

# PARETO MANIFOLD LEARNING: TACKLING MULTIPLE TASKS VIA ENSEMBLES OF SINGLE-TASK MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1        In Multi-Task Learning, tasks may compete and limit the performance achieved  
 2        on each other rather than guiding the optimization trajectory to a common solu-  
 3        tion, superior to its single-task counterparts. There is often not a single solution  
 4        that is optimal for all tasks, leading practitioners to balance tradeoffs between  
 5        tasks’ performance, and to resort to optimality in the Pareto sense. Current Multi-  
 6        Task Learning methodologies either completely neglect this aspect of functional  
 7        diversity, and produce one solution in the Pareto Front predefined by their op-  
 8        timization schemes, or produce diverse but discrete solutions, each requiring a  
 9        separate training run. In this paper, we conjecture that there exist Pareto Sub-  
 10       spaces, i.e., weight subspaces where multiple optimal functional solutions lie. We  
 11       propose *Pareto Manifold Learning*, an ensembling method in weight space that is  
 12       able to discover such a parameterization and produces a continuous Pareto Front  
 13       in a single training run, allowing practitioners to modulate the performance on  
 14       each task during inference on the fly. We validate the proposed method on a di-  
 15       verse set of multi-task learning benchmarks, ranging from image classification to  
 16       tabular datasets and scene understanding, and show that *Pareto Manifold Learning*  
 17       outperforms state-of-the-art algorithms.

[All reviewers]: We have added "notes" such as this one to make the changes more visible. For minor things such as slightly changing plot size, we do not write anything.

## 18 1 INTRODUCTION

19       In Multi-Task Learning (MTL), multiple tasks are learned concurrently within a single model, striv-  
 20       ing towards infusing inductive bias that will help outperform the single-task baselines. Apart from  
 21       the promise of superior performance and some theoretical benefits (Ruder, 2017), such as generaliza-  
 22       tion properties for the learned representation, modeling multiple tasks jointly has practical benefits  
 23       as well, e.g., lower inference times and memory requirements. However, building machine learning  
 24       models presents a multifaceted host of decisions for multiple and often competing objectives, such  
 25       as model complexity, runtime and generalization. Conflicts arise since optimizing for one metric of-  
 26       ten leads to the deterioration of other(s). A single solution satisfying optimally all objectives rarely  
 27       exists and practitioners must balance the inherent trade-offs.

28       In contrary to single-task learning, where one metric governs the comparison between methods (e.g.,  
 29       top-1 accuracy in ImageNet), multiple models can be optimal in Multi-Task Learning; e.g., model  
 30       X yields superior performance on task  $\mathcal{A}$  compared to model Y, but the reverse holds true for task  $\mathcal{B}$ ;  
 31       thus, there is not a single better model among the two. This notion of tradeoffs is formally defined  
 32       as *Pareto optimality*. Intuitively, improvement on an individual task performance can come only at  
 33       the expense of another task. However, there exists no framework addressing the need for efficient  
 34       construction of the Pareto Front, i.e., the set of all Pareto optimal solutions.

35       Recent methods in Multi-Task Learning casted the problem in the lens of multi-objective optimiza-  
 36       tion and introduced the concept of Pareto optimality, resulting in different mechanisms for comput-  
 37       ing the descent direction for the shared parameters. Specifically, Sener & Koltun (2018) produce a  
 38       single solution that lies on the Pareto Front. As an optimization scheme, however, it is biased to-  
 39       wards the task with the smallest gradient magnitude, as argued in Liu et al. (2020). Lin et al. (2019)  
 40       expand this idea and, by imposing additional constraints on the objective space to produce multiple  
 41       solutions on the Pareto Front, each corresponding to a different user-specified tradeoff. Finally, the  
 42       work by Ma et al. (2020) proposes an orthogonal approach that can be applied after training and  
 43       starts with a discrete solution set and produces a continuous set (in weight space) around each so-

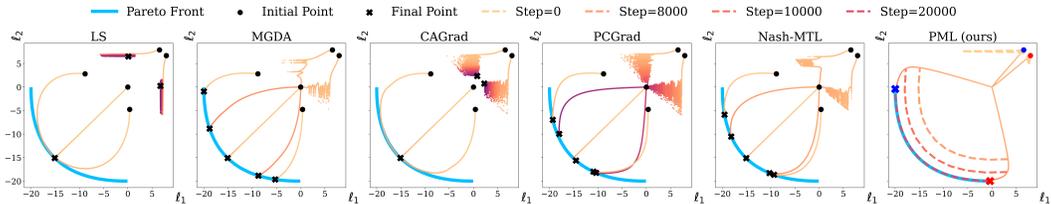


Figure 1: Illustrative example following Yu et al. (2020); Navon et al. (2022). We present the optimization trajectories in loss space starting from different initializations (black bullets) leading to final points (crosses). Color reflects the iteration number when the corresponding value is achieved. To highlight that our method (PML) deals in pairs of models, we use blue and red to differentiate them. Dashed lines show intermediate results of the discovered subspace. While baselines may not reach the Pareto Front or display bias towards specific solutions, PML discovers the entire Pareto Front *in a single run* and shows superior functional diversity.

44 lution, while the overall Pareto Front is continuous *only in objective space* as the union of the local  
 45 (weight-space continuous) Pareto Fronts. As a consequence, the memory requirements grow linearly  
 46 with the number of models stored. Navon et al. (2021); Lin et al. (2021) use hypernetworks to pro-  
 47 duce a Pareto Front in a single training run, but this approach has limited scalability and introduces  
 48 additional design choices.

49 In this paper, we conjecture that we can actually produce a subspace with multiple Pareto stationary  
 50 points in the Multi-Task Learning setting with the hypothesis that local optima (produced by  
 51 different runs or sharing training steps) can be found in close proximity and are connected by  
 52 simple paths. This is motivated by the recent advancements in single task machine learning that  
 53 have explored the geometry of the loss landscape and shown experimentally that local optima are  
 54 connected by simple paths, even linear ones in some cases (Wortsman et al., 2021; Garipov et al.,  
 55 2018; Frankle et al., 2020; Draxler et al., 2018). We assume that, when the problem has multiple  
 56 objectives, it acquires a new dimension relating to the number of tasks. Concretely, there are  
 57 multiple loss landscapes and a solution that satisfies users’ performance requirements must lie in  
 58 the intersection of low loss valleys (for all tasks).

59 Building upon our conjecture, we develop a novel method, *Pareto Manifold Learning*, which casts  
 60 Multi-Task problems as learning an ensemble of single-task predictors by interpolating among (en-  
 61 semble) members during training. By operating in the convex hull of the members’ weight space,  
 62 each single-task model infuses and benefits from representational knowledge to and from the other  
 63 members. During training, the losses are weighted in tandem with the interpolation, i.e., a mono-  
 64 tonic relationship is imposed between the degree of a single-task predictor participation and the  
 65 weight of the corresponding task loss. Consequently, the ensemble as a whole engenders a (weight)  
 66 subspace that explicitly encodes tradeoffs and results in a continuous parameterization of the Pareto  
 67 Front. We identify challenges in guiding the ensemble to such subspaces, designated *Pareto sub-*  
 68 *spaces*, and propose solutions regarding balancing the loss contributions, and regularizing the Pareto  
 69 properties of the subspaces and adapting the interpolation sampling distribution.

70 Experimental results validate that the proposed method is able to discover *Pareto Subspaces*, and out-  
 71 performs baselines on multiple benchmarks. Our training scheme offers two main advantages. First,  
 72 enforcing low loss for all tasks on a linear subspace implicitly penalizes curvature, which has been  
 73 linked to generalization (Chaudhari et al., 2017), benefitting all tasks’ performance. Second, the al-  
 74 gorithm produces a subspace of Pareto Optimal solutions, rather than a single model, enabling prac-  
 75 titioners to handpick during inference the solution that offers the tradeoff that best suits their needs.

76 **2 RELATED WORK**

77 **Multi-Task Learning.** Learning multiple tasks in the Deep Learning setting (Ruder, 2017; Craw-  
 78 shaw, 2020) is usually approached by architectural methodologies (Misra et al., 2016; Ruder et al.,  
 79 2019), where the architectural modules are combined in several layers to govern the joint repre-  
 80 sentation learning, or optimization approaches (Cipolla et al., 2018; Chen et al., 2018), where the  
 81 architecture is standardized to be an encoder-decoder(s), for learning the joint and task-specific rep-

[All reviewers]: Added a short clarification about related work. More details are also provided in “Related work” section.

82 resentations, respectively, and the focus shifts to the descent direction for the shared parameters.  
 83 We focus on the more general track of optimization methodologies fixing the architectural struc-  
 84 ture to Shared-Bottom (Caruana, 1997). The various approaches focus on finding a suitable descent  
 85 direction for the shared parameters. The optimization methods can be broadly categorized into *loss-*  
 86 *balancing* and *gradient-balancing* (Liu et al., 2020). For the former, the goal is to compute an  
 87 appropriate weighting scheme for the losses, e.g., the losses can be weighted via task-dependent ho-  
 88 moscedastic uncertainty (Cipolla et al., 2018), by enforcing task gradient magnitudes to have close  
 89 norms (Chen et al., 2018). The latter class of methodologies manipulate the gradients so that they  
 90 satisfy certain conditions; projecting the gradient of a (random) task on the normal plane of another  
 91 so that gradient conflict is avoided (Yu et al., 2020), enforcing the common descent direction to  
 92 have equal projections for all task gradients (Liu et al., 2020), casting the gradient combination as a  
 93 bargaining game (Navon et al., 2022).

[All reviewers]: Removed discussion about Sener & Koltun (2018) from this point, since the paper is also discussed in the introduction and the next paragraph.

94 **Multi-Task Learning for Pareto Optimality.** The authors in (Sener & Koltun, 2018) were the first  
 95 to view the search for a common descent direction under the Pareto optimality prism and employ  
 96 the Multiple Gradient Descent Algorithm (MGDA) (Désidéri, 2012) in the Deep Learning context.  
 97 However, MGDA did not account for task preferences and the solutions yielded for various initial-  
 98 izations in a synthetic example resulted in similar points in the Pareto Front (Lin et al., 2019). By  
 99 solving a slightly different formulation of the multi-objective problem, they are able to systemati-  
 100 cally introduce task trade-offs and produce a *discrete* Pareto Front. However, this approach requires  
 101 as many training runs as the stated preference combinations and the optimization process for each  
 102 training step of each run introduces a non-negligible overhead. The work in (Ma et al., 2020) pro-  
 103 poses an orthogonal approach for Pareto stationary points; after a model is fitted with any Multi-Task  
 104 Learning method and has converged to a point (seed) in parameter space, a separate phase seeks  
 105 other Pareto stationary points in the vicinity of the seed. The convex hull of these points is guar-  
 106 anteed to lie in the Pareto Front. But training still needs to occur for every seed point, the separate  
 107 phase overhead grows linearly with the number of additional models, and the Pareto Front is not  
 108 continuous across seed points in *parameter space*. Navon et al. (2021) and Lin et al. (2021) employ  
 109 hypernetworks to continuously approximate the Pareto Front in a single run, which introduces ad-  
 110 ditional design choices. Ruchte & Grabocka (2021) address the scalability issues of hypernetworks  
 111 by augmenting the feature space with the preference vector. Raychaudhuri et al. (2022) employ a  
 112 second hypernetwork to also modulate the architecture of the target network addressing.

[All reviewers]: Added prior work.

113 **Ensemble Learning and Mode Connectivity.** Apart from Multi-Task Learning, our algorithm is  
 114 methodologically tied to prior work in the geometry of the neural network optimization landscapes.  
 115 The authors in (Garipov et al., 2018; Draxler et al., 2018) independently and concurrently showed  
 116 that for two local optima  $\theta_1^*, \theta_2^*$  produced by separate training runs (but same initializations) there  
 117 exist nonlinear paths, defined as *connectors* by Wortsman et al. (2021), where the loss remains low.  
 118 The connectivity paths can be extended to include linear in the case of the training runs sharing some  
 119 part of the optimization trajectory (Frankle et al., 2020). These findings can be leveraged to train a  
 120 neural network subspace by enforcing linear connectivity among the subspace endpoints (Wortsman  
 121 et al., 2021). Appendix J discusses more related work regarding ensemble learning and flat minima.

### 122 3 PROBLEM FORMULATION

123 **Notation.** We use bold font for vectors  $\mathbf{x}$ , capital bold for matrices  $\mathbf{X}$  and regular font for scalars  
 124  $x$ .  $T$  is the number of tasks and  $m$  is the number of ensemble members. Each task  $t \in [T]$  has a loss  
 125  $\mathcal{L}_t$ . The overall multi-task loss is  $\mathbf{L} = [\mathcal{L}_1, \dots, \mathcal{L}_T]^\top$ .  $\mathbf{w} \in \Delta_T \subset \mathbb{R}^T$  is the weighting scheme for  
 126 the tasks, i.e., the overall loss is calculated as  $\mathcal{L} = \mathbf{w}^\top \mathbf{L} = \sum_{t=1}^T \alpha_t \mathcal{L}_t$ . Each member  $k \in [m]$  is  
 127 associated with parameters  $\theta_k \in \mathbb{R}^N$  and weighting  $\mathbf{w} \in \Delta_T$ .

128 **Preliminaries.** Our goal lies in solving an unconstrained vector optimization problem of minimiz-  
 129 ing  $\mathbf{L}(\mathbf{y}, \hat{\mathbf{y}}) = [\mathcal{L}_1(y_1, \hat{y}_1), \dots, \mathcal{L}_T(y_T, \hat{y}_T)]^\top$ , where  $\mathcal{L}_i$  corresponds to the objective function for  
 130 the  $i^{\text{th}}$  task, e.g., cross-entropy loss in case of classification. Constructing an optimal solution for all  
 131 tasks is often unattainable due to competing objectives. Hence, an alternative notion of optimality  
 132 is used, as described in Definition 1.

133 **Definition 1** (Pareto Optimality). *A point  $\mathbf{x}$  dominates a point  $\mathbf{y}$  if  $\mathcal{L}_t(\mathbf{x}) \leq \mathcal{L}_t(\mathbf{y})$  for all tasks*  
 134  *$t \in [T]$  and  $\mathbf{L}(\mathbf{x}) \neq \mathbf{L}(\mathbf{y})$ . Then, a point  $\mathbf{x}$  is called Pareto optimal if there exists no point  $\mathbf{y}$  that*  
 135 *dominates it. The set of Pareto optimal points forms the Pareto front  $\mathcal{P}_{\mathbf{L}}$ .*

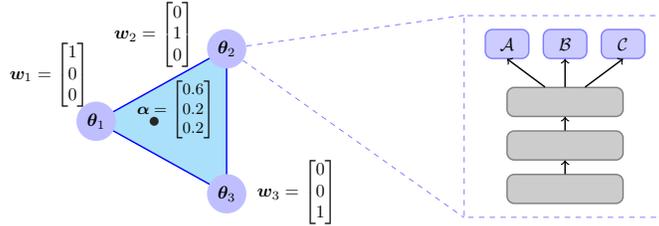


Figure 2: A representation of the encoding in parameter space for  $T = 3$  tasks. Each node corresponds to a tuple of parameters and weighting scheme  $(\theta_v, \mathbf{w}_v) \in \mathbb{R}^N \times \Delta_T$ . The blue dashed frame shows the model, e.g., shared-bottom architecture, implemented by the parameters  $\theta_v$  of each node. For each training step, we sample  $\alpha \in \Delta_T$  and construct the weight combination  $\boldsymbol{\theta} = \alpha^\top \Theta = 0.6 \cdot \theta_1 + 0.2 \cdot \theta_2 + 0.2 \cdot \theta_3$ .

The vector loss function is scalarized by the vector  $\mathbf{w} \in [0, 1]^T$  to form the overall objective  $\mathbf{w}^\top \mathbf{L}$ . Without loss of generality, we assume that  $\mathbf{w}$  lies in the  $T$ -dimensional simplex  $\Delta_T$  by imposing the constraint  $\|\mathbf{w}\| = \sum_{t=1}^T w_t = 1$ . This formulation permits to think of the vector of weights as an encoding of task preferences, e.g., for two tasks letting  $\mathbf{w} = [0.8, 0.2]$  results in attaching more importance to the first task. Overall, the Multi-Task Learning problem can be formulated within the Empirical Risk Minimization (ERM) framework for preference vector  $\mathbf{w}$  and dataset  $\mathcal{D} = \{(x, \mathbf{y})\}_{i=1}$  as:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(x, \mathbf{y}) \sim \mathcal{D}} [\mathbf{L}(\mathbf{y}, \mathbf{f}(x; \boldsymbol{\theta}))] \quad (1)$$

136 Our overall goal is to discover a low-dimensional parameterization in weight space that yields a  
137 (continuous) Pareto Front in functional space. This desideratum leads us to the following definition:

138 **Definition 2** (Pareto Subspace). *Let  $T$  be the number of tasks,  $\mathcal{X}$  the input space,  $\mathcal{Y}$  the multi-*  
139 *task output space,  $\mathcal{R} \subset \mathbb{R}^N$  the parameter space,  $f : \mathcal{X} \times \mathcal{R} \rightarrow \mathcal{Y}$  the function imple-*  
140 *mented by a neural network, and  $\mathbf{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{>0}^T$  be the vector loss. Let  $\{\theta_t \in \mathcal{R} :$   
141  $t \in [T]\}$  be a collection of network parameters and  $\mathcal{S}$  the corresponding convex envelope, i.e.,  
142  $\mathcal{S} = \left\{ \sum_{t=1}^T \alpha_t \theta_t : \sum_{t=1}^T \alpha_t = 1 \text{ and } \alpha_t \geq 0, \forall t \right\}$ . Consider the dataset  $\mathcal{D} = (\mathcal{D}_X, \mathcal{D}_Y)$ . Then,  
143 the subspace  $\mathcal{S}$  is called Pareto if its mapping to functional space via the network architecture  $f$   
144 forms a Pareto Front  $\mathcal{P} = \mathbf{L}(f(\mathcal{D}_X; \mathcal{S}), \mathcal{D}_Y) = \{\mathbf{l} : \mathbf{l} = \mathbf{L}(f(\mathcal{D}_X; \boldsymbol{\theta}), \mathcal{D}_Y), \forall \boldsymbol{\theta} \in \mathcal{S}\}$ .*

## 145 4 METHOD

146 We seek to find a collection of  $m$  neural network models, of identical architecture, whose linear  
147 combination in *weight space* forms a continuous Pareto Front in *objective space*. Model  $i$  corre-  
148 sponds to a tuple of network parameters  $\theta_i$  and task weighting  $w_i$  and implements the function  
149  $\mathbf{f}(\cdot; \theta_i)$ . We impose connectivity among models by modeling the subspace in the convex hull of the  
150 ensemble members. Section 4.1 presents the core of the algorithm, and in Section 4.2 we discuss  
151 various improvements that address Multi-Task Learning challenges.

### 152 4.1 PARETO MANIFOLD LEARNING

Let  $\Theta = [\theta_1, \theta_2, \dots, \theta_m]^\top$  be an  $m \times N$  matrix storing the parameters of all models,  $\mathbf{W} =$   
 $[\mathbf{w}_1, \dots, \mathbf{w}_m]^\top$  be a  $m \times T$  matrix storing the task weighting of ensemble members. By designing  
the subspace as a simplex, the objective now becomes:

$$\mathbb{E}_{(x, \mathbf{y}) \sim \mathcal{D}} [\mathbb{E}_{\alpha \sim \mathcal{P}} [\alpha^\top \mathbf{W} \mathbf{L}(\mathbf{y}, \mathbf{f}(x; \alpha \Theta))]] \quad (2)$$

153 where  $\mathcal{P}$  is the sampling distribution placed upon the simplex. In the case where the en-  
154 semble members are single-task predictors ( $\mathbf{w}$  is one-hot) and the number of tasks coin-  
155 cides with the number of ensemble members ( $m = T$ ), the matrix of task weightings  $\mathbf{W}$   
156 is an identity matrix and Equation 2 simplifies to  $\mathbb{E}_{(x, \mathbf{y}) \sim \mathcal{D}} [\mathbb{E}_{\alpha \sim \mathcal{P}} [\alpha^\top \mathbf{L}(\mathbf{y}, \mathbf{f}(x; \alpha \Theta))]] =$   
157  $\mathbb{E}_{(x, \mathbf{y}) \sim \mathcal{D}} [\mathbb{E}_{\alpha \sim \mathcal{P}} [\sum_{t=1}^T \alpha_t \mathcal{L}_t(\mathbf{y}, \mathbf{f}(x; \sum_{t=1}^T \alpha_t \theta_t))]]$ . By using the same weighting for both



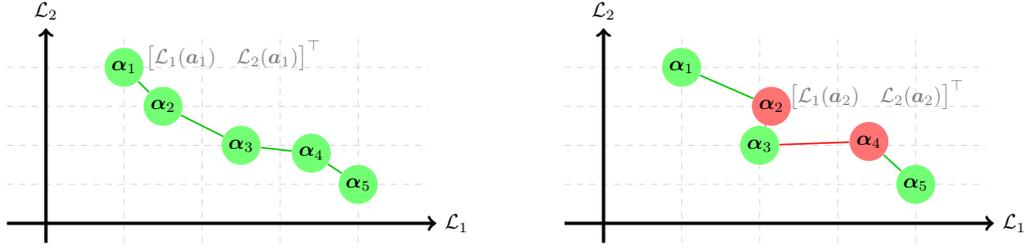


Figure 3: Visual explanation of multiforward regularization, presented in Equation 3. The subfigures depict the loss values for various weightings  $\alpha_i = [\alpha_{i,1}, \alpha_{i,2}]$ . Optimal lies in the origin. We assume that  $\alpha_{1,1} > \dots > \alpha_{5,1}$ . Green color corresponds to Pareto optimality. (Left) all sampled weightings are in the Pareto Front and the regularization term is zero. (Right) The red points are not optimal and, therefore, the regularization term penalizes the violations of the monotonicity constraints for the appropriate task loss:  $\alpha_2$  and  $\alpha_4$  violate the  $\mathcal{L}_1$  and  $\mathcal{L}_2$  orderings w.r.t.  $\alpha_3$ , since  $\alpha_{2,1} > \alpha_{3,1} \not\Rightarrow \mathcal{L}_1(\alpha_2) < \mathcal{L}_1(\alpha_3)$  and  $\alpha_{4,2} > \alpha_{3,2} \not\Rightarrow \mathcal{L}_2(\alpha_4) < \mathcal{L}_2(\alpha_3)$ .

191  $t \in [T]$  w.r.t. the shared parameters. Previously, the update rule occurred with the overall gradient  
 192  $\mathbf{g}_{total} = \alpha^\top \mathbf{G} = \alpha^\top [\mathbf{g}_1 \quad \dots \quad \mathbf{g}_T]$ . We impose a unit  $\ell_2$ -norm for gradients and perform the  
 193 update with  $\tilde{\mathbf{g}}_{total} = \alpha^\top \tilde{\mathbf{G}} = \alpha^\top [\tilde{\mathbf{g}}_1 \quad \dots \quad \tilde{\mathbf{g}}_T]$  where  $\tilde{\mathbf{g}}_t = \frac{\mathbf{g}_t}{\|\mathbf{g}_t\|_2}$ .

194 **Improving stability by Multi-Forward batch regularization.** Consider two different weightings  
 195  $\alpha_1$  and  $\alpha_2 \in \Delta_{T-1}$ . Without loss of generality  $[\alpha_1]_0 = \alpha_1 > [\alpha_2]_0 = \alpha_2$ . Then, ideally, the  
 196 interpolated model closer to the ensemble member for task 1 has the lowest loss on that task, i.e., we  
 197 would want the ordering  $\mathcal{L}_1(\alpha_1) < \mathcal{L}_1(\alpha_2)$ , and, equivalently for the other tasks. Furthermore, if  
 198  $\alpha = [1 - \epsilon, \epsilon/T-1, \dots, \epsilon/T-1]$ , only one member essentially reaps the benefits of the gradient update  
 199 and moves the ensemble towards weight configurations more suitable for one task but, perhaps deleterious  
 200 for the remaining ones. Thus, we propose repeating the forward pass  $W$  times for different  
 201 random weightings  $\{\alpha_i\}_{i \in [W]}$ , allowing the advancement of all ensemble members concurrently in  
 202 a coordinated way. By performing multiple forward passes for various weightings, we achieve a  
 203 lower discrepancy sequence and reduce the variance of such pernicious updates.

We also include a regularization term, which penalizes the wrong orderings and encourages the subspace to have Pareto properties. Let  $\mathcal{V}$  be the set of interpolation weighs sampled in the current batch  $\mathcal{V} = \{\alpha_w = (\alpha_{w,1}, \alpha_{w,2}, \dots, \alpha_{w,T}) \in \Delta_{T-1}\}_{w \in [W]}$ . Then each task defines the *directed* graph  $\mathcal{G}_t = (\mathcal{V}, \mathcal{E}_t)$  where  $\mathcal{E}_t = \{(\alpha_i, \alpha_j) \in \mathcal{V} \times \mathcal{V} : \alpha_{i,t} < \alpha_{j,t}\}$ . The overall loss becomes:

$$\mathcal{L}_{total} = \sum_{i=1}^W \alpha_i^\top \mathbf{L}(\alpha_i) + \lambda \cdot \sum_{t=1}^T \log \left( \frac{1}{|\mathcal{E}_t|} \sum_{(\alpha_i, \alpha_j) \in \mathcal{E}_t} e^{[\mathcal{L}_t(\alpha_i) - \mathcal{L}_t(\alpha_j)]_+} \right) \quad (3)$$

204 The current formulation of the edge set penalizes heavily the connections from vertices with low  
 205 values. For this reason, we only keep one outgoing edge per node, defined by the task lexicographic  
 206 order, resulting in the graph  $\mathcal{G}_t^{\text{LEX}} = (\mathcal{V}, \mathcal{E}_t^{\text{LEX}})$  and  $|\mathcal{E}_t^{\text{LEX}}| = W - 1, \forall t \in [T]$ . Note that the regu-  
 207 larization term is convex as the sum of *log-sum-exp* terms. If no violations occur, the regularization  
 208 term is zero. Figure 3 offers a visual explanation of the regularization term.

209 **The role of sampling.** Another component of Algorithm 1 is the sampling imposed on the convex  
 210 hull parameterization. During training, the sampling distribution dictates the loss weighting used  
 211 and, hence, modulates the degree of task learning. A natural choice is the Dirichlet distribution  
 212  $\text{Dir}(\mathbf{p})$  where  $\mathbf{p} \in \mathbb{R}_{>0}^T$  are the concentration parameters, since its support is the  $T$ -dimensional  
 213 simplex  $\Delta_T$ . For  $\mathbf{p} = p\mathbf{1}_T$ , the distribution is symmetric; for  $p < 1$  the sampling is more  
 214 concentrated near the ensemble members, for  $p > 1$  it is near the centre and for  $p = 1$  it corresponds  
 215 to the uniform distribution. In contrast, for  $p_1 \neq p_2$  the distribution is skewed. In our experiments,  
 216 we use symmetric Dirichlet distributions with  $p \geq 1$  to guide the ensemble to representations best  
 217 suited for Multi-Task Learning.

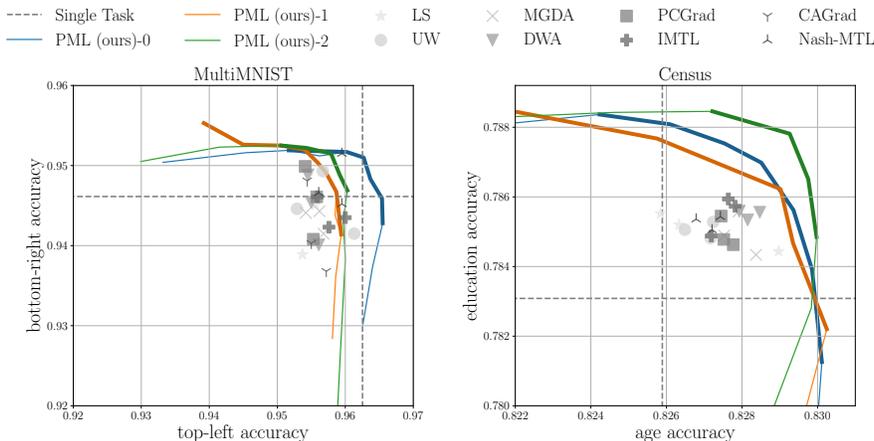


Figure 4: Experimental results on MultiMNIST and Census. Top right is optimal. Three random seeds per method. Solid lines correspond to our method (PML) and thick lines to the Pareto Front. We have used a different color for each seed of PML. Baselines are shown in shades of gray: scatter plot for MTL baselines and dashed lines for single task. In both datasets, Pareto Manifold Learning discovers subspaces with diverse and Pareto-optimal solutions and outperforms the baselines.

218 5 EXPERIMENTS

219 We evaluate our method on several datasets, such as MultiMNIST, Census, MultiMNIST-3,  
 220 UTKFace and CityScapes, and various architectures, ranging from MultiLayer Perceptrons  
 221 (MLPs) to Convolutional Neural Networks (CNNs) and Residual Networks (ResNets). Each  
 222 ensemble member is initialized independently. In all experiments, the learning rate for our method  
 223 is  $m$ -fold the learning rate of the baselines to counteract the fact that the backpropagation step  
 224 scales the gradients by  $m^{-1}$  in expectation. The detailed settings used for each dataset and  
 225 additional experiments are provided in the appendix. Our overarching objective is to construct  
 226 continuous weight subspaces which map to Pareto Fronts in the functional space. However, our  
 227 method produces a continuum of results rather than a single point, rendering tabular presentation  
 228 cumbersome. For this reason, (a) for tables we present the best-of-(sampled)-subspace results, (b)  
 229 we experiment on numerous two-task datasets where plots convey the results succinctly, (c) present  
 230 qualitative results on three-task datasets. The source code will be released after the review process.

231 **Baselines.** We explore various algorithms from the literature: 1. Single-Task Learning (STL),  
 232 2. Linear Scalarization (LS) which minimizes the average loss  $\frac{1}{T} \sum_{t=1}^T \mathcal{L}_t$ , 3. Uncertainty  
 233 Weighting (UW, Cipolla et al. 2018), 4. Multiple-gradient descent algorithm (MGDA, Sener &  
 234 Koltun 2018), 5. Dynamic Weight Averaging (DWA, Liu et al. 2019), 6. Projecting Conflicting  
 235 Gradients (PCGrad, Yu et al. 2020), 7. Impartial Multi-Task Learning (IMTL, Liu et al. 2020),  
 236 8. Conflict-Averse Gradient Descent (CAGrad, Liu et al. 2021) and 9. Bargaining Multi-Task  
 237 Learning (Nash-MTL, Navon et al. 2022).

238 5.1 EXPERIMENTS ON DATASETS WITH TWO CLASSIFICATION TASKS

239 In this section, we focus on datasets with two tasks, both classification. This setting allows for rich  
 240 visualizations that we use to draw insights on the inner workings of the algorithms.

241 **MultiMNIST.** We investigate the effectiveness of Pareto Manifold Learning on digit classification  
 242 using a LeNet model with a shared-bottom architecture. The ensemble consists of two members  
 243 with single task weightings. To gauge the performance of the models lying in the linear segment  
 244 between the nodes, we test the performance on the validation set on the ensemble members as well  
 245 as for 9 models uniformly distributed across the edge, resulting in 11 models in total. We use this  
 246 evaluation/plotting scheme throughout the experiments. We ablate the effect of multi-forward training  
 247 on Appendix D; we use a grid search on window  $W \in \{2, 3, 4, 5\}$  and strength  $\lambda \in \{0, 2, 5, 10\}$   
 248 along with the base case of  $(W, \lambda) = (1, 0)$  and present in the main text the setting that achieves the

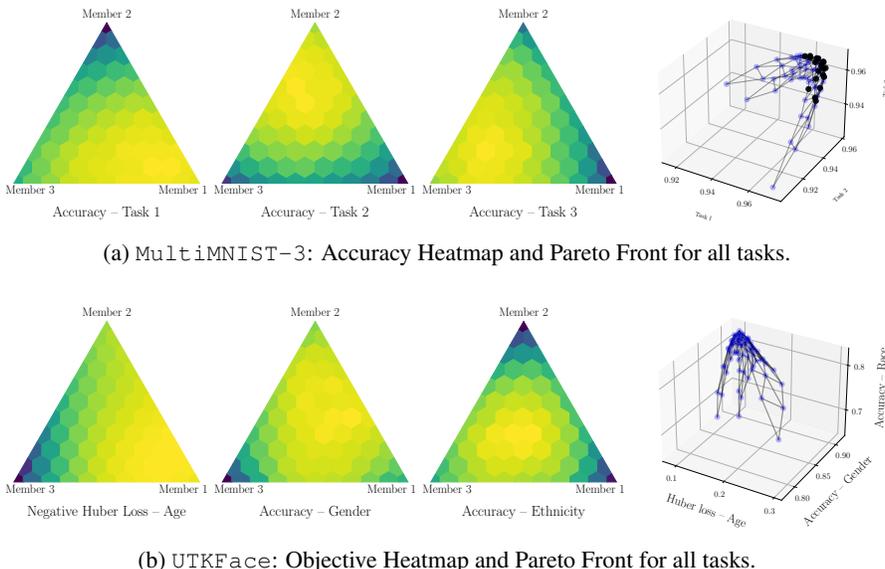


Figure 5: Application of Pareto Manifold Learning on datasets with 3 tasks. Each triangle depicts the performance on a task, using color, as a function of the interpolation weighting, i.e. each hexagon corresponds to a different weighting  $\alpha = [\alpha_1, \alpha_2, \alpha_3] \in \Delta_3$ . The closer the interpolated member is to a single-task predictor, the higher the performance on the corresponding task. The 3D plot, on the right, show the performance of the model in the multi-objective space.

249 highest mean (across seeds) HyperVolume score on the validation set. Figure 4 shows the results on  
 250 MultiMNIST using multi-forward regularization with window  $W = 4$  and strength  $\lambda = 0$ . We  
 251 observe that most baselines are characterized by limited functional diversity; their predefined opti-  
 252 mization schemes lead the differently seeded/initialized training runs to final models with similar  
 253 performance (same markers are clustered in the plots). This lack of functional diversity, as well as  
 254 inability to consistently outperform the Linear Scalarization baseline, are also noted by Kurin et al.  
 255 (2022); Xin et al. (2022). In contrast, all Pareto Manifold Learning seeds find subspaces with diverse  
 256 functional solutions. This statement is quantitatively translated to higher HyperVolume compared  
 257 to the baselines, shown in Table 4 of the appendix, and can be attributed to the observation that  
 258 Equation 2 generalizes the Linear Scalarization method.

[Reviewer NiGo]: added refer-  
 ences [5,6] from the review.

259 **Census.** We explore the method on the tabular dataset Census (Kohavi, 1996) using a Multi-  
 260 Layer Perceptron. We focus on the task combination of predicting age and education level, similar  
 261 to Ma et al. (2020). We perform the same ablation study as before and present the results on Figure 4  
 262 for the best setting ( $W = 3$  and  $\lambda = 10$ ). In the case of MultiMNIST, there exists symmetry  
 263 between the tasks, both digits are drawn from the same distribution and placed in the pixel grid in  
 264 a symmetric way, resulting in equal pace learning. However, in the case of Census, tasks differ in  
 265 statistics and, yet, the proposed method recovers a Pareto subspace with diverse solutions.

266 5.2 BEYOND PAIRS OF CLASSIFICATION TASKS: MULTIMNIST-3 AND UTKFACE

267 We expand the experimental validation to triplets of tasks, consider regression and more complex ar-  
 268 chitectures, graduating from MLPs and CNNs to ResNets (He et al., 2016). For three tasks, we create  
 269 a 2D grid of equidistant points spanning the three single-task predictors. If  $n$  is the number of inter-  
 270 polated points between two (out of three) members, the grid has  $\binom{n+1}{2}$  points. We use  $n = 11$ , result-  
 271 ing in 66 points. For visual purposes, neighboring points are connected. For three tasks, it would be  
 272 visually cluttering to present the discovered subspaces with multiple seeds and baselines. Hence, we  
 273 opt for a more qualitative discussion in this section and present quantitative findings in the appendix.

274 **MultiMNIST-3.** First, we construct an equivalent of MultiMNIST for 3 tasks. Digits are placed  
 275 on top-left, top-right and bottom-centre. Figure 5a shows the results on MultiMNIST-3. As  
 276 argued previously, MNIST variants are characterized by task symmetry and Figure 5a reflects this.  
 277 For this reason, we do not employ any balancing scheme. The 3D plot in conjunction with the

Table 1: Test performance on *CityScapes*. 3 random seeds per method. For Pareto Manifold Learning, we report the mean (across seeds) best results from the final subspace.

	Segmentation		Depth	
	mIoU $\uparrow$	Pix Acc $\uparrow$	Abs Err $\downarrow$	Rel Err $\downarrow$
STL	71.79	92.60	0.0135	32.786
LS	70.94	92.29	0.0192	117.658
UW	70.97	92.24	0.0188	118.168
MGDA	69.23	91.77	0.0138	51.986
DWA	70.87	92.23	0.0190	113.565
PCGrad	71.14	92.32	0.0185	117.797
IMTL	71.54	92.47	0.0151	65.058
CAGrad	70.23	92.06	0.0173	100.162
Nash-MTL	72.07	92.61	0.0148	62.980
PML (ours)	70.28	91.94	0.0140	52.559

278 simplices reveal that the method has the effect of gradual transfer of learned representation from one  
 279 member to the other, and offers a succinct visual confirmation of [Claim 3](#).

280 **UTKFace.** The UTKFace dataset ([Zhang et al., 2017](#)) has more than 20,000 face images and  
 281 three tasks: predicting age (modeled as regression using Huber loss - similar to ([Ma et al., 2020](#))),  
 282 classifying gender and ethnicity. The introduction of a regression task implies that losses have vastly  
 283 different scales, which dictates the use of balancing schemes, as discussed in [Section 4.2](#). We apply  
 284 the proposed gradient-balancing scheme and present the results in [Figure 5b](#). For visual unity and to  
 285 remain in the theme of “higher is better”, the *negative* Huber loss is plotted. Despite the increased  
 286 complexity and the existence of a regression task, the proposed method discovers a *Pareto Subspace*.  
 287 Additional experiments and qualitative results are provided in [Appendix G](#).

### 288 5.3 SCENE UNDERSTANDING

289 We also explore the applicability of Pareto Manifold Learning for *CityScapes* ([Cordts et al.,](#)  
 290 [2016](#)), a scene understanding dataset containing high-resolution images of urban street scenes. Our  
 291 experimental configuration is drawn from [Liu et al. \(2019\)](#); [Yu et al. \(2020\)](#); [Liu et al. \(2021\)](#); [Navon](#)  
 292 [et al. \(2022\)](#) with some modifications. Concretely, we address two tasks: semantic segmentation and  
 293 depth regression. We use a SegNet architecture ([Badrinarayanan et al., 2017](#)) trained for 100 epochs  
 294 with Adam optimizer ([Kingma & Ba, 2015](#)) of initial learning rate  $10^{-4}$ , which is halved after 75  
 295 epochs. The images are resized to  $128 \times 256$  pixels. In the initial training steps any sampling  
 296  $\alpha$  results in a random model, due to initialization, and the algorithm has a warmup period until  
 297 the ensemble members have acquired meaningful representations. Hence, to reduce computational  
 298 overhead and help convergence, the concentration parameter of the Dirichlet distribution is set to  
 299  $p_0 = 5$ . We use gradient balancing, window  $W = 3$  and  $\lambda = 1$ . The results are presented in [Ta-](#)  
 300 [ble 1](#). In Depth Estimation and out of MTL methods, Pareto Manifold Learning is near-optimal with  
 301 MGDA narrowly better. However, the performance compared to the other algorithms is superior. In  
 302 Semantic Segmentation, our method outperforms MGDA, but is worse than other baselines. **Overall**  
 303 **no multi-task method dominates Pareto Manifold Learning**. It is remarkable that, despite our goal of  
 304 discovering *Pareto subspaces*, the proposed method is on par in performance on Semantic Segmen-  
 305 tation with the state-of-the-art algorithms, and better than the vast majority on Depth Estimation.

[Reviewer qexX]: Added short comment addressing weakness 3.

## 306 6 CONCLUSION

307 In this paper, we proposed a weight-ensembling method tailored to Multi-Task Learning; multiple  
 308 single-task predictors are trained in conjunction to produce a subspace formed by their convex hull,  
 309 and endowed with desirable Pareto properties. We experimentally show on a diverse suite of bench-  
 310 marks that the the proposed method is successful in discovering *Pareto subspaces* and outperforms  
 311 some state-of-the-art MTL methods. An interesting future direction is to perform a hierarchical  
 312 weight ensembling, sharing progressively more of the lower layers, given that the features learned  
 313 at low depth are similar across tasks. An alternative exploration venue is to connect our method to  
 314 the challenge of task affinity ([Fifty et al., 2021](#); [Standley et al., 2020](#)) via a geometrical lens of the  
 315 loss landscape.

## 316 REFERENCES

- 317 Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. SegNet: A Deep Convolutional Encoder-  
318 Decoder Architecture for Image Segmentation. *IEEE transactions on pattern analysis and ma-*  
319 *chine intelligence*, (12):2481–2495, 2017.
- 320 Rich Caruana. Multitask learning. *Machine Learning*, (1):41–75, 1997. URL [https://doi.](https://doi.org/10.1023/A:1007379606734)  
321 [org/10.1023/A:1007379606734](https://doi.org/10.1023/A:1007379606734).
- 322 Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian  
323 Borgs, Jennifer T. Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-SGD: Biasing gra-  
324 dient descent into wide valleys. In *International Conference on Learning Representations*. Open-  
325 Review.net, 2017.
- 326 Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. GradNorm: Gradient  
327 Normalization for Adaptive Loss Balancing in Deep Multitask Networks. In *International Con-*  
328 *ference on Machine Learning*, Proceedings of Machine Learning Research, pp. 793–802. PMLR,  
329 2018.
- 330 Roberto Cipolla, Yarin Gal, and Alex Kendall. Multi-task Learning Using Uncertainty to Weigh  
331 Losses for Scene Geometry and Semantics. In *Conference on Computer Vision and Pattern*  
332 *Recognition*), pp. 7482–7491. IEEE, June 2018.
- 333 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo  
334 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes Dataset for Semantic  
335 Urban Scene Understanding. In *Conference on Computer Vision and Pattern Recognition*), 2016.
- 336 Michael Crawshaw. Multi-Task Learning with Deep Neural Networks: A Survey. *arXiv:2009.09796*  
337 *[cs, stat]*, September 2020.
- 338 Jean-Antoine Désidéri. Multiple-gradient descent algorithm (MGDA) for multiobjective optimiza-  
339 tion. *Comptes Rendus Mathématique*, (5-6):313–318, 2012.
- 340 Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize  
341 for deep nets. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th international*  
342 *conference on machine learning, ICML 2017, sydney, NSW, australia, 6-11 august 2017*, vol-  
343 *ume 70 of Proceedings of machine learning research*, pp. 1019–1028. PMLR, 2017. URL [http:](http://proceedings.mlr.press/v70/dinh17b.html)  
344 [//proceedings.mlr.press/v70/dinh17b.html](http://proceedings.mlr.press/v70/dinh17b.html). tex.bibsource: dblp computer sci-  
345 *ence bibliography*, <https://dblp.org> tex.biburl: [https://dblp.org](https://dblp.org/rec/conf/icml/DinhPBB17.bib)  
346 <https://dblp.org/rec/conf/icml/DinhPBB17.bib>  
tex.timestamp: Wed, 29 May 2019 08:41:45 +0200.
- 347 Felix Draxler, Kambis Veschgini, Manfred Salmhofer, and Fred A. Hamprecht. Essentially No Bar-  
348 riers in Neural Network Energy Landscape. In *International Conference on Machine Learning*,  
349 *Proceedings of Machine Learning Research*, pp. 1308–1317. PMLR, 2018.
- 350 Christopher Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, and Chelsea Finn. Efficiently  
351 identifying task groupings for multi-task learning. In *Advances in Neural Information Processing*  
352 *Systems*, 2021.
- 353 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware Minimiza-  
354 tion for Efficiently Improving Generalization. In *9th International Conference on Learning Rep-*  
355 *resentations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- 356 Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. Linear Mode  
357 Connectivity and the Lottery Ticket Hypothesis. In *International Conference on Machine Learn-*  
358 *ing*, Proceedings of Machine Learning Research, pp. 3259–3269. PMLR, 2020.
- 359 Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P. Vetrov, and Andrew Gordon Wil-  
360 *son*. Loss Surfaces, Mode Connectivity, and Fast Ensembling of DNNs. In *Advances in Neural*  
361 *Information Processing Systems*, pp. 8803–8812, 2018.

- 362 Marton Havasi, Rodolphe Jenatton, Stanislav Fort, Jeremiah Zhe Liu, Jasper Snoek, Balaji Lakshmi-  
363 narayanan, Andrew Mingbo Dai, and Dustin Tran. Training independent subnetworks for robust  
364 prediction. In *9th international conference on learning representations, ICLR 2021, virtual event,  
365 austria, may 3-7, 2021*. OpenReview.net, 2021. URL [https://openreview.net/forum?  
366 id=OGg9XnKxFAH](https://openreview.net/forum?id=OGg9XnKxFAH). tex.bibsource: dblp computer science bibliography, <https://dblp.org>  
367 tex.biburl: <https://dblp.org/rec/conf/iclr/HavasiJFLSLDT21.bib> tex.timestamp: Wed, 23 Jun 2021  
368 17:36:39 +0200.
- 369 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
370 nition. In *Conference on Computer Vision and Pattern Recognition*), pp. 770–778, 2016.
- 371 Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry P. Vetrov, and Andrew Gordon Wil-  
372 son. Averaging Weights Leads to Wider Optima and Better Generalization. In *Proceedings of  
373 the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI 2018, Monterey, Cal-  
374 ifornia, USA, August 6-10, 2018*, pp. 876–885. AUAI Press, 2018. URL [http://auai.org/  
375 uai2018/proceedings/papers/313.pdf](http://auai.org/uai2018/proceedings/papers/313.pdf).
- 376 Yiding Jiang\*, Behnam Neyshabur\*, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic  
377 generalization measures and where to find them. Mar 2020. URL [https://openreview.  
378 net/forum?id=SJgIPJBFvH](https://openreview.net/forum?id=SJgIPJBFvH).
- 379 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *3rd Inter-  
380 national Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9,  
381 2015, Conference Track Proceedings*, 2015.
- 382 Ron Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In *Proceeed-  
383 ings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-  
384 96), Portland, Oregon, USA*, pp. 202–207. AAAI Press, 1996.
- 385 Vitaly Kurin, Alessandro De Palma, Ilya Kostrikov, Shimon Whiteson, and M. Pawan Kumar.  
386 In defense of the unitary scalarization for deep multi-task learning. (arXiv:2201.04122), Oct  
387 2022. doi: 10.48550/arXiv.2201.04122. URL <http://arxiv.org/abs/2201.04122>.  
388 arXiv:2201.04122 [cs].
- 389 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and Scalable Predictive  
390 Uncertainty Estimation using Deep Ensembles. In *Advances in Neural Information Processing  
391 Systems*. Curran Associates, Inc., 2017. URL [https://papers.nips.cc/paper/2017/  
392 hash/9ef2ed4b7fd2c810847ffa5fa85bce38-Abstract.html](https://papers.nips.cc/paper/2017/hash/9ef2ed4b7fd2c810847ffa5fa85bce38-Abstract.html).
- 393 Xi Lin, Hui-Ling Zhen, Zhenhua Li, Qingfu Zhang, and Sam Kwong. Pareto Multi-Task Learning.  
394 In *Advances in Neural Information Processing Systems*, pp. 12037–12047, 2019.
- 395 Xi Lin, Zhiyuan Yang, Qingfu Zhang, and Sam Kwong. Controllable pareto multi-task learning.  
396 (arXiv:2010.06313), Feb 2021. doi: 10.48550/arXiv.2010.06313. URL [http://arxiv.org/  
397 abs/2010.06313](http://arxiv.org/abs/2010.06313). arXiv:2010.06313 [cs, stat].
- 398 Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-Averse Gradient Descent  
399 for Multi-task Learning, October 2021.
- 400 Liyang Liu, Yi Li, Zhanghui Kuang, Jing-Hao Xue, Yimin Chen, Wenming Yang, Qingmin Liao, and  
401 Wayne Zhang. Towards Impartial Multi-task Learning. In *International Conference on Learning  
402 Representations*, September 2020.
- 403 Shikun Liu, Edward Johns, and Andrew J. Davison. End-To-End Multi-Task Learning With Atten-  
404 tion. In *Conference on Computer Vision and Pattern Recognition*), pp. 1871–1880. IEEE, June  
405 2019.
- 406 Pingchuan Ma, Tao Du, and Wojciech Matusik. Efficient continuous pareto exploration in multi-task  
407 learning. In *International Conference on Machine Learning*, pp. 6522–6531. PMLR, 2020.
- 408 Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. Cross-Stitch Networks for  
409 Multi-task Learning. In *Conference on Computer Vision and Pattern Recognition*), pp. 3994–  
410 4003. IEEE, June 2016.

- 411 Aviv Navon, Aviv Shamsian, Ethan Fetaya, and Gal Chechik. Learning the Pareto Front with Hy-  
412 pernetworks. In *9th International Conference on Learning Representations, ICLR 2021, Virtual*  
413 *Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- 414 Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik, and  
415 Ethan Fetaya. Multi-task learning as a bargaining game. In *International Conference on Machine*  
416 *Learning, 2022*.
- 417 Dripta S Raychaudhuri, Yumin Suh, Samuel Schuster, Xiang Yu, Masoud Faraki, Amit K Roy-  
418 Chowdhury, and Manmohan Chandraker. Controllable dynamic multi-task architectures. In *Pro-*  
419 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10955–  
420 10964, 2022.
- 421 Michael Ruchte and Josif Grabocka. Scalable pareto front approximation for deep multi-objective  
422 learning. In *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, 2021.
- 423 Sebastian Ruder. An Overview of Multi-Task Learning in Deep Neural Networks. *arXiv:1706.05098*  
424 *[cs, stat]*, June 2017.
- 425 Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. Latent Multi-Task  
426 Architecture Learning. In *AAAI Conference on Artificial Intelligence*, pp. 4822–4829. AAAI  
427 Press, 2019.
- 428 Ozan Sener and Vladlen Koltun. Multi-Task Learning as Multi-Objective Optimization. In *Ad-*  
429 *vances in Neural Information Processing Systems 31: Annual Conference on Neural Information*  
430 *Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 525–536,  
431 2018.
- 432 Trevor Standley, Amir Roshan Zamir, Dawn Chen, Leonidas J. Guibas, Jitendra Malik, and Silvio  
433 Savarese. Which Tasks Should Be Learned Together in Multi-task Learning? In *Proceedings of*  
434 *the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual*  
435 *Event*, Proceedings of Machine Learning Research, pp. 9120–9132. PMLR, 2020.
- 436 Yeming Wen, Dustin Tran, and Jimmy Ba. Batchensemble: an alternative approach to efficient  
437 ensemble and lifelong learning. In *8th international conference on learning representations,*  
438 *ICLR 2020, addis ababa, ethiopia, april 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=Sk1flyrYDr>. tex.bibsource: dblp computer science bibli-  
439 ography, <https://dblp.org> tex.biburl: <https://dblp.org/rec/conf/iclr/WenTB20>. bib tex.timestamp:  
440 Thu, 07 May 2020 17:11:47 +0200.
- 442 Mitchell Wortsman, Maxwell Horton, Carlos Guestrin, Ali Farhadi, and Mohammad Rastegari.  
443 Learning Neural Network Subspaces. In *International Conference on Machine Learning*, Pro-  
444 ceedings of Machine Learning Research, pp. 11217–11227. PMLR, 2021.
- 445 Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo Lopes,  
446 Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Lud-  
447 wig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy  
448 without increasing inference time. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba  
449 Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International conference on machine learning,*  
450 *ICML 2022, 17-23 july 2022, baltimore, maryland, USA*, volume 162 of *Proceedings of ma-*  
451 *chine learning research*, pp. 23965–23998. PMLR, 2022. URL [https://proceedings.](https://proceedings.mlr.press/v162/wortsman22a.html)  
452 [mlr.press/v162/wortsman22a.html](https://proceedings.mlr.press/v162/wortsman22a.html). tex.bibsource: dblp computer science bibli-  
453 ography, <https://dblp.org> tex.biburl: <https://dblp.org/rec/conf/icml/WortsmanIGRLMNF22>. bib  
454 tex.timestamp: Tue, 12 Jul 2022 17:36:52 +0200.
- 455 Derrick Xin, Behrooz Ghorbani, Ankush Garg, Orhan Firat, and Justin Gilmer. Do current multi-task  
456 optimization methods in deep learning even help? (arXiv:2209.11379), Sep 2022. doi: 10.48550/  
457 arXiv.2209.11379. URL <http://arxiv.org/abs/2209.11379>. arXiv:2209.11379 [cs].
- 458 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.  
459 Gradient surgery for multi-task learning. In *Advances in Neural Information Processing Systems*,  
460 pp. 5824–5836. Curran Associates, Inc., 2020.
- 461 Zhang, Yang Zhifei, Song, and Hairong Qi. Age Progression/Regression by conditional adversarial  
462 autoencoder. In *Conference on Computer Vision and Pattern Recognition*). IEEE, 2017.

463 A APPENDIX OVERVIEW

464 As a reference, we provide the following table of contents solely for the appendix.

465	A. Appendix Overview .....	<a href="#">Appendix A</a>
466	B. Details on experimental configurations .....	<a href="#">Appendix B</a>
467	C. Details on Illustrative Example .....	<a href="#">Appendix C</a>
468	D. Ablation on Multi-Forward Regularization .....	<a href="#">Appendix D</a>
469	E. HyperVolume Analysis .....	<a href="#">Appendix E</a>
470	F. Additional experiments on MultiMNIST-3 .....	<a href="#">Appendix F</a>
471	G. Additional experiments on UTKFace .....	<a href="#">Appendix G</a>
472	H. <b>NEW! Details on sampling</b> .....	<a href="#">Appendix H</a>
473	I. <b>NEW! Connection between Pareto Optimality and multiple valley intersections</b>	<a href="#">Appendix I</a>
474	J. <b>NEW! Additional Related Work</b> .....	<a href="#">Appendix J</a>
475	K. <b>NEW! Additional experiments</b> .....	<a href="#">Appendix K</a>

## 476 B EXPERIMENTAL DETAILS

477 **MultiMNIST** MultiMNIST is a synthetic dataset derived from the samples of MNIST. Since  
 478 there is no publicly available version, we create our own by the following procedure. For each  
 479 MultiMNIST image, we sample (with replacement) two MNIST images (of size  $28 \times 28$ ) and place  
 480 them top-left and bottom-right on a  $36 \times 36$  grid. This grid is then resized to  $28 \times 28$  pixels. The  
 481 procedure is repeated 60000 times, 10000 and 10000 times for training, validation and test datasets.  
 482 We use a LeNet shared-bottom architecture. Specifically, the encoder has two convolutional layers  
 483 with 10 and 20 channels and kernel size of 5 followed by Maxpool and a ReLU nonlinearity each.  
 484 The final layer of the encoder is fully connected producing an embedding with 50 features. The  
 485 decoders are fully connected with two layers, one with 50 features and the output layer has 10. We  
 486 use Adam optimizer Kingma & Ba (2015) with learning rate  $10^{-3}$ , no scheduler and the batch size  
 487 is set to 256. Training lasts 10 epochs.

488 **Census** The original version of the Census (Kohavi, 1996) dataset has one task: predicting  
 489 whether a person’s income exceeds \$50000. The dataset becomes suitable for Multi-Task Learning  
 490 by turning one or several features to tasks (Lin et al., 2019). We focus on the task combination of  
 491 predicting age and education level, similar to Ma et al. (2020). The model has a Multi-Layer Percep-  
 492 tron shared-bottom architecture. The encoder has one layer with 256 neurons, followed by a ReLU  
 493 nonlinearity, and two decoders with 2 output neurons each (since the tasks are binary classification).  
 494 Training lasts 10 epochs. We use Adam optimizer learning rate of  $10^{-3}$ .

495 **MultiMNIST-3** The configuration of MultiMNIST is used. Now, the model has three decoders  
 496 and training lasts 20 epochs.

497 **UTKFace** The UTKFace dataset has more than 20,000 face images of dimensions  $200 \times 200$   
 498 pixels and 3 color channels. The dataset has three tasks: predicting age (modeled as regression using  
 499 Huber loss - similar to (Ma et al., 2020)), classifying gender and ethnicity (modeled as classification  
 500 tasks using Cross-Entropy loss). Images are resized to  $64 \times 64$  pixels, age is normalized and a 80/20  
 501 train/test split is used. We use a shared-bottom architecture; the encoder is a ResNet18 (He et al.,  
 502 2016) model without the last fully connected layer. The decoders (task-specific layers) consist of one  
 503 fully-connected layer, where the output dimensions are 1, 2 and 5 for age (modeled as regression),  
 504 gender (binary classification) and ethnicity (classification with 5 classes). Training lasts 100 epochs,  
 505 batch size is 256 and we use Adam optimizer with a learning rate of  $10^{-3}$ . No scheduler is used.

506 **CityScapes** Our experimental configuration is very similar to prior work, namely Liu et al.  
 507 (2019); Yu et al. (2020); Liu et al. (2021); Navon et al. (2022). All images are resized to  $128 \times 256$ .  
 508 The tasks used are coarse semantic segmentation and depth regression. The task of semantic seg-  
 509 mentation has 7 classes, whereas the original has 19. We use a SegNet architecture (Badrinarayanan  
 510 et al., 2017) and train the model for 100 epochs with Adam optimizer (Kingma & Ba, 2015) of an  
 511 initial learning rate  $10^{-4}$ . We employ a scheduler that halves the learning rate after 75 epochs.

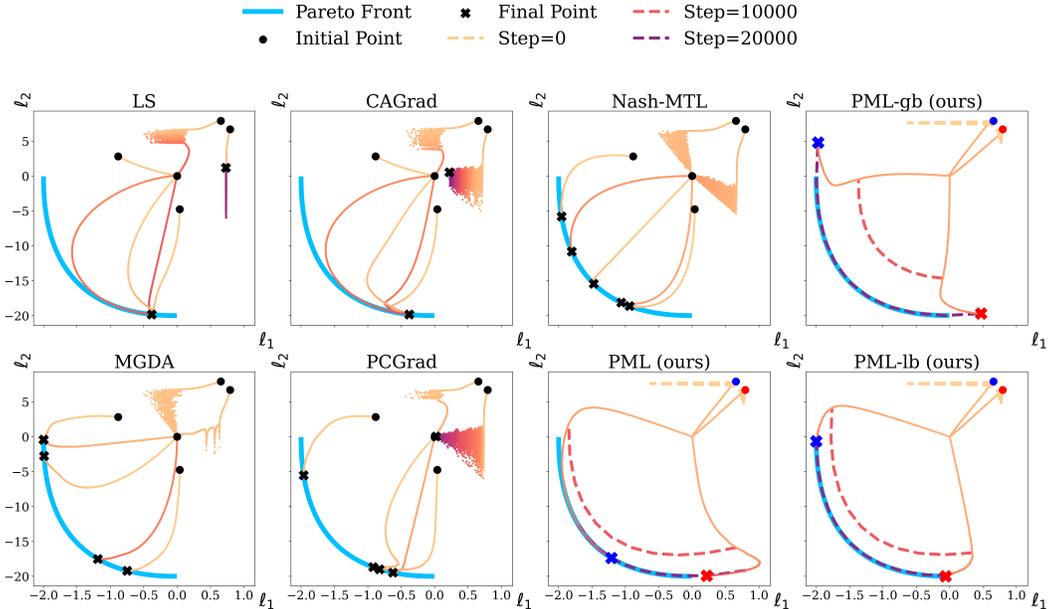


Figure 6: Optimization trajectories in objective space in the case different loss scales. Similar to Figure 1, 5 initializations are shown for baselines and a pair of initializations for Pareto Manifold Learning (PML), in color for clarity. Dashed lines show the evolution of the mapping in loss space for the subspace at the current step. We also show the initial subspace (step= 0). All baselines, except Nash-MTL, and MGDA to a lesser degree, are characterized by trajectories focused on a subset of the Pareto Front, namely minimizing the task with high loss magnitude. The same observation applies to naïvely applying the proposed algorithm PML, because using the same weighting for both the interpolation *and* the losses attaches too much importance on the task with large loss magnitude. However, simple balancing schemes palliate this issue; gradient balancing (PML-gb) discovers a superset of the Pareto Front and loss balancing (PML-lb) discovers the exact Pareto Front.

512 C DETAILS OF THE ILLUSTRATIVE EXAMPLE

513 The details of the illustrative example are provided in this section. We use the configuration pre-  
 514 sented by Navon et al. (2022), which was introduced with slight modifications by Liu et al. (2021)  
 515 and Yu et al. (2020). Specifically, let  $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$  be the parameter vector and  $L = (\ell_1, \ell_2)$   
 516 be the vector objective defined as follows:

$$\tilde{\ell}_1(\theta) = c_1(\theta)f_1(\theta) + c_2(\theta)g_1(\theta) \quad \text{and} \quad \tilde{\ell}_2(\theta) = c_1(\theta)f_2(\theta) + c_2(\theta)g_2(\theta)$$

where

$$\begin{aligned} f_1(\theta) &= \log(\max(|0.5(-\theta_1 - 7) - \tanh(-\theta_2)|, 5e - 6)) + 6, \\ f_2(\theta) &= \log(\max(|0.5(-\theta_1 + 3) - \tanh(-\theta_2) + 2|, 5e - 6)) + 6, \\ g_1(\theta) &= \left( (-\theta_1 + 7)^2 + 0.1 \cdot (-\theta_2 - 8)^2 \right) / 10 - 20, \\ g_2(\theta) &= \left( (-\theta_1 - 7)^2 + 0.1 \cdot (-\theta_2 - 8)^2 \right) / 10 - 20, \\ c_1(\theta) &= \max(\tanh(0.5\theta_2), 0) \quad \text{and} \quad c_2(\theta) = \max(\tanh(-0.5\theta_2), 0) \end{aligned}$$

517 We use the experimental setting outlined by (Navon et al., 2022) with minor modifications, i.e.,  
 518 Adam optimizer with a learning rate of  $2e - 3$  and training lasts for  $50K$  iterations. The  
 519 overall objectives are  $\ell_1 = c \cdot \tilde{\ell}_1$  and  $\ell_2 = \tilde{\ell}_2$  where we explore two configurations for the  
 520 scalar  $c$ , namely  $c \in \{0.1, 1\}$ . For  $c = 1$ , the two tasks have losses at the same scale.  
 521 For  $c = 0.1$ , the difference in loss scales makes the problem more challenging and the algo-

522 rithm used should be characterized by scale invariance in order to find diverse solutions span-  
523 ning the entirety of the Pareto Front. The initialization points are drawn from the following set  
524  $\{(-8.5, 7.5), (0.0, 0.0), (9.0, 9.0), (-7.5, -0.5), (9, -1.0)\}$ . In the case of Pareto Manifold Learn-  
525 ing with two ensemble members there are  $5^2 = 25$  initialization pairs. In the main text we use the  
526 initialization pair with the worst initial objective values.

527 **Figure 6** presents the results for the case of different loss scales, i.e.,  $c = 0.1$ . We plot various  
528 baselines and three versions of the proposed algorithm, Pareto Manifold Learning or PML in short.  
529 We focus on the effect of the balancing schemes, introduced in **Section 4.2**, resulting in the use of no  
530 balancing scheme (denoted as PML), the use of gradient balancing (denoted as PML-gb) and the use  
531 of loss balancing (denoted as PML-lb). We dedicate two figures for each version of the algorithm  
532 and we present all 25 initialization pairs for completeness. **Figure 7** and **Figure 8** correspond to  
533 no balancing scheme in the case of equal loss scales  $c = 1.0$ , i.e., they complement **Figure 1** of  
534 the main text. The subsequent figures focus on the case of unequal loss scales where  $c = 0.1$ ;  
535 **Figure 9** **Figure 10** correspond to no balancing scheme, **Figure 11** and **Figure 12** correspond to the  
536 use of gradient balancing, **Figure 13** and **Figure 14** correspond to the use of loss balancing. The first  
537 figures of each pair show the trajectories for each initialization pair, with markers for initial and final  
538 positions. The other figures of each pair dispense of the visual clutter and focus on the subspace  
539 discovered in the final step of training, which is plotted with dashed lines along with the analytical  
540 Pareto Front in solid light blue. Hence, they provide a succinct overview of whether the method was  
541 able or not to discover the (entire) Pareto Front.

542 For  $c = 1.0$ , the proposed method is able to retrieve the exact Pareto Front with no balancing scheme  
543 for most initialization pairs, as can be seen in **Figure 8**. In three cases (out of 25), the method fails.  
544 In our experiments, we found that allowing longer training times or higher learning rates resolve  
545 the remaining cases. For  $c = 0.1$ , the problem is more challenging and the vanilla version of the  
546 algorithm results in a subset of the analytical Pareto Front. **Figure 10** shows that this subset is  
547 consistent across initialization pairs, excluding the ones the method fails, and focuses on the task  
548 with higher loss magnitude. Applying gradient balancing, shown in **Figure 11** and **Figure 12**, allows  
549 the method to retrieve (a superset of) the Pareto Front for all initialization pairs. Similarly, loss  
550 balancing, shown in **Figure 13** and **Figure 14**, results in the exact Pareto Front. Hence, the inclusion  
551 of balancing schemes endows scale invariance in the proposed algorithm. Balancing schemes are  
552 used for the more challenging datasets, such as `CityScapes`.

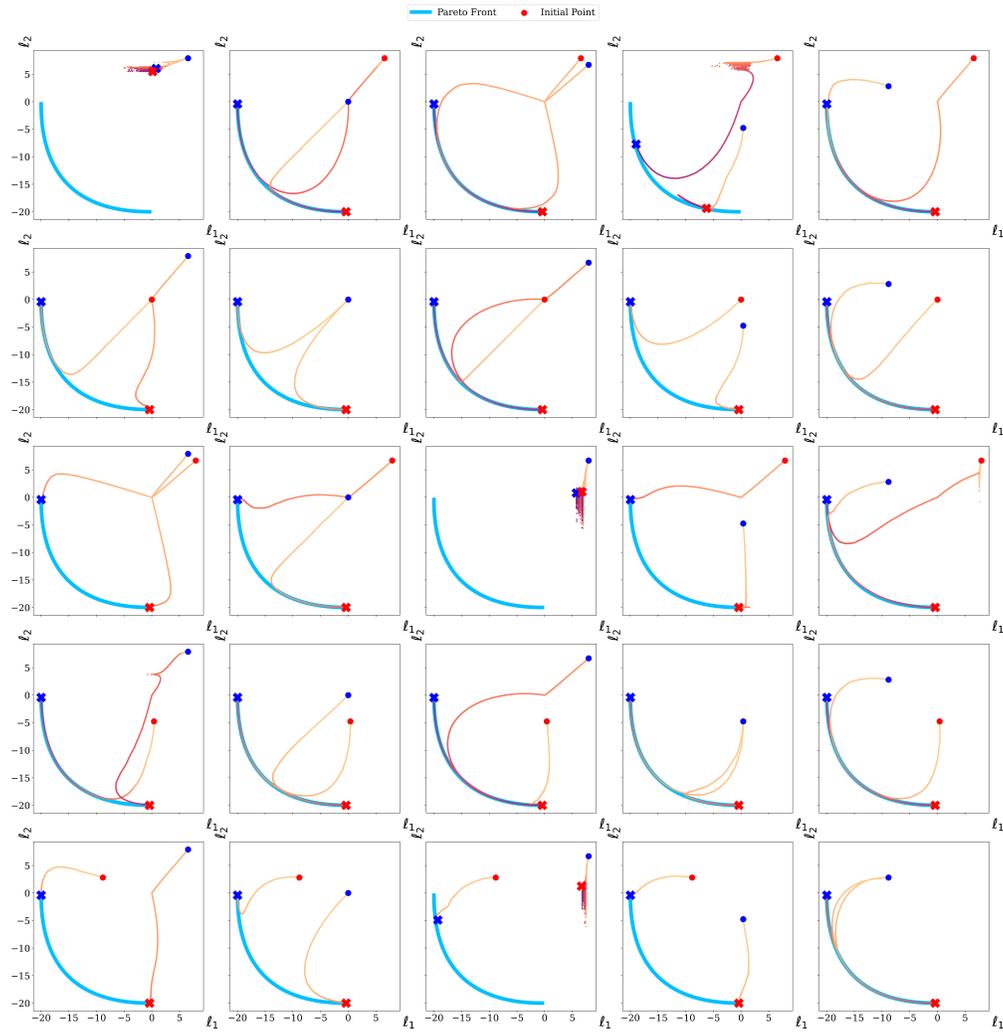


Figure 7: *Illustrative example*. Optimization trajectories in objective space for all initialization pairs in the case of equal loss scales ( $c = 1.0$ ) and application of the proposed method with no balancing scheme. Blue and red markers show each ensemble member’s loss value, dots and “X”s correspond to the initial and final step, accordingly. In all but four cases, Pareto Manifold Learning retrieves the entirety of the Pareto Front (can be seen clearly in Figure 8). Allowing longer training times or higher learning rates solves the remaining initialization pairs.

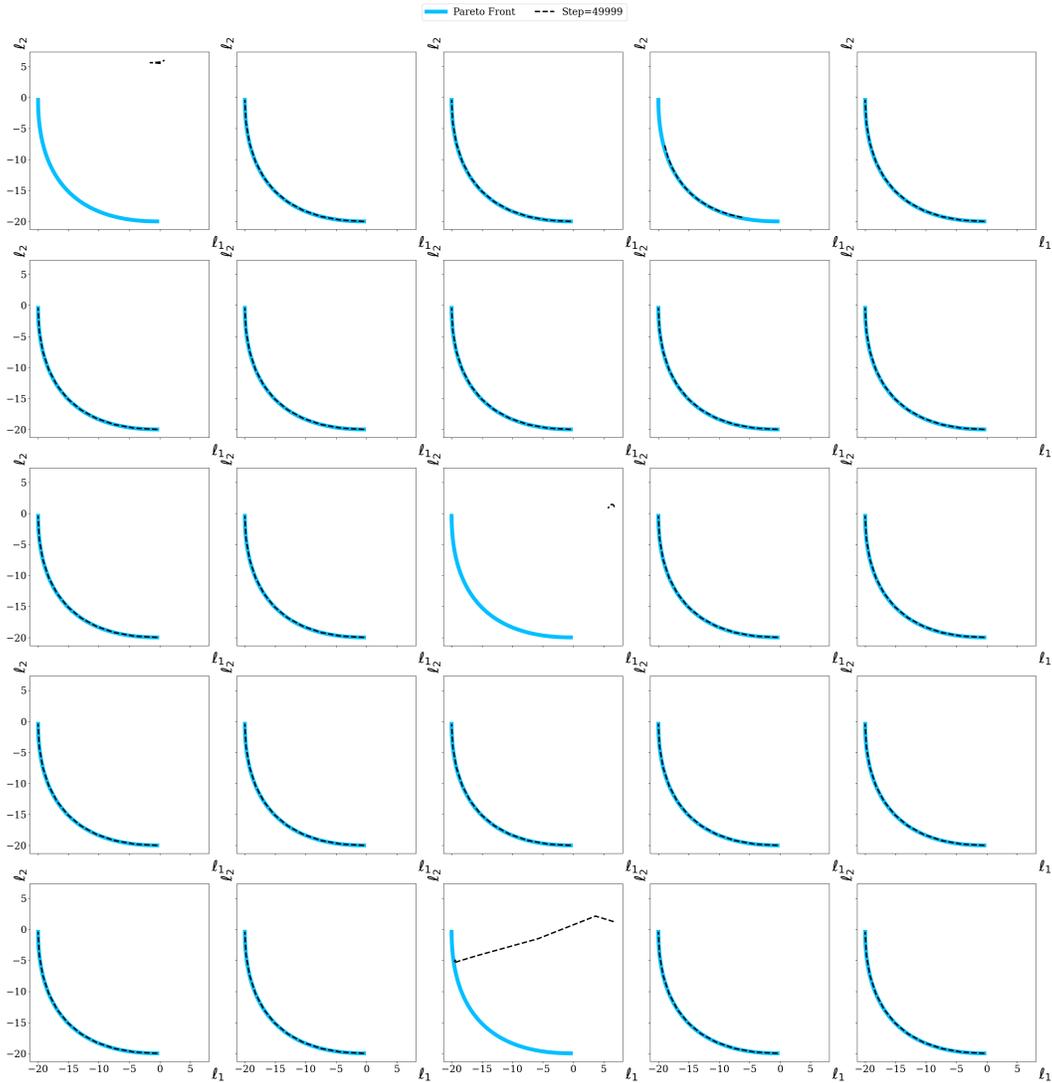


Figure 8: *Illustrative example.* Mapping in objective space of the weight subspace discovered by the proposed method with no balancing scheme, in the case of equal loss scales ( $c = 1.0$ ). The analytic Pareto Front is plotted in light blue. In all but four cases, the dashed line (our method) coincides with the full analytic Pareto Front.

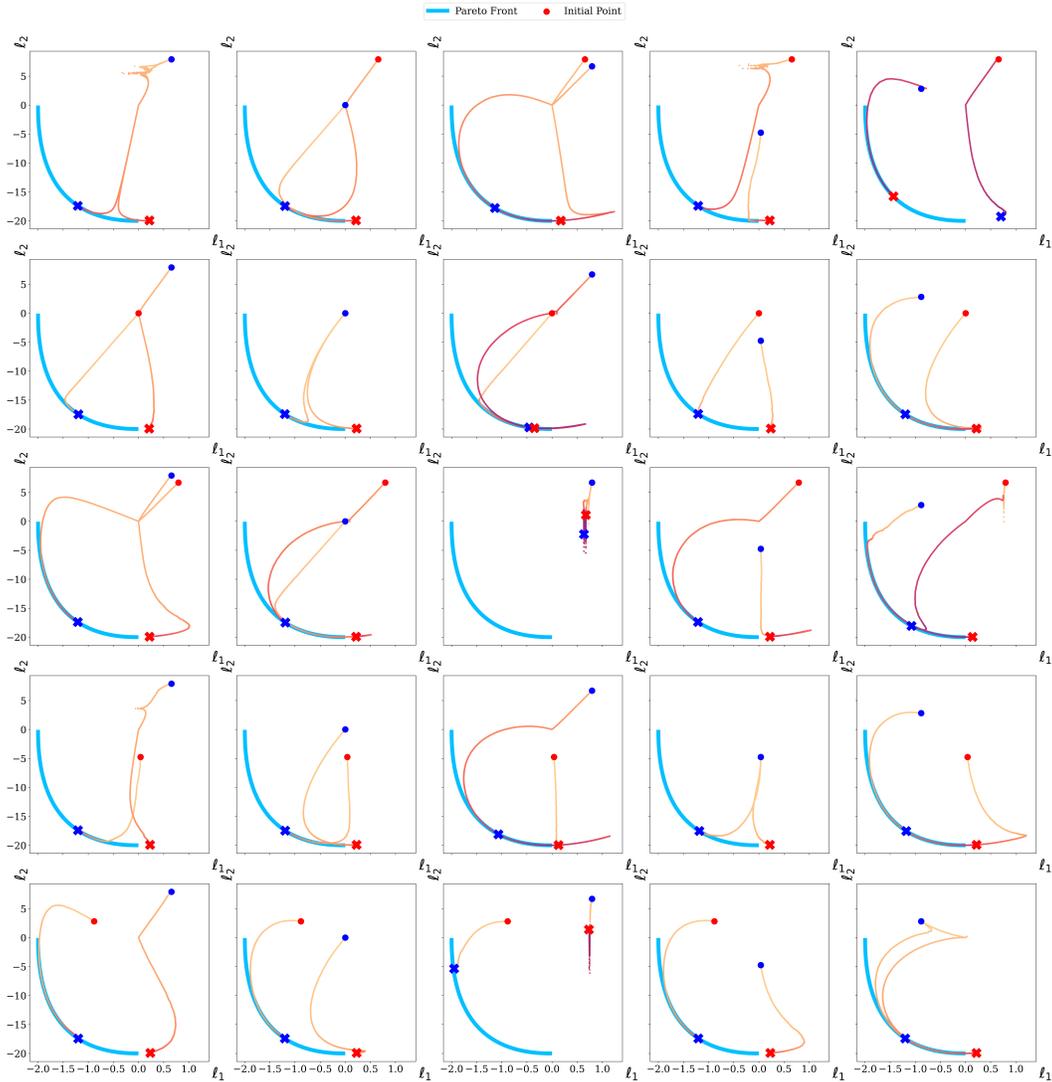


Figure 9: *Illustrative example.* Optimization trajectories in objective space for all initialization pairs in the case of unequal loss scales ( $c = 0.1$ ) and application of the proposed method with no balancing scheme. Blue and red markers show each ensemble member’s loss value, dots and “X”s correspond to the initial and final step, accordingly. For the vast majority of initialization pairs, the lack of balancing scheme guides the ensemble to a subset of the Pareto Front, influenced by the task with higher loss magnitude (can be seen clearly in Figure 10).

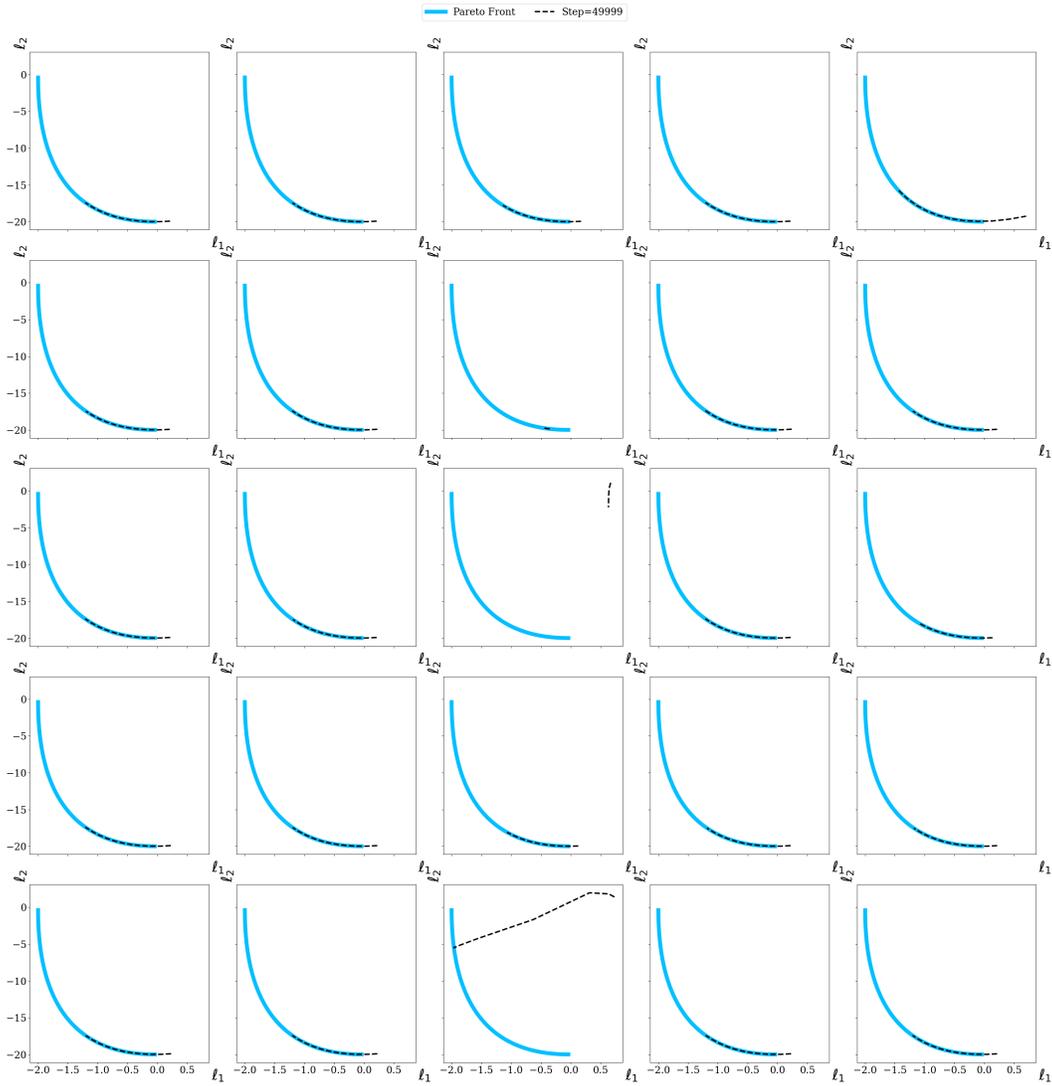


Figure 10: *Illustrative example*. Mapping in objective space of the weight subspace discovered by the proposed method with no balancing scheme, in the case of unequal loss scales ( $c = 0.1$ ). The analytic Pareto Front is plotted in light blue. The lack of balancing scheme renders optimization difficult; the method either completely fails or retrieves a narrow subset of the analytic Pareto Front. Applying balancing schemes resolve these issues.

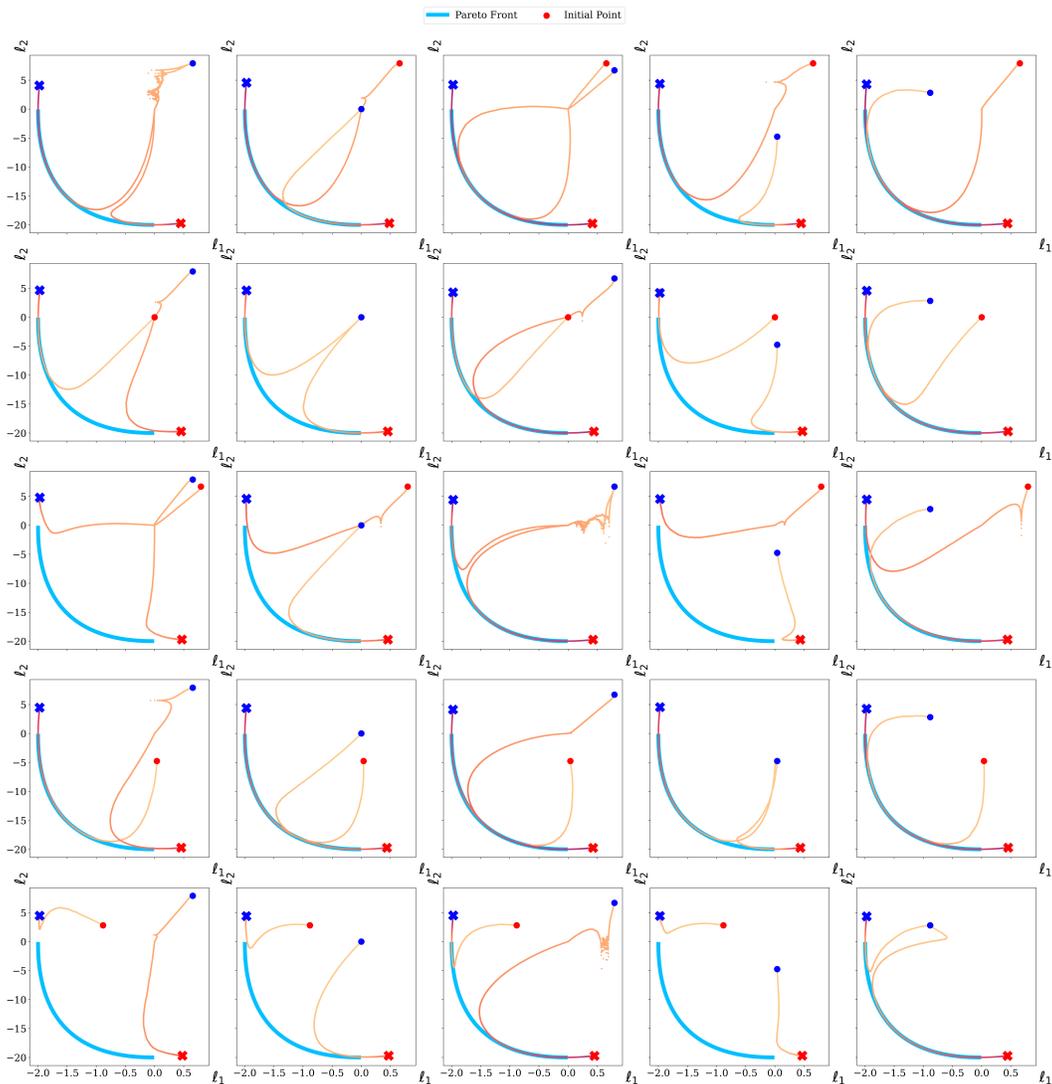


Figure 11: *Illustrative example.* Optimization trajectories in objective space for all initialization pairs in the case of unequal loss scales ( $c = 0.1$ ) and application of the proposed method with gradient balancing scheme. Blue and red markers show each ensemble member’s loss value, dots and “X”s correspond to the initial and final step, accordingly. The proposed method discovers a subspace whose mapping in objective space results in a superset of the Pareto Front. This can be clearly seen in [Figure 12](#).

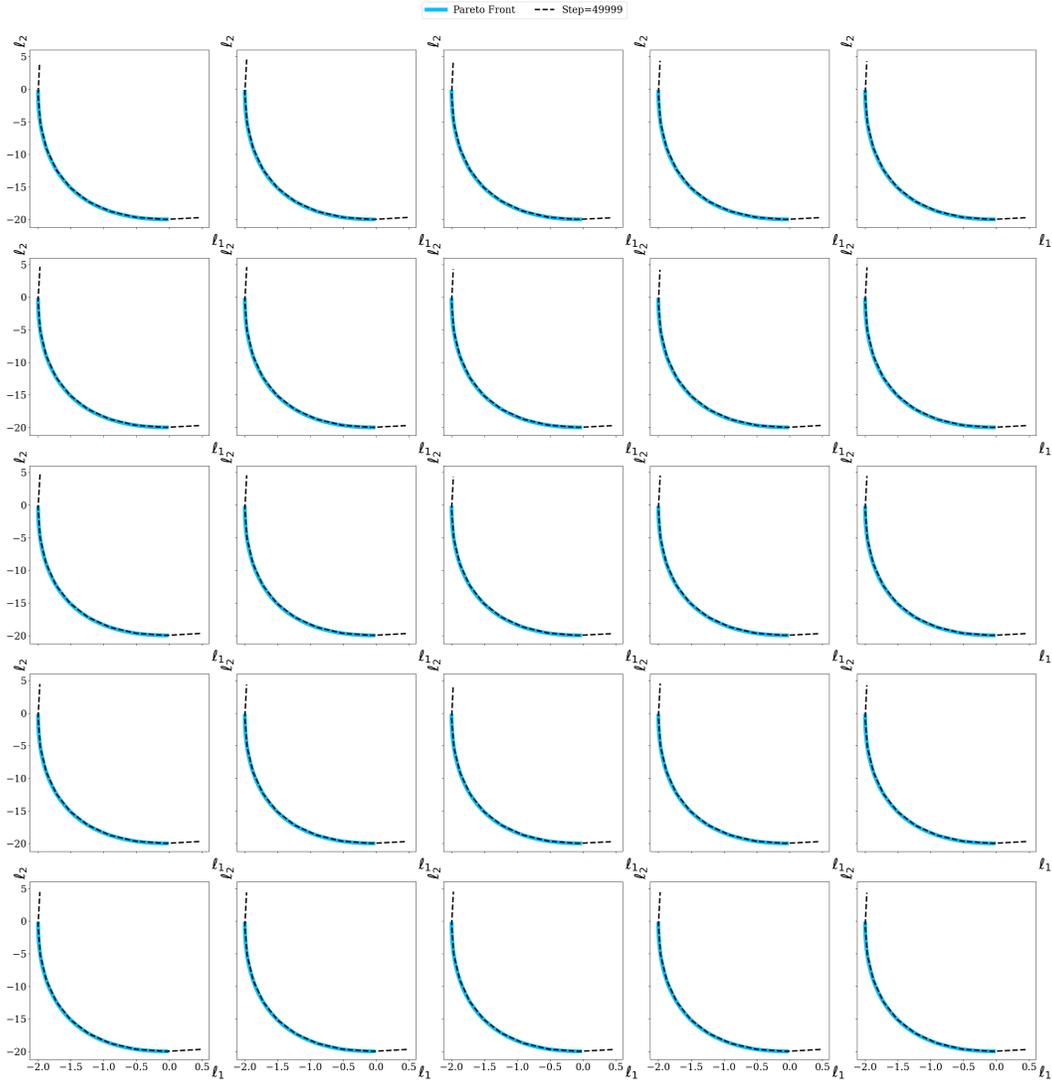


Figure 12: *Illustrative example.* Mapping in objective space of the weight subspace discovered by the proposed method with gradient balancing scheme, in the case of unequal loss scales ( $c = 0.1$ ). The analytic Pareto Front is plotted in light blue. The proposed method consistently finds the same subspace, which is a superset of the analytic Pareto Front.

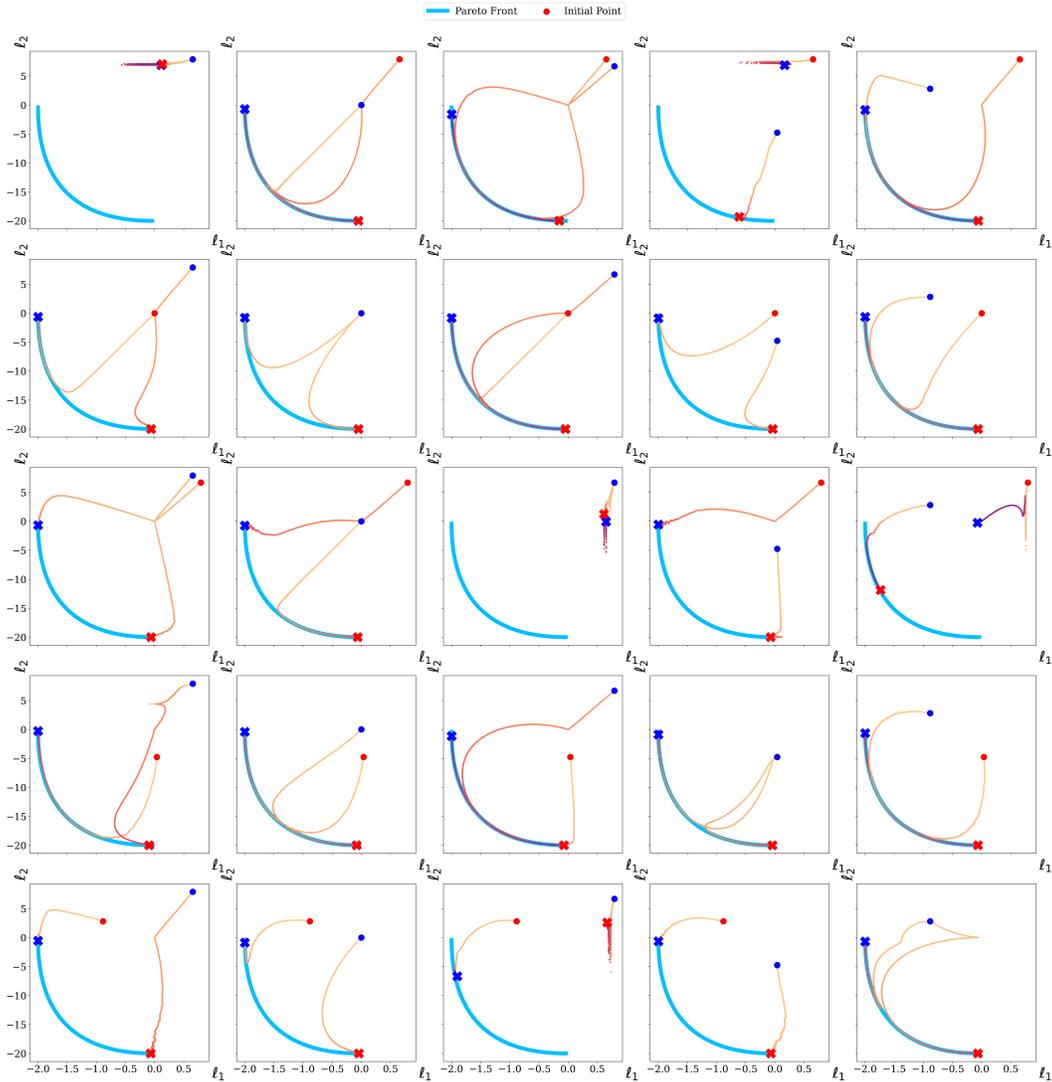


Figure 13: *Illustrative example.* Optimization trajectories in objective space for all initialization pairs in the case of unequal loss scales ( $c = 0.1$ ) and application of the proposed method with loss balancing scheme. Blue and red markers show each ensemble member’s loss value, dots and “X”s correspond to the initial and final step, accordingly. For all but five cases, the proposed method discovers a subspace whose mapping in objective space results in the exact Pareto Front. This can be clearly seen in [Figure 14](#).

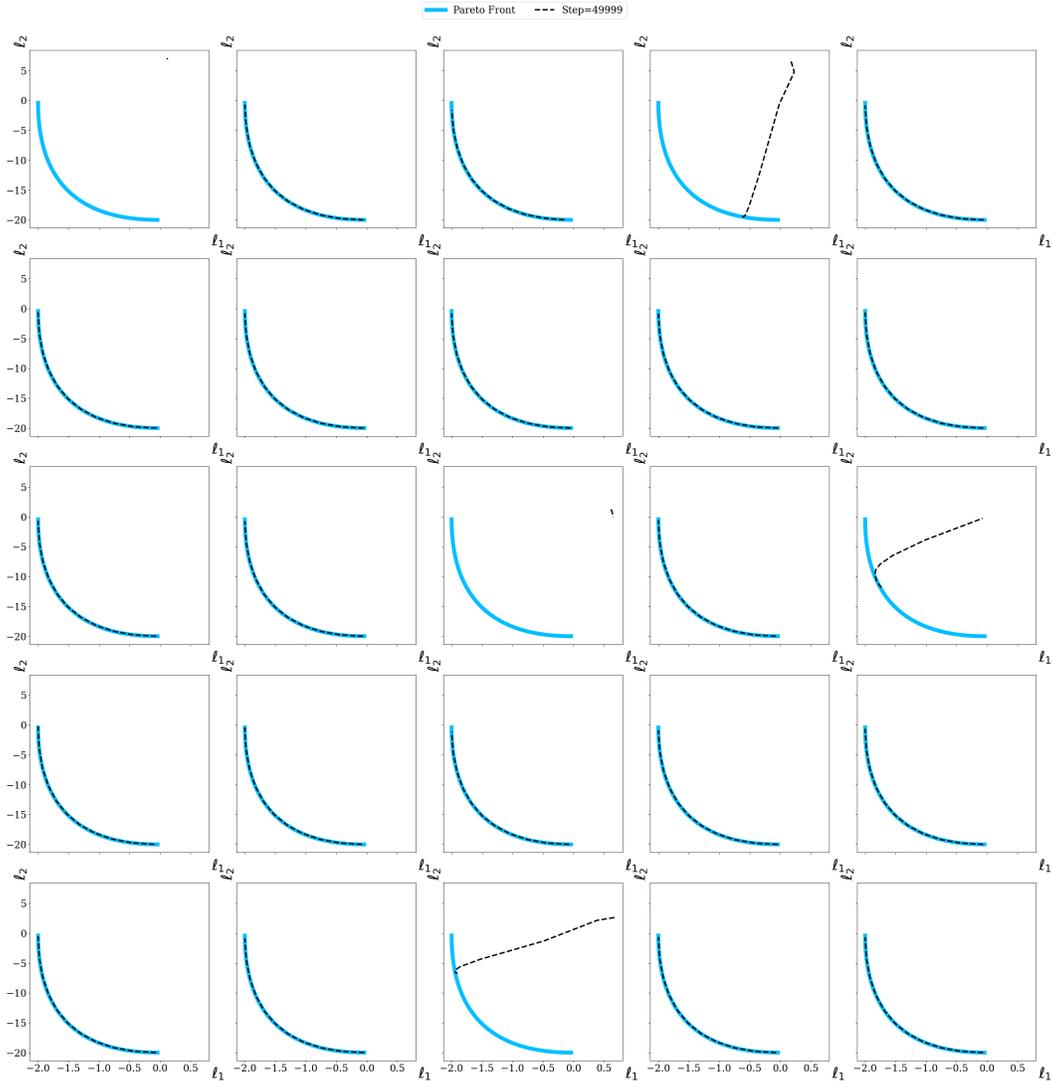


Figure 14: *Illustrative example.* Mapping in objective space of the weight subspace discovered by the proposed method with loss balancing scheme, in the case of unequal loss scales ( $c = 0.1$ ). The analytic Pareto Front is plotted in light blue. Using loss balancing endows scale invariance and the solutions are more functionally diverse, in comparison with no balancing scheme in Figure 10. However, the same initialization pairs continue to be problematic as in the case of equal loss scales (see Figure 8). Allowing for longer training or higher learning rates solves the remaining initialization pairs.

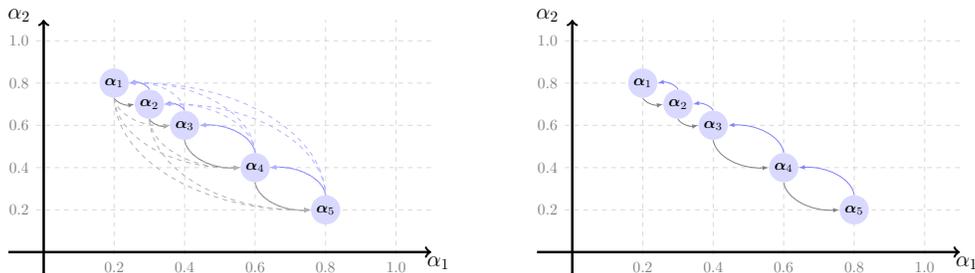


Figure 15: Multi-Forward Graph: case of two tasks. We assume a window of  $W = 5$ . The nodes lie in the line segment  $\alpha_2 + \alpha_1 = 1$ ,  $\alpha_1, \alpha_2 \in [0, 1]$ . (Left) Full graph and dashed edges will be removed. (Right) Final graph.

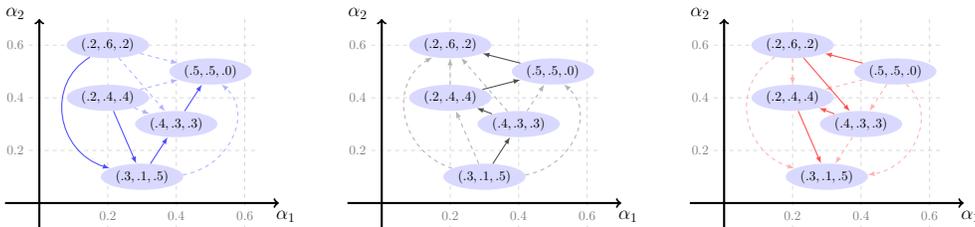


Figure 16: Multi-Forward Graph for three tasks. Left, middle and right present the case of the first, second and third task, respectively. Each node is noted by its weighting, summing up to 1. Edges are drawn if the two nodes obey the total ordering imposed by the task. Dashed edges are omitted from the final graph.

553 D ABLATION ON MULTI-FORWARD REGULARIZATION

554 Multi-Forward regularization, introduced in Section 4.2, penalizes the ensemble if the interpolated  
 555 models’ losses (sampled within a batch) are not in accordance with the tradeoff imposed by the  
 556 corresponding interpolation weights. Simply put, the closer we sample to the member corresponding  
 557 to task 1, the lower the loss should be on task 1. The same applies to the other tasks. Equation 3  
 558 in the main text presents the case of two tasks, where the idea of the regularization is outlined in  
 559 loss space. For completeness, we present the underlying graph construction for the cases of two and  
 560 three tasks in Figure 15 and Figure 16, respectively. The nodes of the graphs are associated with  
 561 the sampled weightings and the edges for the graph  $\mathcal{G}_t$  of task  $t$  are drawn w.r.t. the corresponding  
 562 partial ordering. If the loss ordering is violated for a given edge, a penalty term is added.

563 We ablate the effect multi-forward training and the corresponding regularization have on perfor-  
 564 mance. We explore the MultiMNIST and Census datasets using the same experimental configu-  
 565 rations as in the main text. We are interested in two parameters:

- 566 •  $W$ : number of  $\alpha$  re-samplings per batch. This parameter is also referred as *window*.
- 567 •  $\lambda$ : the regularization strength as presented in Algorithm 1. For  $\lambda = 0$ , no regularization is  
 568 applied but the subspace is still sampled  $W$  times and the total loss takes into account all  
 569 the respective interpolated models.

570 Figure 17 and Table 2 present the results for MultiMNIST. Figure 18 and Table 3 present the re-  
 571 sults for Census. It is important to note that MultiMNIST is symmetric, while Census is not.  
 572 As a result, the features learned for each single-task predictor are helpful to one another and the case  
 573 of  $\lambda = 0$ , i.e., no regularization and only multi-forward training, is beneficial for MultiMNIST but  
 574 not for Census. Intuitively, both digit classification tasks have the same difficulty and posterior  
 575 distribution, which produces few violations of monotonicity constraints and renders the regulariza-  
 576 tion less applicable. On the other hand, severe regularization such as  $\lambda = 10$  can be harmful and  
 577 hinder training. More details in table and figure captions.

[Reviewer qexX]: expanded commentary on  $\lambda = 0$  for MultiMNIST.

Table 2: `MULTIMNIST`: Ablation on multi-forward training and regularization, presented in Section 4.2. Validation performance in terms of HyperVolume (HV) metric. Higher is better, except for standard deviation (std). The visual complement of the table appears in Figure 17. For each configuration, we track the Hypervolume across three random seeds and present Mean HV, max HV and standard deviation. We annotate with bold the best per column. In the main text, we report the best result in terms of mean HV, i.e.,  $W = 4$  and  $\lambda = 0$ .

		Seed - 0	Seed - 1	Seed - 2	Mean HV	Max HV	std
$W = 2$	$\lambda = 0$	0.9205	0.9083	0.9100	0.9129	0.9205	0.0054
	$\lambda = 2$	0.9121	0.9105	0.9037	0.9088	0.9121	0.0036
	$\lambda = 5$	0.9132	0.9016	0.8979	0.9043	0.9132	0.0065
	$\lambda = 10$	0.8766	0.8932	0.8470	0.8723	0.8932	0.0191
$W = 3$	$\lambda = 0$	0.9215	0.9141	0.9111	0.9156	0.9215	0.0044
	$\lambda = 2$	0.9176	0.9150	0.9122	0.9149	0.9176	0.0022
	$\lambda = 5$	0.9155	0.9138	0.9140	0.9144	0.9155	<b>0.0008</b>
	$\lambda = 10$	0.9122	0.9050	0.8962	0.9045	0.9122	0.0066
$W = 4$	$\lambda = 0$	<b>0.9220</b>	0.9187	0.9143	<b>0.9184</b>	<b>0.9220</b>	0.0032
	$\lambda = 2$	0.9213	0.9149	<b>0.9157</b>	0.9173	0.9213	0.0028
	$\lambda = 5$	0.9158	0.9139	0.9132	0.9143	0.9158	0.0011
	$\lambda = 10$	0.9177	0.9022	0.9102	0.9100	0.9177	0.0063
$W = 5$	$\lambda = 0$	0.9131	0.9180	0.9156	0.9156	0.9180	0.0020
	$\lambda = 2$	0.9158	<b>0.9203</b>	0.9146	0.9169	0.9203	0.0024
	$\lambda = 5$	0.9138	0.9082	0.9140	0.9120	0.9140	0.0027
	$\lambda = 10$	0.9165	0.9158	0.9121	0.9148	0.9165	0.0019

Table 3: `Census`: Ablation on multi-forward training and regularization, presented in Section 4.2. Validation performance in terms of HyperVolume (HV) metric. Higher is better, except for standard deviation (std). The visual complement of the table appears in Figure 18. For each configuration, we track the Hypervolume across three random seeds and present Mean HV, max HV and standard deviation. We annotate with bold the best per column. In the main text, we report the best result in terms of mean HV, i.e.,  $W = 2$  and  $\lambda = 5$ .

		Seed - 0	Seed - 1	Seed - 2	Mean HV	Max HV	std
$W = 2$	$\lambda = 0$	0.6517	0.6530	0.6532	0.6526	0.6532	0.0006
	$\lambda = 2$	0.6575	0.6564	0.6560	0.6566	0.6575	0.0006
	$\lambda = 5$	<b>0.6577</b>	<b>0.6574</b>	<b>0.6590</b>	<b>0.6581</b>	<b>0.6590</b>	0.0007
	$\lambda = 10$	0.6548	0.6557	0.6554	0.6553	0.6557	<b>0.0004</b>
$W = 3$	$\lambda = 0$	0.6517	0.6496	0.6501	0.6505	0.6517	0.0009
	$\lambda = 2$	0.6540	0.6523	0.6544	0.6536	0.6544	0.0009
	$\lambda = 5$	0.6552	0.6539	0.6536	0.6542	0.6552	0.0007
	$\lambda = 10$	0.6574	0.6567	0.6566	0.6569	0.6574	<b>0.0004</b>
$W = 4$	$\lambda = 0$	0.6488	0.6516	0.6504	0.6503	0.6516	0.0011
	$\lambda = 2$	0.6492	0.6522	0.6504	0.6506	0.6522	0.0012
	$\lambda = 5$	0.6499	0.6514	0.6525	0.6513	0.6525	0.0011
	$\lambda = 10$	0.6529	0.6549	0.6558	0.6545	0.6558	0.0012
$W = 5$	$\lambda = 0$	0.6497	0.6502	0.6484	0.6494	0.6502	0.0008
	$\lambda = 2$	0.6478	0.6497	0.6495	0.6490	0.6497	0.0009
	$\lambda = 5$	0.6492	0.6509	0.6489	0.6497	0.6509	0.0009
	$\lambda = 10$	0.6507	0.6538	0.6508	0.6518	0.6538	0.0014

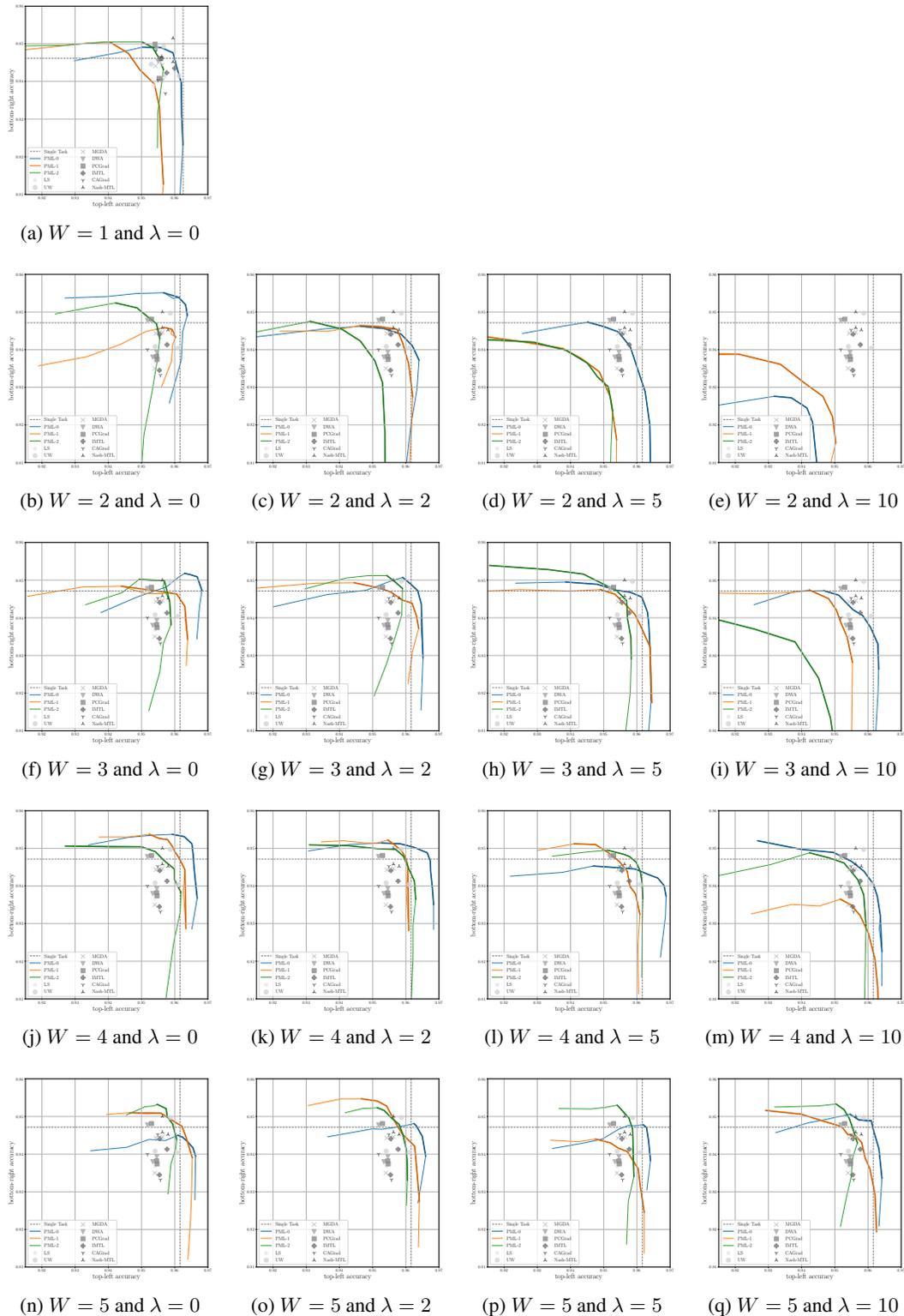


Figure 17: MultiMNIST: Effect of multi-forward on the window  $W$  and the regularization coefficient  $\lambda$  on the validation dataset. The case of no multi-forward ( $W = 1$ ) is presented in the first row. Multi-forward regularization for higher  $W$  values is beneficial. Intuitively, attaching serious weight on the regularization  $\lambda \in \{5, 10\}$  while sampling few times  $W \in \{2, 3\}$  leads to suboptimal performance since the update step focuses on an uninformed regularization term. The accompanying quantitative analysis appears in Table 2.



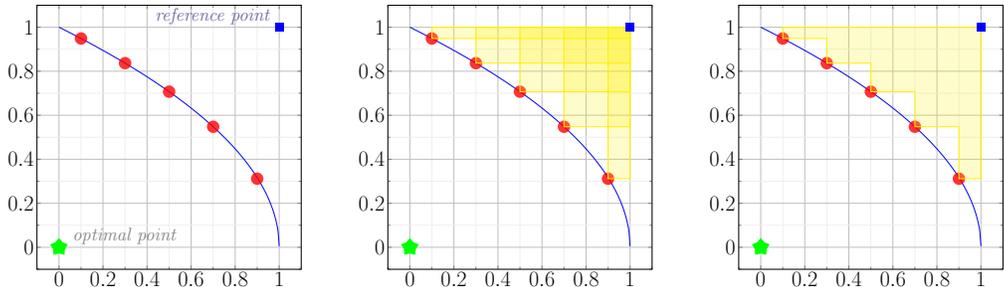


Figure 19: Visual Explanation of Hypervolume. The metric captures the union of axis-aligned rectangles defined by the reference point (star) and the corresponding sample points (red circles). This example showcases loss and the perfect oracle lies in the origin. The point (1, 1) is used for reference. Hence, higher hypervolume implies that the objective space is better explored/covered.

578 E HYPERVOLUME ANALYSIS ON MULTIMNIST AND CENSUS

579 HyperVolume is a metric widely used in multi-objective optimization that captures the quality of  
 580 exploration. A visual explanation of the metric is given in Figure 19. Table 4 presents the results  
 581 of Figure 4 of the main text in a tabular form. We present the best three results per column (higher  
 582 is better) to succinctly and visually show that all Pareto Manifold Learning seeds outperform the  
 583 baselines.

Table 4: Tabular complement to Figure 4. Classification accuracy for both tasks and HyperVolume (HV) metric (higher is better). Three random seeds per method. For baselines, we show the mean accuracy and HV (across seeds). For PML, we show the results per seed; HV and max accuracies for the subspace yielded by that seed. We use underlined bold, solely bold and solely underlined font for the best, second best and third best results. We observe that the best results are concentrated in the rows concerning the proposed method (PML). Note that the use of three decimals leads to ties.

	MultimNIST			Census		
	Task 1	Task 2	HV	Task 1	Task 2	HV
LS	0.955	0.944	0.907	0.827	0.785	0.651
UW	0.957	0.945	0.913	0.827	0.785	0.650
MGDA	0.956	0.943	0.904	0.828	0.785	0.651
DWA	0.955	0.945	0.907	0.828	0.785	0.651
PCGrad	0.955	0.946	0.908	0.828	0.785	0.650
IMTL	0.958	0.944	0.908	0.828	0.786	0.651
Nash-MTL	0.958	0.948	0.913	0.827	0.785	0.650
PML - 0	<b><u>0.968</u></b>	<b><u>0.951</u></b>	<b><u>0.92</u></b>	<b><u>0.830</u></b>	<b><u>0.789</u></b>	<b><u>0.655</u></b>
PML - 1	<b><u>0.961</u></b>	<b><u>0.953</u></b>	<b><u>0.916</u></b>	<b><u>0.830</u></b>	<b><u>0.789</u></b>	<b><u>0.655</u></b>
PML - 2	<b><u>0.964</u></b>	<b><u>0.953</u></b>	<b><u>0.919</u></b>	<b><u>0.829</u></b>	<b><u>0.788</u></b>	<b><u>0.653</u></b>

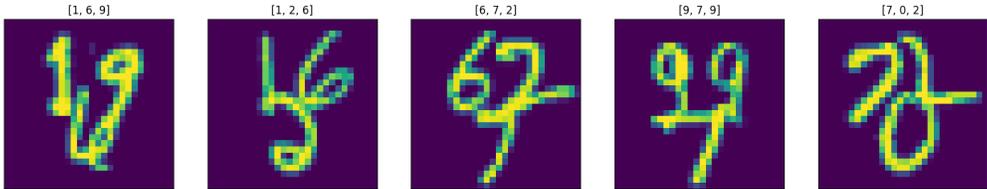


Figure 20: Examples of samples and corresponding labels for the MultiMNIST-3 dataset.

## 584 F MULTIMNIST-3 ADDITIONAL RESULTS

585 This section serves as supplementary to Section 5.2 of the main text. MultiMNIST-3 is  
 586 a synthetic dataset generated by MNIST samples in a manner similar to the creation of the  
 587 MultiMNIST dataset, which is ubiquitous in the Multi-Task Learning literature. Specifically, each  
 588 MultiMNIST-3 sample is created with the following procedure. Three randomly sampled digits  
 589 of size  $28 \times 28$  are placed in the top-left, top-right and bottom middle pixels of a  $42 \times 42$  grid.  
 590 For the pixels where the initial digits overlap, the maximum value is selected. Finally, the image is  
 591 resized to  $28 \times 28$  pixels. Figure 20 shows some examples of the dataset, which consists of three  
 592 digit classification tasks.

593 Section 5 compares the performance of baselines and the proposed method while Figure 21 presents  
 594 visually the performance achieved on the discovered subspace.

Table 5: MultiMNIST-3: Mean Accuracy and standard deviation of accuracy (over 3 random seeds). For the proposed method (PML), we report the mean and standard deviation of the best performance from the interpolated models in the sampled subspace. No balancing schemes and regularization are applied. **Bold is used for the best performing multi-task method.**

	Task 1	Task 2	Task 3
STL	$96.97 \pm 0.06$	$96.10 \pm 0.17$	$96.40 \pm 0.22$
LS	$96.26 \pm 0.20$	$95.48 \pm 0.14$	$95.87 \pm 0.37$
UW	$96.48 \pm 0.08$	$95.42 \pm 0.30$	$95.77 \pm 0.06$
MGDA	$96.50 \pm 0.20$	$94.80 \pm 0.22$	$95.71 \pm 0.08$
DWA	$96.42 \pm 0.26$	$95.26 \pm 0.29$	$95.75 \pm 0.08$
PCGrad	$96.45 \pm 0.06$	$95.39 \pm 0.15$	$95.88 \pm 0.01$
IMTL	$96.58 \pm 0.22$	$95.18 \pm 0.12$	$96.08 \pm 0.31$
CAGrad	$96.70 \pm 0.13$	$95.20 \pm 0.26$	$95.66 \pm 0.06$
Nash-MTL	<b><math>96.85 \pm 0.08</math></b>	$95.25 \pm 0.23$	$96.18 \pm 0.13$
PML (ours)	<b><math>96.85 \pm 0.43</math></b>	<b><math>95.72 \pm 0.22</math></b>	<b><math>96.27 \pm 0.32</math></b>

[Reviewer qexX]: Table update.

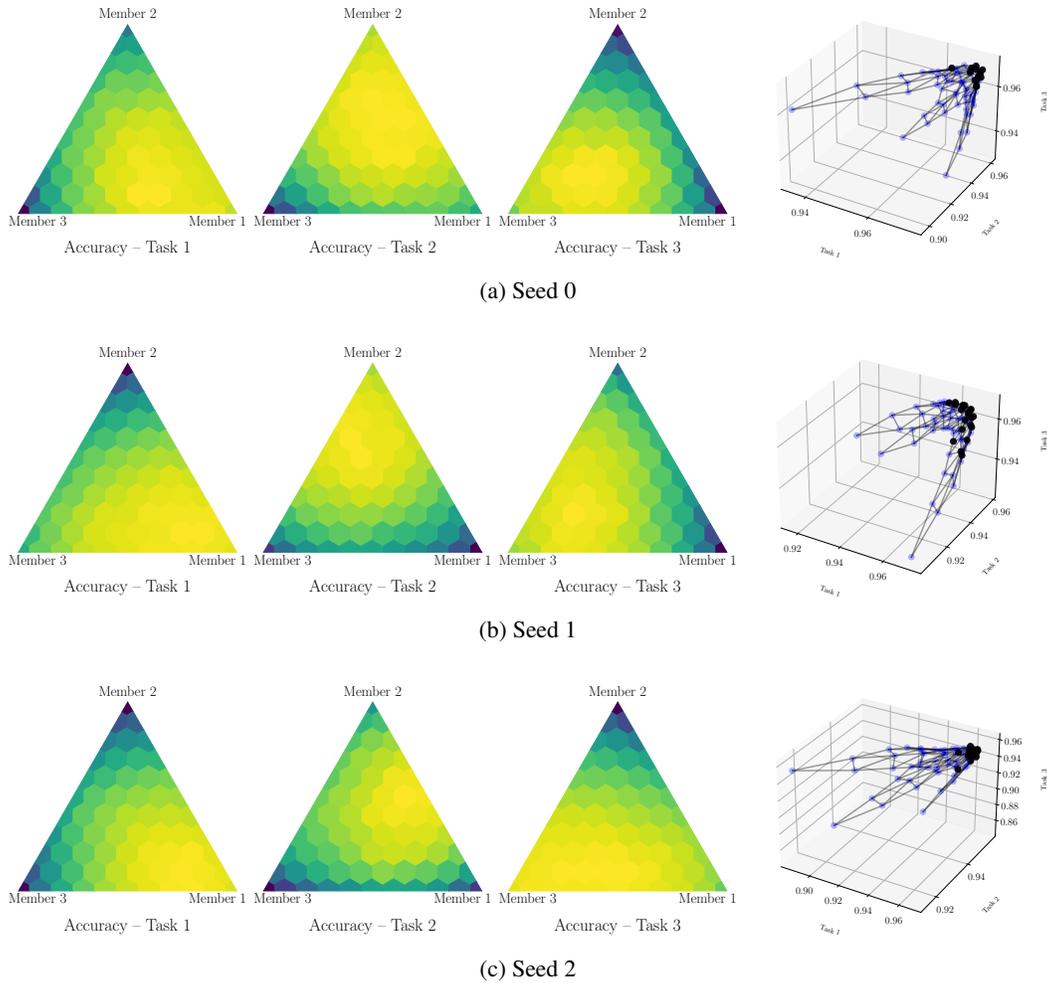


Figure 21: MultiMNIST-3 results for all three seeds. Each triangle shows the 66 points in the convex hull and color is used for the performance on the associated task. The 3d plot shows the mapping of the subspace to the multi-objective space. No balancing scheme is used.

## 595 G UTKFACE ADDITIONAL RESULTS

596 This section serves as supplementary to Section 5.2. Section 6 compares the performance of the  
 597 baselines and the proposed method. We experiment without balancing schemes and with gradient-  
 598 balancing, and present the results in Figure 22 and Figure 23, respectively. Together with the quan-  
 599 titative results, we observe that for datasets with varying task difficulties, scales, etc. the lack of  
 600 balancing can be impeding. On the other hand, its inclusion makes the subspace functionally di-  
 601 verse and boosts overall performance. For instance, Huber loss on the task of age prediction is  
 602 significantly improved.

Table 6: UTKFace: Mean Accuracy and standard deviation of accuracy (over 3 random seeds). For the proposed method (PML), we report the mean and standard deviation of the best performance from the interpolated models in the sampled subspace. No multi-forward training is applied. We present Pareto Manifold Learning with no balancing scheme and with gradient balancing, denoted as *gb*. **Bold is used for the best performing multi-task method.**

	Age ↓	Gender ↑	Ethnicity ↑
STL	0.081 ± 0.005	90.79 ± 0.55	82.38 ± 0.40
LS	0.086 ± 0.003	91.66 ± 0.55	82.78 ± 0.60
UW	0.093 ± 0.007	91.86 ± 0.75	83.62 ± 0.02
MGDA	<b>0.075 ± 0.003</b>	91.17 ± 0.59	74.06 ± 2.66
DWA	0.093 ± 0.008	91.65 ± 0.46	82.85 ± 0.20
PCGrad	0.101 ± 0.018	91.85 ± 0.90	83.57 ± 0.43
IMTL	0.091 ± 0.004	91.24 ± 0.34	82.52 ± 1.15
CAGrad	0.083 ± 0.002	<b>91.93 ± 0.53</b>	<b>83.71 ± 0.33</b>
Nash-MTL	0.095 ± 0.001	90.40 ± 0.16	79.59 ± 0.92
PML (ours)	0.096 ± 0.002	90.97 ± 0.63	81.78 ± 0.14
PML- <i>gb</i> (ours)	0.086 ± 0.003	91.61 ± 0.52	81.77 ± 0.86

[Reviewer qexX]: Table update.

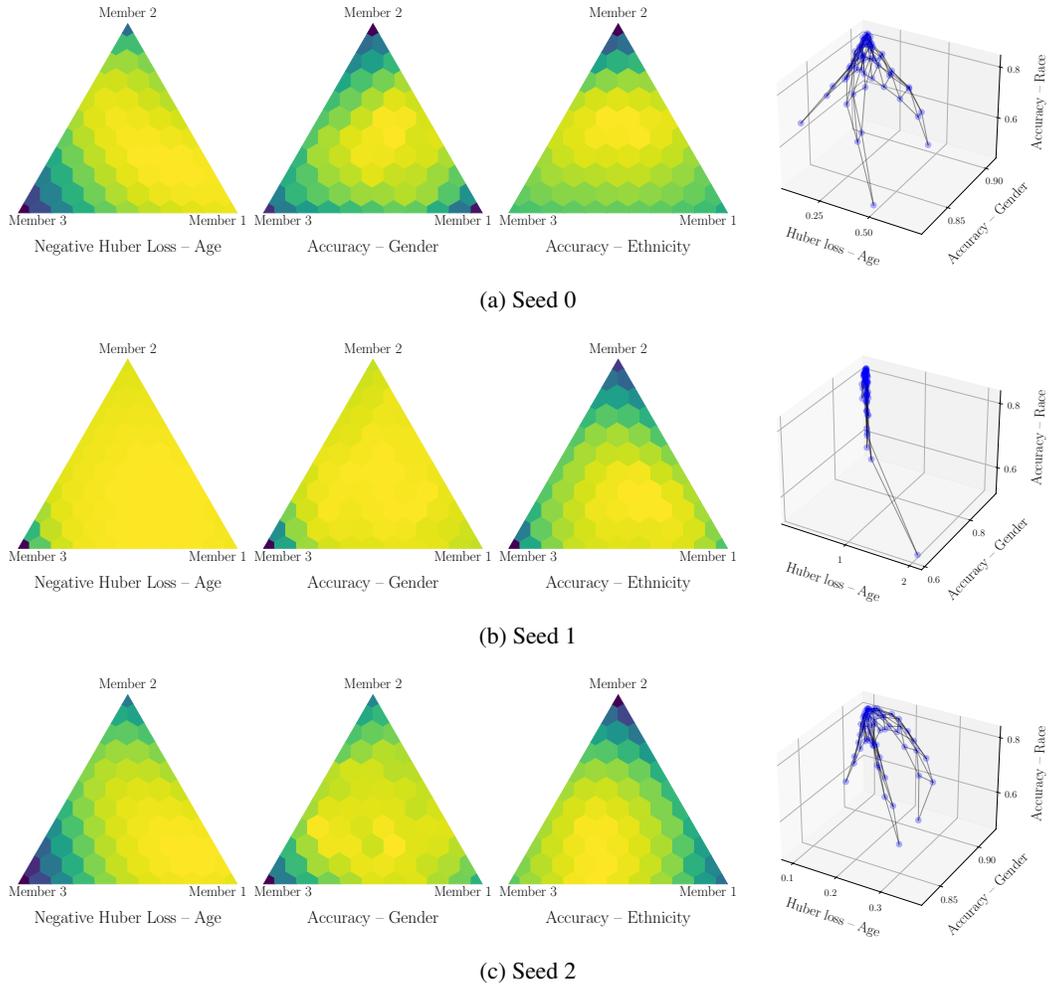


Figure 22: UTKFace results with Linear Scalarization for all three seeds. Each triangle shows the 66 points in the convex hull and color is used for the performance on the associated task. The 3d plot shows the mapping of the subspace to the multi-objective space. Applying no balancing scheme for datasets with different loss scales, e.g., regression and classification tasks, may lead to limited functional diversity, such as for seed 1.

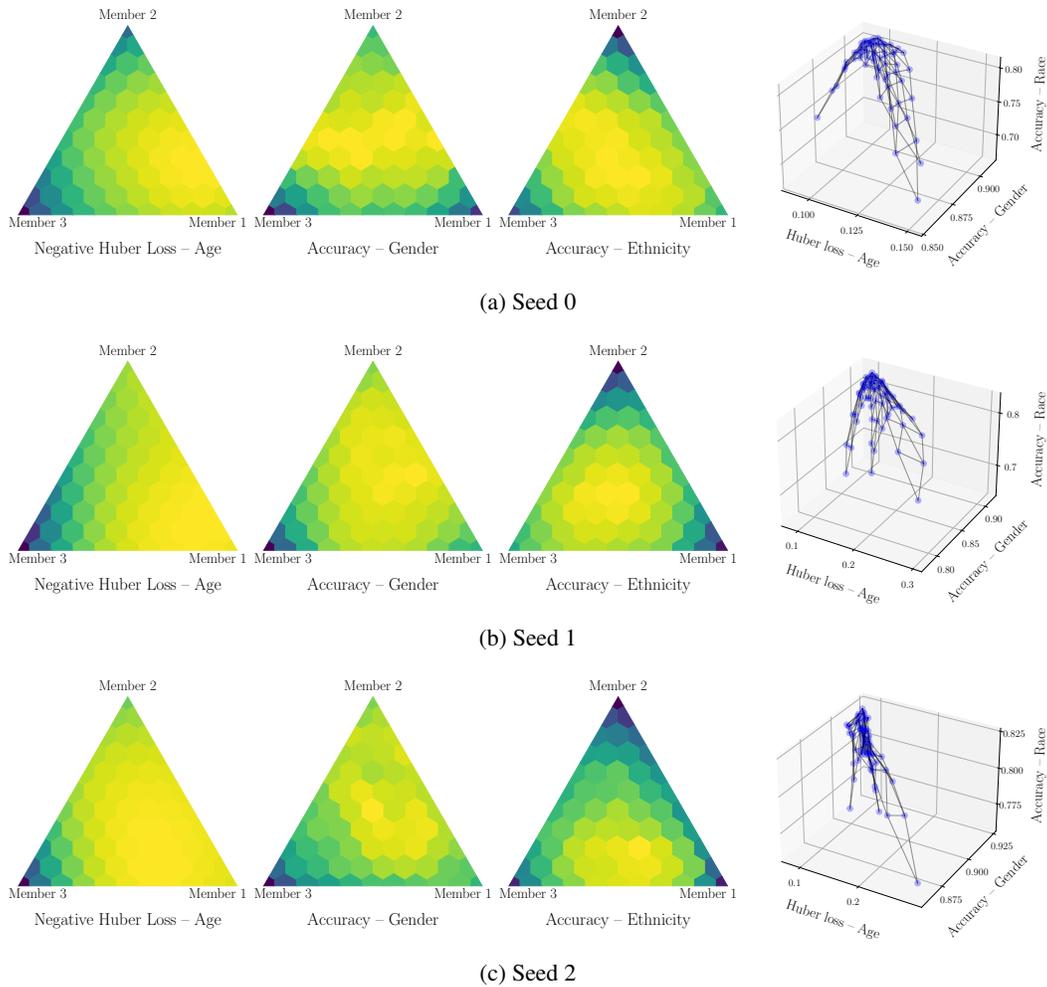


Figure 23: UTKFace results with Gradient-Balancing Scheme for all three seeds. Each triangle shows the 66 points in the convex hull and color is used for the performance on the associated task. The 3d plot shows the mapping of the subspace to the multi-objective space. For datasets with tasks of varying loss scales, applying gradient balancing improves functional diversity and performance, as shown in Section 6.

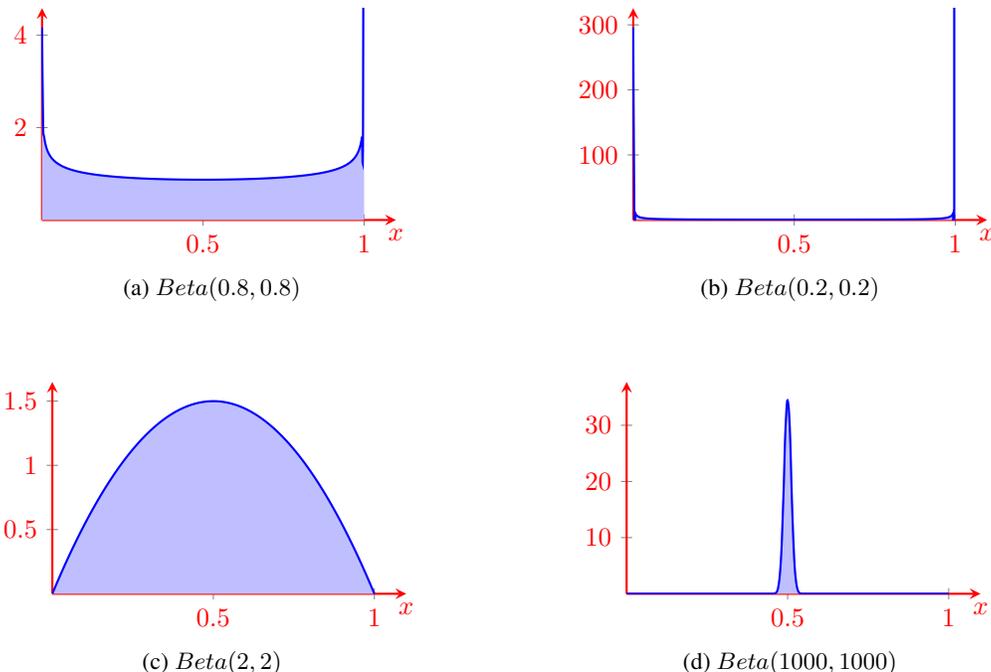


Figure 24: Dirichlet distribution in the case of two tasks. Top row:  $p < 1$  and the distribution is more concentrated towards the ensemble members. Bottom row:  $p > 1$  and the distribution focuses more on the midpoint which corresponds to all tasks having the same weight. Right column: extreme choices  $p \rightarrow 0$  or  $p \rightarrow \infty$ . Left column: milder choices.

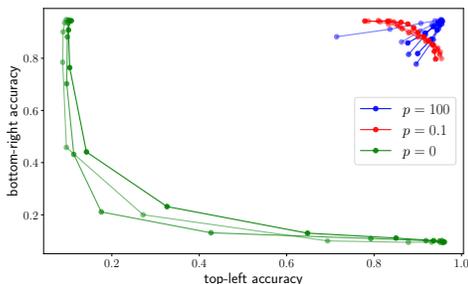
## 603 H DETAILS ON SAMPLING

604

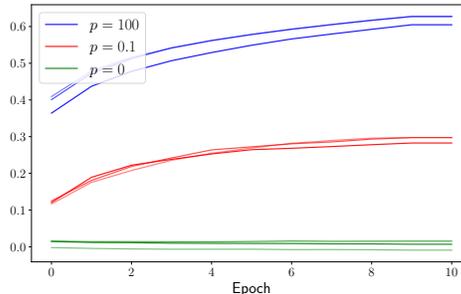
605 This appendix expands on Section 4.2 and, specifically, presents in greater detail the intuition behind  
 606 the sampling distribution’s parameters. Let  $\mathbf{p} \in \mathbb{R}_+^T$  be the parameters of the Dirichlet distribution.  
 607 Assuming no prior knowledge on the tasks, e.g., task difficulties or affinities, a symmetric distribu-  
 608 tion is used by setting  $\mathbf{p} = p\mathbf{1}_T$ . This design choice results in three cases:

- 609 •  $p = 1$ : the distribution is uniform on the simplex. Intuitively this means that all tasks are  
 610 equally important and we care about the diversity of solutions for all tradeoffs (reflected in  
 611 the linear scalarization weights)
- 612 •  $p \in (0, 1)$ : the distribution is more concentrated towards the ensemble members, as in the  
 613 top row of Figure 24. Assume an extreme case of two tasks and  $p \rightarrow 0$ . Then the distribution  
 614 degenerates to a Bernoulli distribution. Effectively, at each iteration one of the ensemble  
 615 members is selected and its weights are updated, which will result in two separate and inde-  
 616 pendent single-task predictors with no common representation infused about the other task.  
 617 Then, linearly interpolating in weight space will result in models with random predictions  
 618 for both tasks, since the training procedure has not focused in retrieving a Pareto Subspace.  
 619 For milder cases (e.g.  $p = 0.7$ ), we observed that the models in the middle of the linear  
 620 interpolation suffered in performance which can be attributed to the fact that the sampling  
 621 focused more on single-task rather than multi-task representations and performance.
- 622 •  $p > 1$ . Then the distribution is more concentrated towards the midpoint of the simplex, as in  
 623 the bottom row of Figure 24. Assume an extreme case of two tasks and  $p \rightarrow \infty$ . Then, the  
 624 distribution becomes deterministic and outputs equal weights for all tasks. The randomly  
 625 and independently initialized ensemble members will collapse to each other, resulting in  
 626 duplicate ensemble members. Similarly, for very large values (e.g.  $p = 100$ ), the functional  
 627 diversity of the ensemble will suffer since the weights produced by the distribution will be  
 628 almost equal for all tasks, resulting in a milder version of the aforementioned phenomenon.  
 629 In contrast, we found that small values such as  $p = 2$  or  $p = 3$  can help convergence since

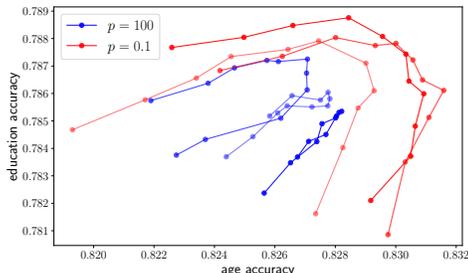
[Reviewer SqFR]: Added ap-  
 pendix regarding sampling.



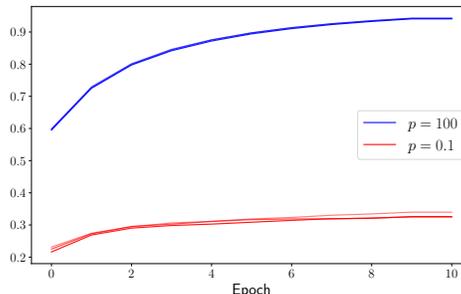
(a) MultiMNIST: Experimental results using three random seeds per method.



(b) MultiMNIST: Cosine similarities of ensemble members.



(c) Census: Experimental results using three random seeds per method.



(d) Census: Cosine similarities of ensemble members.

Figure 25: Experimental results on MultiMNIST and Census varying the concentration parameters  $p = p\mathbf{1}_T$  of the sampling distribution. Three seeds depicted in shades of the same colors for the various  $p$ .

630 they put more emphasis towards common representation (compared to  $p = 1$ ), but may  
 631 limit functional diversity.

632 Figure 25 presents experimental results on MultiMNIST and Census for various concentration  
 633 parameters  $p \in \{0, 0.1, 100\}$  of the Dirichlet distribution. Let  $\theta_1$  and  $\theta_2$  be the parameters of  
 634 the ensemble members. For  $p = 0$ , the ensemble consists of two single-task predictors with no  
 635 multitask learning representational knowledge, since their interpolation meets a low accuracy/high  
 636 loss barrier. We omit the case of  $p = 0$  for Census for visual clarity. This lack of common  
 637 representation is evident in the cosine similarities as well, where for  $p = 0$   $\cos(\theta_1, \theta_2) \approx 0$ . On the  
 638 other hand, for  $p = 0.1$ , common representations are infused into the ensemble and the experimental  
 639 results show that the test performance is characterized by diversity. However, this comes at the  
 640 expense of the interpolated models at the middle of the line segment, where the performance is  
 641 suboptimal compared to  $p = 100$  for MultiMNIST. This behavior is also illustrated in the cosine  
 642 similarities, where for  $p = 100$  the ensemble weights  $\alpha$  are in an  $\epsilon$ -ball around the midpoint causing  
 643 the independently initialized models to progressively collapse. For Census, we also observe that  
 644 this collapsing leads to very high cosine similarity  $\cos(\theta_1, \theta_2) > 0.9$  and the ensemble is suboptimal  
 645 compared to  $p = 0.1$ .

## I CONNECTION BETWEEN PARETO OPTIMALITY AND MULTIPLE VALLEY INTERSECTIONS

646  
647  
648

649 In this section, we investigate the connection between the intersection of multiple loss landscapes, pareto optimality and the effect of the proposed algorithm Pareto Manifold Learning. We use the illustrative example, presented in Figure 1. Let  $\Theta$  be the parameter space of the model and  $\mathcal{L}_t : \Theta \rightarrow \mathbb{R}, t \in \{1, 2\}$ , be the losses of the problem. For  $\alpha \in [0, 1]$  and  $\theta \in \Theta$ , the overall objective is  $\mathcal{L}(\theta, \alpha) = \alpha \mathcal{L}_1(\theta) + (1 - \alpha) \mathcal{L}_2(\theta)$ .

654 Figure 26 and the accompanying Figure 27 present the overall loss objective as  $\alpha$  varies from 0 to 1. For the extreme values of the range, the loss landscape is inherently single-task. The subspace discovered by the method is depicted in blue, while a black 'x' is used for the corresponding interpolated model, i.e., it corresponds to  $\mathcal{L}(\alpha \theta_1 + (1 - \alpha) \theta_2, \alpha)$ . Figure 28 presents the overall losses on the subspace by fixing as a function of one of the parameters. In other words, the proposed method tracks the optimum in parameter space as the overall objective evolves and the various loss landscapes are weighted accordingly. While an acceptable multi-task solution lies in the intersection of low loss landscapes, Pareto Manifold Learning focuses on the aforementioned dynamic scenario of loss weighting.

[Reviewer qexX]: Added appendix regarding weakness 1. For clarity during the rebuttal, this discussion has been added as a standalone appendix. It will be incorporated in the appendix regarding the illustrative example.

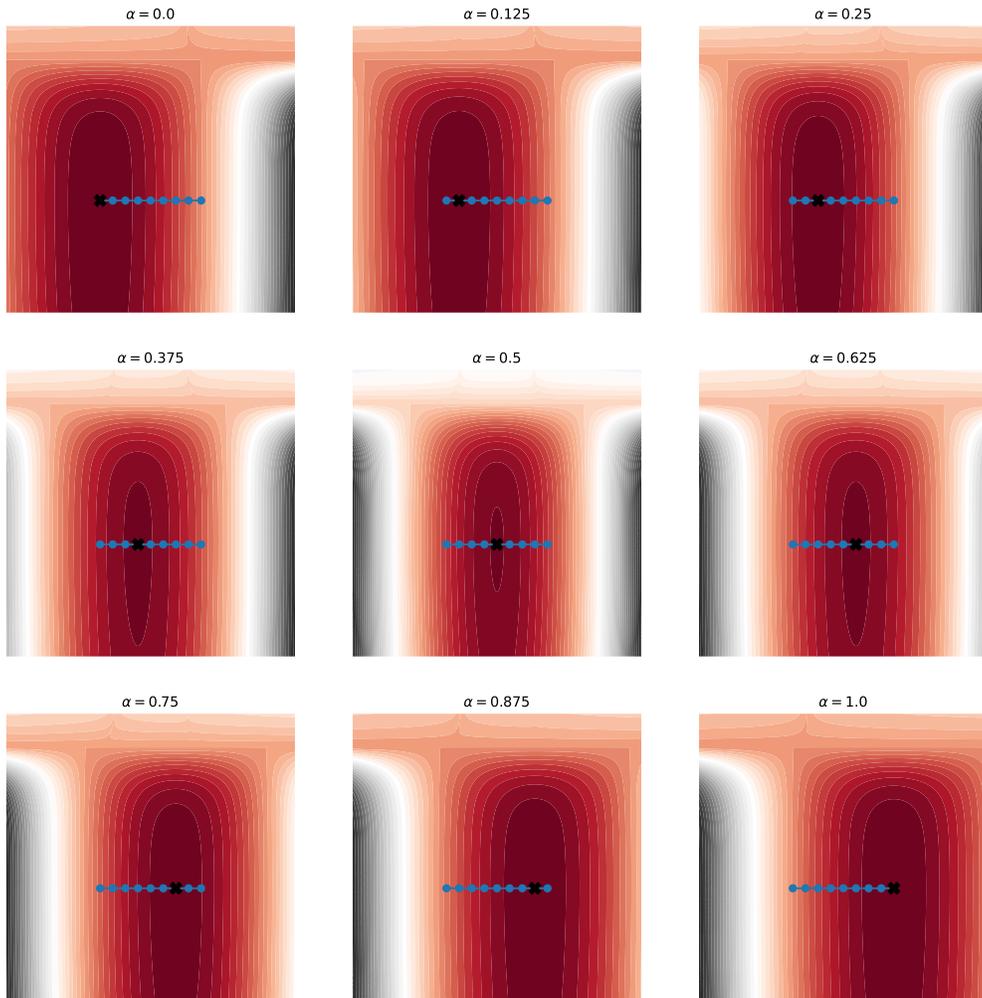


Figure 26: *Illustrative example:* (Overall) loss surface as a function of the model’s weights. The overall objective is  $\mathcal{L}(\theta, \alpha) = \alpha \mathcal{L}_1(\theta) + (1 - \alpha) \mathcal{L}_2(\theta)$  and is shown for various values of  $\alpha$ . The Pareto subspace discovered by the proposed method is depicted in blue. ‘X’ shows the solution of the method for the corresponding  $\alpha$ .

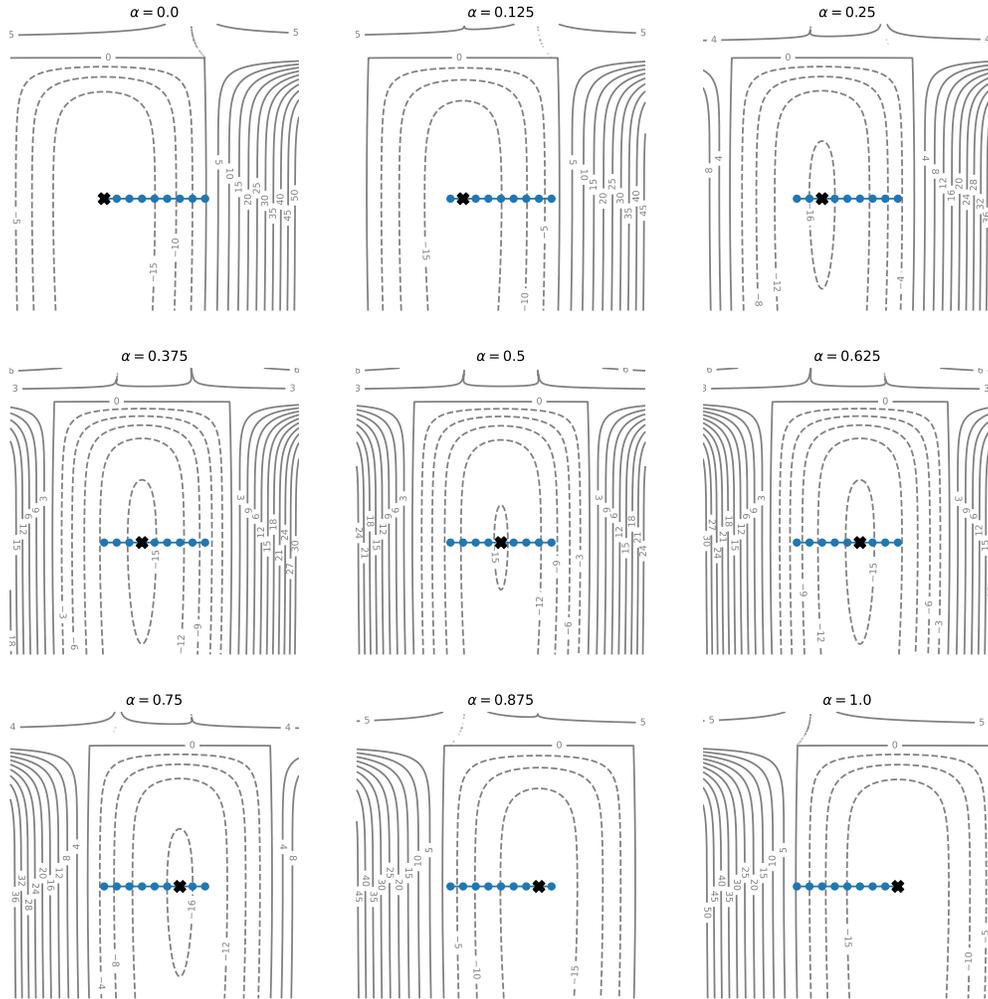


Figure 27: *Illustrative example*: Alternate view of Figure 26. Refer to the text for details.

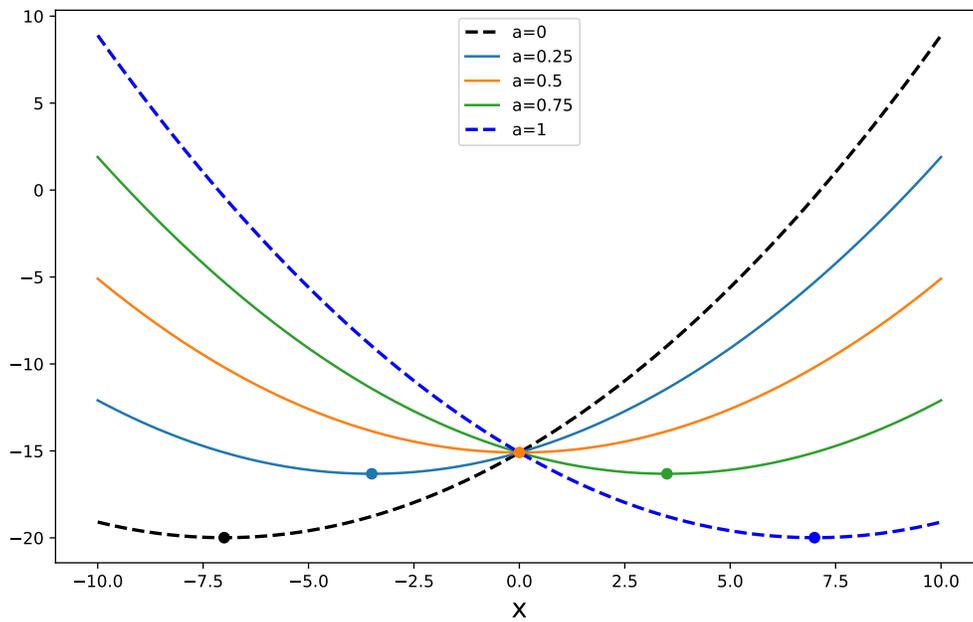


Figure 28: *Illustrative example*: Overall loss for various weightings  $\alpha \in [0, 1]$  as a function of one of the parameters, denoted as  $x$ . Points corresponds to the loss achieved by the parameter vector  $\theta(\alpha) = \alpha\theta_1 + (1 - \alpha)\theta_2$ . The subspace discovered by the model spans the range  $[-7, 7]$ .

## J ADDITIONAL RELATED WORK

663  
664

665 In this appendix, we further expand on prior work. Linear mode connectivity, as in (Wortsman et al.,  
666 2021), encourages flatness and, therefore, is linked with methods explicitly enforcing flat minima  
667 (Chaudhari et al., 2017; Foret et al., 2021; Dinh et al., 2017; Jiang\* et al., 2020). These approaches  
668 are applicable when designing a single objective, e.g. average of losses in Multi-Task Learning,  
669 but do not allow for the infusion of Pareto properties and the inclusion of tradeoffs. Izmailov et al.  
670 (2018) produce flat minima by averaging multiple weight vectors discovered during the optimization  
671 trajectory, so that the final model lies in the middle of the low-loss basin. Wortsman et al. (2022)  
672 perform weight ensembling with fine-tuned models produced via different hyperparameter config-  
673 urations. Apart from the recent weight ensembling works, output ensembling has been one of the  
674 staples of machine learning literature. Lakshminarayanan et al. (2017) utilize deep ensembles for  
675 uncertainty prediction but inference scales linearly with the number of ensemble members. Wen  
676 et al. (2020) improve on the computational complexity of output ensembles by sharing the bulk of  
677 the parameters among members and differentiating them via rank-1 matrices, while Havasi et al.  
678 (2021) employ a multi-input multi-output network by accommodating independent subnetworks for  
679 each ensemble and allowing a single-forward pass ensemble prediction. However, this results in  
680 subnetworks with incompatible architecture which does not allow for a continuous approximation  
681 of the Pareto Front.

[Reviewer NFGo]: Added ap-  
pendix discussing additional  
related work.[Reviewer NFGo]: added  
works on "flat minima"[Reviewer NFGo]: added prior  
work on ensemble learning.

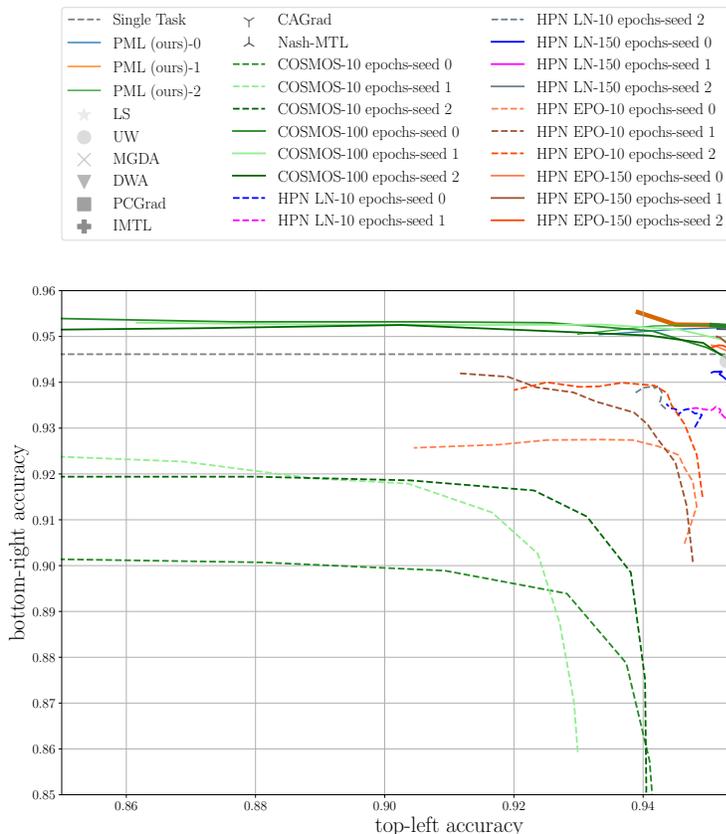


Figure 29: MultiMNIST: Figure 4 with additional baselines.

## 682 K ADDITIONAL EXPERIMENTS

683

684 In this section, we supplement our experimental findings on MultiMNIST with additional base-  
 685 lines, namely HPN-LN and HPN-EPO (Navon et al., 2021) and COSMOS (Ruchte & Grabocka,  
 686 2021)<sup>1</sup>. We use hyperparameters of Ruchte & Grabocka (2021) for both methods. We provide two  
 687 experimental settings:

- 688 • *Setting I*: 10 epochs and no learning rate scheduler, i.e., the setting used for all other meth-  
 689 ods in Figure 4,
- 690 • *Setting II*: the experimental setting used by (Ruchte & Grabocka, 2021), i.e., 100 epochs for  
 691 COSMOS and 150 epochs for HPN-LN/HPN-EPO with multi-step learning scheduler.

692 Figure 29 presents the results with the additional baselines, using three seeds each. We use dashed  
 693 lines for *setting I* and solid lines for *setting II* and group the three methods in various color shades  
 694 (blue, green, red) for visual clarity. We observe that in the original setting of 10 epochs, all new base-  
 695 lines are suboptimal compared to all methodologies. For *setting II*, the hypernetwork methodologies  
 696 are competitive with some baselines but are suboptimal compared to the proposed method. For  
 697 COSMOS, only one seed is competitive with the proposed method. Moreover, HPN-LN, HPN-EPO  
 698 employ a hypernetwork of 1.6m parameters, while the target network has < 50k parameters.

<sup>1</sup>We use the open source implementation provided by Ruchte & Grabocka (2021) making minimal changes. Our implementation of the MultiMNIST dataset has images of size  $28 \times 28$  rather than  $36 \times 36$  resulting in slightly different models.

[All reviewers]: Added appendix with additional baselines.