRODUCTION

038 Learning-enabled systems are rapidly transforming various sectors, with applications in autonomous vehicles, medical diagnosis, and financial prediction. These systems often rely on ML models that are 040 trained on vast amounts of data. However, the inherent complexity of the environments in which these 041 systems operate often presents critical challenges, e.g., dealing with corner cases and rare scenarios 042 that deviate from the norm. Additionally, real-world scenarios continuously evolve, presenting new 043 challenges and requiring the system to adapt. These necessitate an iterative development process 044 where models are constantly refined and improved based on new data. Continuously updating the model has become a norm especially in the era of large foundation models, e.g., ChatGPT has experienced several cycles of development from GPT3.5 to GPT4 and GPT40 and recent GPT01. 046

However, this iterative model development process raises a significant issue, i.e., the model development for improving the existing capabilities or acquiring new capabilities may inadvertently lose the previously acquired capabilities of the old model. This issue has been widely observed and documented as catastrophic forgetting when models are trained to learn a sequence of contents (McCloskey & Cohen, 1989). Tremendous studies have been conducted to mitigate the forgetting problem in continual learning literature (Zhou et al., 2022; Rolnick et al., 2019; Shin et al., 2017; Li & Hoiem, 2016; Kirkpatrick et al., 2017). However, these works primarily focus on mitigating the catastrophic forgetting problem, by trading off performance on previous tasks and new tasks to have good av-

054Figure 1: Performance of recog-<br/>nizing 6 weather conditions for au-<br/>tonomous driving with two rounds056of model development using new 3<br/>data. The Round 1 development tar-<br/>gets at *overcast* and Round 2 aims<br/>to improve recognizing *foggy*. Base<br/>refers to the CLIP model finetuned<br/>on BDD100K data.

063

094

096

098

099

100

101

102



erage performance (Wang et al., 2024), but do not strictly retain existing abilities (i.e., ensuring zero forgetting) while learning new tasks. Ensuring zero forgetting is crucial for many applications especially in safety-critical domains, as failure to ensure strict preservation of the model's original capabilities not only introduces safety risks and uncertainties but also imposes substantial expenses in the re-improving and re-validation of existing measures, such as in autonomous driving, where validation and verification are challenging and could cost billions of dollar (Rajabli et al., 2020; Koopman & Wagner, 2016; Company, 2023). This presents a significant challenge for iterative model development process.

To address this challenge, this paper formally introduces model developmental safety (MDS) as a 072 guarantee of a learning system such that in the model development process the new model should 073 strictly retain the existing protected capabilities of the old model while improving its performance on 074 target tasks. This concept subtly differs from trading off performance between previous tasks and new 075 tasks to have good average performance of existing continual learning approaches. Moreover, MDS 076 cannot be achieved by the naive weighting method that optimizes a weighted loss via combining 077 the losses of protected tasks and target tasks and tuning the weight to preserve existing protected 078 capabilities. This approach does not necessarily preserve the performance of the model on all 079 protected tasks even if the weight is large enough, as shown in Table 3, and will yield no improvement 080 on target tasks if the weight is too large. A better algorithm is required to enable an efficient search of a model that not only retains the performance on protected tasks but also improves the performance 081 on target tasks. To the best of our knowledge, no such algorithm currently exists. 082

083 This paper aims to address this critical gap by introducing a novel retention-centric framework to 084 ensure MDS. We propose to formulate the MDS as data-dependent constraints, which offers statistical 085 guarantee for strict preservation of performance for all protected tasks. With this framework, we explore developing a pretrained CLIP model for acquiring new capabilities or improving existing ones in image classification. We propose an efficient constrained optimization algorithm with theoretical 087 guarantee. With insights from theoretical analysis, we finetune the CLIP model with task-dependent 880 heads to facilitate MDS. Finally, we demonstrate the efficacy of our approach through experiments on 089 enhancing vision-based perception capabilities in autonomous driving dataset and scene recognition 090 dataset, highlighting the practical importance of MDS in real-world scenarios. Our contributions are 091 summarized below: 092

- We introduce a retention-centric framework by formulating the MDS as data-dependent constraints, which offer statistical guarantee for strictly preserving performance for every protected task.
- We propose an efficient constrained optimization algorithm with theoretical guarantee to develop a pretrained vision-language model for acquiring new capabilities or improving existing capabilities of image classification.
- We conduct comprehensive experiments to study the proposed algorithm and compare our approach with existing baselines to demonstrate its effectiveness. An experimental result for ensuring MDS in improving vision-based perception capabilities of autonomous driving is shown in Figure 1.
- 2 RELATED WORK

Continual learning. This work is closely related to Continual learning (CL), also known as lifelong
learning, yet it exhibits nuanced differences. Continual learning usually refers to learning a sequence
of tasks one by one and accumulating knowledge like human instead of substituting knowledge (Wang
et al., 2024; Qu et al., 2021). There is a vast literature of CL of deep neural networks (DNNs) (Aljundi
et al., 2018; Lopez-Paz & Ranzato, 2017a; Farajtabar et al., 2019; Lee et al., 2017; Guo et al.,
2020; Parisi et al., 2018). The core issue in CL is known as catastrophic forgetting (McCloskey &

108 Cohen, 1989), i.e., the learning of the later tasks may **significantly** degrade the performance of the 109 model for the earlier tasks. Different approaches have been investigated to mitigate catastrophic 110 forgetting, including regularization-based approaches (Castro et al., 2018; Kirkpatrick et al., 2017; 111 Zenke et al., 2017; Li & Hoiem, 2016), expansion-based approaches (Zhou et al., 2022; Li et al., 112 2019; Rusu et al., 2016; van de Ven et al., 2020), and memory-based approaches (Rolnick et al., 2019; Cha et al., 2021; Guo et al., 2020; Lopez-Paz & Ranzato, 2017a; Chaudhry et al., 2019). 113 In the era of large language models(LLMs), another type of continual learning method, known as 114 knowledge/representation editing(Meng et al., 2022; Zou et al., 2023; Luo et al., 2024; Huang et al., 115 2024; Liu et al., 2024), emerges to efficiently modify the behavior of LLMs with minimal impact on 116 unrelated inputs (Wang et al., 2023b), such as to update stale facts, eliminate unintended biases, or 117 reduce undesired hallucinations. 118

The framework proposed in this work is similar to conventional memory-based approaches in the 119 sense that both use examples of existing tasks to regulate learning. However, the key difference 120 is that most existing continual learning focuses on the trade-off between learning plasticity and 121 memory stability and aims to find a proper balance between performance on previous tasks and new 122 tasks (Wang et al., 2024). Hence, they do not provide a guarantee for MDS. A recent work (Peng 123 et al., 2023) has proposed an ideal continual learner that never forgets by assuming that all tasks 124 share the same optimal solution. However, it is not practical and not implementable for deep learning 125 problems. Besides, existing continual learning studies usually highlight resource efficiency when 126 accumulating knowledge by reducing the number of samples of previous tasks. In contrast, this work 127 tends to utilize more examples to construct constraints for protected tasks to facilitate MDS.

128 Constrained Learning. Our work is also related to constrained learning. While most traditional 129 constrained optimization works focus on convex objectives or convex constraints, the research interest 130 recently has been directed to non-convex optimization (Boob et al., 2023; Facchinei et al., 2021; Li 131 et al., 2024; Chamon et al., 2022; Alacaoglu & Wright, 2024), due to its increasing importance in 132 modern machine learning problems, such as in applications concerned with fairness (Cotter et al., 133 2019), robustness (Robey et al., 2021), and safety (Paternain et al., 2019b) problems. Nevertheless, 134 none of existing algorithms can be directly applied to our large-scale deep learning problem (4), 135 due to either prohibitive running cost or failure to handle biased stochastic gradients caused by compositional structure. We include more discussion in Appendix B. 136

# <sup>137</sup> 3 NOTATIONS AND PRELIMINARIES

139 **Notations.** We consider developing a model w to improve its capabilities on a target task  $\mathbb{T}_o$  while preserving its performance on a set of protected tasks denoted by  $\mathbb{T}_1, \ldots, \mathbb{T}_m$ . A task can be as 140 simple as predicting a class for multi-class classification or as complicated as coding ability of LLMs. 141 In the paper, we focus on classification using CLIP models and each task refers to one class. For 142 example, we can consider tasks of predicting different weather conditions in autonomous driving, 143 e.g., foggy, overcast, cloudy, clear, rainy, etc. We assume that each task is associated with a data 144 distribution denoted by  $\mathfrak{D}_k$ . Let  $(\mathbf{x}, y) \sim \mathfrak{D}_k$  denote random data of task  $\mathbb{T}_k$  with input  $\mathbf{x} \in \mathcal{X}$  (e.g., 145 an image) and output  $y \in \mathcal{Y}$  (e.g., its class label). We assume that each protected task has a set 146 of examples denoted by  $\mathcal{D}_k = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_k}$ , sampled from  $\mathfrak{D}_k$ . Let  $\ell_k(\mathbf{w}, \mathbf{x}, y) = \ell_k(s(\mathbf{x}; \mathbf{w}), y)$ 147 denote a loss function that measures the loss of the model's prediction  $s(\mathbf{x}; \mathbf{w})$  with respect to the 148 groundtruth y for task k. For classification, the loss could be zero-one loss  $\ell_{0-1}$  that measures the classification error or the cross-entropy loss  $\ell_{ce}$  that is differentiable for learning. We will define these losses shortly for using CLIP models. We denote by  $\mathcal{L}_k(\mathbf{w}, \mathfrak{D}_k) = \mathbb{E}_{\mathbf{x}, y \sim \mathfrak{D}_k} \ell_k(\mathbf{w}, \mathbf{x}, y)$  as the expected loss, and by  $\mathcal{L}(\mathbf{w}, \mathcal{D}_k) = \frac{1}{n_k} \sum_{(\mathbf{x}_i, y_i) \sim \mathcal{D}_k} \ell_k(\mathbf{w}, \mathbf{x}_i, y_i)$  as the empirical loss for task k. 149 150 151 152

The CLIP model and Contrastive Loss. The contrastive loss has been successfully applied to learning the CLIP model (Radford et al., 2021a), which has exhibited remarkable performance for classifying images. We consider optimizing a two-way contrastive loss for each image-text pair  $(\mathbf{x}_i, \mathbf{t}_i)$  following Yuan et al. (2022):

$$L_{\text{ctr}}(\mathbf{w}; \mathbf{x}_i, \mathbf{t}_i, \mathcal{T}_i^-, \mathcal{I}_i^-) := -\tau \log \frac{\exp(E_1(\mathbf{w}, \mathbf{x}_i)^\top E_2(\mathbf{w}, \mathbf{t}_i)/\tau)}{\sum_{\mathbf{t}_j \in \mathcal{T}_i^-} \exp(E_1(\mathbf{w}, \mathbf{x}_i)^\top E_2(\mathbf{w}, \mathbf{t}_j)/\tau)}$$
(1)  
$$-\tau \log \frac{\exp(E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_i)/\tau)}{\sum_{\mathbf{x}_j \in \mathcal{I}_i^-} \exp(E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_j)/\tau)},$$

157 158 159

1

where  $E_1(\mathbf{w}, \mathbf{x})$  and  $E_2(\mathbf{w}, \mathbf{t})$  denotes a (normalized) encoded representation of a image  $\mathbf{x}$ , and a text  $\mathbf{t}$ , respectively,  $\mathcal{T}_i^-$  denotes the set of all texts to be contrasted with respect to (w.r.t)  $\mathbf{x}_i$  (including itself) and  $\mathcal{I}_i^-$  denotes the set of all images to be contrasted w.r.t  $\mathbf{t}_i$  (including itself).

To utilize a CLIP model for multi-class classification with classes  $C = \{c_1, \ldots, c_K\}$ , we will convert 166 a class  $c_k$ , e.g., "rainy", into a text description of  $c_k$ , denoted by  $\hat{\mathbf{t}}_k$ , e.g., "the weather is rainy", similar 167 to the zero-shot classification scheme of the well-known CLIP model (Radford et al., 2021a). Hence, 168 a prediction score (i.e., a logit) for an image x and a text description  $\hat{\mathbf{t}}_k$  of class  $c_k$  is calculated by 169  $s_k(\mathbf{x}; \mathbf{w}) = E_1(\mathbf{w}, \mathbf{x})^\top E_2(\mathbf{w}, \hat{\mathbf{t}}_k)$ . The predicted class label is given by  $\hat{y} = \arg \max_{c_k \in \mathcal{C}} s_k(\mathbf{x}; \mathbf{w})$ . 170 Hence, given the true class  $y \in C$ , the zero-one loss is given by  $\ell_{0,1}(\mathbf{w}, \mathbf{x}, y) = \mathbb{I}(\hat{y} \neq y)$ , and 171 the cross-entropy loss is given by  $\ell_{ce}(\mathbf{w}, \mathbf{x}, y) = -\log \frac{\exp(s_y(\mathbf{x}; \mathbf{w})/\tau_0)}{\sum_{\ell=1}^{K} \exp(s_\ell(\mathbf{x}; \mathbf{w})/\tau_0)}$ , where  $\tau_0 > 0$  is a 172 temperature parameter that controls the balance between the approximation error of the zero-one loss 173 and the smoothness of the function. In particular, a smaller  $\tau_0$  gives a smaller approximation error 174 and a larger  $\tau_0$  indicates a smaller gradient Lipschitz constant of the loss function in terms of logits. 175

#### 4 A RETENTION-CENTRIC FRAMEWORK

178 4.1 MODEL DEVELOPMENTAL SAFETY

Wnew

To measure the model developmental safety, it is necessary to evaluate how the performance of the model changes in protected tasks from the old model  $w_{old}$  to a new model  $w_{new}$ . We introduce the formal definition of model developmental safety (MDS) in Definition 1, which ensures the new model strictly preserves performance on each individual protected task.

**Definition 1 (Model Developmental Safety (MDS))** In model development process, the model developmental safety is satisfied if  $\mathcal{L}_k(\mathbf{w}_{new}, \mathfrak{D}_k) \leq \mathcal{L}_k(\mathbf{w}_{old}, \mathfrak{D}_k), \forall k \in \{1, \ldots, m\}$ , where  $\mathcal{L}_k(\mathbf{w}, \mathfrak{D}_k) = \mathbb{E}_{\mathbf{x}, y \sim \mathfrak{D}_k} \ell_k(\mathbf{w}, \mathbf{x}, y).$ 

187 188

192

184

185

176

177

In practice, the developmental safety will be measured using a set of examples  $S_j \sim \mathfrak{D}_j$  for each protected task. Hence, we define the empirical developmental safety metric, corresponding to Definition 1, for evaluation:

$$DevSafety = \min_{k \in \{1, \dots, m\}} \left( \mathcal{L}_k(\mathbf{w}_{old}, \mathcal{S}_k) - \mathcal{L}_k(\mathbf{w}_{new}, \mathcal{S}_k) \right).$$
(2)

When we use the zero-one loss  $\ell_{0-1}$  in the above definitions, we refer to the above developmental safety metric as DevSafety(acc).

## 4.2 A RETENTION-CENTRIC APPROACH FOR MODEL DEVELOPMENTAL SAFETY

197 The key to our retention-centric framework is to utilize examples of protected tasks to define empirical 198 retention constraints when updating the model on a target task. In order to develop the model for 199 improving the performance on a target task  $\mathbb{T}_o$ , we assume that a set of data  $\mathcal{D}$  for  $\mathbb{T}_o$  is constructed 200 and a proper objective is given based on application, denoted by  $F(\mathbf{w}, \mathcal{D})$ . Then, our retention-centric 201 approach for model development is imposed by solving the following problem:

202 203

$$F = \underset{\mathbf{w}}{\operatorname{arg\,min}} F(\mathbf{w}, \mathcal{D})$$
  
s.t.  $\mathcal{L}_k(\mathbf{w}, \mathcal{D}_k) - \mathcal{L}_k(\mathbf{w}_{\operatorname{old}}, \mathcal{D}_k) \le 0, \ k = 1, \cdots, m.$  (3)

We will propose an algorithm to directly solve this data-dependent constrained optimization problem with a contrastive objective in the context of developing a CLIP model in next section.

**Generalization Analysis.** Since we can only use empirical data  $\mathcal{D}_1, \ldots, \mathcal{D}_m$  in (3), there exist 207 generalization errors between the retention constraints in (3) and the MDS we want to ensure 208 in Definition 1. The lemma below uses a standard tool of statistical error analysis to bound the 209 generalization error of retention. For simplicity, we assume each protected task is associated with 210 the same loss function, namely,  $\ell_k = \ell$  for  $k = 1, \ldots, m$ . In the analysis, we use the Rademacher 211 complexity of the loss class  $\mathcal{H} = \{\ell(\mathbf{w}, \cdot, \cdot) : \mathcal{X} \times \mathcal{Y} \to [0, 1] | \mathbf{w} \in \mathbb{R}^d\}$  induced by the model  $\mathbf{w}$ 212 on n data points, which is denoted by  $R_n(\mathcal{H})$ . We assume that  $R_n(\mathcal{H}) \leq Cn^{-\alpha}$  for some  $C \geq 0$ 213 and  $\alpha \leq 0.5$ . We note that  $\alpha = 0.5$  in the vast majority of model and loss families, including linear 214 models (Kakade et al., 2008), deep neural networks (Bartlett & Mendelson, 2002), and model families with bounded VC dimension (Bartlett & Mendelson, 2002). 215

219 220 221

222

223

224

225

244 245

246

266 267

216 Lemma 1 (Generalization Error of Rentention) Suppose that  $R_n(\mathcal{H}) \leq Cn^{-\alpha}$  for some  $C \geq 0$ 217 and  $\alpha \leq 0.5$ . Then, with probability at least  $1 - \delta$ , it holds that

$$\mathcal{L}_{k}(\mathbf{w}_{new},\mathfrak{D}_{k}) - \mathcal{L}_{k}(\mathbf{w}_{old},\mathfrak{D}_{k}) \leq \mathcal{L}_{k}(\mathbf{w}_{new},\mathcal{D}_{k}) - \mathcal{L}_{k}(\mathbf{w}_{old},\mathcal{D}_{k}) + \frac{4C}{n_{k}^{\alpha}} + 2\sqrt{\frac{\ln(2m/\delta)}{2n_{k}}}, \forall k.$$

**Remark:** The lemma indicates that as long as the empirical retention constraints are satisfied, i.e.,  $\mathcal{L}_k(\mathbf{w}_{\text{new}}, \mathcal{D}_k) - \mathcal{L}_k(\mathbf{w}_{\text{old}}, \mathcal{D}_k) \leq 0$ , the model developmental safety is ensured up to a statistical error in the order of  $O(n^{-\alpha})$ , where  $n = \min_k n_k$ . Hence, the more examples used to construct the constraints, the more likely the new model meets MDS requirement. The proof is given in C.1.

#### <sup>226</sup> 5 **RETENTION-CENTRIC DEVELOPMENT OF CLIP MODELS**

227 Based on the proposed framework above, in this section, we present an efficient algorithm for 228 improving a pretrained CLIP model on a target task while ensuring MDS on a set of protected tasks. 229 The CLIP model is of particular interest because (i) it is a foundation model that has been used 230 extensively in many applications; and (ii) can adapt to the open-world for handling new classes using 231 languages. However, existing studies have shown that directly applying a pretrained CLIP model (e.g., OpenAI's CLIP model) to a certain downstream application yields varying performance across 232 different classes (Parashar et al., 2024). Rare concepts (e.g., foggy) usually has worse performance 233 than frequent concepts (e.g., clear), making it necessary to continuously update. 234

Suppose a CLIP model  $\mathbf{w}_{old}$  has been trained. We aim to improve it for a target task  $\mathbb{T}_o$  (e.g., classifying foggy). To this end, we collect a set of image-text pairs related to the target task, denoted by  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{t}_i)\}_{i=1}^{n_o}$ . As labeled data for rare scenarios (e.g., *foggy*) are usually limited in practice, we consider augmenting the dataset  $\mathcal{D}$  by using a query prompt to search for target-related image-text pairs from the internet (detailed in Appendix A.2). For each image-text pair, a set of negative texts has been collected to be contrasted w.r.t.  $\mathbf{x}_i$ , which together with  $\mathbf{t}_i$  form  $\mathcal{T}_i^-$ , and a set of negative images has been also collected to be contrasted w.r.t.  $\mathbf{t}_i$ , which together with  $\mathbf{x}_i$  form  $\mathcal{I}_i^-$ .

To develop the CLIP model in our retention-centric framework, we instantiate (3) as: 1 - 1

$$\min_{\mathbf{w}} F(\mathbf{w}, \mathcal{D}) := \frac{1}{n_o} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}} L_{\text{ctr}}(\mathbf{w}; \mathbf{x}_i, \mathbf{t}_i, \mathcal{T}_i^-, \mathcal{I}_i^-)$$
s.t.  $h_k(\mathbf{w}) := \mathcal{L}_k(\mathbf{w}, \mathcal{D}_k) - \mathcal{L}_k(\mathbf{w}_{\text{old}}, \mathcal{D}_k) \le 0, \ k = 1, \cdots, m.$ 
(4)

#### 247 5.1 EFFICIENT OPTIMIZATION AND CONVERGENCE ANALYSIS

The optimization problem in (4) is challenging for multiple reasons. First, this problem involves a non-convex objective and non-convex constraints, so finding a global optimal solution is intractable in general. Second, the objective and constraint functions are formulated using a large dataset, so we need to sample from the dataset in order to construct stochastic gradients of the functions to update the solution. Lastly, (4) may contain a large number of constraints, so updating the solutions using the gradients of all constraints may be prohibited. Given these challenges, we need to develop a stochastic optimization for (4) based on advanced techniques and constraint sampling.

255 Our method is motivated by the stochastic quadratic penalty method in (Alacaoglu & Wright, 2024), 256 which first converts (4) into an unconstrained problem by adding a quadratic penalty on the constraints 257 violation to the objective function and then solves the unconstrained problem using a variance-reduced 258 stochastic gradient method. Unfortunately, their method can not be directly applied to (4) because (i) they only consider equality constraints while (4) involves inequality constraints and (ii) they require 259 an unbiased stochastic gradients for each update while the stochastic gradients for (4) will be biased 260 due to the compositional structure. Note that an augmented Lagrangian algorithm (ALA) is also 261 studied by Alacaoglu & Wright (2024), which has the same issue as their penalty method. We only 262 consider quadratic penalty method for (4) because it has the same complexity as the ALA but is more 263 intuitive and easier to implement. 264

A quadratic penalty method converts (4) into the following unconstrained problem:

$$\min_{\mathbf{w}} \Phi(\mathbf{w}) := F(\mathbf{w}, \mathcal{D}) + \frac{1}{m} \sum_{k=1}^{m} \frac{\beta}{2} ([h_k(\mathbf{w})]_+)^2$$
(5)

where  $[\cdot]_{+} = \max\{\cdot, 0\}$  and  $\beta \ge 0$  is the penalty parameter. Under mild conditions(Bertsekas, 2014), a large enough  $\beta$  will ensure the optimal solution to (5) is also an optimal solution to (4). In the following, we introduce an efficient stochastic algorithm to solve (5). It is notable that both terms are of finite-sum coupled compositional structure (Wang & Yang, 2022), i.e.,  $\sum_i f(g_i(\mathbf{w}))$ , where f is non-linear.

We discuss how to approximate the gradient of two terms of the objective using mini-batch samples below. Define  $g_{1i}(\mathbf{w}) = \frac{1}{|\mathcal{T}_i^-|} \sum_{\mathbf{t}_j \in \mathcal{T}_i^-} \exp\left(E_1(\mathbf{w}, \mathbf{x}_i)^\top E_2(\mathbf{w}, \mathbf{t}_j) - E_1(\mathbf{w}, \mathbf{x}_i)^\top E_2(\mathbf{w}, \mathbf{t}_i)/\tau\right)$ and  $g_{2i}(\mathbf{w}) = \frac{1}{|\mathcal{I}_i^-|} \sum_{\mathbf{x}_j \in \mathcal{I}_i^-} \exp\left(E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_j) - E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_i)/\tau\right)$ . Then,  $E(\mathbf{w}, \mathcal{D}) = \frac{1}{2} \sum_{\mathbf{x}_j \in \mathcal{I}_i^-} \exp\left(E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_j) - E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_i)/\tau\right)$ .

 $F(\mathbf{w}, \mathcal{D}) = \frac{1}{n_o} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}} \tau \log g_{1i}(\mathbf{w}) + \tau \log g_{2i}(\mathbf{w}) \text{ and its gradient is given by}$ 

$$abla F(\mathbf{w}, \mathcal{D}) = rac{ au}{n_o} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}} \left( rac{
abla g_{1i}(\mathbf{w})}{g_{1i}(\mathbf{w})} + rac{
abla g_{2i}}{g_{2i}} 
ight)$$

The major cost of computing  $\nabla F(\mathbf{w}; \mathcal{D})$  lies on calculating  $g_{1i}(\mathbf{w})$  and  $g_{2i}(\mathbf{w})$  and their gradient for each pair, as it involves all the samples in  $\mathcal{T}_i^-$  and  $\mathcal{I}_i^-$ . Directly approximating  $g_{1i}$  and  $g_{2i}$  by a mini-batch of samples from  $\mathcal{T}_i^-$  and  $\mathcal{I}_i^-$  will reduce the computational cost but lead to a biased stochastic gradient of  $\nabla F(\mathbf{w}, \mathcal{D})$  due to the non-linear dependence of  $\nabla F(\mathbf{w}; \mathcal{D})$  on  $g_{1i}$  and  $g_{2i}$ , which will cause the issue of requiring a large batch size in order to converge.

To address this issue, we employ the moving average estimators for estimating  $g_{1i}$  and  $g_{2i}$  which gradually reduces the aforementioned biases to zero (Yuan et al., 2022). More specifically, let  $\mathbf{w}^t$ be the solution at iteration t. We randomly sample a mini batch  $\mathcal{B} \subset \mathcal{D}$ , and construct mini-batch negatives  $\mathcal{B}_{1,i} \subset \mathcal{T}_i^-$ ,  $\mathcal{B}_{2,i} \subset \mathcal{I}_i^-$  for each data  $(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{B}$  and construct the following stochastic estimations of  $g_{1i}(\mathbf{w}^t)$  and  $g_{2i}(\mathbf{w}^t)$ :

292 293 294

295 296

297 298 299

300 301

302

303

278 279

$$\hat{g}_{1i}(\mathbf{w}^t) := \frac{1}{|\mathcal{B}_{1,i}|} \sum_{\mathbf{t}_j \in \mathcal{B}_{1,i}} \exp((E_1(\mathbf{w}, \mathbf{x}_i)^\top E_2(\mathbf{w}, \mathbf{t}_j) - E_1(\mathbf{w}, \mathbf{x}_i)^\top E_2(\mathbf{w}, \mathbf{t}_i))/\tau)$$
$$\hat{g}_{2i}(\mathbf{w}^t) := \frac{1}{|\mathcal{B}_{2,i}|} \sum_{\mathbf{x}_j \in \mathcal{B}_{2,i}} \exp((E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_j) - E_2(\mathbf{w}, \mathbf{t}_i)^\top E_1(\mathbf{w}, \mathbf{x}_i))/\tau).$$

The moving averaging estimators of  $g_{1i}(\mathbf{w}^t)$  and  $g_{2i}(\mathbf{w}^t)$  denoted by  $u_{1i}^t$  and  $u_{2i}^t$  are updated by:

$$u_{1i}^{t+1} = (1 - \gamma_1)u_{1i}^t + \gamma_1 \hat{g}_{1i} \left(\mathbf{w}^t\right), \ u_{2i}^{t+1} = (1 - \gamma_1)u_{2i}^t + \gamma_1 \hat{g}_{2i} \left(\mathbf{w}^t\right), \tag{6}$$

where  $\gamma_1 \in (0, 1)$  is a hyper-parameter. The gradient estimator of  $F(\mathbf{w}^t, \mathcal{D})$  is computed by

$$G_{1}^{t} = \frac{\tau}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \left( \nabla \hat{g}_{1i} \left( \mathbf{w}^{t} \right) / u_{1i}^{t} + \nabla \hat{g}_{2i} \left( \mathbf{w}^{t} \right) / u_{2i}^{t} \right).$$
(7)

The gradient of the quadratic penalized term at  $\mathbf{w}^t$  can be approximated similarly by

$$G_2^t = \frac{1}{|\mathcal{B}_c|} \sum_{k \in \mathcal{B}_c} \beta[u_k^t]_+ \nabla \hat{h}_k(\mathbf{w}^t), \tag{8}$$

where  $\mathcal{B}_c$  denotes a sampled subset of protected tasks,  $\hat{h}_k(\mathbf{w}^t)$  denotes a mini-batch estimator of  $h_k(\mathbf{w}^t)$  using mini-batch  $\mathcal{B}_k \subset \mathcal{D}_k$ , and  $u_k^t$  is the moving average estimator of  $h_k(\mathbf{w}^t)$  computed by

$$u_{k}^{t+1} = (1 - \gamma_{2})u_{k}^{t} + \gamma_{2}\hat{h}_{k}(\mathbf{w}^{t}), \ \hat{h}_{k}(\mathbf{w}^{t}) = \frac{1}{|\mathcal{B}_{k}|} \sum_{j \in \mathcal{B}_{k}} \ell_{ce}(\mathbf{w}, \mathbf{x}_{j}, y_{j}) - \ell_{ce}(\mathbf{w}_{old}, \mathbf{x}_{j}, y_{j}).$$
<sup>(9)</sup>

We emphasize that the gradient estimator in (8) related to the protected tasks, where each protected task has an effective weight  $\beta[u_k^t]_+$  that is dynamically changing in the learning process, is the key difference from the native weighting method mentioned at the beginning.

The key steps are presented in Algorithm 1. For analysis, we make the following assumptions.

**Assumption 1** (a)  $g_1(\cdot)$  and  $g_2(\cdot)$  are  $L_g$ -Lipschitz continuous and  $L_{\nabla g}$ -smooth. (b) There exist  $C_g > 0$  and  $c_g > 0$  such that  $c_g \le \min\{g_1(\cdot), g_2(\cdot)\}$  and  $\max\{g_1(\cdot), g_2(\cdot)\} \le C_g$ . (c)  $h_k(\cdot)$  is  $L_h$ -Lipschitz continuous and  $L_{\nabla h}$ -smooth for  $k = 1, \cdots, m$ .

Assumption 2 There exists  $\mathbf{w}^0$  such that  $h_k(\mathbf{w}^0) \leq 0$  for  $k = 1, \cdots, m$ .

Assumption 3 (a)  $\mathbb{E}[\|\hat{g}_{1i}(\mathbf{w}) - g_{1i}(\mathbf{w})\|^2] \leq \sigma_g^2/|\mathcal{B}_{1i}|, \mathbb{E}[\|\hat{g}_{2i}(\mathbf{w}) - g_{2i}(\mathbf{w})\|^2] \leq \sigma_g^2/|\mathcal{B}_{2i}|;$ (b)  $\mathbb{E}[\|\nabla\hat{g}_1(\mathbf{w}) - \nabla g_{1i}(\mathbf{w})\|^2] \leq \sigma_{\nabla g}^2/|\mathcal{B}_{1i}|, \mathbb{E}[\|\nabla\hat{g}_{2i}(\mathbf{w}) - \nabla g_{2i}(\mathbf{w})\|^2] \leq \sigma_{\nabla g}^2/|\mathcal{B}_{2i}|;$  (c)  $\mathbb{E}[\|\nabla\hat{h}_k(\mathbf{w}) - \nabla h_k(\mathbf{w})\|^2] \leq \sigma_{\nabla h}^2/|\mathcal{B}_k|;$  (d)  $\mathbb{E}[\|\hat{h}_k(\mathbf{w}) - h_k(\mathbf{w})\|^2] \leq \sigma_h^2/|\mathcal{B}_k|$  for  $k = 1, \cdots, m$ . Assumption 4 There exists a constant  $\delta > 0$  such that  $\|\nabla h(\mathbf{w}^t)[h(\mathbf{w}^t)]_+\| \geq \delta \|[h(\mathbf{w}^t)]_+\|$  for

- Assumption 4 There exists a constant  $\delta > 0$  such that  $\|\nabla h(\mathbf{w}^t)[h(\mathbf{w}^t)]_+\| \ge \delta \|[h(\mathbf{w}^t)]_+\|$  for t = 0, ..., T, where  $h(\mathbf{w}) = [h_1(\mathbf{w}), \dots, h_m(\mathbf{w})]^\top$  and  $\nabla h(\mathbf{w}) = [\nabla h_1(\mathbf{w}), \dots, \nabla h_m(\mathbf{w})]$ .
- **Remark:** Assumption 1 has been justified in the earlier work (Yuan et al., 2022; Qiu et al., 2023) for optimizing a global contrastive loss. Assumption 2 is easily satisfied with  $\mathbf{w}^0 = \mathbf{w}_{old}$ . Assumption 3

324 Algorithm 1 Algorithm for solving (4) 325 1: **Initialization:** choose  $\mathbf{w}^0$ ,  $\beta$ ,  $\gamma_1$ ,  $\gamma_2$ ,  $\theta$  and step sizes  $\eta$ . 326 2: for  $t = 0, 1, \dots, T - 1$  do 327 Sample image-text pairs  $\mathcal{B}$  from  $\mathcal{D}$  and protected tasks  $\mathcal{B}_c$  from  $\{1, \dots, m\}$ . 3: 328 4: for each  $(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{B}$  do Update  $u_{1i}^t$  and  $u_{2i}^t$  by Eqn. (6) 5: 330 6: end for 331 Update the estimator of gradient  $\nabla F(\mathbf{w}^t, \mathcal{D})$  by  $G_1^t$  as in Eqn. (7) 7: 332 8: for each  $k \in \mathcal{B}_c$  do 333 9: Sample a minibatch of data from  $\mathcal{D}_k$  denoted by  $\mathcal{B}_k$ . Update the estimators of  $h_k$  by Eqn. (9). 10: 334 11: end for 335 Compute the stochastic gradient estimator  $G_2^t$  as in Eqn. (8) 12: 336 Update Gradient Estimator  $v^{t+1} = (1-\theta)v^t + \theta(G_1^t + G_2^t)$ Update w by  $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta v^{t+1}$ . 13: 337 14: 338 15: end for 339 is a standard one that bounds the variance of mini-batch estimators. Assumption 4 is also made in 340 many existing works on optimization with non-convex constraints (Sahin et al., 2019; Xie & Wright, 341 2021; Alacaoglu & Wright, 2024; Lin et al., 2022; Li et al., 2024). This assumption is equivalent to 342 that the quadratic penalty term  $H(\mathbf{w}) := \frac{\beta}{2m} \|[\mathbf{h}(\mathbf{w})]_+\|^2$  satisfies the Polyak-Lojasiewicz inequality 343 at  $\mathbf{w} = \mathbf{w}^t$ , meaning that there exists  $\delta \ge 0$  such that  $\|\nabla H(\mathbf{w}^t)\|^2 \ge \frac{2\delta^2\beta}{m}H(\mathbf{w}^t)$ . Without this assumption, (4) may be intractable because there may exist an iterate  $\mathbf{w}^t$  such that  $H(\mathbf{w}^t) > 0$  but 344 345  $\nabla H(\mathbf{w}^t) = 0$ , meaning that  $\mathbf{w}^t$  is infeasible but at a flat location of  $H(\mathbf{w})$  so  $\mathbf{w}^t$  may get trapped at 346 347 this location forever. We will show later that a small  $\delta$  in Assumption 4 will increase the complexity of our algorithm. Hence, we will present an approach in next subsection to increase  $\delta$ . 348 349 For a non-convex optimization problem like (4), finding a globally optimal solution is intractable, so 350 almost all numerical algorithms for non-convex problems can only guarantee a Karush-Kuhn-Tucker 351 (KKT) solution defined below. 352 **Definition 2** A solution w is a KKT solution to (4) if there exist  $\lambda = (\lambda_1, \dots, \lambda_m)^\top \in \mathbb{R}^m_+$  such 353 that  $\nabla F(\mathbf{w}, \mathcal{D}) + \nabla h(\mathbf{w}) \lambda = 0$ ,  $h(\mathbf{w}) \leq 0$  and  $\lambda_k h_k(\mathbf{w}) = 0$  for  $k = 1, \dots, m$ . 354 355 We present the convergence theorem of Algorithm 1 as follows, which shows the iteration complexity 356 of Algorithm 1 for finding an  $\epsilon$ -KKT solution, i.e., a solution satisfying the three conditions in 357 Definition 2 up to  $\epsilon$  precision. The proof of the theorem is presented in Appendix C.3. **Theorem 1** Suppose Assumptions 1, 2, 3 and 4 hold. Also, suppose, in Algorithm 1, set  $\beta = \frac{1}{\epsilon\delta}$ ,  $\theta = \min\{\frac{\epsilon^4\delta^2\min\{|\mathcal{B}_c|,|\mathcal{B}_k|\}}{672(\sigma_{\nabla h}^2+L_h^2)}, \frac{\epsilon^2\min\{|\mathcal{B}|,|\mathcal{B}_{1i}|,|\mathcal{B}_{2i}|\}}{1344L_f^2(\sigma_{\nabla g}^2+L_g^2)}\}, \gamma_1 = \gamma_2 = \min\{\frac{5n_0\theta}{3|\mathcal{B}|}, \frac{5m\theta}{3|\mathcal{B}_c|}, \frac{\epsilon^4\delta^2|\mathcal{B}_k|}{26880\sigma_h^2\tilde{C}_{\nabla h}^2}\}$  and  $\eta = \min\{\frac{1}{12(L_F+\beta L_H)}, \frac{\theta}{8\sqrt{3}L_F}, \frac{\theta}{8\sqrt{3}L_H\beta}, \frac{\gamma_1|\mathcal{B}|}{40\sqrt{6}L_gL_f\tilde{C}_{\nabla g}n_0}, \frac{\gamma_2|\mathcal{B}_c|}{40\sqrt{6}\beta L_h\tilde{C}_{\nabla h}m}\}, where \tilde{C}_{\nabla g} := \sigma_{\nabla g} + L_g, \tilde{C}_{\nabla h} := \sigma_{\nabla h} + L_h, L_f := \frac{\tau}{c_g}, L_{\nabla f} := \frac{\tau}{c_g^2}, L_F := 2(L_{\nabla g}L_f + L_{\nabla f}L_g^2)$  and  $L_H := 2L_{\nabla h} + L_h^2$ . Then there exists  $\lambda \in \mathbb{R}^m_+$  such that after  $T = O(\epsilon^{-7}\delta^{-3})$  iterations Algorithm 1 satisfies 358 359 360 361 362 363 364 365  $\mathbb{E}\left[\|\nabla F(\mathbf{w}^{\hat{t}}, \mathcal{D}) + \nabla \boldsymbol{h}(\mathbf{w}^{\hat{t}})\boldsymbol{\lambda}\|\right] \leq \epsilon, \quad \mathbb{E}[\|[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}\|] \leq \epsilon, \quad \mathbb{E}[\boldsymbol{\lambda}^{\top}[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}] \leq \epsilon$ 366 where  $\hat{t}$  selected uniformly at random from  $\{1, \dots, T\}$ . 367 368 **Remark:** It is notable that the order of complexity in terms of  $\epsilon$  is higher than that of standard 369 learning (i.e.,  $O(\epsilon^{-4})$ ). While the complexity for a stochastic constrained optimization could be 370 inherently higher than unconstrained optimization (Alacaoglu & Wright, 2024), we note that the 371 above complexity is also weaker than the state-of-the-art complexity of stochastic constrained 372 optimization (Alacaoglu & Wright, 2024). We remark that this is a limitation of the present work 373 due to two reasons: (i) we use the moving average gradient estimator for sake of implementation; 374 in contrast, they use the advanced variance reduced gradient estimator (STORM), which incurs

additional overhead; (ii) we use a constant  $\beta$  and they use an increasing  $\beta$ . In our experiments shown in ablation studies, we find that using a constant  $\beta$  is generally better than using an increasing  $\beta$ . Additionally, the dependence on  $\delta$  could also slow down the convergence. We mitigate this issue by utilizing task-dependent heads for CLIP models justified below.

#### 378 5.2PROMOTING DEVELOPMENTAL SAFETY VIA TASK-DEPENDENT HEADS 379

Below, we present a way to design the text encoder of the CLIP model such that the value of  $\delta$ 380 could be larger. Without causing confusion, we denote by w the parameter of the text encoder, 381 which consists of two components u and W such that the text embedding  $E_2(\mathbf{w}, \mathbf{t}) \in \mathbb{R}^{d_2}$  can 382 be represented as  $E_2(\mathbf{w}, \mathbf{t}) = W \cdot \overline{E}_2(\mathbf{u}, \mathbf{t})$ , where  $\overline{E}_2(\mathbf{u}, \cdot) \in \mathbb{R}^{d_1}$  is a backbone encoder while  $W \in \mathbb{R}^{d_2 \times d_1}$  is called the head. The idea of task-dependent heads is to let each task k have 384 its own head  $W_k = W + U_k V_k^{\top}$  using low rank matrices  $U_k \in \mathbb{R}^{d_2 \times r}$  and  $V_k \in \mathbb{R}^{d_1 \times r}$ , where 385  $r < \min(d_1, d_2)$  is the rank chosen as a hyper-parameter. The output of this class-specific text encoder 386 for task k is  $E_2(\mathbf{u}, W, U_k, V_k, \mathbf{t}_k) = (W + U_k V_k^{\top}) \cdot \overline{E}_2(\mathbf{u}, \mathbf{t}_k)$ . Note that  $\|\nabla \mathbf{h}(\mathbf{w}^t)^{\top} [\mathbf{h}(\mathbf{w}^t)]_+ \|^2 \ge |\nabla \mathbf{h}(\mathbf{w}^t)|^2$  $\lambda_{\min}(\nabla \mathbf{h}(\mathbf{w}^t)^\top \nabla \mathbf{h}(\mathbf{w}^t)) \| [\mathbf{h}(\mathbf{w}^t)]_+ \|^2$ , where  $\lambda_{\min}(\cdot)$  represents the smallest eigenvalue of a matrix. This means  $\min_t \lambda_{\min}(\nabla \mathbf{h}(\mathbf{w}^t)^\top \nabla \mathbf{h}(\mathbf{w}^t))$  is a lower bound of  $\delta$  in Assumption 4. The following 387 388 lemma shows that, after expanding  $\mathbf{w}$  with  $U_k$  and  $V_k$ ,  $\lambda_{\min}(\nabla \mathbf{h}(\mathbf{w}^t)^\top \nabla \hat{\mathbf{h}}(\mathbf{w}^t))$  may increase at 389 some  $U_k$  and  $V_k$ , providing some insight on why the task-dependent heads help to increase the 390 parameter  $\delta$  in Assumption 4, reducing the total complexity of our algorithm according to Theorem 1. 391 392 **Lemma 2** Let  $\mathbf{U} = (U_1, ..., U_m)$  and  $\mathbf{V} = (V_1, ..., V_m)$ . Let  $\mathbf{w} = (W, \mathbf{u})$ ,  $\hat{\mathbf{w}} = (W, \mathbf{u}, \mathbf{U}, \mathbf{V})$ ,  $h_k(\mathbf{w}) = h_k(W, \mathbf{u})$ , and  $\hat{h}_k(\hat{\mathbf{w}}) = h_k(W + U_k V_k^{\top}, \mathbf{u})$ . Suppose  $U_k V_k^{\top} = \mathbf{0}$  for all k's. We have 393 394  $\lambda_{\min}\left(\nabla \widehat{\boldsymbol{h}}(\widehat{\mathbf{w}})^{\top} \nabla \widehat{\boldsymbol{h}}(\widehat{\mathbf{w}})\right) \geq \lambda_{\min}\left(\nabla \boldsymbol{h}(\mathbf{w})^{\top} \nabla \boldsymbol{h}(\mathbf{w})\right) + \min_{k} \left\{ \|\nabla_{W} h_{k}(\mathbf{w}) V_{k}\|_{F}^{2}, \left\|\nabla_{W} h_{k}(\mathbf{w})^{\top} U_{k}\right\|_{F}^{2} \right\},$ 395 396 where  $\hat{\boldsymbol{h}}(\hat{\mathbf{w}}) = [\hat{h}_1(\hat{\mathbf{w}}), \dots, \hat{h}_m(\hat{\mathbf{w}})]^\top$  and  $\nabla \hat{\boldsymbol{h}}(\hat{\mathbf{w}}) = [\nabla \hat{h}_1(\hat{\mathbf{w}}), \dots, \nabla \hat{h}_m(\hat{\mathbf{w}})].$ 397 398

Following this lemma, in our experiments, we employ the task-dependent heads by setting the initial value of  $U_k$  to zero so  $U_k V_k^{\top} = 0$ . The proof of the above lemma is given in Appendix C.2 399

#### 6 EXPERIMENTS 401

400

In this section, we conduct extensive experiments to understand our proposed method, including 402 an overview of how our approach works, performance comparison with other strong baselines and 403 potential of our method in multi-round model development. Detailed ablation studies about design 404 choices are included in Appendix A.6. 405

406 **Dataset.** We experiment on the large-scale diverse driving image dataset, namely BDD100K (Seita, 2018). This dataset involves classification of six weather conditions, i.e., *clear, overcast, snowy, rainy,* 407 partly cloudy, foggy, and of six scene types, i.e., highway, residential area, city street, parking lot, 408 gas station, tunnel. We consider three settings with foggy, overcast and tunnel as the target class 409 separately and other weather conditions or scenes as protected tasks. Moreover, we experiment on 410 the scene recognition dataset, Places365 which has 365 classes (Zhou et al., 2017), to verify the 411 effectiveness of the proposed method in handling a large number of constraints. We consider *dressing* 412 room as the target class, as it has fewest samples in the dataset. More experimental settings are in A.1. 413

**Evaluation Metrics.** We measure improvement on target task with  $\Delta Acc(Target) =$ 414  $Acc(Target, \mathbf{w}_{new}) - Acc(Target, \mathbf{w}_{old})$ . Besides, we utilize "DevSafety(acc)" (i.e., Eqn. 2) to mea-415 sure the empirical MDS. As optimization involves randomness, we run all the experiments with five 416 different random seeds then calculate the average target accuracy and the percentage of times that 417 DevSafety(acc) is non-negative, denoted as **Rentention Ratio**, to measure the possibility of strictly 418 preserving the performance on protected tasks. (e.g., the Retention Ratio is 60% if 3 out of 5 runs of 419 the method preserve previous performance for all protected tasks.) 420

**Baselines.** To verify the effectiveness of our algorithms, we compare our proposed algorithm with 421 the following baseline methods: (1) FLYP (Goyal et al., 2023), a state-of-the-art CLIP finetuning 422 method that optimizes a contrastive loss on all available data including those used in our objective 423 and constraints. In our experiments, we utilize the same global contrastive loss (GCL) (Yuan et al., 424 2022) instead of mini-batch contrastive loss; (2) Weighted Combination of Contrastive Losses 425 (WCCL), which utilizes a weight to combine GCL losses on protected tasks and the target task to 426 control the tradeoff between them to achieve model developmental safety; (3) GEM (Lopez-Paz & 427 Ranzato, 2017b), which is a strong CL baseline motivated by a similar idea utilizing data of previous 428 tasks for constraints; (4)  $Co^2L$  (Cha et al., 2021), which is a recent SOTA contrastive continual learning baseline; (5) Regularization Method (RM), as commonly adopted in continual learning 429 literature (Rebuffi et al., 2017; Castro et al., 2018), directly takes the constraints in Eqn. (4) as a 430 regularization term by adding it to the objective function with a regularization weight  $\alpha$ . All methods 431 start from the same CLIP model. More details about baselines are presented in Appendix A.1.4.



Figure 2: Visualization of the learning trajectory. Each dot denotes a solution with lighter color being earlier iterations and darker being later iterations.



Figure 3: Performance Comparison with Baselines. Dot lines represent the performance of the base model on the target task. Detailed numbers are presented in Table 4, 5, 6.

#### 6.1 VISUALIZATION OF LEARNING PROCESS

444

445

446

447

448

453

454

455

To provide a direct understanding of why and how the proposed algorithm works, we present the 456 learning trajectory of the algorithm in Figure 2. Each dot in this figure represents a solution during 457 the learning processing, with lighter colors indicating earlier stages and darker colors representing 458 later stages. From the top four figures for training sets, we can observe a common trend that solutions 459 start from the lower left and move toward the upper right, indicating the algorithm endeavors to 460 enhance the performance of the targeted task while improving developmental safety on protected 461 tasks. Similarly, this trend extends to the validation sets, shown in the bottom row, demonstrating the 462 generalization capability of the proposed algorithm. It is striking to see that, when targeting *Dressing* 463 Room in Places365 dataset with all other 364 classes as protected tasks, our method are still able to 464 achieve developmental safety in training set and generalize to validation set. These observations can 465 also be found in separate views of DevSafety vs epochs and  $\Delta Acc(Target)$  vs epochs shown in Fig 5.

#### 466 467 6.2 Comparison with Baselines for Model developmental safety

In this part, we compare the proposed method with baselines to demonstrate the superiority. Specifi-468 cally, we focus on two metrics, i.e., Rention Ratio for measuring the possibility of strictly preserving 469 the performance on all protected tasks and accuracy on the target task. The details of hyperparameter 470 tuning is presented in Appendix A.1.3. On autonomous driving BDD100K dataset, we conduct 471 experiments with different numbers of data for constraints, i,e., 100, 1k, 2k, 4k from each task. The 472 comparison results are presented in Figure 3. The figure illustrates that improving the base model on 473 the target tasks is not challenging, as nearly all methods accomplish this effortlessly. However, all 474 baselines, including the strong continual learning baselines GEM and Co<sup>2</sup>L, exhibit a zero Rention 475 Ratio across almost all settings, showing the insufficiency of existing methods for ensuring zero 476 forgetting on protected tasks. In contrast, our method begins to ensure developmental safety with 1k 477 samples per protected task and even 100 samples for the target class *tunnel*. Besides, the Rention Ratio increases when using more data for constraints, consistent with the result obtained in Lemma 1 478 (Refer to Table 7 for more results). Notably, our method achieves a 100% Rention Ratio with 4k 479 samples per protected task in all three settings, while improving accuracy on the target class. We also 480 see that the target overcast is most difficult to improve as the base model already has 73.6% accuracy. 481

From Figure 3, we notice that the baseline RM fails to achieve MDS, even though it has a tunable weight parameter  $\alpha$  for protected tasks. Comparison with RM can directly verify the advantage of our method as the only difference between the two methods is how to handle the protected tasks. From Eqn. (8), we can see that our algorithm has an effective weight  $\beta[u_k^t]_+$  for each protected task. It is adaptively adjusted during learning, and depends on the degree of violation of constraints, i.e., the



Figure 4: (Left) Adaptive weight adjustments for each protected task during training (Targeting *Foggy*). Weights shown are averaged over every 600 iterations for visualization. (Right) Performance comparison with baseline RM when targeting *Dresssing Room* on Places365 Dataset, with 2k samples per constraint. Red line denotes base model's performance, green diamonds denote the target class. RM baseline shown is for weight  $\alpha = 1000$  and more plots for other weights are presented in 6.

499 larger the violation, the larger the weight. Figure 4 (left) shows that these effective weights gradually 500 decrease to zero during the learning of our algorithm, which allows the model to learn from the target task while satisfying constraints. This mechanism plays a big role in not only achieving MDS but 501 also improving the performance on the target task. In contrast, RM uses a constant weight  $\alpha$  for 502 every protected task. Simply increasing  $\alpha$  may not ensure MDS, due to varied learning difficulty between protected tasks. Besides, too large  $\alpha$  will also harm the performance of the target task. We 504 further investigate the phenomenon in Appendix A.4 and find that, with a uniform weight for all the 505 protected tasks, it might preserve previous performance on some of the protected tasks but fail to 506 achieve MDS for all the protected tasks, even with a very high weight  $\alpha$ . 507

To further verify the effectiveness of our method in handling a large number of constraints, we experiment on the Place365 dataset, compared with RM, targeting *Dressing room* class and protecting the other 364 tasks in Figure 4 (right). It shows that even with hundreds of protected tasks, our method is still effective in preserving their performance, whereas RM causes performance drops in around 30 classes, failing to ensure MDS.

513 We include more ablation studies in Appendix A.6 to (i) demonstrate the benefit of augmenting target 514 dataset with retrieved target-related image-text pairs; (ii) verify the benefit of task-dependent heads; 515 (iii) verify theoretical result Lemma 2; (iv) compare constant  $\beta$  vs increasing  $\beta$ .

516 6.3 PERFORMANCE WITH MULTIPLE ROUNDS OF MODEL DEVELOPMENT

517 Finally, to demonstrate the effectiveness of the proposed retention-centric framework in iterative 518 model development process, we conduct two consecutive rounds of development on recognizing 519 weather conditions. Specifically, we first target at overcast task, taking all the other five weather 520 conditions as protected tasks, then with one selected improved model, we successively improve 521 the model, targeting at improving the performance of the *foggy* task. As shown in Fig. 1, our 522 method notably improves the performance of the *overcast* task in the first round while ensuring 523 the performance of other tasks does not decrease. In the second round, it continues to enhance the performance of the *foggy* task. Simultaneously, it preserves the performance, if not boosts it, across 524 other tasks, with only a slight decrease on the *snowy* task, showing the effectiveness of the proposed 525 framework for maintaining the model developmental safety. 526

7 CONCLUSION

527

528 In this paper, we introduced the concept of "model developmental safety" to ensure that model devel-529 opment not only acquires new capabilities but also strictly preserves those already owns, addressing 530 the critical developmental safety oversight in existing ML/AI studies. To ensure model developmental 531 safety, we proposed a retention-centric framework by formulating the model developmental safety 532 as data-dependent constraints. We proposed an efficient constrained optimization algorithm with 533 theoretical guarantees to develop a pretrained vision-language model (CLIP model) for improving 534 existing image classification capabilities. Comprehensive experiments demonstrate the effectiveness of the algorithm in enhancing vision-based perception capabilities in autonomous driving and scene recognition, showing its practical value in real-world scenarios. As the proposed framework in this paper is a generic retention-centric optimization framework, it can be potentially extended to various scenarios or models, such as finetuning LLMs or enhancing object detection systems and motion 538 prediction tasks for autonomous driving. We hope our work can inspire researchers in safety-critical application domain for more exploration.

# 540 REFERENCES

547

559

561

566

569

570

571

579

- Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In International conference on machine learning, pp. 22–31. PMLR, 2017.
- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna M. Wallach. A reductions approach to fair classification. *CoRR*, abs/1803.02453, 2018. URL http://arxiv.org/abs/1803.02453.
- Ahmet Alacaoglu and Stephen J Wright. Complexity of single loop algorithms for nonlinear
   programming with stochastic objective and constraints. In *International Conference on Artificial Intelligence and Statistics*, pp. 4627–4635. PMLR, 2024.
- Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning. *CoRR*, abs/1812.03596, 2018. URL http://arxiv.org/abs/1812.03596.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul F. Christiano, John Schulman, and Dan Mané.
   Concrete problems in AI safety. *CoRR*, abs/1606.06565, 2016. URL http://arxiv.org/ abs/1606.06565.
- Aleksandr Y Aravkin, James V Burke, Dmitry Drusvyatskiy, Michael P Friedlander, and Scott Roy.
   Level-set methods for convex optimization. *Mathematical Programming*, 174:359–390, 2019.
  - Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Dimitri P Bertsekas. Constrained optimization and Lagrange multiplier methods. Academic press, 2014.
- Digvijay Boob, Qi Deng, and Guanghui Lan. Stochastic first-order methods for convex and nonconvex
   functional constrained optimization. *Mathematical Programming*, 197(1):215–279, 2023.
- 567 Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification: A survey of some recent advances. *ESAIM: probability and statistics*, 9:323–375, 2005.
  - Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. arXiv preprint arXiv:2303.12712, 2023.
- Archana Bura, Aria HasanzadeZonuzy, Dileep Kalathil, Srinivas Shakkottai, and Jean-Francois
   Chamberland. Dope: Doubly optimistic and pessimistic exploration for safe reinforcement
   *Advances in neural information processing systems*, 35:1047–1059, 2022.
- Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari.
   End-to-end incremental learning. In *Proceedings of the European conference on computer vision* (*ECCV*), pp. 233–248, 2018.
  - Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *Proceedings of the IEEE/CVF International conference on computer vision*, pp. 9516–9525, 2021.
- Luiz FO Chamon, Santiago Paternain, Miguel Calvo-Fullana, and Alejandro Ribeiro. Constrained
   learning with non-convex losses. *IEEE Transactions on Information Theory*, 69(3):1739–1760, 2022.
- Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient
   lifelong learning with a-GEM. In *International Conference on Learning Representations*, 2019.
   URL https://openreview.net/forum?id=Hkf2\_sC5FX.
- Xinlei Chen, Abhinav Shrivastava, and Abhinav Gupta. Neil: Extracting visual knowledge from web data. In *Proceedings of the IEEE international conference on computer vision*, pp. 1409–1416, 2013.
- 592 McKinsey & Company. The future of autonomous driving. https://www. 593 mckinsey.com/industries/automotive-and-assembly/our-insights/ autonomous-drivings-future-convenient-and-connected, 2023.

594 595 596	Andrew Cotter, Heinrich Jiang, Maya Gupta, Serena Wang, Taman Narayan, Seungil You, and Karthik Sridharan. Optimization with non-differentiable constraints with applications to fairness, recall, churn, and other goals. <i>Journal of Machine Learning Research</i> , 20(172):1–59, 2019.
597 598 599 600	Dongsheng Ding, Xiaohan Wei, Zhuoran Yang, Zhaoran Wang, and Mihailo Jovanovic. Provably efficient safe exploration via primal-dual policy optimization. In <i>International conference on artificial intelligence and statistics</i> , pp. 3304–3312. PMLR, 2021.
601 602 603	Francisco Facchinei, Vyacheslav Kungurtsev, Lorenzo Lampariello, and Gesualdo Scutari. Ghost penalties in nonconvex constrained optimization: Diminishing stepsizes and iteration complexity. <i>Mathematics of Operations Research</i> , 46(2):595–627, 2021.
604 605 606	Mehrdad Farajtabar, Navid Azizan, Alex Mott, and Ang Li. Orthogonal gradient descent for continual learning. <i>CoRR</i> , abs/1910.07104, 2019. URL http://arxiv.org/abs/1910.07104.
607 608 609	Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like you pretrain: Improved finetuning of zero-shot vision models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 19338–19347, 2023.
610 611 612 613 614 615	Yunhui Guo, Mingrui Liu, Tianbao Yang, and Tajana Rosing. Improved schemes for episodic memory-based lifelong learning. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria- Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/ hash/0b5e29aalacf8bdc5d8935d7036fa4f5-Abstract.html.
616 617 618	Erfan Yazdandoost Hamedani and Necdet Serhat Aybat. A primal-dual algorithm with line search for general convex-concave saddle point problems. <i>SIAM Journal on Optimization</i> , 31(2):1299–1329, 2021.
619 620 621	Han Huang, Haitian Zhong, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. Kebench: A benchmark on knowledge editing for large vision-language models. <i>arXiv preprint arXiv:2403.07350</i> , 2024.
622 623 624	Sham M Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. <i>Advances in neural information processing systems</i> , 21, 2008.
625 626 627 628 629 630	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. <i>Proceedings of the National Academy of Sciences</i> , 114(13):3521–3526, 2017. doi: 10.1073/pnas.1611835114. URL https://www.pnas.org/doi/abs/10.1073/pnas.1611835114.
631 632 633 634	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35: 22199–22213, 2022.
635 636	Philip Koopman and Michael Wagner. Challenges in autonomous vehicle testing and validation. <i>SAE International Journal of Transportation Safety</i> , 4(1):15–24, 2016.
637 638 639	Guanghui Lan and Renato DC Monteiro. Iteration-complexity of first-order penalty methods for convex programming. <i>Mathematical Programming</i> , 138(1):115–139, 2013.
640 641	Guanghui Lan and Renato DC Monteiro. Iteration-complexity of first-order augmented lagrangian methods for convex programming. <i>Mathematical Programming</i> , 155(1):511–547, 2016.
643 644	Guanghui Lan and Zhiqiang Zhou. Algorithms for stochastic optimization with expectation con- straints. <i>arXiv preprint arXiv:1604.03887</i> , 2016.
645 646 647	Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang. Overcoming catastrophic forgetting by incremental moment matching. In <i>Proceedings of the 31st International Conference on Neural Information Processing Systems</i> , NIPS'17, pp. 4655–4665, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.

657

664

- Kilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong. Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting. In *International Conference on Machine Learning*, 2019.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. CoRR, abs/1606.09282, 2016. URL
   http://arxiv.org/abs/1606.09282.
- Zichong Li, Pin-Yu Chen, Sijia Liu, Songtao Lu, and Yangyang Xu. Stochastic inexact augmented lagrangian method for nonconvex expectation constrained optimization. *Computational Optimization and Applications*, 87(1):117–147, 2024.
- Mingfu Liang, Jong-Chyi Su, Samuel Schulter, Sparsh Garg, Shiyu Zhao, Ying Wu, and Manmohan
   Chandraker. Aide: An automatic data engine for object detection in autonomous driving. In
   *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14695–14706, 2024.
- Qihang Lin, Runchao Ma, and Tianbao Yang. Level-set methods for finite-sum constrained convex
   optimization. In *International conference on machine learning*, pp. 3112–3121. PMLR, 2018a.
- Qihang Lin, Selvaprabu Nadarajah, and Negar Soheili. A level-set method for convex optimization with a feasible solution path. *SIAM Journal on Optimization*, 28(4):3290–3311, 2018b.
- Qihang Lin, Runchao Ma, and Yangyang Xu. Complexity of an inexact proximal-point penalty method
   for constrained smooth non-convex optimization. *Computational optimization and applications*,
   82(1):175–224, 2022.
- Sheng Liu, Haotian Ye, and James Zou. Reducing hallucinations in vision-language models via latent space steering. *arXiv preprint arXiv:2410.15778*, 2024.
- David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning.
  In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pp. 6470–6479, Red Hook, NY, USA, 2017a. Curran Associates Inc. ISBN 9781510860964.
- 677
   678
   679
   David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. Advances in neural information processing systems, 30, 2017b.
- Jinqi Luo, Tianjiao Ding, Kwan Ho Ryan Chan, Darshan Thaker, Aditya Chattopadhyay, Chris
   Callison-Burch, and René Vidal. Pace: Parsimonious concept engineering for large language
   models. arXiv preprint arXiv:2406.04331, 2024.
- Runchao Ma, Qihang Lin, and Tianbao Yang. Quadratically regularized subgradient methods for weakly convex optimization with weakly convex constraints. In *International Conference on Machine Learning*, pp. 6554–6564. PMLR, 2020.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
   Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing
   memory in a transformer. *arXiv preprint arXiv:2210.07229*, 2022.
- Tom Mitchell, William Cohen, Estevam Hruschka, Partha Talukdar, Bishan Yang, Justin Betteridge,
   Andrew Carlson, Bhavana Dalvi, Matt Gardner, Bryan Kisiel, et al. Never-ending learning.
   *Communications of the ACM*, 61(5):103–115, 2018.
- Arkadi Nemirovski. Prox-method with rate of convergence o (1/t) for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004.

702 703 704	Masahiro Ono, Marco Pavone, Yoshiaki Kuwata, and J Balaram. Chance-constrained dynamic programming with application to risk-aware robotic space exploration. <i>Autonomous Robots</i> , 39: 555–571, 2015.
705 706	Shubham Parashar, Zhiqiu Lin, Tian Liu, Xiangjue Dong, Yanan Li, Deva Ramanan, James Caverlee, and Shu Kong. The neglected tails of vision-language models 2024
707 708 709 710	German Ignacio Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. <i>CoRR</i> , abs/1802.07569, 2018. URL
711 712 713	<ul> <li>Santiago Paternain, Miguel Calvo-Fullana, Luiz FO Chamon, and Alejandro Ribeiro. Learning safe policies via primal-dual methods. In 2019 IEEE 58th Conference on Decision and Control (CDC), pp. 6491–6497. IEEE, 2019a.</li> </ul>
714 715 716 717	Santiago Paternain, Luiz Chamon, Miguel Calvo-Fullana, and Alejandro Ribeiro. Constrained reinforcement learning has zero duality gap. <i>Advances in Neural Information Processing Systems</i> , 32, 2019b.
718 719	Liangzu Peng, Paris Giampouras, and René Vidal. The ideal continual learner: An agent that never forgets. In <i>International Conference on Machine Learning</i> , pp. 27585–27610. PMLR, 2023.
720 721 722 723	Samuel Pfrommer, Tanmay Gautam, Alec Zhou, and Somayeh Sojoudi. Safe reinforcement learn- ing with chance-constrained model predictive control. In <i>Learning for Dynamics and Control</i> <i>Conference</i> , pp. 291–303. PMLR, 2022.
724 725 726	Tu-Hoa Pham, Giovanni De Magistris, and Ryuki Tachibana. Optlayer-practical constrained opti- mization for deep reinforcement learning in the real world. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pp. 6236–6243. IEEE, 2018.
727 728	VT Polyak and NV Tret'yakov. The method of penalty estimates for conditional extremum problems. USSR Computational Mathematics and Mathematical Physics, 13(1):42–58, 1973.
730 731 732 733 734 735 736	Zi-Hao Qiu, Quanqi Hu, Zhuoning Yuan, Denny Zhou, Lijun Zhang, and Tianbao Yang. Not all semantics are created equal: Contrastive self-supervised learning with automatic temperature individualization. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pp. 28389–28421. PMLR, 2023. URL https://proceedings.mlr.press/v202/qiu23a.html.
737 738	Haoxuan Qu, Hossein Rahmani, Li Xu, Bryan Williams, and Jun Liu. Recent advances of continual learning in computer vision: An overview. <i>arXiv preprint arXiv:2109.11369</i> , 2021.
739 740 741 742 743 744 745	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event, volume 139 of Proceedings of Machine Learning Research, pp. 8748–8763. PMLR, 2021a. URL http://proceedings.mlr.press/v139/ radford21a.html.
746 747 748 749	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021b.
750 751 752 753	Nijat Rajabli, Francesco Flammini, Roberto Nardone, and Valeria Vittorini. Software verification and validation of safe autonomous cars: A systematic literature review. <i>IEEE Access</i> , 9:4797–4819, 2020.
754 755	Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In <i>Proceedings of the IEEE conference on</i> <i>Computer Vision and Pattern Recognition</i> , pp. 2001–2010, 2017.

756 757 758 759	Alexander Robey, Luiz Chamon, George J Pappas, Hamed Hassani, and Alejandro Ribeiro. Adversar- ial robustness with semi-infinite constrained learning. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 34:6198–6215, 2021.
760 761 762 763 764	David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 32. Cur- ran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper/2019/ file/fa7cdfadla5aaf8370ebeda47alfflc3-Paper.pdf.
765 766 767 768	Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. <i>CoRR</i> , abs/1606.04671, 2016. URL http://arxiv.org/abs/1606.04671.
769 770 771	Mehmet Fatih Sahin, Ahmet Alacaoglu, Fabian Latorre, Volkan Cevher, et al. An inexact augmented lagrangian framework for nonconvex optimization with nonlinear constraints. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
772 773 774	Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. <i>arXiv preprint arXiv:2111.02114</i> , 2021.
776 777 777	Daniel Seita. Bdd100k: A large-scale diverse driving video database. <i>The Berkeley Artificial Intelligence Research Blog. Version</i> , 511:41, 2018.
779 780	Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. <i>arXiv preprint arXiv:1610.03295</i> , 2016.
781 782 783 784 785	Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran As- sociates, Inc., 2017. URL https://proceedings.neurips.cc/paper/2017/file/ 0efbe98067c6c73dba1250d2beaa81f9-Paper.pdf.
786 787 788	Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. <i>arXiv</i> preprint arXiv:1805.11074, 2018.
789 790 791 792	Brijen Thananjeyan, Ashwin Balakrishna, Suraj Nair, Michael Luo, Krishnan Srinivasan, Minho Hwang, Joseph E Gonzalez, Julian Ibarz, Chelsea Finn, and Ken Goldberg. Recovery rl: Safe reinforcement learning with learned recovery zones. <i>IEEE Robotics and Automation Letters</i> , 6(3): 4915–4922, 2021.
793 794 795	Garrett Thomas, Yuping Luo, and Tengyu Ma. Safe reinforcement learning by imagining the near future. <i>Advances in Neural Information Processing Systems</i> , 34:13859–13869, 2021.
796 797 798 799	Matteo Turchetta, Andrey Kolobov, Shital Shah, Andreas Krause, and Alekh Agarwal. Safe rein- forcement learning via curriculum induction. <i>Advances in Neural Information Processing Systems</i> , 33:12151–12162, 2020.
800 801	Gido M. van de Ven, Hava T. Siegelmann, and Andreas Savas Tolias. Brain-inspired replay for continual learning with artificial neural networks. <i>Nature Communications</i> , 11, 2020.
802 803 804	Akifumi Wachi, Xun Shen, and Yanan Sui. A survey of constraint formulations in safe reinforcement learning, 2024.
805 806 807 808 809	Bokun Wang and Tianbao Yang. Finite-sum coupled compositional stochastic optimization: Theory and applications. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), <i>International Conference on Machine Learning, ICML 2022, 17- 23 July 2022, Baltimore, Maryland, USA</i> , volume 162 of <i>Proceedings of Machine Learning</i> <i>Research</i> , pp. 23292–23317. PMLR, 2022. URL https://proceedings.mlr.press/ v162/wang22ak.html.

010	Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Cheijan Xu,
811	Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of
812	trustworthiness in gpt models. In NeurIPS, 2023a.

- Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: theory, method and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan
  Cheng, Kangwei Liu, Guozhou Zheng, et al. Easyedit: An easy-to-use knowledge editing framework for large language models. *arXiv preprint arXiv:2308.07269*, 2023b.
- Yixuan Wang, Simon Sinong Zhan, Ruochen Jiao, Zhilu Wang, Wanxin Jin, Zhuoran Yang, Zhaoran Wang, Chao Huang, and Qi Zhu. Enforcing hard constraints with soft barriers: Safe reinforcement learning in unknown stochastic environments. In *International Conference on Machine Learning*, pp. 36593–36604. PMLR, 2023c.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
   Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
   *arXiv preprint arXiv:2206.07682*, 2022.
- Yue Xie and Stephen J Wright. Complexity of proximal augmented lagrangian for nonconvex optimization with nonlinear equality constraints. *Journal of Scientific Computing*, 86:1–30, 2021.
- Yangyang Xu. Iteration complexity of inexact augmented lagrangian methods for constrained convex
   programming. *Mathematical Programming*, 185:199–244, 2021.
- Zhuoning Yuan, Yuexin Wu, Zi-Hao Qiu, Xianzhi Du, Lijun Zhang, Denny Zhou, and Tianbao
   Yang. Provable stochastic optimization for global contrastive learning: Small batch does not harm
   performance. In *International Conference on Machine Learning*, pp. 25760–25782. PMLR, 2022.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, pp. 3987–3995. JMLR.org, 2017.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10
  million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- Ba-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards memory-efficient class-incremental learning. *arXiv preprint arXiv:2205.13218*, 2022.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,
   Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A
   top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
- 849 850

813

- 851 852
- 853
- 854 855
- 856

857

- 858 859
- 860
- 861

862

# A More Experimental Details and Results

866 A.1 EXPERIMENTAL DETAILS.

All experiments in our paper are run on two High Performance Research Computing platforms. One
 contains 117 GPU nodes, each with two A100 40GB GPUs. Another contains 100 GPU nodes, each
 with four A40 48GB GPUs.

870 871 A.1.1 DATASET.

We choose the large-scale diverse driving image dataset, namely BDD100K (Seita, 2018), for part of our experiments. This dataset involves six weather conditions, i.e., *clear, overcast, snowy, rainy, partly cloudy, foggy*, and six scene types, i.e., *highway, residential area, city street, parking lot, gas station, tunnel*. Since the labels of the official testing dataset are not released, we utilize the official validation set for testing and partition the training dataset into training and validation sets using an 80%/20% ratio.

Moreover, we experiment on a scene recognition dataset, Places365 (Zhou et al., 2017), to verify the effectiveness of the proposed method in handling a large number of constraints. We utilize the standard version of the dataset (i.e., Places365-Standard), with 1.8 million training and 36500 validation images from 365 scene classes. The number of examples for each class varies between 3,068 and 5,000 in the training set. We merge the training dataset and validation dataset and randomly split the whole set into training set, validation set and test set with an 60%/20%/20% ratio.

Table 1: Datasets Statistics for BDD100K Dataset

Training	Validation	Testing
29865	7479	5346
4445	1104	769
4119	951	738
3992	959	738
7043	1727	1239
57	43	43
Training	Validation	Testing
13952	3427	2499
6458	1616	1253
34862	8654	6112
297	80	49
62	47	47
	Training           29865           4445           4119           3992           7043           57           Training           13952           6458           34862           297           62	Training         Validation           29865         7479           4445         1104           4119         951           3992         959           7043         1727           57         43           Training         Validation           13952         3427           6458         1616           34862         8654           297         80           62         47

899 900 901

884

885

A.1.2 EXPERIMENTAL SETTINGS.

We employ the CLIP ViT-B/16 (Radford et al., 2021b) as the backbone network in all our experiments.

For BDD100K dataset, we obtain a base model by fine-tuning the pretrained CILP model, following 904 the method proposed in Yuan et al. (2022), on the BDD100K training dataset without foggy and 905 tunnel data. Subsequently, we undertake secondary development to improve the performance of 906 a target class separately. We consider three settings with *foggy*, *overcast* and *tunnel* as the target 907 class. For targeting *foggy*, we consider other weather conditions as protected tasks, for targeting 908 *overcast* we consider other weather conditions except for foggy as protected tasks due to that there 909 is a lack of foggy data in BDD100k for defining a significant constraint. For the same reason, we 910 consider other scence types except gas station as protected tasks for targeting *tunnel*. The image-text 911 pairs for the objective function are from the training set of BDD100K and the external LAION400M 912 (Schuhmann et al., 2021) dataset. Specifically, for each target class, we use a query prompt (detailed 913 in Appendix A.2) to search for target-related image-text pairs in LAION400M to augment the set 914  $\mathcal{D}$ . Additionally, we randomly sample a set of image-text pairs from LAION400M that is 10 times 915 larger than target-related pairs as negative data for contrasting. The data of protected tasks used for developmental safety constraints are sampled from the BDD100K training set with varying sizes. 916 Statistics for BDD100K in our experiments are shown in Appendix Table 1. The text templates used 917 for BDD100K dataset are "the weather is [Weather]" and "the scene is a [Scene]".

For Places365 dataset, we directly utilize the pretrained CLIP model released by Radford et al. (2021b) as the base mode. Then we conduct continual development to improve the performance of *dressing room* class, which has the fewest samples in the dataset, and consider all the other 364 classes as protected tasks. Similar to the setting for BDD100K dataset, we also use a query prompt (detailed in Appendix A.2) to search for target-related image-text pairs in LAION400M to augment the set  $\mathcal{D}$ . The data of protected tasks used for developmental safety constraints are sampled from the Places365 training set. The text templates used for Places365 dataset are "*the scene is a(n) [Scene]*".

925 A.1.3 HYPERPARAMETER TUNING. 926

For all methods in our experiments, we tune the learning rate in {1e-5, 1e-6} with Cosine scheduler and AdamW optimizer, using a weight decay of 0.1.

929 For BDD100K dataset, we set temperature  $\tau_0$  as 0.05. We run each method for a total of 40 epochs with a batch size of 256 and 600 iterations per epoch, except for GEM whose total epochs are tuned 930 in  $\{1,2,5\}$  with a batch size of 64 since more iterations lead to exacerbated catastrophic forgetting 931 problems as shown in their paper. For our method, we tune  $\beta$  in {100, 200, 400},  $\gamma_2$  in {0.4, 0.6, 932 0.8} and set r = 32,  $|\mathcal{B}_c| = m$ ,  $|\mathcal{B}_k| = 10$ . We set  $\gamma_1$  to 0.8,  $\tau$  to 0.05 in FLYP, WCCL, RM, and our 933 method. For WCCL, we vary the weight parameter  $\alpha$  in {0.5,0.9,0.99}. For GEM, we tune their 934 small constant  $\gamma$  in {0.5, 1.0}. For Co<sup>2</sup>L, we tune their  $\tau$  in {0.05, 0.1},  $\kappa$  in {0.1, 0.2},  $\kappa^*$  in {0.01, 935 0.1,  $\lambda$  in {0.1, 1, 10}. For RM, we tune regularization weight  $\alpha$  in {0.1, 1, 10}. In hyper-parameters 936 selection for all methods, we prioritize larger **retention ratio** first and consider larger  $\Delta Acc$  (Target) 937 if there is a tie in terms of safety ratio, as we look for models that maximize  $\Delta Acc$  (Target) while 938 satisfying DevSafety  $\geq 0$ .

For Places365 dataset, the temperature  $\tau_0$  is set as 0.01. Since there are as many as 364 constraints, we set  $|\mathcal{B}_c| = 240$ ,  $|\mathcal{B}_k| = 2$ . We tune  $\beta$  in {600, 1000, 4000} for our method and regularization weight  $\alpha$  in {1, 10, 100, 1000, 10000} for RM. We run each method five times for a total of 40 epochs with 1400 iterations per epoch, with a batch size of 64.

#### 944 A.1.4 DETAILS ABOUT BASELINES

FLYP. In the original FLYP paper (Goyal et al., 2023), the author presents extensive experiments demonstrating the superiority of employing the contrastive loss used during pre-training instead of the typical cross-entropy for finetuning image-text models for zero-shot vision classification. As the local contrastive loss, defined over the mini-batch samples, utilized in their paper requires a very large mini-batch size to converge, we follow Yuan et al. (2022) to employ a global constrastive loss (GCL) as indicated in Eqn. 10 to address this issue:

$$\min_{\mathbf{w}} \quad \frac{1}{n_{all}} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}_{all}} L_{ctr}(\mathbf{w}; \mathbf{x}_i, \mathbf{t}_i, \mathcal{D}_{all}, \mathcal{D}_{all})$$
(10)

w  $n_{all} \succeq (\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}_{all}$  and  $(\mathbf{y}, \mathbf{y}, \mathbf{y}, \mathbf{y}, \mathbf{u}, \mathbf{t}_{all})$  and  $(\mathbf{y}, \mathbf{y}, \mathbf{y}, \mathbf{u}, \mathbf{t}_{all})$ where  $\mathcal{D}_{all} = \mathcal{D} \cup \mathcal{D}_- \cup \mathcal{D}_1 \cup \cdots \cup \mathcal{D}_m, n_{all} = n_o + 10 * n_o + n_1 + \cdots + n_m, \mathcal{D}_-$  is the negative data collected form LAION400M as discussed in AppendixA.2. All available data, including those used in our objective and constraints, are utilized for fine-tuning. The simple text prompts for the labeled BDD100k dataset are the same as those used for our method, i.e., "the weather is [Weather]" and "the scene is a [Scene]".

WCCL. Weighted Combination of Contrastive Losses(WCCL) is a straightforward baseline that utilizes a weight to combine GCL losses on protected tasks and the target task to balance protected tasks and the target task and achieve model developmental safety. Specifically, the objective can be formulated as:

$$\min_{\mathbf{w}} \quad \alpha \left( \frac{1}{m} \sum_{k=1}^{m} \frac{1}{n_k} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}_k} L_{\text{ctr}}(\mathbf{w}; \mathbf{x}_i, \mathbf{t}_i, \mathcal{T}_{ik}^-, \mathcal{I}_{ik}^-) \right) \\
+ (1 - \alpha) \left( \frac{1}{n_o} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}_o} L_{\text{ctr}}(\mathbf{w}; \mathbf{x}_i, \mathbf{t}_i, \mathcal{T}_{io}^-, \mathcal{I}_{io}^-) \right)$$
(11)

963 964 965

962

951

where  $\mathcal{T}_{ik}^{-} = \{\mathbf{t}_j : (\mathbf{x}_j, \mathbf{t}_j) \in \mathcal{D}_{all} \setminus \mathcal{D}_k\} \cup \{\mathbf{t}_i\}, \mathcal{I}_{ik}^{-} = \{\mathbf{x}_j : (\mathbf{x}_j, \mathbf{t}_j) \in \mathcal{D}_{all} \setminus \mathcal{D}_k\} \cup \{\mathbf{x}_i\}, \mathcal{D}_{all} \setminus \mathcal{D}_k$  denotes all training samples excluding samples from  $\mathcal{D}_k$ . Similarly,  $\mathcal{T}_{io}^{-} = \{\mathbf{t}_j : (\mathbf{x}_j, \mathbf{t}_j) \in \mathcal{D}_{all} \setminus \mathcal{D}_o\} \cup \{\mathbf{t}_i\}, \mathcal{I}_{io}^{-} = \{\mathbf{x}_j : (\mathbf{x}_j, \mathbf{t}_j) \in \mathcal{D}_{all} \setminus \mathcal{D}_o\} \cup \{\mathbf{t}_i\}, \mathcal{I}_{io}^{-} = \{\mathbf{x}_j : (\mathbf{x}_j, \mathbf{t}_j) \in \mathcal{D}_{all} \setminus \mathcal{D}_o\} \cup \{\mathbf{x}_i\}.$  Consistent with other methods, the simple text prompts for this baseline are also "the weather is [Weather]" and "the scene is a [Scene]".

972 GEM. GEM (Lopez-Paz & Ranzato, 2017b) is a strong continual learning baseline which motivated 973 by a similar idea, utilizing data of previous tasks for constraints. But it doesn't solve the constrained 974 optimization problem directly but project gradients to reduce the increase in the loss of previous 975 tasks. For GEM, we start from pretrained image encoder of the same CLIP model and initialize the linear classification heads  $W \in \mathbb{R}^{d \times (m+1)}$  with the representations outputted by the text encoder 976 with input "the weather is [Weather]" or "the scene is a [Scene]". For each task k, cross entropy loss 977 is employed  $\mathcal{L}_k(\mathbf{w}, W, \mathcal{D}_k) = \frac{1}{n_k} \sum_{(\mathbf{x}_i, y_i) \sim \mathcal{D}_k} -\log \frac{\exp(W_k^\top E_1(\mathbf{w}, \mathbf{x}_i)/\tau_0)}{\sum_{\ell=1}^{m+1} \exp(W_l^\top E_1(\mathbf{w}, \mathbf{x}_i)/\tau_0)}$ , where  $\tau_0 > 0$  is a 978 979 temperature parameter,  $W_k$ ,  $W_l$  denoted the  $k_{th}$ ,  $l_{th}$  column vector of W respectively, and  $E_1(\mathbf{w}, \mathbf{x}_i)$ 980 is the normalized image representation of  $\mathbf{x}_i$ . For consistency,  $\tau_0$  is fixed to 0.05 as the one used 981 in our method. In each iteration, 10 examples are drawn from each protected task to calculate the 982 corresponding loss gradient vector for each task. 983

RM. In continual learning literature, adding explicit regularization terms is a widely used approach to balance old and new tasks, exploiting a frozen copy of previously-learned model to help prevent catastrophic forgetting (Rebuffi et al., 2017; Castro et al., 2018). Similarly, the Regularization Method(RM) baseline incorporates the constraints from Eqn. (4) as a regularization term, adding it to the objective function with an associated regularization weight:

991 992 993

994

 $\min_{\mathbf{w}} \frac{1}{n_o} \sum_{(\mathbf{x}_i, \mathbf{t}_i) \in \mathcal{D}_o} L_{\text{ctr}}(\mathbf{w}; \mathbf{x}_i, \mathbf{t}_i, \mathcal{T}_{io}^-, \mathcal{I}_{io}^-) + \alpha \left( \frac{1}{m} \sum_{k=1}^m \frac{1}{n_k} \sum_{(\mathbf{x}, y) \in \mathcal{D}_k} \ell_{ce}(\mathbf{w}, \mathbf{x}, y) \right)$ (12)

#### A.2 RETRIEVING EXTERNAL DATA FROM LIAON400M

To improve the performance of the CLIP model on the target task, we retrieve related image-text 995 pairs from an external dataset. Specifically, for each target task, we retrieve task-related image-text 996 pairs from Laion400M (Schuhmann et al., 2021) to improve target performance, by going through 997 the dataset and retrieving the image-text pairs with text containing the specific target task names, 998 e.g., 'foggy', 'overcast', 'tunnel', 'dressing room'. Similar approaches have been used in (Liang 999 et al., 2024; Mitchell et al., 2018; Chen et al., 2013), where Liang et al. (2024) used this approach 1000 to improve the detection of rare or unseen categories in object detection for autonomous driving 1001 systems. However, their study is different from ours in the sense that they do not provide guaratnee 1002 on the model developmental safety.

1003 Moreover, we refine the retrieved datasets. Let's take the task 'tunnel' as an example. For task 1004 'tunnel', the retrieved data contained excessive noise, including numerous image-text pairs unrelated 1005 to tunnels, but contained 'tunnel' in the text. Therefore, we employed the GPT-40 API to filter the retrieved data with prompt "Determine whether the following caption mentions a tunnel or related context. First provide reasoning for your answer, and then respond with 'True' if it mentions a tunnel, 1008 or 'False' if it does not.", thereby decreasing the noise of our retrieved data. The statistics of obtained 1009 task-related image-text pairs are presented in the Table 2. Additionally, for each target class, we randomly sample a set of image-text pairs from LAION400M that is 10 times larger than the positive 1010 set as negative data for contrasting. 1011

1012

1013 1014

1015 1016 1017

Table 2: Statistics of Data Collected from LIAON400M

Task	Foggy	Overcast	Tunnel	Dressing room
Size	11415	4134	23484	6786

1018 1019

#### 1020 1021 A.3 VISUALIZATION OF MODELS' LEARNING CURVES

1022 Along with the learning trajectory in the main paper, we present the training and validation curves 1023 in Fig. 5 to further illustrate the learning process of the algorithm. From the figure, we can see 1024 that the DevSafety(acc) fluctuates along the safety line while  $\Delta Acc(Target)$  continues to increase, 1025 which shows the model is striving to improve the model's performance while satisfying the safety 1026 requirements.



Figure 5: Models' Training and Validation Curves

Table 3: Detailed performance comparison between our method and baseline RM on targeting *Foggy* with 4k samples for each protected task. Bold numbers highlight the performance decrease over the base model.

	Clear	Overcast	Protected Tasks Snowy	Rainy	Partly cloudy	Target Task Foggy	Average
Base	0.8938	0.7014	0.7503	0.7195	0.6734	0.3953	0.6889
Ours	+0.0115(0.0054)	+0.0831(0.0228)	+0.0120(0.0079)	+0.0230(0.0081)	+0.1047(0.0168)	0.0326(0.0316)	+0.0430(0.0027)
$\begin{array}{c c} \operatorname{RM} \alpha = 0.1 \\ \operatorname{RM} \alpha = 1 \\ \operatorname{RM} \alpha = 10 \end{array}$	-0.0189(0.0039) -0.0129(0.0055) -0.0106(0.0085)	+0.0667(0.0392) +0.0910(0.0102) +0.1131(0.0068)	+0.0328(0.0113) +0.0666(0.0139) +0.0656(0.0302)	+0.0081(0.0074) +0.0217(0.0215) +0.0163(0.0182)	+0.1253(0.0227) +0.1168(0.0112) +0.0830(0.0201)	+0.0559(0.0617) -0.0604(0.0634) -0.1674(0.0174)	+0.0450(0.0071) +0.0372(0.0114) +0.0167(0.0050)

1050 A.4 DEFICIENCY OF WEIGHTING METHODS

1042

1043

1044

1051 As observed in Figure 3, the naive weighting approach RM fail to achieve model developmental 1052 safety, even though they tradeoff the performance on the target task and protected tasks with weight 1053 parameter  $\alpha$ . To have a close look at why this happens, we show the detailed performance RM when 1054 targeting foggy with 4k samples for each protected task in Table 3. We find that, with a uniform 1055 weight for all the protected tasks, the method might preserve previous performance on some of the 1056 protected tasks but fail to achieve MDS for all the protected tasks, even with a very high  $\alpha$ . Moreover, 1057 with the weight  $\alpha$  getting larger, the performance on the target task drops dramatically while the 1058 decrease gap goes smaller, e.g., Clear tasks for RM. In contrast, our proposed method is able to preserve all the protected tasks' performance and improve the target task, as the mechanism of our 1059 1060 algorithm is very different from using the uniform weight. In our method, weights for constraints depend on the loss of those tasks, i.e., the larger the violation, the larger the weight. As shown in 1061 Figure 4, the weight for each protected task is adaptively adjusted during learning and once one 1062 protected task constraint is satisfied, it will not be penalized (weight becomes zero). This mechanism 1063 plays a big role in enabling the model to find feasible solutions to ensure zero-forgetting on all the 1064 protected tasks.

To further demonstrate the deficiency of the weighting method, we compare RM with our method on the Place365 dataset, targeting *Dressing room* class and protecting the other 364 tasks in Figure 6. With  $\alpha = 1, 10, 100, 1000, 10000$ , RM causes performance drops in 50, 35, 33, 32, and 35 classes, respectively. Although larger weights reduce the number of classes where performance drops, RM still cannot ensure MDS for all protected tasks. In contrast, we can see that even with hundreds of protected tasks, our method is still effective in preserving their performance whiling improving the target task.

#### 1073 A.5 DETAILED PERFORMANCE COMPARISON WITH BASELINES

In this part, we present a detailed performance comparison with baselines. Specifically, we include the DevSafety(acc) numbers for each method in Table 4, 5, 6, which directly show the largest decrease over all the protected tasks. We can see that baselines usually lead to 3-10 percent decrease when targeting Tunnel, 1.5-7 percent decrease when targeting Foggy, 3-30 percent decrease when targeting Overcast. In contrast, our method demonstrates a smaller performance drop when there is insufficient data for constraints and ensures zero forgetting on the protected task when sufficient constraint data is available.

#### Table 4: Detailed Performance Comparison on Targeting Tunnel

1085	Method	Measures	100	1k	2k	4k
1086	Base	Rentention Ratio//DevSafety(acc) Target Tunnel	100%//0.00(0.0000) 0.1064(0.0000)	100%//0.00(0.0000) 0.1064(0.0000)	100%//0.00(0.0000) 0.1064(0.0000)	100%//0.00(0.0000) 0.1064(0.0000)
1088	FLYP	Rentention Ratio//DevSafety(acc) Target Tunnel	0.00%//-0.0398(0.0067 0.9361(0.0330)	0.00%//-0.0660(0.0126) 0.9702(0.0318)	0.00%//-0.0647(0.0123) 0.9915(0.0170)	0.00%//-0.0774(0.0069) 0.9659(0.0170)
1089	WCCL	Rentention Ratio//DevSafety(acc) Target Tunnel	0.00%//-0.0836(0.0164) 0.9957(0.0085)	0.00%//-0.0756(0.0090) 0.6000(0.1002)	0.00%//-0.0673(0.0103) 0.6553(0.0282)	0.00%//-0.0893(0.0089) 0.6383(0.0485)
1090	GEM	Rentention Ratio//DevSafety(acc) Target Tunnel	0.00%//-0.1019(0.0267) 0.8255(0.1214)	0.00%//-0.1034(0.0153) 0.5915(0.2020)	0.00%//-0.1301(0.0169) 0.6085(0.0768)	0.00%//-0.0873(0.0231) 0.3915(0.1819)
1092	RM	Rentention Ratio//DevSafety(acc) Target Tunnel	0.00%//-0.1021(0.0022) 0.9574(0.0233)	0.00%//-0.0969(0.0036) 0.8894(0.0340)	0.00%//-0.0955(0.0057) 0.8808(0.0170)	0.00%//-0.0897(0.0068) 0.8681(0.0085)
1093	Ours	Rentention Ratio//DevSafety(acc) Target Tunnel	40.00%//-0.0050(0.0076) 0.9362(0.0699)	60.00%//-0.0001(0.0043) 0.8723(0.0233)	100.00%//0.0105(0.0053) 0.9106(0.0159)	100.00%//0.0186(0.0058) 0.8723(0.0233)
1034						

#### Table 5: Detailed Performance Comparison on Targeting Foggy

Method	Measures	100	1k	2k	4k
Base	Rentention Ratio//DevSafety(acc)	100%//0.00(0.0000)	100%//0.00(0.0000)	100%//0.00(0.0000)	100%//0.00(0.0000)
	Target Foggy	0.3953(0.0000)	0.3953(0.0000)	0.3953(0.0000)	0.3953(0.0000)
FLYP	Rentention Ratio//DevSafety(acc)	0.00%//-0.0590(0.0140)	20.00%//-0.0281(0.0167)	0.00%//-0.0254(0.0101)	0.00%//-0.0201(0.0105)
	Target Foggy	0.5721(0.0315)	0.5209(0.0581)	0.5302(0.0228)	0.4977(0.0186)
WCCL	Rentention Ratio//DevSafety(acc)	0.00%//-0.0504(0.0123)	0.00%//-0.0259(0.0080)	20.00%//-0.0141(0.0111)	0.00%//-0.0132(0.0076)
	Target Foggy	0.3395(0.0865)	0.2186(0.0186)	0.2093(0.0208)	0.2000(0.0114)
GEM	Rentention Ratio//DevSafety(acc)	0.00%//-0.0695(0.0099)	0.00%//-0.0339(0.0053)	0.00%//-0.0424(0.0060)	0.00%//-0.0424(0.0060)
	Target Foggy	0.3349(0.0865)	0.2837(0.0271)	0.2558(0.0000)	0.2558(0.0000)
RM	Rentention Ratio//DevSafety(acc)	0.00%//-0.0418(0.0062)	0.00%//-0.0173(0.0054)	0.00%//-0.0159(0.0034)	20.00%//-0.0124(0.0091)
	Target Foggy	0.5674(0.0378)	0.5023(0.0186)	0.4419(0.0658)	0.2279(0.0174)
Ours	Rentention Ratio//DevSafety(acc)	0.00%//-0.0241(0.0082)	60.00%//-0.0009(0.0044)	100.00%//0.0044(0.0033)	100.00%//0.0061(0.0047)
	Target Foggy	0.5721(0.0406)	0.4930(0.0174)	0.4326(0.0186)	0.4279(0.0316)

#### Table 6: Detailed Performance Comparison on Targeting Overcast

			•		
Method	Measures	100	1k	2k	4k
Base	Rentention Ratio//DevSafety(acc)	100%//0.00(0.0000)	100%//0.00(0.0000)	100%//0.00(0.0000)	100%//0.00(0.0000
	Target Overcast	0.7361(0.0000)	0.7361(0.0000)	0.7361(0.0000)	0.7361(0.0000)
FLYP	Rentention Ratio//DevSafety(acc)	0.00%//-0.0749(0.0049)	0.00%//-0.0449(0.0140)	0.00%//-0.0434(0.0095)	0.00%//-0.0314(0.011
	Target Overcast	0.9143(0.0111)	0.8559(0.0241)	0.8412(0.0294)	0.8247(0.0255)
WCCL	Rentention Ratio//DevSafety(acc)	0.00%//-0.1192(0.0294)	0.00%//-0.0716(0.0053)	0.00%//-0.0424(0.0091)	0.00%//-0.0414(0.010
	Target Overcast	0.9315(0.0112)	0.9296(0.0092)	0.9207(0.0022)	0.9172(0.0064)
GEM	Rentention Ratio//DevSafety(acc)	0.00%//-0.0677(0.0042)	0.00%//-0.0711(0.0050)	0.00%//-0.0807(0.0128)	0.00%//-0.0634(0.004
	Target Overcast	0.9282(0.0051)	0.9233(0.0037)	0.9149(0.0088)	0.9165(0.0049)
RM	Rentention Ratio//DevSafety(acc)	0.00%//-0.2932(0.0365)	0.00%//-0.3016(0.0228)	0.00%//-0.2444(0.0120)	0.00%//-0.2634(0.010
	Target Overcast	0.9787(0.0050)	0.9730(0.0028)	0.9588(0.0041)	0.9647(0.0023)
Ours	SafetyRatio//DevSafety(acc)	0.00%//-0.0655(0.0249)	20.00%//-0.0043(0.0037)	60.00%//0.0012(0.0029)	100.00%//0.0046(0.00
	Target Overcast	0.8789(0.0464)	0.7827(0.0225)	0.7562(0.0167)	0.7525(0.0366)



Figure 6: Performance comparison between our method and baseline RM when targeting Dresssing 1157 Room on Places365 Dataset, with 2k samples per constraint. Red line denotes base model's perfor-1158 mance, green diamonds denote the target class. RM baseline shown is with weight  $\alpha = 1$  (Top Left), 1159 weight  $\alpha = 10$  (Top Right), weight  $\alpha = 100$  (Middle Left), weight  $\alpha = 100$  (Middle Right), weight 1160  $\alpha = 10000$ (Bottom). 1161

1164 Table 7: Effect of the Number of Samples for Constraints. Numbers in parentheses denote standard 1165 deviation.

1166							
1167	Target	Measures	Base model	100	1k	2k	4k
1168 1169	Tunnel	DevSafety(acc) Rentention Ratio Target Acc	0.00(0.0000) 100.00% 0.1064(0.0000)	-0.0050(0.0076) 40.00% 0.9362(0.0699)	-0.0001(0.0043) 60.00% 0.8723(0.0233)	0.0105(0.0053) 100.00% 0.9106(0.0159)	0.0186(0.0058) 100.00% 0.8723(0.0233)
1170 1171 1172	Foggy	DevSafety(acc) Rentention Ratio Target Acc	0.00(0.0000) 100.00% 0.3953(0.0000)	-0.0241(0.0082) 0.00% 0.5721(0.0406)	-0.0009(0.0044) 60.00% 0.4930(0.0174)	0.0044(0.0033) 100.00% 0.4326(0.0186)	0.0061(0.0047) 100.00% 0.4279(0.0316)
1173 1174	Overcast	DevSafety(acc) Rentention Ratio Target Acc	0.00(0.0000) 100.00% 0.7361(0.0000)	-0.0655(0.0249) 0.00% 0.8789(0.0464)	-0.0043(0.0037) 20.00% 0.7827(0.0225)	0.0012(0.0029) 60.00% 0.7562(0.0167)	0.0046(0.0016) 100.00% 0.7525(0.0366)

1175 1176

1162 1163

> **DETAILED ABLATION STUDIES** A.6

#### 1177 A.6.1 THE EFFECT OF DIFFERENT NUMBER OF SAMPLES USED FOR CONSTRAINTS.

1178 In our framework, we propose to formulate the model developmental safety as data-dependent 1179 constraints. As discussed in Section 4.2, since we only have access to a finite set of empirical samples, 1180 there exists generalization errors between the safety constraints in Eqn. 3. In this part, we examine 1181 the impact of varying the amount of data used for constraints. Specifically, we conduct experiments 1182 with different numbers of data for constraints, i,e., 100, 1k, 2k, 4k from each task. The results are 1183 summarized in Tab. 7. From the table, we can observe that DevSafety(acc) and Safety Ratio increase 1184 when using more data for constraints. This is consistent with results obtained in Lemma 1 which 1185 shows a larger number of data for constraints leads to a higher probability of being safe. One the other hand, we found that, with the likelihood of developmental safety increasing, the improvement 1186 of the targeted task decreases, which indicates there may still exist a tradeoff between enhancing the 1187 targeted task's performance and satisfying the developmental safety requirements.



# 1188Table 8: The Effect of External Image-text Pairs from LIAON400M. Numbers in parentheses denote1189std.



# A.6.2 IMPORTANCE OF THE EXTERNAL DATA FROM LAION400M

1205 We conduct experiments on targeting *foggy* to investigate the benefits of the external data retrieved from LAION400M dataset. In detail, we vary the number of retrieved target-related image-text pairs 1206 utilized in the objective function, i.e.,  $\{0, 2k, 5k, 11k\}$ , with 1k samples from each protected task as 1207 constraints. From Tab. 8, we can see that, with only 57 foggy samples from BDD100k dataset (i.e., 0 1208 samples from the external data), the model does not improve the target accuracy at all. However, with 1209 more and more retrieved image-text pairs utilized to augment the dataset  $\mathcal{D}$ , the improvement on the 1210 targeted task appears and becomes significant, showing the advantages of incorporating the retrieved 1211 target-related image-text pairs for boosting target task accuracy. Regarding safety ratios, we don't 1212 observe a clear correlation between the amount of retrieved data and the safety ratios. 1213

#### A.6.3 IMPORTANCE OF TASK-DEPENDENT HEADS

1215 As introduced in Section 5.2, to reduce the total complexity of our algorithm, we propose task-1216 dependent heads to increase the parameter  $\delta$  in Assumption 4, avoiding getting trapped at a flat location where  $\mathbf{w}^t$  is infeasible but  $\nabla H(\mathbf{w}^t) = 0$ . To verify the effectiveness of the design, we 1217 1218 experiment on targeting *overcast* and *foggy* tasks with varying numbers of data for constraints. The results are presented in Figure 7. The results show that models equipped with task-dependent heads 1219 almost consistently exhibit both higher safety ratio and higher accuracy on the target task. Besides, 1220 without task-dependent heads, models may have trouble achieving 100% developmental safety, 1221 such as targeting *Overcast*, demonstrating the importance of task-dependent heads for promoting 1222 developmental safety. 1223

- 1224 A.6.4 VERIFICATION OF LEMMA 2
- 1225

To verify the theoretical result in Lemma 2, we empirically calculate  $\nabla \hat{\mathbf{h}}(\hat{\mathbf{w}})$  and  $\nabla \mathbf{h}(\mathbf{w})$  with CLIP models. Specifically, for targeting *overcast*, we compute the minimal singular values of  $\nabla \hat{\mathbf{h}}(\hat{\mathbf{w}})$  and  $\nabla \mathbf{h}(\mathbf{w})$  on the base model and two trained models, with 1k samples for each protected task. The initial value of  $U_k$  is set to zero so  $U_k V_k^{\top} = \mathbf{0}$ . From the results presented in Table 9, we can observe that, on the initial model, the minimal singular value of  $\nabla \hat{\mathbf{h}}(\hat{\mathbf{w}})$  is slightly larger than that of  $\nabla \mathbf{h}(\mathbf{w})$ and the gap become much significant after training, which is consistent with the theoretical result in Lemma 2 and also provides some insight on the empirical results in Figure 7.

**1233** A.6.5 CONSTANT  $\beta$  vs Increasing  $\beta$ 

1234 In theory, an increasing penalty parameter  $\beta$  may help reduce the complexity of constrained problems 1235 as shown in Alacaoglu & Wright (2024), but in our empirical experiments, we find that using a 1236 constant  $\beta$  is generally behave better than using an increasing  $\beta$ . As shown in Fig. 8 for target task 1237 foggy, models with a constant  $\beta$  are able to achieve 100% safety ratio with 2k or 4k sampler per constraint. On the contrary, models using a cosine increasing  $\beta$  obtain both lower safety ratio and lower accuracy on the target task compared to models with constant  $\beta$ . We conjecture that this is 1239 because models with an increasing  $\beta$  might leave the feasible developmental safety region too far 1240 in the initial stages as they have a relatively small penalty weight  $\beta$  at this time. Given the high 1241 non-convexity and complexity of the model space, it becomes increasingly challenging in the later



#### Table 9: Minimal Singular Values of $\nabla \mathbf{h}(\mathbf{w})$ and $\nabla \widehat{\mathbf{h}}(\widehat{\mathbf{w}})$



stages to return to a feasible solution that satisfies developmental safety constraints while significantly
 improving target accuracy.

## 1264 B MORE RELATED WORK

1242

1260

1265 **Continual learning.** This work is closely related to Continual learning (CL), also known as lifelong 1266 learning, yet it exhibits nuanced differences. Continual learning usually refers to learning a sequence 1267 of tasks one by one and accumulating knowledge like human instead of substituting knowledge (Wang 1268 et al., 2024; Qu et al., 2021). There is a vast literature of CL of deep neural networks (DNNs) (Aljundi 1269 et al., 2018; Lopez-Paz & Ranzato, 2017a; Farajtabar et al., 2019; Lee et al., 2017; Guo et al., 2020; 1270 Parisi et al., 2018). The core issue in CL is known as catastrophic forgetting (McCloskey & Cohen, 1989), i.e., the learning of the later tasks may **significantly** degrade the performance of the models 1271 learned for the earlier tasks. Different approaches have been investigated to mitigate catastrophic 1272 forgetting, including regularization-based approaches (Kirkpatrick et al., 2017; Zenke et al., 2017; Li 1273 & Hoiem, 2016), expansion-based approaches (Zhou et al., 2022; Li et al., 2019; Rusu et al., 2016; 1274 van de Ven et al., 2020), and memory-based approaches (Rolnick et al., 2019; Shin et al., 2017; 1275 Guo et al., 2020; Lopez-Paz & Ranzato, 2017a; Chaudhry et al., 2019). In the era of large language 1276 models(LLMs), another type of continual learning method, known as knowledge/representation 1277 editing(Liu et al., 2024; Zou et al., 2023; Meng et al., 2022), emerges to efficiently modify the 1278 behavior of LLMs with minimal impact on unrelated inputs (Wang et al., 2023b), such as to update 1279 stale facts, eliminate unintended biases, and reduce undesired hallucinations.

1280 The framework proposed in this work is similar to memory-based approaches in the sense that both 1281 use examples of existing tasks to regulate learning. However, the key difference is that most existing 1282 continual learning focuses on the trade-off between learning plasticity and memory stability and aims 1283 to find a proper balance between performance on previous tasks and new tasks (Wang et al., 2024). 1284 Hence, they do not provide a guarantee for MDS. A recent work (Peng et al., 2023) has proposed an 1285 ideal continual learner that never forgets by assuming that all tasks share the same optimal solution. 1286 However, it is not implementable for deep learning problems. Besides, existing continual learning 1287 studies usually highlight resource efficiency when accumulating knowledge by reducing the number of samples of previous tasks. In contrast, this work tends to utilize more examples to construct 1288 developmental safety constraints for protected tasks to facilitate MDS. 1289

AI Safety. Our notion of model developmental safety should not be confused with AI safety. The
 latter is a field concerned with mitigating risks associated with AI, whose surge in attention stems
 from the growing capabilities of AI systems, particularly large foundation models (Kojima et al.,
 2022; Wei et al., 2022; Bubeck et al., 2023; Radford et al., 2021b). As these models become more
 adept at complex tasks, concerns around potential misuse, bias, and unintended consequences rise
 proportionally. Amodei et al. (2016) presents several practical research problems related to AI safety,
 including avoiding side effects, avoiding reward hacking, scalable oversight, safe exploration, and

1296 robustness to distributional shift. More recently, Wang et al. (2023a) elaborate on eight different 1297 perspectives to evaluate the trustworthiness of LLMs, including toxicity, stereotype bias, adversarial 1298 robustness, out-of-distribution robustness, robustness on adversarial demonstrations, privacy, machine 1299 ethics, and fairness. These AI safety issues arise in the usage of AI models, and they are distinctive 1300 from model developmental safety studied in this work, which arises in the development of AI models. Note that the term "safety" in model developmental safety is to underline that it is important and 1301 must be enforced in practice. Therefore, this work provides another dimension for consideration in 1302 AI safety, i.e., retention of safety. Any safety features of an AI system that have been acquired and 1303 validated should be retained safely in continuous development. 1304

1305 SafeRL. This work is partially related to SafeRL (Safe Reinforcement Learning), which focuses on 1306 developing algorithms and techniques to ensure safety (avoid harmful actions) of RL agents, such as in autonomous driving (Shalev-Shwartz et al., 2016), robotics areas (Pham et al., 2018). Many 1307 studies have been conducted in SafeRL domain. A popular approach in SafeRL is to maximize 1308 the expected cumulative reward subject to specific safety constraints (Wachi et al., 2024), such as 1309 expected cumulative safety constraint (Ding et al., 2021; Bura et al., 2022; Tessler et al., 2018; 1310 Achiam et al., 2017), state constraint (Thomas et al., 2021; Turchetta et al., 2020; Wang et al., 2023c; 1311 Thananjeyan et al., 2021), joint chance constraint (Ono et al., 2015; Pfrommer et al., 2022), etc. 1312 However, as SafeRL heavily relies on the special structure of policy optimization for RL, it is different 1313 from our work that study a generic developmental safety in model development process. Hence, 1314 although sharing the similarity of solving a constrained problem, the algorithms for SafeRL are not 1315 applicable to our problem.

1316 Constrained Learning. Constrained learning has attracted significant attention in the literature. 1317 Traditional works for constrained optimization include three primary categories: 1) primal methods 1318 which do not involve the Lagrange multipliers, e.g., cooperative subgradient methods (Lan & Zhou, 1319 2016; Polyak & Tret'yakov, 1973) and level-set methods (Aravkin et al., 2019; Lin et al., 2018a;b); 2) primal-dual methods which reformulate constrained optimization problems as saddle point prob-1321 lems (Hamedani & Aybat, 2021; Nemirovski, 2004); 3) penalty-based approaches which incorporate 1322 constraints by adding a penalty term to the objective function (Xu, 2021; Lan & Monteiro, 2013; 1323 2016). However, most of these works are limited to convex objectives or convex constraints. In recent years, due to its increasing importance in modern machine learning problems, such as in applications 1324 concerned with fairness (Cotter et al., 2019; Agarwal et al., 2018), robustness (Robey et al., 2021; 1325 Madry et al., 2017), and safety (Paternain et al., 2019b;a) problems, the research interest has been 1326 directed to developing efficient algorithms for non-convex optimization (non-convex objective and 1327 non-convex constraint) (Boob et al., 2023; Facchinei et al., 2021; Ma et al., 2020; Li et al., 2024; 1328 Chamon et al., 2022; Alacaoglu & Wright, 2024). Among these, Chamon et al. (2022) studies how to 1329 solve constrained learning learning with expected non-convex loss and expected non-convex con-1330 straints by using empirical data to ensure the PAC learnability, and proposed a primal-dual algorithm 1331 to solve constrained optimization problems in the empirical dual domain. However, their algorithm 1332 requires solving the primal problem up to a certain accuracy, which is theoretically not feasible 1333 for general non-convex problems. Boob et al. (2023) introduces a new proximal point method that 1334 transforms a non-convex problem into a sequence of convex problems by adding quadratic terms to both the objective and constraints. For solving non-convex optimization problems with equality 1335 constraints, Alacaoglu & Wright (2024) propose single-loop quadratic penalty and augmented La-1336 grangian algorithms with variance reduction techniques to improve the complexity. Nevertheless, none of these algorithms can be directly applied to our large-scale deep learning problem (4), due to 1338 either prohibitive running cost or failure to handle biased stochastic gradients caused by compositional 1339 structure. 1340

1341 C PROOFS

1343 C.1 PROOF OF LEMMA 1

**Proof** Consider task k. Recall that  $D_k$  contains  $n_k$  data points. According to Theorem 3.2 in Boucheron et al. (2005), we have with probability at least  $1 - \delta/m$ , for all w,

$$|\mathcal{L}_k(\mathbf{w}, \mathfrak{D}_k) - \mathcal{L}_k(\mathbf{w}, \mathcal{D}_k)| \le 2R_{n_k}(\mathcal{H}) + \sqrt{\frac{\ln(2m/\delta)}{2n_k}} \le \frac{2C}{n_k^{\alpha}} + \sqrt{\frac{\ln(2m/\delta)}{2n_k}},$$

1348 1349

1346 1347

where the second inequality is by the assumption on  $R_n(\mathcal{H})$ . Combining the inequalities above with  $\mathbf{w} = \mathbf{w}_{new}$  and  $\mathbf{w} = \mathbf{w}_{old}$ , we have with probability at least  $1 - \delta/m$ 

$$\mathcal{L}_{k}(\mathbf{w}_{new},\mathfrak{D}_{k}) - \mathcal{L}_{k}(\mathbf{w}_{old},\mathfrak{D}_{k}) \leq \mathcal{L}_{k}(\mathbf{w}_{new},\mathcal{D}_{k}) - \mathcal{L}_{k}(\mathbf{w}_{old},\mathcal{D}_{k}) + \frac{4C}{n_{k}^{\alpha}} + 2\sqrt{\frac{\ln(2m/\delta)}{2n_{k}}}$$

Applying the union bound with the events above for k = 1, ..., m leads to the conclusion of this lemma. 

C.2 PROOF OF LEMMA 2 

**Proof** Recall that w has two component u and W. The gradient of  $h_k(\mathbf{w})$  with respect to W and u are denoted by  $\nabla_W h_k(\mathbf{w})$  and  $\nabla_{\mathbf{u}} h_k(\mathbf{w})$ , respectively. Hence, 

 $\nabla h_k(\mathbf{w}) = (\nabla_{\mathbf{u}} h_k(\mathbf{w}), \nabla_W h_k(\mathbf{w}))$ 

for  $k = 1, \ldots, m$ . Similarly, after adding the task-dependent heads,  $\hat{\mathbf{w}}$  has four component  $\mathbf{u}, W, \mathbf{U}$ and V. The gradients  $\nabla_{\mathbf{u}} h_k(\hat{\mathbf{w}})$ ,  $\nabla_W h_k(\hat{\mathbf{w}}) \nabla_U h_k(\hat{\mathbf{w}})$  and  $\nabla_V h_k(\hat{\mathbf{w}})$  are defined correspondingly, and

$$\nabla \hat{h}_k(\hat{\mathbf{w}}) = \left( \nabla_{\mathbf{u}} \hat{h}_k(\hat{\mathbf{w}}), \nabla_W \hat{h}_k(\hat{\mathbf{w}}), \nabla_{\mathbf{U}} \hat{h}_k(\hat{\mathbf{w}}), \nabla_{\mathbf{V}} \hat{h}_k(\hat{\mathbf{w}}) \right).$$

Recall that 

$$\hat{h}_k(\hat{\mathbf{w}}) = h_k(W + U_k V_k^{\top}, \mathbf{u})$$
 for  $k = 1, \dots, m_k$ 

Therefore, 

and

$$\nabla_{\mathbf{u}} \hat{h}_k(\hat{\mathbf{w}}) = \nabla_{\mathbf{u}} h_k(W + U_k V_k^{\top}, \mathbf{u}), \qquad \nabla_W \hat{h}_k(\hat{\mathbf{w}}) = \nabla_W h_k(W + U_k V_k^{\top}, \mathbf{u})$$

$$\nabla_{\mathbf{U}}\hat{h}_k(\hat{\mathbf{w}}) = \left(\mathbf{0}, \dots, \mathbf{0}, \underbrace{\nabla_W h_k(\mathbf{W} + U_k V_k^{\top}, \mathbf{u}) V_k}_{The \ kth \ block}, \mathbf{0}, \dots, \mathbf{0}\right)$$

$$\nabla_{\mathbf{V}} \hat{h}_k(\hat{\mathbf{w}}) = \left(\mathbf{0}, \dots, \mathbf{0}, \underbrace{\nabla_W h_k(\mathbf{W} + U_k V_k^\top, \mathbf{u})^\top U_k}_{The \ kth \ block}, \mathbf{0}, \dots, \mathbf{0}\right)^\top$$

where the sparsity patterns of  $\nabla_{\mathbf{U}} \hat{h}_k(\hat{\mathbf{w}})$  and  $\nabla_{\mathbf{V}} \hat{h}_k(\hat{\mathbf{w}})$  are because  $\hat{h}_k$  does not depend on  $U_j$  and  $V_j$  with  $j \neq k$ . 

1379  
1379  
1380  
Suppose 
$$U_k V_k^{\top} = \mathbf{0}$$
 for all  $k$ . It holds that  $h_k(\mathbf{w}) = \hat{h}_k(\hat{\mathbf{w}})$  and  
 $\nabla h_k(\mathbf{w}) = (\nabla_{\mathbf{u}} h_k(\mathbf{w}), \nabla_W h_k(\mathbf{w})) = \left(\nabla_{\mathbf{u}} \hat{h}_k(\hat{\mathbf{w}}), \nabla_W \hat{h}_k(\hat{\mathbf{w}})\right)$ .  
1382  
1384  
1384  
 $\lambda_{\min} \left( \left[ \nabla \hat{h}_1(\hat{\mathbf{w}}), \dots, \nabla \hat{h}_m(\hat{\mathbf{w}}) \right]^{\top} \left[ \nabla \hat{h}_1(\hat{\mathbf{w}}), \dots, \nabla \hat{h}_m(\hat{\mathbf{w}}) \right] \right)$   
1385  
1386  
1386  
1389  
1389  
1389  
1389  
1390  

$$= \min_{\boldsymbol{\alpha}, s.t, \|\boldsymbol{\alpha}\|=1} \left( \left\| \sum_{k=1}^m \alpha_k \nabla_{\mathbf{u}} \hat{h}_k(\hat{\mathbf{w}}) \right\|^2 + \left\| \sum_{k=1}^m \alpha_k \nabla_W \hat{h}_k(\hat{\mathbf{w}}) \right\|^2 + \left\| \sum_{k=1}^m \alpha_k \nabla_{\mathbf{U}} \hat{h}_k(\hat{\mathbf{w}}) \right\|^2 + \left\| \sum_{k=1}^m \alpha_k \nabla_{\mathbf{V}} \hat{h}_k(\hat{\mathbf{w}}) \right\|^2$$

$$\begin{aligned}
& \lim_{\mathbf{\alpha},s,t,\|\mathbf{\alpha}\|=1} \left( \left\| \sum_{k=1}^{\infty} \alpha_k \nabla h_k(\mathbf{w}) \right\| + \sum_{k=1}^{\infty} \alpha_k^2 \left\| \nabla_W h_k(\mathbf{w}) V_k \right\|_F^2 + \sum_{k=1}^{\infty} \alpha_k^2 \left\| \nabla_W h_k(\mathbf{w})^\top U_k \right\|_F^2 \right) \\
& \lim_{\mathbf{\alpha},s,t,\|\mathbf{\alpha}\|=1} \left( \left[ \nabla h_1(\mathbf{w}), \dots, \nabla h_m(\mathbf{w}) \right]^\top \left[ \nabla h_1(\mathbf{w}), \dots, \nabla h_m(\mathbf{w}) \right] \right) \\
& \lim_{k=1}^{\infty} \left\| \nabla_W h_k(\mathbf{w}) V_k \right\|_F^2 + \min_{k} \left\| \nabla_W h_k(\mathbf{w})^\top U_k \right\|_F^2,
\end{aligned}$$

where the first two equalities are by definitions and the third equality is because  $U_k V_k^{\top} = \mathbf{0}$  for all k.

C.3 PROOF OF THEOREM 1 

In this section, we present the proof of the Theorem 1. Recall that the problem is formulated as  $n_0$ 

1402  
1403 
$$\min_{\mathbf{w}} F(\mathbf{w}, \mathcal{D}) := \frac{1}{n_0} \sum_{i=1}^{n_0} \left( f(g_{1i}(\mathbf{w})) + f(g_{2i}(\mathbf{w})) \right) \quad \text{s.t.} \quad \frac{1}{m} h_k(\mathbf{w}; \mathcal{D}_k) \le 0, \ k = 1, \cdots, m.$$
(13)

w

with  $f(\cdot) = \tau \log(\cdot)$ . With the quadratic penalty method, the problem is converted to 

$$\min_{\mathbf{w}} \Phi(\mathbf{w}) := F(\mathbf{w}, \mathcal{D}) + \underbrace{\frac{1}{m} \sum_{k=1}^{m} \frac{\beta}{2} ([h_k(\mathbf{w}; \mathcal{D}_k)]_+)^2}_{H(\mathbf{w})}.$$
(14)

By Assumptions 1, we can get f is  $L_f$ -Lipschitz continuous and  $L_{\nabla f}$ -smooth with  $L_f = \frac{\tau}{c_a}$  and  $L_{\nabla f} = \frac{\tau}{c_{a}^{2}}$ . By noticing that  $\ell_{ce}$  is a cross entropy loss, we find that  $|h_{k}(\cdot)|$  can be bounded by a constant  $C_h$  with  $C_h = 2$ . Then, we can get  $\Phi(\mathbf{w})$  is  $L_\beta$ -smooth with  $L_\beta := L_F + \beta L_H$  where  $L_F := 2(L_{\nabla g}L_f + L_{\nabla f}L_g^2)$  and  $L_H := L_{\nabla h}C_h + L_h^2$ . We also define  $\tilde{C}_{\nabla g} := \sigma_{\nabla g} + L_g$  and  $\tilde{C}_{\nabla h} := \sigma_{\nabla h} + L_h$ . To facilitate our discussion, we let 

$$v_1^t = (1 - \theta)v_1^{t-1} + \theta G_1^t,$$
  
$$v_2^t = (1 - \theta)v_2^{t-1} + \theta G_2^t,$$

$$v^t = v^t_1 + v^t_2.$$

To prove our main theorem, we need following lemmas. 

**Lemma 3** If  $\theta \leq \frac{1}{3}$ , the gradient variance  $\Delta_1^t := \|v_1^t - \nabla F(\mathbf{w}^t, \mathcal{D})\|^2$  can be bounded as  $2L^2$ 

with  $\Xi_1^{t+1} := \frac{1}{n_0} \| \boldsymbol{u}_1^{t+1} - \boldsymbol{g}_1(\mathbf{w}^{t+1}) \|^2 = \frac{1}{n_0} \sum_{i=1}^{n_0} \| u_{1i}^{t+1} - g_{1i}(\mathbf{w}^{t+1}) \|^2$  and  $\Xi_2^{t+1} := \frac{1}{n_0} \| \boldsymbol{u}_2^{t+1} - \boldsymbol{g}_{2i}(\mathbf{w}^{t+1}) \|^2$  $\boldsymbol{g}_2(\mathbf{w}^{t+1}) \|^2 = \frac{1}{n_0} \sum_{i=1}^{n_0} \| u_{2i}^{t+1} - g_{2i}(\mathbf{w}^{t+1}) \|^2.$ 

Proof

$$\Delta_1^{t+1} = \|v_1^{t+1} - \nabla F(\mathbf{w}^{t+1})\|^2 = \|(1-\theta)v_1^t + \theta G_1^t - \nabla F(\mathbf{w}^{t+1})\|^2$$
$$= \left\| \widehat{(I)} + \widehat{(2)} + \widehat{(3)} + \widehat{(4)} \right\|^2,$$

$$\begin{aligned} & \text{where } (\widehat{I}, (\widehat{Q}), \widehat{3}, (\widehat{q} \text{ are defined as} \\ & \widehat{I} = (1 - \theta)(v_1^t - \nabla F(\mathbf{w}^t)), \quad (\widehat{Q}) = (1 - \theta)(\nabla F(\mathbf{w}^t) - \nabla F(\mathbf{w}^{t+1})), \\ & \text{if } 3 = \frac{\theta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}^{t+1}} \nabla \hat{g}_{1i}(\mathbf{w}^{t+1}) \left( \nabla f(u_{1i}^t) - \nabla f(g_{1i}(\mathbf{w}^{t+1})) \right) + \nabla \hat{g}_{2i}(\mathbf{w}^{t+1}) \left( \nabla f(u_{2i}^t) - \nabla f(g_{2i}(\mathbf{w}^{t+1})) \right), \\ & \text{if } 439 \\ & \text{if } 440 \\ & \widehat{I} = \frac{\theta}{|\mathcal{B}|} \sum_{i \in \mathcal{B}^{t+1}} \nabla \hat{g}_{1i}(\mathbf{w}^{t+1}) \nabla f(g_{1i}(\mathbf{w}^{t+1})) + \nabla \hat{g}_{2i}(\mathbf{w}^{t+1}) \nabla f(g_{2i}(\mathbf{w}^{t+1})) - \nabla F(\mathbf{w}^{t+1}). \\ & \text{if } 441 \\ & \text{if } 442 \\ & \text{if } 442 \\ & \text{if } E_t \left[ \left| (\widehat{I} + (\widehat{Q}), \widehat{\mathcal{A}} \right| \right] = \mathbb{E}_t [\langle (\widehat{Q}), \widehat{\mathcal{A}} \rangle] = 0. \text{ Then, by the Young's inequality, we can get} \\ & \mathbb{E}_t \left[ \left\| (\widehat{I} + (\widehat{Q}), \widehat{\mathcal{A}} \right\|^2 \right] \\ & = \left\| (\widehat{I}) \right\|^2 + \left\| (\widehat{Q}) \right\|^2 + \mathbb{E}_t \left\| (\widehat{\mathcal{A}}) \right\|^2 + \mathbb{E}_t \left\| (\widehat{\mathcal{A}}) \right\|^2 + 2 \left\langle (\widehat{I}, (\widehat{Q}) \right\rangle + 2\mathbb{E}_t [\langle (\widehat{I}), (\widehat{\mathcal{A}}) \rangle] + 2\mathbb{E}_t [\langle (\widehat{\mathcal{Q}}), \widehat{\mathcal{A}} \rangle] \\ & \leq (1 + \theta) \left\| (\widehat{I}) \right\|^2 + 2 \left( 1 + \frac{1}{\theta} \right) \left\| (\widehat{\mathcal{Q}}) \right\|^2 + \frac{2 + 3\theta}{\theta} \mathbb{E}_t \left\| (\widehat{\mathcal{A}}) \right\|^2 + 2\mathbb{E}_t \left\| (\widehat{\mathcal{A}}) \right\|^2 . \\ & \text{We can also get} \end{aligned}$$

$$(1+\theta)\|\widehat{U}\|^{2} = (1+\theta)(1-\theta)^{2}\|v_{1}^{t} - \nabla F(\mathbf{w}^{t})\|^{2} \le (1-\theta)\|v_{1}^{t} - \nabla F(\mathbf{w}^{t})\|^{2}$$
$$2\left(1+\frac{1}{\theta}\right)\left\|\widehat{\mathcal{Q}}\right\|^{2} = 2\left(1+\frac{1}{\theta}\right)(1-\theta)^{2}\|\nabla F(\mathbf{w}^{t}) - \nabla F(\mathbf{w}^{t+1})\|^{2} \le \frac{2L_{F}^{2}}{\theta}\|\mathbf{w}^{t+1} - \mathbf{w}^{t}\|^{2}$$

We first bound the first term  $\frac{(2+3\theta)\theta}{|\mathcal{B}|} \mathbb{E}_{t} \sum_{i=ot+1} \left\| \nabla \hat{g}_{1i}(\mathbf{w}^{t+1}) \right\|^{2} \left\| \nabla f(u_{1i}^{t}) - \nabla f(g_{1i}(\mathbf{w}^{t+1})) \right\|^{2}$  $\leq \frac{(2+3\theta)\theta L_{f}^{2}}{|\mathcal{B}|} \mathbb{E}_{t} \left| \sum_{1=sk+1} \left\| \nabla \hat{g}_{1i}(\mathbf{w}^{t+1}) \right\|^{2} \left\| u_{1i}^{t} - g_{1i}(\mathbf{w}^{t+1}) \right\|^{2} \right|$  $= (2+3\theta)\theta L_f^2 \mathbb{E}_t \left\| \frac{1}{|\mathcal{B}|} \sum_{\mathbf{k} \in \mathbf{M}^{t+1}} \mathbb{E}_t \left[ \left\| \nabla \hat{g}_{1i}(\mathbf{w}^{k+1}) \right\|^2 | i \in \mathcal{B}^{t+1} \right] \left\| u_{1i}^t - g_{1i}(\mathbf{w}^{k+1}) \right\|^2 \right\|$  $\leq (2+3\theta)\theta L_f^2 \tilde{C}_{\nabla g}^2 \mathbb{E}_t \left[ \frac{1}{|\mathcal{B}|} \sum_{1 \leq i \neq t+1} \left\| u_{1i}^t - g_{1i}(\mathbf{w}^{t+1}) \right\|^2 \right]$  $\leq (2+3\theta)\theta L_{f}^{2}\tilde{C}_{\nabla g}^{2}\left((1+\delta)\mathbb{E}_{t}\left|\frac{1}{n_{0}}\sum_{i=1}^{n_{0}}\left\|u_{1i}^{t+1}-g_{1i}(\mathbf{w}^{t+1})\right\|^{2}\right|+(1+1/\delta)\mathbb{E}_{t}\left|\frac{1}{n_{0}}\sum_{i=1}^{n_{0}}\left\|u_{1i}^{t+1}-u_{1i}^{t}\right\|^{2}\right|\right)$  $= (2+3\theta)\theta L_{f}^{2}\tilde{C}_{\nabla g}^{2}\left((1+\delta)\mathbb{E}_{t}\left[\frac{1}{n_{0}}\sum_{i=1}^{n_{0}}\left\|u_{1i}^{t+1}-g_{i}(\mathbf{w}^{t+1})\right\|^{2}\right] + (1+1/\delta)\mathbb{E}_{t}\left[\frac{1}{n_{0}}\sum_{i\in\mathcal{R}^{t+1}}\left\|u_{1i}^{t+1}-u_{1i}^{t}\right\|^{2}\right]\right)$ If  $\theta \leq \frac{1}{3}$  and  $\delta = \frac{3\theta}{2}$ , we have  $(2+3\theta)\theta(1+\delta) \leq 5\theta$  and  $(2+3\theta)\theta(1+1/\delta) \leq 3$ . And similarly, we can get the bound for the second term. Then, by combining them, we can get  $\frac{2+3\theta}{\theta} \mathbb{E}\left[\|\widehat{\mathcal{B}}\|^2\right] \leq 5\theta L_f^2 \tilde{C}_{\nabla g}^2 \mathbb{E}[\Xi_1^{t+1} + \Xi_2^{t+1}] + 3L_f^2 \tilde{C}_{\nabla g}^2 \mathbb{E}\left|\frac{1}{n_0} \sum_{1 \leq n \leq t+1} \left(\left\|u_{1i}^{t+1} - u_{1i}^t\right\|^2 + \left\|u_{2i}^{t+1} - u_{2i}^t\right\|^2\right)\right|.$  $\mathbb{E}_t \left\| \widehat{\mathcal{A}} \right\|^2$  $=\theta^{2}\mathbb{E}_{t}\left[\left\|\frac{1}{|\mathcal{B}|}\sum_{i\in\mathbf{w}^{t+1}}\nabla\hat{g}_{1i}(\mathbf{w}^{k+1})\nabla f(g_{1i}(\mathbf{w}^{k+1})) - \frac{1}{n_{0}}\sum_{i=1}^{n_{0}}\nabla g_{1i}(\mathbf{w}^{t+1})\nabla f(g_{1i}(\mathbf{w}^{t+1}))\right\|^{2}\right]$  $+ \theta^{2} \mathbb{E}_{t} \left\| \left\| \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{D}^{t+1}} \nabla \hat{g}_{2i}(\mathbf{w}^{k+1}) \nabla f(g_{2i}(\mathbf{w}^{t+1})) - \frac{1}{n_{0}} \sum_{i=1}^{n_{0}} \nabla g_{2i}(\mathbf{w}^{t+1}) \nabla f(g_{2i}(\mathbf{w}^{t+1})) \right\| \right\|$  $=\theta^{2}\mathbb{E}_{t}\left[\left\|\frac{1}{|\mathcal{B}|}\sum_{i\in\mathcal{B}^{t+1}}\nabla\hat{g}_{1i}(\mathbf{w}^{t+1})\nabla f(g_{1i}(\mathbf{w}^{t+1})) - \frac{1}{|\mathcal{B}|}\sum_{i\in\mathcal{B}^{t+1}}\nabla g_{1i}(\mathbf{w}^{t+1})\nabla f(g_{1i}(\mathbf{w}^{t+1}))\right\|^{2}\right]$  $+ \theta^{2} \mathbb{E}_{t} \left\| \left\| \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{D}^{t+1}} \nabla g_{1i}(\mathbf{w}^{t+1}) \nabla f(g_{1i}(\mathbf{w}^{t+1})) - \frac{1}{n_{0}} \sum_{i=1}^{n_{0}} \nabla g_{1i}(\mathbf{w}^{t+1}) \nabla f(g_{1i}(\mathbf{w}^{t+1})) \right\|^{2} \right\|$  $+\theta^{2}\mathbb{E}_{t}\left\|\left\|\frac{1}{|\mathcal{B}|}\sum_{i=1,2,1,1}\nabla\hat{g}_{2i}(\mathbf{w}^{t+1})\nabla f(g_{2i}(\mathbf{w}^{t+1}))-\frac{1}{|\mathcal{B}|}\sum_{i=1,2,1,1}\nabla g_{2i}(\mathbf{w}^{t+1})\nabla f(g_{2i}(\mathbf{w}^{t+1}))\right\|^{2}\right\|$  $+\theta^{2}\mathbb{E}_{t}\left\|\frac{1}{|\mathcal{B}|}\sum_{i=n+1}\nabla g_{2i}(\mathbf{w}^{t+1})\nabla f(g_{2i}(\mathbf{w}^{t+1})) - \frac{1}{n_{0}}\sum_{i=1}^{n_{0}}\nabla g_{2i}(\mathbf{w}^{t+1})\nabla f(g_{2i}(\mathbf{w}^{t+1}))\right\|^{2}\right\|$  $\leq \frac{2\theta^2 L_f^2(\sigma_{\nabla g}^2 + L_g^2)}{\min\{|\mathcal{B}|, |\mathcal{B}_{1i}|, |\mathcal{B}_{2i}|\}}$ Therefore, we can get  $\mathbb{E}[\Delta_1^{t+1}] \le (1-\theta) \mathbb{E}[\Delta_1^t] + \frac{2L_F^2}{\theta} \mathbb{E}[\|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2] + 5\theta L_f^2 \tilde{C}_{\nabla g}^2 \mathbb{E}[\Xi_1^{t+1} + \Xi_2^{t+1}]$  $+3L_{f}^{2}\tilde{C}_{\nabla g}^{2}\mathbb{E}\left|\frac{1}{n_{0}}\sum_{i,j\in\mathbb{N}^{t+1}}\left(\left\|u_{1i}^{t+1}-u_{1i}^{t}\right\|^{2}+\left\|u_{2i}^{t+1}-u_{2i}^{t}\right\|^{2}\right)\right|+\frac{2\theta^{2}L_{f}^{2}(\sigma_{\nabla g}^{2}+L_{g}^{2})}{\min\{|\mathcal{B}_{1i}|,|\mathcal{B}_{2i}|\}}$ 

$$\begin{split} & \text{Lemma 4 } If \gamma_1 \leq 1/5, \text{ function value variance } \mathbb{E}_1^1 := \frac{1}{n_0} \|u_1^1 - g_1(\mathbf{w}^{1+1})\|^2 \text{ can be bounded as} \\ & \mathbb{E}[\mathbb{E}_1^{t+1}] \leq \left(1 - \frac{\gamma_1[B]}{4n_0}\right) \mathbb{E}\left[\mathbb{E}_1^t\right] + \frac{5n_0L_2^t\mathbb{E}[\||\mathbf{w}^{t+1} - \mathbf{w}^t||^2]}{\gamma_1[B]} + \frac{2\gamma_1^2\sigma_2^2[B]}{n_0[B_{11}]} - \frac{1}{4n_0} \mathbb{E}\left[\sum_{l \in B^{t+1}} \|u_{1^{t+1}}^{l+1} - u_{1^{t}}^t\|^2\right], \\ & (16) \\ & \text{Lemma 5 } If \gamma_1 \leq 1/5, \text{ function value variance } \mathbb{E}_2^t := \frac{1}{n_0} \|u_2^t - g_2(\mathbf{w}^t)\|^2 \text{ can be bounded as} \\ & \mathbb{E}[\mathbb{E}_2^{t+1}] \leq \left(1 - \frac{\gamma_1[B]}{4n_0}\right) \mathbb{E}\left[\mathbb{E}_2^t\right] + \frac{5n_0L_2^{t+0}[\|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2]}{\gamma_1[B]} + \frac{2\gamma_1^2\sigma_2^2[B]}{n_0[B_{22}]} - \frac{1}{4n_0} \mathbb{E}\left[\sum_{l \in B^{t+1}} \|u_{2^{t+1}}^{t+1} - u_{2^{t}}^t\|^2\right]. \\ & (17) \\ & \text{Since the proof of Lemma 4 and Lemma 5 are almost the same, we only presents the proof of \\ Lemma 4 as follows. \\ & \textbf{Proof Define } \phi_1^t(\mathbf{u}_1) = \frac{1}{2}\|\mathbf{u}_1 - g_1(\mathbf{w}^{t+1})\|^2 = \frac{1}{2}\mathbb{E}_{2^{t+1}}^{t+1}\|u_1^{t+1} - g_1(\mathbf{w}^{t+1}) + \frac{1}{2}\|u_1^{t+1} - u_1^t(\mathbf{w}^t) - g_1(\mathbf{w}^{t+1})\|^2 + \frac{1}{2}\|u_1^{t+1} - u_1^t\|^2 \\ & = \frac{1}{2}\||u_1^{t+1} - g_1(\mathbf{w}^{t+1})|\|^2 + \frac{1}{2}(\mathbb{E}_{2^{t+1}}^{t+1}) + \frac{1}{2}(\mathbb{E}_{2^{t+1}}^{t+1} - u_1^t)\|^2 \\ & = \frac{1}{2}\|u_1^{t+1} - g_1(\mathbf{w}^{t+1}) - g_1(\mathbf{w}^{t+1}) + g_1(\mathbf{w}^{t+1}) + u_{11}^{t+1} - u_{11}^t) \\ & + \sum_{e \in B^{t+1}} (\hat{g}_{1e}(\mathbf{w}^{t+1}) - g_{1e}(\mathbf{w}^{t+1}), u_{11}^{t+1} - u_{11}^t) + \frac{1}{2}(\mathbb{E}_{2^{t+1}}^{t+1} - g_{1e}(\mathbf{w}^{t+1}), u_{11}^{t+1} - u_{11}^t) \|^2 \\ & = \sum_{i \in B^{t+1}} (u_{1i}^{t-1} - \hat{g}_{1i}(\mathbf{w}^{t+1}), g_{1i}(\mathbf{w}^{t+1}) - u_{1i}^t) + \frac{1}{2}(\mathbb{E}_{2^{t+1}}^{t+1} + u_{1i}^t) + \frac{1}{2}(\mathbf{w}^{t+1}) +$$

$$\begin{aligned} & \text{Note that } \frac{1}{2\gamma_{1}} \sum_{i \notin \mathcal{B}^{t+1}} \|u_{1i}^{t} - g_{1i}(\mathbf{w}^{t+1})\|^{2} = \frac{1}{2\gamma_{1}} \sum_{i \notin \mathcal{B}^{t+1}} \|u_{1i}^{t+1} - g_{1i}(\mathbf{w}^{t+1})\|^{2}, \text{ which implies} \\ & \frac{1}{2\gamma_{1}} \sum_{i \in \mathcal{B}^{t+1}} (\|u_{1i}^{t} - g_{1i}(\mathbf{w}^{t+1})\|^{2} - \|u_{1i}^{t+1} - g_{1i}(\mathbf{w}^{t+1})\|^{2}) = \frac{1}{2\gamma_{1}} (\|u_{1}^{t} - g_{1}(\mathbf{w}^{t+1})\|^{2} - \|u_{1}^{t+1} - g_{1}(\mathbf{w}^{t+1})\|^{2}) \\ & \text{Esides, we also have } \mathbb{E} \left[\sum_{i \in \mathcal{B}^{t+1}} \|\hat{g}_{1i}(\mathbf{w}^{t+1}) - u_{1i}^{t}\right] \right] = \frac{|B|}{n_{0}} \sum_{i=1}^{n_{0}} \langle u_{1i}^{t} - g_{1i}(\mathbf{w}^{t+1}) - u_{1i}^{t} \rangle \\ & = -\frac{|B|}{n_{0}} |u_{1}^{t} - g_{1i}(\mathbf{w}^{t+1}), g_{1i}(\mathbf{w}^{t+1}) - u_{1i}^{t} \rangle \right] \\ & = -\frac{|B|}{n_{0}} \|u_{1}^{t} - g_{1i}(\mathbf{w}^{t+1}), g_{1i}(\mathbf{w}^{t+1}) - u_{1i}^{t} \rangle \\ & = -\frac{|B|}{n_{0}} \|u_{1}^{t} - g_{1i}(\mathbf{w}^{t+1})\|^{2}. \end{aligned}$$

$$Then we can obtain \\ & \left(\frac{1}{2} + \frac{1}{2\gamma_{1}}\right) \mathbb{E} \left[\|u_{1}^{t} + g_{1}(\mathbf{w}^{t+1})\|^{2}\right] + \frac{\gamma_{1}|B|\sigma_{g}^{2}}{|B_{1i}|} - \frac{\gamma_{1} + 1}{8\gamma_{1}} \mathbb{E} \left[\sum_{i \in \mathcal{B}^{t+1}} \|u_{1i}^{t+1} - u_{1i}^{t}\|^{2}\right]. \end{aligned}$$

$$Divide both sides by \frac{\gamma_{1}+1}{\gamma_{1}} we can get \\ \mathbb{E} \left[\|u_{1}^{t+1} - g_{1}(\mathbf{w}^{t+1})\|^{2}\right] \leq \frac{\gamma_{1} + 1 - 2\gamma_{1}|B|}{\gamma_{1} + 1} \mathbb{E} \left[\|u_{1}^{t} - g_{1}(\mathbf{w}^{t+1})\|^{2}\right] + \frac{1}{2} \frac{\gamma_{1}|B|\sigma_{g}^{2}}{2\gamma_{1} + 1} - \frac{\gamma_{1}|B|}{|B_{1i}|} \right] \\ - \frac{1}{4} \mathbb{E} \left[\sum_{i \in \mathcal{B}^{t+1}} \|u_{1i}^{t+1} - u_{1i}^{t}\|^{2}\right]. \end{aligned}$$

$$Note that \frac{\gamma_{1}+1-2\gamma_{1}\frac{m_{0}}{m_{0}}}{\gamma_{1}+1} \leq \frac{\gamma_{1}+1-2\gamma_{1}|B|}{(\gamma_{1}+1)m_{0}} \leq (1 - \frac{\gamma_{1}|B|}{(\gamma_{1}+1)m_{0}} \leq 1 - \frac{\gamma_{1}|B|}{2m_{0}} \|ad \frac{1}{\gamma_{1}+1} \leq 1 for \gamma_{1} \in (0, 1].$$

$$Besides, we have \|u_{1}^{t} - g_{1}(\mathbf{w}^{t+1})\|^{2} \leq (1 + \frac{\gamma_{1}|B|}{(\gamma_{1}+1)m_{0}} \leq (1 - \frac{\gamma_{1}|B|}{4m_{0}})(1 - \frac{\gamma_{1}|B|}{4m_{0}})(1 - \frac{\gamma_{1}|B|}{2m_{0}}) = (1 - \frac{\gamma_{1}|B|}{2m_{0}}) = (1 - \frac{\gamma_{1}|B|}{m_{0}}) = \left(\frac{1}{m_{0}} \|u_{1}^{t} - g_{1}(\mathbf{w}^{t+1})\|^{2}\right]$$

$$= \left(1 - \frac{\gamma_{1}|B|}{m_{0}} \|E \left[\frac{1}{m_{0}} \|u_{1}^{t} - g_{1}(\mathbf{w}^{t+1})\|^{2}\right] + \frac{5m_{0}L_{g}^{2}|B|}{\gamma_{1}|B|} + \frac{2\gamma_{1}^{2}\tau_{0}^{2}|B|}{m_{0}|B|} - \frac{1}{m_{0}|B|} = \frac{1}{2m_{0}} \left[\frac{1}{m_{0}} \left[\frac{1}{$$

$$\begin{array}{ll} \text{1606}\\ \text{1607}\\ \text{1608}\\ \text{1608}\\ \text{1608}\\ \mathbb{E}[\Delta_{2}^{t+1}] \leq (1-\theta)\mathbb{E}[\Delta_{2}^{t}] + \frac{2\beta^{2}L_{H}^{2}}{\theta}\mathbb{E}\left[\|\mathbf{w}^{t+1} - \mathbf{w}^{t}\|^{2}\right] + 5\theta\beta^{2}\tilde{C}_{\nabla h}^{2}\mathbb{E}[\Gamma_{t+1}]\\ \text{1609}\\ \text{1610}\\ \text{1610}\\ \text{1611}\\ \text{1612}\\ \text{1612}\\ \text{1612}\\ \text{1613}\\ \text{with } \Gamma_{t+1} := \frac{1}{m}\|\boldsymbol{u}^{t+1} - \boldsymbol{h}(\mathbf{w}^{t+1})\|^{2}. \end{array}$$
(19)

Proof  
1618  
1619  

$$\Delta_2^{t+1} = \|v_2^{t+1} - \nabla H(\mathbf{w}^{t+1})\|^2 = \|(1-\theta)v_2^t + \theta G_2^t - \nabla H(\mathbf{w}^{t+1})\|^2$$

$$\|(\widehat{I} + \widehat{Q}) + \widehat{Q} + \widehat{Q}\|^2,$$

where (1), (2), (3) and (4) are defined as  $(1-\theta)(v_2^t - \nabla H(\mathbf{w}^t)), \quad (2) = (1-\theta)(\nabla H(\mathbf{w}^t) - \nabla H(\mathbf{w}^{t+1})),$  $(\mathfrak{J}) = \frac{\theta}{|\mathcal{B}_c|} \beta \sum_{k \in \mathbf{v}^{t+1}} \left( [u_k^t]_+ \nabla \hat{h}_k(\mathbf{w}^{t+1}) - [h_k(\mathbf{w}^{t+1})]_+ \nabla \hat{h}_k(\mathbf{w}^{t+1}) \right)$  $(\widehat{\mathcal{A}}) = \theta \left( \frac{1}{|\mathcal{B}_c|} \beta \sum_{\substack{l \in \mathbf{v}^{t+1}}} [h_k(\mathbf{w}^{t+1})]_+ \nabla \hat{h}_k(\mathbf{w}^{t+1}) - \nabla H(\mathbf{w}^{t+1}) \right)$ Note that  $\mathbb{E}_t[\langle (\underline{I}, \underline{4}) \rangle] = \mathbb{E}_t[\langle (\underline{2}, \underline{4}) \rangle] = 0$ . Then, by the Young's inequality, we can get  $\mathbb{E}_{t}\left[\|(1)+(2)+(3)+(4)\|^{2}\right]$  $= \|\widehat{(I)}\|^{2} + \|\widehat{(2)}\|^{2} + \mathbb{E}_{t}\|\widehat{(3)}\|^{2} + \mathbb{E}_{t}\|\widehat{(4)}\|^{2} + 2\langle \widehat{(I)}, \widehat{(2)}\rangle + 2\mathbb{E}_{t}[\langle \widehat{(I)}, \widehat{(3)}\rangle] + 2\mathbb{E}_{t}[\langle \widehat{(2)}, \widehat{(3)}\rangle] + 2\mathbb{E}_{t}[\langle \widehat{(3)}, \widehat{(4)}\rangle]$  $\leq (1+\theta) \|\widehat{(1)}\|^2 + 2\left(1+\frac{1}{\theta}\right) \|\widehat{(2)}\|^2 + \frac{2+3\theta}{\theta} \mathbb{E}_t \|\widehat{(3)}\|^2 + 2\mathbb{E}_t \|\widehat{(4)}\|^2.$ We can also get  $(1+\theta)\|\widehat{(I)}\|^2 = (1+\theta)(1-\theta)^2\|v_2^t - \nabla H(\mathbf{w}^t)\|^2 \le (1-\theta)\|v_2^t - \nabla H(\mathbf{w}^t)\|^2$  $2\left(1+\frac{1}{\theta}\right)\|\widehat{\mathcal{Q}}\|^2 = 2\left(1+\frac{1}{\theta}\right)(1-\theta)^2\|\nabla H(\mathbf{w}^t) - \nabla H(\mathbf{w}^{t+1})\|^2$  $\leq \frac{2}{\theta} \left\| \frac{1}{m} \sum_{k=1}^{m} \beta \left( \nabla h_k(\mathbf{w}^{t+1})^\top [h_k(\mathbf{w}^{t+1})]_+ - \nabla h_k(\mathbf{w}^t)^\top [h_k(\mathbf{w}^t)]_+ \right) \right\|^2$  $\leq \frac{2\beta^2 L_H^2}{\rho} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2$  $\frac{2+3\theta}{\theta} \|\widehat{\mathcal{J}}\|^2 \le \frac{2+3\theta}{\theta} \frac{\theta^2 \beta^2}{|\mathcal{B}_c|} \sum_{\substack{k \in \mathcal{B}^{t+1}}} \|\nabla \hat{h}_k(\mathbf{w}^{t+1})\|^2 \|[u_k^t]_+ - [h_k(\mathbf{w}^{t+1})]_+\|^2$  $\leq \frac{(2+3\theta)\theta\beta^2}{|\mathcal{B}_c|} \sum_{k=nt+1} \|\nabla \hat{h}_k(\mathbf{w}^{t+1})\|^2 \|u_k^t - h_k(\mathbf{w}^{t+1})\|^2$ Consider that  $\mathbf{w}^{t+1}$  and  $u_k^t$  do not depend on either  $\mathcal{B}_c^{t+1}$  or  $\mathcal{B}_k$ , we have  $(2+3\theta)\theta\beta^{2}\mathbb{E}_{t}\left[\frac{1}{|\mathcal{B}_{c}|}\sum_{k=n^{t+1}}\|\nabla\hat{h}_{k}(\mathbf{w}^{t+1})\|^{2}\|u_{k}^{t}-h_{k}(\mathbf{w}^{t+1})\|^{2}\right]$  $= (2+3\theta)\theta\beta^{2}\mathbb{E}_{t} \left| \frac{1}{|\mathcal{B}_{c}|} \sum_{k, c \neq t+1} \mathbb{E}_{t} \left[ \|\nabla \hat{h}_{k}(\mathbf{w}^{t+1})\|^{2} | k \in \mathcal{B}_{c}^{t+1} \right] \|u_{k}^{t} - h_{k}(\mathbf{w}^{t+1})\|^{2} \right|$  $\leq (2+3\theta)\theta\beta^2 \tilde{C}_{\nabla h}^2 \mathbb{E}_t \left[ \frac{1}{|\mathcal{B}_c|} \sum_{\mathbf{w} = t+1} \|u_k^t - h_k(\mathbf{w}^{t+1})\|^2 \right]$  $\leq \frac{(2+3\theta)\theta(1+\delta)\beta^{2}\tilde{C}_{\nabla h}^{2}}{m} \sum_{k \in [m]} \mathbb{E}_{t} \left[ \|u_{k}^{t+1} - h_{k}(\mathbf{w}^{t+1})\|^{2} \right] + \frac{(2+3\theta)\theta(1+1/\delta)\beta^{2}\tilde{C}_{\nabla h}^{2}}{m} \mathbb{E}_{t} \left[ \sum_{k \in [m]} \|u_{k}^{t+1} - u_{k}^{t}\|^{2} \right]$  $=\frac{(2+3\theta)\theta(1+\delta)\beta^{2}\tilde{C}_{\nabla h}^{2}}{m}\sum_{k=1,\dots,k=1}\mathbb{E}_{t}\left[\|u_{k}^{t+1}-h_{k}(\mathbf{w}^{t+1})\|^{2}\right]+\frac{(2+3\theta)\theta(1+1/\delta)\beta^{2}\tilde{C}_{\nabla h}^{2}}{m}\mathbb{E}_{t}\left[\sum_{k=1,\dots,k=1}\|u_{k}^{t+1}-u_{k}^{t}\|^{2}\right]$ where the last equation holds by noting that  $u_k^{t+1} = u_k^t$  for all  $i \notin \mathcal{B}_c^{t+1}$ . If  $\theta \leq \frac{1}{3}$  and  $\delta = \frac{3\theta}{2}$ , we have  $(2+3\beta)\beta(1+\delta) \leq 5\theta$  and  $(2+3\beta)\beta(1+1/\delta) \leq 3$ . Therefore, we can get  $\mathbb{E}\left[\frac{2+3\theta}{\theta}\|\widehat{\mathcal{I}}\|^2\right] \leq 5\theta\beta^2 \tilde{C}_{\nabla h}^2 \mathbb{E}[\Gamma_{t+1}] + \frac{3\beta^2 \tilde{C}_{\nabla h}^2}{m} \mathbb{E}\left|\sum_{k=p^{t+1}} \|u_k^{t+1} - u_k^t\|^2\right|$ 

*Next, we give the upper bound of*  $\mathbb{E}_t || \widehat{(4)} ||^2$ .  $\mathbb{E}_{t} \|\widehat{\mathcal{A}}\|^{2} = \theta^{2} \beta^{2} \mathbb{E}_{k} \left\| \left\| \frac{1}{|\mathcal{B}_{c}|} \sum_{k \in \mathcal{B}^{t+1}} [h_{k}(\mathbf{w}^{t+1})]_{+} \nabla \hat{h}_{k}(\mathbf{w}^{t+1}) - \frac{1}{m} \sum_{k=1}^{m} [h_{k}(\mathbf{w}^{t+1})]_{+} \nabla h_{k}(\mathbf{w}^{t+1}) \right\|^{2} \right\|^{2}$  $\leq \theta^2 \beta^2 \mathbb{E}_t \left[ \left\| \frac{1}{|\mathcal{B}_c|} \sum_{k \in \mathcal{B}^{t+1}} [h_k(\mathbf{w}^{t+1})]_+ \nabla \hat{h}_k(\mathbf{w}^{t+1}) - \frac{1}{|\mathcal{B}_c|} \sum_{k \in \mathcal{B}^{t+1}} [h_k(\mathbf{w}^{t+1})]_+ \nabla h_k(\mathbf{w}^{t+1}) \right\|^2 \right]$  $+\theta^{2}\beta^{2}\mathbb{E}_{t}\left[\left\|\frac{1}{|\mathcal{B}_{c}|}\sum_{k\in\mathcal{B}_{c}^{t+1}}[h_{k}(\mathbf{w}^{t+1})]_{+}\nabla h_{k}(\mathbf{w}^{t+1})-\frac{1}{m}\sum_{k=1}^{m}[h_{k}(\mathbf{w}^{t+1})]_{+}\nabla h_{k}(\mathbf{w}^{t+1})\right\|^{2}\right]$  $\mathcal{I} \theta^2 \beta^2 C_h^2 (\sigma_{\nabla h}^2 + L_h^2)$ 

$$\leq \min\{|\mathcal{B}_c|, |\mathcal{B}_k|\}$$

Combine above inequalities, we can get

$$\mathbb{E}[\Delta_{2}^{t+1}] \leq (1-\theta)\mathbb{E}[\Delta_{2}^{t}] + \frac{2\beta^{2}L_{H}^{2}}{\theta}\mathbb{E}\left[\|\mathbf{w}^{t+1} - \mathbf{w}^{t}\|^{2}\right] + 5\theta\beta^{2}\tilde{C}_{\nabla h}^{2}\mathbb{E}[\Gamma_{t+1}] \\ + \frac{3\beta^{2}\tilde{C}_{\nabla h}^{2}}{m}\mathbb{E}\left[\sum_{k\in\mathcal{B}_{c}^{t+1}}\|u_{k}^{t+1} - u_{k}^{t}\|^{2}\right] + \frac{\theta^{2}\beta^{2}C_{h}^{2}(\sigma_{\nabla h}^{2} + L_{h}^{2})}{\min\{|\mathcal{B}_{c}|, |\mathcal{B}_{k}|\}}.$$

$$\begin{array}{ll} \text{1698} & \text{Lemma 7 } If \gamma_2 \leq 1/5, \text{ function value variance } \Gamma_t := \frac{1}{m} \| \boldsymbol{u}^t - \boldsymbol{h}(\mathbf{w}^t) \|^2 \text{ can be bounded as} \\ \text{1699} & \\ \text{1700} & \\ \mathbb{E}[\Gamma_{t+1}] \leq \left(1 - \frac{\gamma_2 |\mathcal{B}_c|}{4m}\right) \mathbb{E}[\Gamma_t] + \frac{5mL_h^2 \mathbb{E}[\| \mathbf{w}^{t+1} - \mathbf{w}^t \|^2]}{\gamma |\mathcal{B}_c|} + \frac{2\gamma_2^2 \sigma_h^2 |\mathcal{B}_c|}{m |\mathcal{B}_k|} - \frac{1}{4m} \mathbb{E}\left[\sum_{k \in \mathcal{B}_c^{t+1}} \| \boldsymbol{u}_k^{t+1} - \boldsymbol{u}_k^t \|^2\right] \\ \text{1702} & \\ \text{(20)} \end{array}$$

$$\begin{array}{l} \text{Proof Define } \psi_{k}(\boldsymbol{u}) = \frac{1}{2} \|\boldsymbol{u} - \boldsymbol{h}(\mathbf{w}^{t})\|^{2} = \frac{1}{2} \sum_{k=1}^{m} \|\boldsymbol{u}_{k} - \boldsymbol{h}_{k}(\mathbf{w}^{t})\|^{2}, \text{ which is 1-strongly convex.} \\ \psi_{t+1}(\boldsymbol{u}^{t+1}) = \frac{1}{2} \|\boldsymbol{u}^{t+1} - \boldsymbol{h}(\mathbf{w}^{t+1})\|^{2} = \frac{1}{2} \|\boldsymbol{u}^{t} - \boldsymbol{h}(\mathbf{w}^{t+1})\|^{2} + \langle \boldsymbol{u}^{t} - \boldsymbol{h}(\mathbf{w}^{t+1}), \boldsymbol{u}^{t+1} - \boldsymbol{u}^{t} \rangle + \frac{1}{2} \|\boldsymbol{u}^{t+1} - \boldsymbol{u}^{t}\|^{2} \\ = \frac{1}{2} \|\boldsymbol{u}^{t} - \boldsymbol{h}(\mathbf{w}^{t+1})\|^{2} + \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle \boldsymbol{u}_{k}^{t} - \hat{\boldsymbol{h}}_{k}(\mathbf{w}^{t+1}), \boldsymbol{u}_{k}^{t+1} - \boldsymbol{u}_{k}^{t} \rangle + \frac{1}{2} \sum_{k \in \mathcal{B}_{c}^{t+1}} \|\boldsymbol{u}_{k}^{t+1} - \boldsymbol{u}_{k}^{t}\|^{2} \\ + \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle \hat{\boldsymbol{h}}_{k}(\mathbf{w}^{t+1}) - \boldsymbol{h}_{k}(\mathbf{w}^{t+1}), \boldsymbol{u}_{k}^{t+1} - \boldsymbol{u}_{k}^{t} \rangle \end{array}$$

$$\tag{21}$$

$$Note that u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}) = (q_{i}^{k} - q_{i}^{k+1})/\gamma_{2} and 2\langle b - a, a - c \rangle \leq \|b - c\|^{2} - \|a - b\|^{2} - \|a - c\|^{2}.$$

$$\sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}), u_{k}^{t+1} - u_{k}^{t} \rangle$$

$$= \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}), h_{k}(\mathbf{w}^{t+1}) - u_{k}^{t} \rangle + \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}), u_{k}^{t+1} - h_{k}(\mathbf{w}^{t+1}) \rangle$$

$$= \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}), h_{k}(\mathbf{w}^{t+1}) - u_{k}^{t} \rangle + \frac{1}{\gamma_{2}} \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - u_{k}^{t+1}, u_{k}^{t+1} - h_{k}(\mathbf{w}^{t+1}) \rangle$$

$$= \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}), h_{k}(\mathbf{w}^{t+1}) - u_{k}^{t} \rangle + \frac{1}{\gamma_{2}} \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - u_{k}^{t+1}, u_{k}^{t+1} - h_{k}(\mathbf{w}^{t+1}) \rangle$$

$$= \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - \hat{h}_{k}(\mathbf{w}^{t+1}), h_{k}(\mathbf{w}^{t+1}) - u_{k}^{t} \rangle + \frac{1}{\gamma_{2}} \sum_{k \in \mathcal{B}_{c}^{t+1}} \langle u_{k}^{t} - h_{k}(\mathbf{w}^{t+1}) | 2 - \|u_{k}^{t+1} - u_{k}^{t}\|^{2} - \|u_{k}^{t+1} - h_{k}(\mathbf{w}^{t+1})\|^{2} \rangle$$

$$\begin{aligned} & |I| \gamma_{2} \leq \frac{1}{3}, \text{ we have} \\ & - \frac{1}{2} \left( \frac{1}{\gamma_{2}} - 1 - \frac{\gamma_{2} + 1}{4\gamma_{2}} \right) \sum_{k \in \mathbb{R}^{d+1}} \|h_{k}^{k+1} - u_{k}^{k}\|^{2} + \sum_{k \in \mathbb{R}^{d+1}} \left( \hat{h}_{k} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}), u_{k}^{t+1} - u_{k}^{k} \right) \\ & \leq -\frac{1}{4\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|h_{k}^{k+1} - u_{k}^{k}\|^{2} + \gamma_{2} \sum_{k \in \mathbb{R}^{d+1}} \|h_{1} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}) \|^{2} + \frac{1}{4\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k+1} - u_{k}^{k}\|^{2} \\ & = \gamma_{2} \sum_{k \in \mathbb{R}^{d+1}} \|h_{k}^{k} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}) \|^{2} + \frac{1}{2\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k} - h_{k} (\mathbf{w}^{t+1}) \|^{2} - \frac{1}{2\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k+1} - h_{k} (\mathbf{w}^{t+1}) \|^{2} \\ & + \gamma_{2} \sum_{k \in \mathbb{R}^{d+1}} \|h_{k}^{k} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}) \|^{2} - \frac{1}{2\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k+1} - u_{k}^{k}\|^{2} \\ & + \gamma_{2} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k} - h_{k} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}) \|^{2} - \frac{1}{2\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k+1} - u_{k}^{k}\|^{2} \\ & + \gamma_{2} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k} - h_{k} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}) - u_{k}^{k} \right). \end{aligned}$$
Note that  $\frac{1}{2\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k} - h_{k} (\mathbf{w}^{t+1}) \|^{2} = \frac{1}{2\gamma_{2}} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k+1} - h_{k} (\mathbf{w}^{t+1}) \|^{2} \right) = \frac{1}{2\gamma_{2}} \left( \|u^{k} - h_{k} (\mathbf{w}^{t+1}) \|^{2} - \|u^{k+1} - h_{k} (\mathbf{w}^{t+1}) \|^{2} \right) \\ \\ Besides, we also have \mathbb{E} \left[ \sum_{k \in \mathbb{R}^{d+1}} \|h_{k} (\mathbf{w}^{t+1}) - u_{k}^{k} \right] \right] = \frac{|B_{k}|}{m} \sum_{k \in \mathbb{R}^{d+1}} \|u_{k}^{k} - h_{k} (\mathbf{w}^{t+1}) - h_{k} (\mathbf{w}^{t+1}) \|^{2} \right] \\ \\ & = \left[ \frac{1}{2\gamma_{2}} \left\|u^{k} - h_{k} (\mathbf{w}^{t+1}) \right] \right] \mathbb{E} \left[ \|u^{k+1} - h_{k} (\mathbf{w}^{t+1}) \|^{2} \right] \\ \\ & = \left[ \frac{1}{2} \frac{1}{2\gamma_{2}} \sum \left[ u^{k+1} - u^{k} \|^{2} \right] \right] \\ \\ Besides, we also have \mathbb{E} \left[ \sum_{k \in \mathbb{R}^{d+1}} \|h_{k} (\mathbf{w}^{t+1}) - u_{k}^{k} \right] \right] \\ \\ & = \left[ \frac{1}{2} \frac{1}{2\gamma_{2}} \sum \left[ \frac{1}{2} \frac{1}{2\gamma_{2}$ 

$$\begin{aligned} & \text{final} \\ \text{final} \\ \text{final} \\ \text{function} \\ \text{for the transform} \\ \text{final} \\ \text{function} \\ \text{final} \\ \text{final}$$

We state the main theorem again for convenience and present the proof. 

**Theorem 2** Suppose Assumptions 1, 2, 3 and 4 hold, and set  $\beta = \frac{1}{\epsilon\delta}$ ,  $\theta = \min\{\frac{\epsilon^4\delta^2\min\{|\mathcal{B}_k|,|\mathcal{B}_c|\}}{672(\sigma_{\nabla h}^2 + L_h^2)}, \frac{\epsilon^2\min\{|\mathcal{B}|,|\mathcal{B}_{1i}|,|\mathcal{B}_{2i}|\}}{1344L_f^2(\sigma_{\nabla g}^2 + L_g^2)}\}, \gamma_1 = \gamma_2 = \min\{\frac{5n_0\theta}{3|\mathcal{B}|}, \frac{5m\theta}{3|\mathcal{B}_c|}, \frac{\epsilon^4\delta^2|\mathcal{B}_k|}{26880\sigma_h^2\tilde{C}_{\nabla h}^2}\} and$   $\eta = \min\{\frac{1}{12(L_F + \beta L_H)}, \frac{\theta}{8\sqrt{3}L_F}, \frac{\theta}{8\sqrt{3}L_H\beta}, \frac{\gamma_1|\mathcal{B}|}{40\sqrt{6}L_gL_f\tilde{C}_{\nabla g}n_0}, \frac{\gamma_2|\mathcal{B}_c|}{40\sqrt{6}\beta L_h\tilde{C}_{\nabla h}m}\}.$  Then there exists  $\lambda$  such that such that  $\mathbb{E}\left[ \|\nabla F(\mathbf{w}^{\hat{t}}) + \nabla \boldsymbol{h}(\mathbf{w}^{\hat{t}})\boldsymbol{\lambda})\| \right] \leq \epsilon$  $\mathbb{E}[\|[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_+\|] < \epsilon$  $\mathbb{E}[\boldsymbol{\lambda}^{\top}[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}] \leq \epsilon$ 

with number of iterations T of Algorithm 1 bounded by  $O(\epsilon^{-7}\delta^{-3})$  and  $\hat{t}$  selected uniformly at random from  $\{1, \cdots, T\}$ . 

**Proof** Since  $\Phi(\mathbf{w})$  is  $L_{\beta}$ -smooth with  $L_{\beta} = L_F + \beta L_H$  where  $L_F := 2(L_{\nabla g}L_f + L_{\nabla f}L_g^2)$  and  $L_H := L_{\nabla h}C_h + L_hC_{\nabla h}$ , we have 

$$\begin{aligned} & \Phi(\mathbf{w}^{t+1}) \leq \Phi(\mathbf{w}^t) + \langle \nabla \Phi(\mathbf{w}^t), \mathbf{w}^{t+1} - \mathbf{w}^t \rangle + \frac{L_{\beta}}{2} \| \mathbf{w}^{t+1} - \mathbf{w}^t \|^2 \\ & = \Phi(\mathbf{w}^t) + \langle v^t, \mathbf{w}^{t+1} - \mathbf{w}^t \rangle + \langle \nabla \Phi(\mathbf{w}^t) - v^t, \mathbf{w}^{t+1} - \mathbf{w}^t \rangle + \frac{L_{\beta}}{2} \| \mathbf{w}^{t+1} - \mathbf{w}^t \|^2 \\ & = \Phi(\mathbf{w}^t) + \langle v^t, \mathbf{w}^{t+1} - \mathbf{w}^t \rangle + \left( \frac{L_{\beta}}{2} + \frac{1}{4\eta} \right) \| \mathbf{w}^{t+1} - \mathbf{w}^t \|^2 + \eta \| \nabla \Phi(\mathbf{w}^t) - v^t \|^2. \end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

(22)

Since  $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta v^t$ , which is equivalent to  $\mathbf{w}^{t+1} = \arg\min_{\mathbf{w}} \langle v^t, \mathbf{w} \rangle + \frac{1}{2\eta} \|\mathbf{w} - \mathbf{w}^t\|^2$ , we have ) 

$$\langle v^t, \mathbf{w}^{t+1} - \mathbf{w}^t \rangle \le -\frac{1}{2\eta} \| \mathbf{w}^{t+1} - \mathbf{w}^t \|^2.$$
 (23)

Then we can get 

$$\Phi(\mathbf{w}^{t+1}) \le \Phi(\mathbf{w}^t) + \left(\frac{L_{\beta}}{2} - \frac{1}{4\eta}\right) \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 + \eta \|\nabla\Phi(\mathbf{w}^t) - v^t\|^2$$
(24)

$$\Phi(\mathbf{w}^{t+1}) \le \Phi(\mathbf{w}^{t}) + \left(\frac{L_{\beta}}{2} - \frac{1}{4\eta}\right) \|\mathbf{w}^{t+1} - \mathbf{w}^{t}\|^{2} + 2\eta \|\nabla\Phi(\mathbf{w}^{t}) - v^{t}\|^{2} - \eta \|\nabla\Phi(\mathbf{w}^{t}) - v^{t}\|^{2}$$
(25)

$$\eta \|\nabla \Phi(\mathbf{w}^{t}) - v^{t}\|^{2} \leq \Phi(\mathbf{w}^{t}) - \Phi(\mathbf{w}^{t+1}) + \left(\frac{L_{\beta}}{2} - \frac{1}{4\eta}\right) \|\mathbf{w}^{t+1} - \mathbf{w}^{t}\|^{2} + 2\eta \|\nabla \Phi(\mathbf{w}^{t}) - v^{t}\|^{2}.$$
(26)

Then we want to bound  $\mathbb{E} \| \nabla \Phi(\mathbf{w}^t) - v^t \|^2$ . 

$$\begin{aligned} \|\nabla\Phi(\mathbf{w}^{t}) - v^{t}\|^{2} &= \|(1-\theta)(v_{1}^{t-1} + v_{2}^{t-1}) + \theta(G_{1}^{t} + G_{2}^{t}) - \nabla\Phi(\mathbf{w}^{t})\|^{2} \\ &= \|(1-\theta)v_{1}^{t-1} + \theta G_{1}^{t} - \nabla F(\mathbf{w}^{t}) + (1-\theta)v_{2}^{t-1} + \theta G_{2}^{t} - \nabla H(\mathbf{w}^{t})\|^{2} \end{aligned}$$
(27)

$$= \|v_1^t - \nabla F(\mathbf{w}^t) + v_2^t - \nabla H(\mathbf{w}^t)\|^2$$

$$\begin{aligned} & \text{Since } \mathbb{E}_{k} [ \nabla \phi(\mathbf{w}^{+}), \mathbf{v}_{t+1} ]^{2} = \mathbb{E}_{t} | \mathbf{v}_{t}^{k+1} - \nabla F(\mathbf{w}^{t+1}) |^{2} + \mathbb{E}_{t} | \mathbf{v}_{t}^{k+1} - \nabla H(\mathbf{w}^{t+1}) | ^{2} = (28) \\ & \text{Summing } (15), \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \tilde{C}_{t}^{2} m^{2}}{\gamma_{1} | \mathbf{S} |} (16) \text{ and } \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \tilde{C}_{t}^{2} m^{2}}{\gamma_{1} | \mathbf{S} |} \mathbb{E} [ \| \mathbf{u}_{t}^{k+1} - \nabla F(\mathbf{w}^{t+1}) \|^{2} \| \mathbb{E} \| \| \mathbf{u}_{t}^{k+1} - \mathbf{w}^{t} \|^{2} ] \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \tilde{C}_{t}^{2} m^{2}}{\gamma_{1} | \mathbf{S} |} \left( 1 - \frac{2}{H} \frac{2}{H} + \frac{100 L_{t}^{2} \tilde{L}_{t}^{2} \tilde{C}_{t}^{2} m^{2} \tilde{C}_{t}^{2} m^{2} }{\gamma_{1} | \mathbf{S} |} \right) \mathbb{E} [ \frac{1}{n_{t}} \sum_{e \in \mathbb{R}^{k+1}} \| \mathbf{u}_{t+1}^{t+1} - \mathbf{u}_{t}^{t} \|^{2} + \| \mathbf{u}_{2}^{t+1} - \mathbf{u}_{2}^{t} \|^{2} \end{bmatrix} \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \tilde{C}_{t}^{2} m^{2} (1 - \gamma_{1} | \mathbf{S} |)}{\gamma_{1} | \mathbf{S} |} \right) \mathbb{E} [ \frac{1}{n_{t}} \sum_{e \in \mathbb{R}^{k+1}} \| \mathbf{u}_{t+1}^{t+1} - \mathbf{u}_{t}^{t} \|^{2} + \| \mathbf{u}_{2}^{t+1} - \mathbf{u}_{2}^{t} \|^{2} \end{bmatrix} \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \tilde{C}_{t}^{2} m^{2} L_{t}^{2} \\ m(|\mathbf{S}|_{t+1}|| \mathbf{S}_{2} || \mathbf{S} \| + \frac{100 L_{t}^{2} L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{m_{t}} \| \mathbf{u}_{t+1}^{t+1} - \mathbf{u}_{t}^{t} \|^{2} \end{bmatrix} \\ & \quad + \frac{20^{2} L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}^{2} + \frac{200 r_{t}^{2} L_{t}^{2} (2 \tilde{C}_{t}^{2} + L_{t}^{2})}{m(|\mathbf{S}|_{t+1}||} \mathbb{E} [ \mathbf{I}_{t}^{t+1} - \mathbf{u}_{t}^{t} \|^{2} \end{bmatrix} \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}^{2} + \frac{200 r_{t}^{2} L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{\eta_{t}} \| \mathbb{E} | \mathbf{I}_{t} - \Gamma_{t+1} | \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}^{2} + \frac{200 r_{t}^{2} L_{t}^{2} L_{t}^{2} \sigma_{t}^{2}}{\eta_{t}} \| \mathbb{E} [ \mathbf{I}_{t}^{t+1} - \mathbf{u}_{t}^{t} \| \|^{2} \end{bmatrix} \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{\eta_{t}} \| \mathbb{E} | \mathbf{I}_{t} - \Gamma_{t+1} | \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{\eta_{t}} \| \mathbb{E} | \mathbf{I}_{t} - \mathbf{U}_{t} \| \|^{2} \end{bmatrix} \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{\eta_{t}} \| \mathbb{E} | \mathbf{I}_{t} - \Gamma_{t+1} \| \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{\eta_{t}} \| \mathbb{E} | \mathbf{I}_{t} - \Gamma_{t+1} \| \\ & \quad + \frac{200 L_{t}^{2} \tilde{C}_{t}^{2} \sigma_{t}}{\eta_{t}} \| \mathbb{E} | \| \mathbb{E} [ \mathbf{I}_{t} - \mathbf{I}_{t} \| \mathbb{E} ] \| \| \mathbb{$$

Dividing both sides by  $\frac{\eta}{24}$  and taking the average over T we can get  $\frac{1}{T} \sum_{t=1}^{t-1} \mathbb{E} \left[ \eta^{-2} \| \mathbf{w}^{t+1} - \mathbf{w}^t \|^2 + \| \nabla \Phi(\mathbf{w}^t) - v^t \|^2 \right]$  $\leq \frac{24\mathbb{E}[Y_0]}{nT} + \frac{192\theta L_f^2(\sigma_{\nabla g}^2 + L_g^2)}{\min\{|\mathcal{B}|_i|, |\mathcal{B}_{2i}|\}} + \frac{7680\gamma\sigma_g^2 L_f^2 \tilde{C}_{\nabla g}^2}{\min\{|\mathcal{B}_{1i}|, |\mathcal{B}_{2i}|\}} + \frac{96\theta\beta^2(\sigma_{\nabla h}^2 + L_h^2)}{\min\{|\mathcal{B}_{c}|, |\mathcal{B}_{k}|\}} + \frac{3840\gamma\beta^2\sigma_h^2 \tilde{C}_{\nabla h}^2}{|\mathcal{B}_{k}|}$ (33) $Y_{0} = \Phi(\mathbf{w}^{0}) + \frac{4\eta}{4} \|\nabla\Phi(\mathbf{w}^{0}) - v^{0}\|^{2} + \frac{80\eta L_{f}^{2} C_{\nabla g}^{2} n_{0}}{2|\mathcal{B}|} (\Xi_{1}^{0} + \Xi_{2}^{0}) + \frac{80\eta \beta^{2} \tilde{C}_{\nabla h}^{2} m}{2|\mathcal{B}|} \Gamma_{0}$  $= \Phi(\mathbf{w}^{0}) + \frac{4\eta}{\theta} \|\nabla\Phi(\mathbf{w}^{0}) - v^{0}\|^{2} + \frac{80\eta L_{f}^{2}C_{\nabla g}^{2}}{\gamma|\mathcal{B}|} (\|\boldsymbol{u}_{1}^{0} - \boldsymbol{g}_{1}(\mathbf{w}^{0})\|^{2} + \|\boldsymbol{u}_{2}^{0} - \boldsymbol{g}_{2}(\mathbf{w}^{0})\|^{2})$  $+\frac{80\eta\beta^2C_{\nabla h}^2}{\gamma^2|\mathcal{B}|}\|\boldsymbol{u}^0-\boldsymbol{h}(\mathbf{w}^0)\|^2.$ Since  $\mathbf{w}^0$  is a feasible solution, we have  $\Phi(\mathbf{w}^0) = \frac{1}{n_0} \sum_{i=1}^{n_0} f(g_i(\mathbf{w}^0))$ . Since g is bounded by Assumption 1 and f is Lipschitz continuous, we can show that there exists a constant  $C_F :=$  $\max\{\tau | \log(c_g^2)|, \tau | \log(C_g^2)| \} \text{ such that } |F(\mathbf{w}, \mathcal{D})| \le C_F. \text{ We assume that } u_k^0 = h_k(\mathbf{w}^0), \ \hat{u}_{1i}^0 = \hat{g}_{1i}(\mathbf{w}^0), \ u_{2i}^0 = \hat{g}_{2i}(\mathbf{w}^0) \text{ and } v^0 = \frac{1}{n_0} \sum_{i=1}^{n_0} (\nabla \hat{g}_{1i}(\mathbf{w}^0)^\top \nabla f(\hat{g}_{1i}(\mathbf{w}^0) + \nabla \hat{g}_{2i}(\mathbf{w}^0)^\top \nabla f(\hat{g}_{2i}(\mathbf{w}^0)),$ we can get  $\mathbb{E}[\|\nabla \Phi(\mathbf{w}^0) - v^0\|^2] \le 2(C_{\nabla a}^2 + \sigma_{\nabla a}^2)C_{\nabla f}^2$  $\mathbb{E}[\|\boldsymbol{u}_1^0 - \boldsymbol{g}_1(\mathbf{w}^0)\|^2] \le \sigma_a^2$ (34) $\mathbb{E}[\|\boldsymbol{u}_2^0 - \boldsymbol{g}_2(\mathbf{w}^0)\|^2] \le \sigma_a^2$  $\mathbb{E}[\|\boldsymbol{u}^0 - \boldsymbol{h}(\mathbf{w}^0)\|^2] = 0$ Therefore, we can get  $\frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}\left[\eta^{-2} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 + \|\nabla \Phi(\mathbf{w}^t) - v^t\|^2\right]$  $\leq \frac{24C_F}{nT} + \frac{192(C_{\nabla g}^2 + \sigma_{\nabla g}^2)C_{\nabla f}^2}{\theta T} + \frac{1920L_f^2 \tilde{C}_{\nabla g}^2 \sigma_g^2}{|\mathcal{B}| \sim T}$ (35) $+\frac{192\theta L_{f}^{2}(\sigma_{\nabla g}^{2}+L_{g}^{2})}{\min\{|\mathcal{B}_{1:}|,|\mathcal{B}_{2:}|\}}+\frac{7680\gamma\sigma_{g}^{2}L_{f}^{2}\tilde{C}_{\nabla g}^{2}}{\min\{|\mathcal{B}_{1:}|,|\mathcal{B}_{2:}|\}}+\frac{96\theta\beta^{2}(\sigma_{\nabla h}^{2}+L_{h}^{2})}{\min\{|\mathcal{B}_{c}|,|\mathcal{B}_{k}|\}}+\frac{3840\gamma\beta^{2}\sigma_{h}^{2}\tilde{C}_{\nabla h}^{2}}{|\mathcal{B}_{k}|}.$  $Let \beta = \frac{1}{\epsilon\delta}, \theta = \min\left\{\frac{\epsilon^4\delta^2 \min\{|\mathcal{B}_k|, |\mathcal{B}_c|\}}{672(\sigma_{\nabla h}^2 + L_h^2)}, \frac{\epsilon^2 \min\{|\mathcal{B}|, |\mathcal{B}_{1i}|, |\mathcal{B}_{2i}|\}}{1344L_f^2(\sigma_{\nabla g}^2 + L_g^2)}\right\} = O(\epsilon^4\delta^2),$   $\gamma_1 = \gamma_2 = \gamma \le \min\left\{\frac{5n_0\theta}{3|\mathcal{B}|}, \frac{5m\theta}{3|\mathcal{B}_c|}, \frac{\epsilon^4\delta^2|\mathcal{B}_k|}{26880\sigma_h^2\tilde{C}_{\nabla h}^2}\right\} = O(\epsilon^4\delta^2),$   $\eta = \min\left\{\frac{1}{12(L_F + \beta L_H)}, \frac{\theta}{8\sqrt{3}L_F}, \frac{\theta}{8\sqrt{3}L_H\beta}, \frac{\gamma_1|\mathcal{B}|}{40\sqrt{6}L_g L_f \tilde{C}_{\nabla g}n}, \frac{\gamma_2|\mathcal{B}_c|}{40\sqrt{6}\beta L_h \tilde{C}_{\nabla h}m}\right\} = O(\epsilon^5\delta^3) \text{ and}$   $T = O(\epsilon^{-7}\delta^{-3}) \text{ we here}$  $T = O(\epsilon^{-7}\delta^{-3})$ , we have  $\frac{1}{T}\sum_{t=1}^{T-1} \mathbb{E}\left[\eta^{-2} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 + \|\nabla\Phi(\mathbf{w}^t) - v^t\|^2\right] \le O(\epsilon^2)$ By the definition of  $\mathbf{w}^{t+1}$ , we have  $\mathbf{w}^{t+1} - \mathbf{w}^t + nv^t = 0$  $\Leftrightarrow \eta^{-1}(\mathbf{w}^t - \mathbf{w}^{t+1}) + (\nabla \Phi(\mathbf{w}^t) - v^t) + (\nabla \Phi(\mathbf{w}^{t+1}) - \nabla \Phi(\mathbf{w}^t)) = \nabla \Phi(\mathbf{w}^{t+1})$  $\Leftrightarrow \eta^{-1}(\mathbf{w}^t - \mathbf{w}^{t+1}) + (\nabla \Phi(\mathbf{w}^t) - v^t) + (\nabla \Phi(\mathbf{w}^{t+1}) - \nabla \Phi(\mathbf{w}^t))$  $= \nabla F(\mathbf{w}^{t+1}) + \frac{\beta}{m} \nabla \boldsymbol{h}(\mathbf{w}^{t+1}) [\boldsymbol{h}(\mathbf{w}^{t+1})]_{+}$ 

 $\begin{array}{ll} \begin{array}{ll} & \text{This gives} \\ & \text{1945} \\ & \text{1945} \\ & \text{1946} \\ & \text{1947} \\ & \text{1947} \\ & \text{1948} \\ & \text{1949} \\ & \text{1949} \\ & \text{1950} \\ & \text{1950} \\ & \text{1951} \end{array} \end{array} \xrightarrow{} \begin{array}{l} & \left\| \nabla F(\mathbf{w}^{t+1}) + \frac{\beta}{m} \nabla \boldsymbol{h}(\mathbf{w}^{t+1})^\top [\boldsymbol{h}(\mathbf{w}^{t+1})]_+ \right\|^2 \\ & \text{1950} \\ & \text{1951} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \\ & \text{100} \end{array} \xrightarrow{} \begin{array}{l} & \text{100} \end{array}$ 

$$\leq 6(\eta^{-2} \| \mathbf{w}^t - \mathbf{w}^{t+1} \|^2 + \| \nabla \Phi(\mathbf{w}^t) - v^t \|^2)$$

1953 Therefore, we can achieve that

$$\frac{1}{T}\sum_{\substack{t=0\\w_{t}}}^{T-1} \mathbb{E}\left[\left\|\nabla F(\mathbf{w}^{t+1}) + \frac{\beta}{m}\nabla \boldsymbol{h}(\mathbf{w}^{t+1})[\boldsymbol{h}(\mathbf{w}^{t+1})]_{+}\right\|^{2}\right] \le O(\epsilon^{2})$$
(36)

1957 By Jensen's inequality, we can get

$$\mathbb{E}\left[\left\|\nabla F(\mathbf{w}^{\hat{t}}) + \frac{\beta}{m}\nabla \boldsymbol{h}(\mathbf{w}^{\hat{t}})[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}\right\|\right] \le O(\epsilon),$$
(37)

1960 with  $\hat{t}$  selected uniformly at random from  $\{1, \dots, T\}$ .

Then, with the full rank assumption on the Jacobian, which is  $\|\nabla h(\mathbf{w}^t)[h(\mathbf{w}^t)]_+\| \ge \delta \|[h(\mathbf{w}^t)]_+\|$ as in Assumption 4, we can get

$$\begin{aligned} \|[\boldsymbol{h}(\mathbf{w}^{t+1})]_{+}\|^{2} &\leq \frac{1}{\delta^{2}} \|\nabla \boldsymbol{h}(\mathbf{w}^{t+1})[\boldsymbol{h}(\mathbf{w}^{t+1})]_{+}\|^{2} \\ &= \frac{m^{2}}{\beta^{2}\delta^{2}} \|\nabla F(\mathbf{w}^{t+1}) + \frac{\beta}{m} \nabla \boldsymbol{h}(\mathbf{w}^{t+1})[\boldsymbol{h}(\mathbf{w}^{t+1})]_{+} - \nabla F(\mathbf{w}^{t+1})\|^{2} \\ &\leq \frac{2m^{2}}{\beta^{2}\delta^{2}} \left[ \|\nabla F(\mathbf{w}^{t+1})\|^{2} + \left\|\nabla F(\mathbf{w}^{t+1}) + \frac{\beta}{m} \nabla \boldsymbol{h}(\mathbf{w}^{t+1})[\boldsymbol{h}(\mathbf{w}^{t+1})]_{+}\right\|^{2} \right] \end{aligned}$$
(38)

1971 Taking the average over T, we can get

$$\frac{1972}{1973} = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|[\boldsymbol{h}(\mathbf{w}^{t+1})]_{+}\|^{2} \leq \frac{1}{T} \sum_{t=0}^{T-1} \frac{2m^{2}}{\beta^{2}\delta^{2}} \mathbb{E} \left[ \left\| \nabla F(\mathbf{w}^{t+1}) \right\|^{2} + \left\| \nabla F(\mathbf{w}^{t+1}) + \frac{\beta}{m} \nabla \boldsymbol{h}(\mathbf{w}^{t+1}) [\boldsymbol{h}(\mathbf{w}^{t+1})]_{+} \right\|^{2} \right] \leq O(\epsilon^{2})$$

$$(39)$$

1977 and using  $\lambda = \frac{\beta}{m} [h(\mathbf{w}^{\hat{t}})]_+$ . By Jensen's inequality, we can get 

$$\mathbb{E}\|[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}\| \le O(\epsilon) \tag{40}$$

 $\leq O(\epsilon).$ 

$$\mathbb{E}|\boldsymbol{\lambda}^{\top}[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}| = \mathbb{E}\left|\frac{\beta}{m}[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}^{\top}[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}\right| = \frac{\beta}{m}\mathbb{E}\|[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}\|^{2}$$
$$= \frac{1}{m\delta\epsilon}\mathbb{E}\|[\boldsymbol{h}(\mathbf{w}^{\hat{t}})]_{+}\|^{2}$$
(41)