

---

# RotoGrad: Gradient Homogenization in Multi-Task Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Multi-task learning is being increasingly adopted in applications domains like  
2 computer vision and reinforcement learning. However, optimally exploiting its ad-  
3 vantages remains a major challenge due to the effect of negative transfer. Previous  
4 works have tracked down this issue to the disparities in gradient magnitudes and  
5 directions across tasks, when optimizing the shared network parameters. While  
6 recent work has acknowledged that negative transfer is a two-fold problem, exist-  
7 ing approaches fall short as they focus only on either homogenizing the gradient  
8 magnitude across tasks; or greedily change the gradient directions, overlooking  
9 future conflicts. In this work, we introduce RotoGrad, an algorithm that tackles  
10 negative transfer as a whole: it jointly homogenizes gradient magnitudes and direc-  
11 tions, while ensuring training convergence. We show that RotoGrad outperforms  
12 competing methods in complex problems, including multi-label classification in  
13 CelebA and computer vision tasks in the NYUv2 dataset.

## 14 1 Introduction

15 As neural network architectures get larger in order to solve increasingly more complex tasks, the  
16 idea of jointly learning multiple tasks (for example, depth estimation and semantic segmentation in  
17 computer vision) with a single network is becoming more and more appealing. This is precisely the  
18 idea of multi-task learning (MTL) [3], which promises higher performance in the individual tasks  
19 and better generalization to unseen data, while drastically reducing the number of parameters [27].

20 Unfortunately, sharing parameters between tasks may also lead to difficulties during training as  
21 tasks compete for shared resources, often resulting in poorer results than solving individual tasks, a  
22 phenomenon known as *negative transfer* [27]. Previous works have tracked down this issue to the  
23 two types of differences between task gradients. First, *differences in magnitude* across tasks can make  
24 some tasks dominate the others during the learning process. Several methods have been proposed to  
25 homogenize gradient magnitudes such as MGDA [28], GradNorm [6], or IMTL-G [18]. However,  
26 little attention has been put towards the second source of the problem: *conflicting directions* of the  
27 gradients for different tasks. Due to the way gradients are added up, gradients of different tasks may  
28 cancel each other out if they point to opposite directions of the parameter space, thus leading to a poor  
29 update direction for a subset or even all tasks. Only very recently a handful of works have started to  
30 propose methods to mitigate the conflicting gradients problem, for example, by removing conflicting  
31 parts of the gradients [33], or randomly ‘dropping’ some elements of the gradient vector [7].

32 In this work we propose RotoGrad, an algorithm that tackles negative transfer as a whole by ho-  
33 mogenizing both gradient magnitudes and directions across tasks. RotoGrad addresses the gradient  
34 magnitude discrepancies by re-weighting task gradients at each step of the learning, while encourag-  
35 ing learning those tasks that have converged the least thus far. In that way, it makes sure that no task is  
36 overlooked during training. Additionally, instead of directly modifying gradient directions, RotoGrad

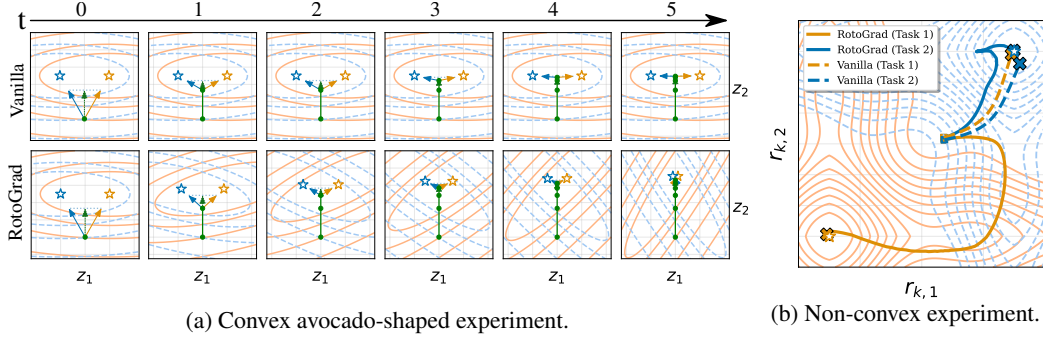


Figure 1: Level plots showing the evolution of two regression MTL problems with/without RotoGrad, see Section 4. RotoGrad is able to reach the optimum ( $\star$ ) for both tasks. (a) In the space of  $z$ , RotoGrad rotates the function-spaces to align task gradients (blue/orange arrows), finding shared features  $z$  (green arrow) closer to the (matched) optima. (b) In the space of  $r_k$ , RotoGrad rotates the shared feature  $z$ , providing per-task features  $r_k$  that better fit each task.

37 smoothly rotates the shared feature space differently for each task, seamlessly aligning gradients in  
 38 the long run. As shown by our theoretical insights, the cooperation between gradient magnitude-  
 39 and direction-homogenization ensures the stability of the overall learning process. Finally, we run  
 40 extensive experiments to empirically demonstrate that RotoGrad leads to stable (convergent) learning,  
 41 scales up to complex network architectures, and outperforms competing methods in multi-label  
 42 classification settings in CIFAR10 and CelebA, as well as in computer vision tasks using the NYUv2  
 43 dataset. Alongside this paper, we will provide a simple-to-use library to include RotoGrad in any  
 44 Pytorch pipeline with a few lines of code.

## 45 2 Multi-task learning and negative transfer

46 The goal of MTL is to simultaneously learn  $K$  different tasks, that is, finding  $K$  mappings from a  
 47 common input dataset  $\mathbf{X} \in \mathbb{R}^{N \times D}$  to a task-specific set of labels  $\mathbf{Y}_k \in \mathbb{Y}_k^N$ . Most settings consider  
 48 a hard-parameter sharing architecture, which is characterized by two components: the *backbone* and  
 49 *heads* networks. The backbone uses a set of shared parameters,  $\theta$ , to transform each input  $x \in \mathbf{X}$   
 50 into a shared intermediate representation  $z = f(x; \theta) \in \mathbb{R}^d$ , where  $d$  is the dimensionality of  $z$ .  
 51 Additionally, each task  $k = 1, 2, \dots, K$  has a head network  $h_k$ , with exclusive parameters  $\phi_k$ , that  
 52 takes this intermediate feature  $z$  and outputs the prediction  $h_k(x) = h_k(z; \phi_k)$  for the corresponding  
 53 task. This architecture is illustrated in Figure 2, where we have added task-specific rotation matrices  
 54  $\mathbf{R}_k$  that will be necessary for the proposed approach, RotoGrad. Note that the general architecture  
 55 described above is equivalent to the one in Figure 2 when all rotations  $\mathbf{R}_k$  correspond to identity  
 56 matrices, such that  $r_k = z$  for all  $k$ .

57 MTL aims to learn the architecture parameters  
 58  $\theta, \phi_1, \phi_2, \dots, \phi_K$  by simultaneously minimizing  
 59 all task losses, that is,  $L_k(h_k(x), \mathbf{y}_k)$  for  
 60  $k = 1, \dots, K$ . Although this is a priori a multi-  
 61 objective optimization problem [28], in practice  
 62 a single surrogate loss consisting of a linear combination  
 63 of the task losses,  $L = \sum_k \omega_k L_k$ , is opti-  
 64 mized. While this approach leads to a simpler  
 65 optimization problem, it may also trigger *negat-  
 66 ive transfer* between tasks, hurting the overall  
 67 MTL performance due to an imbalanced competition among tasks for the shared parameters [27].

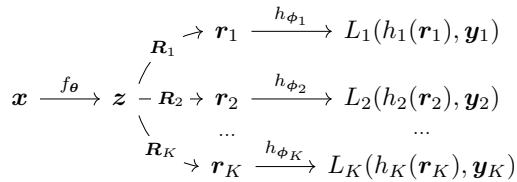


Figure 2: Hard-parameter sharing architecture including the rotation matrices  $\mathbf{R}_k$  of RotoGrad.

68 The negative transfer problem can be studied through the updates of the shared parameters  $\theta$ . At each  
 69 training step,  $\theta$  is updated according to a linear combination of task gradients,  $\nabla_\theta L = \sum_k \omega_k \nabla_\theta L_k$ ,  
 70 which may suffer from two problems. First, **magnitude differences** of the gradients across tasks  
 71 may lead to a subset of tasks dominating the total gradient, and therefore to the model prioritizing  
 72 them over the others. Second, **conflicting directions** of the gradients across tasks may lead to update

73 directions that do not improve any of the tasks. Figure 1 shows an example of poor direction updates  
 74 (left) as well as magnitude dominance (right).

75 In this work, we tackle negative transfer as a whole by homogenizing tasks gradients both in magnitude  
 76 and direction. Note that homogenizing gradients with respect to  $\theta$  is equivalent to homogenizing  
 77 gradients with respect to the shared feature  $z$  due to the chain rule,  $\nabla_{\theta} L_k = \nabla_{\theta} z \cdot \nabla_z L_k$ . Thus,  
 78 from now on we focus on homogenizing the feature-level task gradients  $\nabla_z L_k$ .

### 79 3 RotoGrad

80 In this section we introduce RotoGrad, a novel algorithm that addresses the negative transfer problem  
 81 as a whole. RotoGrad consists of two building blocks which, respectively, homogenize task-gradient  
 82 magnitudes and directions. Moreover, these blocks complement each other and provide convergence  
 83 guarantees of the network training. Next, we detail each of these building blocks and show how they  
 84 are combined towards an effective MTL learning process.

#### 85 3.1 Gradient-magnitude homogenization

86 As discussed in Section 2, we aim to homogenize gradient magnitudes across tasks, as large magnitude  
 87 disparities can lead to a subset of tasks dominating the learning process. Thus, the first goal of  
 88 RotoGrad is to homogenize the magnitude of the gradients across tasks at each step of the training.

89 Let us denote the feature-level task gradient of the  $k$ -th task for the  $n$ -th datapoint, at iteration  $t$ , by  
 90  $\mathbf{g}_{n,k} := \nabla_z L_k(h_k(\mathbf{x}_n), \mathbf{y}_{n,k})$ , and its batch versions by  $\mathbf{G}_k^{\top} := [\mathbf{g}_{1,k}, \mathbf{g}_{2,k}, \dots, \mathbf{g}_{B,k}]$ , where  $B$   
 91 is the batch size. Then, equalizing gradient magnitudes amounts to finding weights  $\omega_k$  that normalize  
 92 and scale each gradient  $\mathbf{G}_k$ , that is,

$$\|\omega_k \mathbf{G}_k\| = \|\omega_i \mathbf{G}_i\| \quad \forall i \iff \omega_k \mathbf{G}_k = \frac{C}{\|\mathbf{G}_k\|} \mathbf{G}_k = C \mathbf{U}_k \quad \forall k, \quad (1)$$

93 where  $\mathbf{U}_k := \frac{\mathbf{G}_k}{\|\mathbf{G}_k\|}$  denotes the normalized task gradient and  $C$  is the target magnitude for all tasks.  
 94 Note that, in the above expression,  $C$  is a free parameter that we need to select.

95 In RotoGrad, we select  $C$  such that all tasks converge at a similar rate. We motivate this choice  
 96 by the fact that, by scaling all gradients, we change their individual step size, interfering with the  
 97 convergence guarantees provided by their Lipschitz-smoothness (for an introduction to non-convex  
 98 optimization see, for example, [25]). Therefore, we seek for the value of  $C$  providing the best  
 99 step-size for those tasks that have converged the least up to iteration  $t$ . Specifically, we set  $C$  to be a  
 100 convex combination of the task-wise gradient magnitudes,  $C := \sum_k \alpha_k \|\mathbf{G}_k\|$ , where the weights  
 101  $\alpha_1, \alpha_2, \dots, \alpha_K$  measure the relative convergence of each task and sum up to one, that is,

$$\alpha_k = \frac{\|\mathbf{G}_k\| / \|\mathbf{G}_k^0\|}{\sum_i \|\mathbf{G}_i\| / \|\mathbf{G}_i^0\|}, \quad (2)$$

102 with  $\mathbf{G}_k^0$  being the initial gradient of the  $k$ -th task, i.e., the gradient at iteration  $t = 0$  of the training.

103 As a result, we obtain a (hyper)parameter-free approach that equalizes the gradient magnitude across  
 104 tasks to encourage learning slow-converging tasks. Note that the resulting approach resembles  
 105 Normalized Gradient Descent (NGD) [8] for single-task learning, which has been proved to quickly  
 106 escape saddle points during optimization [24]. Thus, we expect a similar behavior for RotoGrad,  
 107 where slow-converging tasks will force quick-converging tasks to escape from saddle points.

108 The resulting training algorithm may however diverge as a consequence of constantly oscillating  
 109 between (slow-converging) tasks. For example, in scenarios where one task improves, there is always  
 110 another task(s) that deteriorates. Fortunately, as shown in the following result (proof in Appendix A),  
 111 such a phenomenon does not appear in the absence of conflicting gradients.

112 **Proposition 3.1.** *Let  $\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_K$  be the task gradients with respect to  $\mathbf{Z}$  as defined above. If  
 113  $K = 2$ ; or  $\cos\_sim(\mathbf{G}_i, \mathbf{G}_j) \geq 0$  pairwise; then there exists a small-enough step size  $\varepsilon > 0$  such  
 114 that, for all tasks, we have that  $L_k(h_k(\mathbf{Z} - \varepsilon \cdot C \sum_k \mathbf{U}_k; \phi_k); \mathbf{Y}_k) < L_k(h_k(\mathbf{Z}; \phi_k); \mathbf{Y}_k)$ .*

115 In other words, Proposition 3.1 shows that, when gradients do not conflict in direction with each  
 116 other, following the feature-level gradient  $C \sum_k \mathbf{U}_k$  improves all (lower-bounded) task losses for

117 the given batch. This result, while restricted to the given batch and to the gradient with respect to  
 118 the shared representation  $\mathbf{Z}$ , still provides useful insights in favor of having as *desideratum* of an  
 119 efficient MTL pipeline the absence of conflicting gradients.

### 120 3.2 Gradient-direction homogenization

121 In the previous subsection, we have shown that avoiding conflicting gradients may not only be  
 122 necessary to avoid negative transfer, but also to ensure the stability of the training. In this section  
 123 we introduce the second building block of RotoGrad, an algorithm that homogenizes task-gradient  
 124 directions. The main idea of this approach is to smoothly rotate the feature-space  $\mathbf{z}$  in order to reduce  
 125 the gradient conflict between tasks—in following iterations—of the training by bringing (local)  
 126 optima for different tasks closer to each other (in the parameter space). As a result, it complements  
 127 the previous magnitude-scaling approach and reduces the likelihood of the training to diverge.

128 In order to homogenize gradients, for each task  $k = 1, \dots, K$ , RotoGrad introduces a matrix  $\mathbf{R}_k$   
 129 so that, instead of optimizing  $L_k(\mathbf{z})$  with  $\mathbf{z}$  being the last shared representation, we optimize an  
 130 equivalent loss function  $L_k(\mathbf{R}_k \mathbf{z})$ . As we are only interested in changing directions (not the gradient  
 131 magnitudes), we choose  $\mathbf{R}_k \in SO(d)$  to be a rotation matrix<sup>1</sup> leading to per-task representations  
 132  $\mathbf{r}_k := \mathbf{R}_k \mathbf{z}$ . RotoGrad thus extends the standard MTL architecture by adding task-specific rotations  
 133 before each head, as depicted in Figure 2.

134 Unlike all other network parameters, matrices  $\mathbf{R}_k$  do not seek to reduce their task’s loss. Instead,  
 135 these additional parameters are optimized to reduce the direction conflict of the gradients across  
 136 tasks. To this end, for each task we optimize  $\mathbf{R}_k$  to maximize the batch-wise cosine similarity or,  
 137 equivalently, to minimize

$$\mathcal{L}_{\text{rot}}^k := - \sum_n \langle \mathbf{R}_k^\top \tilde{\mathbf{g}}_{n,k}, \mathbf{v}_n \rangle, \quad (3)$$

138 where  $\tilde{\mathbf{g}}_{n,k} := \nabla_{\mathbf{r}_k} L_k(h_k(\mathbf{x}_n), \mathbf{y}_{n,k})$  (which holds that  $\mathbf{g}_{n,k} = \mathbf{R}_k^\top \tilde{\mathbf{g}}_{n,k}$ ) and  $\mathbf{v}_n$  is the target vector  
 139 that we want all task gradients pointing towards. We set the target vector  $\mathbf{v}_n$  to be the gradient we  
 140 would have followed if all task gradients weighted the same, that is,  $\mathbf{v}_n := \frac{1}{K} \sum_k \mathbf{u}_{n,k}$ , where  $\mathbf{u}_{n,k}$   
 141 is a row vector of the normalized batch gradient matrix  $\mathbf{U}_k$ , as defined before.

142 As a result, in each training step of RotoGrad we simultaneously optimize the following two problems:

$$\text{Network: minimize}_{\boldsymbol{\theta}, \{\phi\}_k} \sum_k \omega_k L_k, \quad \text{Rotation: minimize}_{\{\mathbf{R}_k\}_k} \sum_k \mathcal{L}_{\text{rot}}^k \quad (4)$$

144 The above problem can be interpreted as a Stackelberg game: a two player-game in which *leader*  
 145 and *follower* alternately make moves in order to minimize their respective losses,  $L_l$  and  $L_f$ , and the  
 146 leader knows what will be the follower’s response to their moves. Such an interpretation allows us to  
 147 derive simple guidelines to guarantee training convergence—that is, that the network loss does not  
 148 oscillate as a result of optimizing the two different objectives in Equation 4. Specifically, following  
 149 Fiez et al. [10], we can ensure that problem 4 converges as long as the rotations’ optimizer (leader)  
 150 is a slow-learner compared with the network optimizer (follower). That is, as long as we make the  
 151 rotations’ learning rate decrease faster than that of the network, we know that RotoGrad will converge  
 152 to a local optimum for both objectives. A more extensive discussion can be found in Appendix B.

### 153 3.3 RotoGrad: the full picture

154 After the two main building blocks of RotoGrad, we can now summarize the overall proposed  
 155 approach in Algorithm 1. At each step, RotoGrad first homogenizes the gradient magnitudes such  
 156 that there is no dominant task and the step size is set by the slow-converging tasks. Additionally,  
 157 RotoGrad smoothly updates the rotation matrices—using the local information given by the task  
 158 gradients—to seamlessly align task gradients in the following steps, thus reducing direction conflicts.

### 159 3.4 Practical considerations

160 In this section, we discuss the main practical considerations to account for when implementing  
 161 RotoGrad and propose efficient solutions.

<sup>1</sup>The special orthogonal group,  $SO(d)$ , denotes the set of all (proper) rotation matrices of dimension  $d$ .

---

**Algorithm 1** Training step with RotoGrad

---

**Input** input samples  $\mathbf{X}$ , task labels  $\{\mathbf{Y}_k\}$ , network’s (RotoGrad’s) learning rate  $\eta$  ( $\eta_{\text{roto}}$ )

**Output** backbone (heads) parameters  $\theta$  ( $\{\phi_k\}$ ), RotoGrad’s parameters  $\{\mathbf{R}_k\}$

```
1: compute shared feature  $\mathbf{Z} = f(\mathbf{X}; \theta)$ 
2: for  $k = 1, 2, \dots, K$  do
3:   compute task-specific loss  $L_k = \sum_n L_k(h_k(\mathbf{R}_k \mathbf{z}_n; \phi_k), \mathbf{y}_{n,k})$ 
4:   compute gradient of shared feature  $\mathbf{G}_k = \nabla_{\mathbf{z}} L_k$ 
5:   compute gradient of task-specific feature  $\tilde{\mathbf{G}}_k = \mathbf{R}_k \mathbf{G}_k$   $\triangleright$  Treated as constant w.r.t.  $\mathbf{R}_k$ .
6:   compute unitary gradients  $\mathbf{U}_k = \mathbf{G}_k / \|\mathbf{G}_k\|$ 
7:   compute relative task convergence  $\alpha_k = \|\mathbf{G}_k\| / \|\mathbf{G}_k^0\|$ 
8: end for
9: make  $\{\alpha_k\}$  sum up to one  $[\alpha_1, \alpha_2, \dots, \alpha_K] = [\alpha_1, \alpha_2, \dots, \alpha_K] / \sum_k \alpha_k$ 
10: compute shared magnitude  $C = \sum_k \alpha_k \|\mathbf{G}_k\|$ 
11: update backbone parameters  $\theta = \theta - \eta C \sum_k \mathbf{U}_k$ 
12: compute target vector  $\mathbf{V} = \frac{1}{K} \sum_k \mathbf{U}_k$ 
13: for  $k = 1, 2, \dots, K$  do
14:   compute RotoGrad’s loss  $L_k^{\text{roto}} = -\sum_n \langle \mathbf{R}_k^\top \tilde{\mathbf{g}}_{n,k}, \mathbf{v}_n \rangle$ 
15:   update RotoGrad’s parameters  $\mathbf{R}_k = \mathbf{R}_k - \eta_{\text{roto}} \nabla_{\mathbf{R}_k} L_k^{\text{roto}}$ 
16:   update head’s parameters  $\phi_k = \phi_k - \eta \nabla_{\phi_k} L_k$ 
17: end for
```

---

162 **Unconstrained optimization.** As previously discussed, parameters  $\mathbf{R}_k$  are defined as rotation  
163 matrices, and thus the *Rotation* optimization in problem 4 is a constrained problem. While this would  
164 typically imply using expensive algorithms like Riemannian gradient descent [1], we can leverage  
165 recent work on manifold parametrization [5] and, instead, apply unconstrained optimization methods  
166 by automatically<sup>2</sup> parametrizing  $\mathbf{R}_k$  via exponential maps on the Lie algebra of  $SO(d)$ .

167 **Memory efficiency and time complexity.** Second, as we need one rotation matrix per task, we have  
168 to store  $O(Kd^2)$  additional parameters. In practice, we only need  $Kd(d-1)/2$  parameters due to the  
169 aforementioned parametrization and, in most cases, this amounts to a small part of the total number  
170 of parameters. Moreover, as described by Casado et al. [5], parametrizing  $\mathbf{R}_k$  enables efficient  
171 computations compared with traditional methods, with a time complexity of  $O(d^3)$  independently of  
172 the batch size. In our case, the time complexity is of  $O(Kd^3)$ , which scales better with respect to the  
173 number of tasks than existing methods (for example,  $O(K^2d)$  for PCGrad [33]). Moreover, caching  
174  $\mathbf{R}_k$  in the forward pass and GPU parallelization can further reduce training time.

175 **Scaling-up RotoGrad.** Even though we can efficiently compute and optimize the rotation matrix  $\mathbf{R}_k$ ,  
176 in some application domains, like computer vision, in which the size  $d$  of the shared representation  $\mathbf{z}$   
177 is large, the time complexity for updating the rotation matrix may become comparable to the one of  
178 the network updates. In those cases, we propose to only rotate a subspace of the feature space, that  
179 is, rotate only  $m \ll d$  dimensions of  $\mathbf{z}$ . Then, we can simply apply a transformation of the form  
180  $\mathbf{r}_k = [\mathbf{R}_k \mathbf{z}_{1:m}, \mathbf{z}_{m+1:d}]$ , where  $\mathbf{z}_{a:b}$  denotes the elements of  $\mathbf{z}$  with indexes  $a, a+1, \dots, b$ . While  
181 there exist other possible solutions, such as using block-diagonal rotation matrices  $\mathbf{R}_k$ , we defer them  
182 to future work.

## 183 4 Illustrative examples

184 In this section, we illustrate the behavior of RotoGrad in two synthetic scenarios, providing clean  
185 qualitative results about its effect on the optimization process. Appendix C.1 provides a detailed  
186 description of the experimental setups.

187 To this end, we propose two different multi-task regression problems of the form

$$L(\mathbf{x}) = L_1(\mathbf{x}) + L_2(\mathbf{x}) = \varphi(\mathbf{R}_1 f(\mathbf{x}; \theta), 0) + \varphi(\mathbf{R}_2 f(\mathbf{x}; \theta), 1), \quad (5)$$

188 where  $\varphi$  is a test function with a single global optimum whose position is parametrized by the second  
189 argument, that is, both tasks are identical (and thus related) up to a translation. We use a single input

---

<sup>2</sup>For example, Geotorch [4] makes this transparent to the user.



190  $\mathbf{x} \in \mathbb{R}^2$  and drop task-specific network parameters. As backbone, we take a simple network of the  
191 form  $\mathbf{z} = \mathbf{W}_2 \max(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1, 0) + \mathbf{b}_2$  with  $\mathbf{b}_1 \in \mathbb{R}^{10}$ ,  $\mathbf{b}_2 \in \mathbb{R}^2$ , and  $\mathbf{W}_1, \mathbf{W}_2^\top \in \mathbb{R}^{10 \times 2}$ .

192 For the first experiment we choose a simple (avocado-shaped) convex objective function and, for  
193 the second one, we opt for a non-convex function with several local optima and a single global  
194 optimum. Figure 1 shows the training trajectories in the presence (and absence) of RotoGrad in both  
195 experiments, depicted as level plots in the space of  $\mathbf{z}$  and  $r_k$ , respectively. We can observe that in  
196 the first experiment (Figure 1a), RotoGrad finds both optima—which is in stark contrast to the vanilla  
197 case—by rotating the feature space and matching the (unique) local optima of the tasks. Similarly,  
198 the second experiment (Figure 1b) shows that, as we have two symmetric tasks and a non-equidistant  
199 starting point, in the vanilla case the optimization is dominated by the task with an optimum closest to  
200 the starting point. RotoGrad avoids this behavior by equalizing gradients and, by aligning gradients,  
201 is able to find the optima of both functions.

## 202 5 Related Work

203 Understanding and improving the interaction between tasks is one of the most fundamental problems  
204 of MTL, since any improvement in this regard would translate to all MTL systems. Consequently,  
205 several approaches to address this problem have been adopted in the literature. Among the different  
206 lines of work, the one most related to the present work is gradient homogenization.

207 **Gradient homogenization.** Since the problem is two-fold, there are two main lines of work. On  
208 the one hand, we have task-weighting approaches that focus on alleviating magnitude differences.  
209 Similar to us, GradNorm [6] attempts to learn all tasks at a similar rate, yet they propose to learn  
210 these weights as parameters. Instead, we provide a closed-form solution in Equation 1, and so does  
211 IMTL-G [18]. However, IMTL-G scales all task gradients such that all projections of  $\mathbf{G}$  onto  $\mathbf{G}_k$  are  
212 equal. MGDA [28], instead, adopts an iterative method based on the Frank-Wolfe algorithm in order  
213 to find the set of weights  $\{\omega_k\}$  (with  $\sum_k \omega_k = 1$ ) such that  $\sum_k \omega_k \mathbf{G}_k$  has minimum norm. On the  
214 other hand, recent works have started to put attention on the conflicting direction problem. Maninis  
215 et al. [22] first proposed adversarial training to make task gradients statistically indistinguishable  
216 as part of a bigger image-tailored architecture. More recently, PCGrad [33] proposed to drop the  
217 projection of one task gradient onto another if they are in conflict, whereas GradDrop [7] randomly  
218 drops elements of the task gradients based on a sign-purity score.

219 In the literature, we can also find other approaches which, while orthogonal to the gradient homoge-  
220 nization, are **complementary to our work** and thus could be used along with RotoGrad. Next, we  
221 provide a brief overview of them.

222 A prominent approach for MTL is task clustering, that is, selecting which tasks should be learned  
223 together. This approach dates back to the original task-clustering algorithm [31], but new work in  
224 this direction keeps coming out [29, 35]. Alternative approaches, for example, scale the loss of each  
225 task differently based on different criteria such as task uncertainty [14], task prioritization [11], or  
226 similar loss magnitudes [18]. Moreover, while most models fall into the hard-parameter sharing  
227 umbrella, there exists other architectures in the literature. Soft-parameter sharing architectures [27],  
228 for example, do not have shared parameters but instead impose some kind of shared restrictions to the  
229 entire set of parameters. An interesting approach consists in letting the model itself learn which parts  
230 of the architecture should be used for each of the tasks [12, 23, 30, 32]. Other architectures, such  
231 as MTAN [19], make use of task-specific attention to select relevant features for each task. Finally,  
232 problems triggered by the differences between task gradients (in magnitude and direction) have also  
233 been studied in other domains like meta-learning [34] and continual learning [21].

## 234 6 Experiments

235 In this section we assess the performance of RotoGrad on a wide range of datasets and MTL  
236 architectures. First, we check the effect of the learning rates of the rotation and network updates  
237 on the stability of the learning process of RotoGrad. Then, with the goal of applying RotoGrad  
238 in scenarios with extremely large sizes of  $\mathbf{z}$ , we explore the effect of rotating a subspace of  $\mathbf{z}$   
239 instead of the whole shared representation. Finally, we compare our approach with competing MTL  
240 solutions in the literature, showing that RotoGrad consistently outperforms all other methods. Refer  
241 to Appendix C for a more detailed description of the experimental setups and additional results.

242 **Relative task improvement.** Since MTL uses different metrics for different tasks, throughout this  
 243 section we group results by means of the relative task improvement, first introduced in [22]. Given a  
 244 task  $k$ , and the metrics obtained during test time by our model,  $M_k$ , and by a baseline model,  $S_k$ ,  
 245 which consists of  $K$  networks trained on each task individually, the relative task improvement for the  
 246  $k$ -th task is defined as

$$\Delta_k := 100 \cdot (-1)^{l_k} \frac{M_k - S_k}{S_k}, \quad (6)$$

247 where  $l_k = 1$  if  $M_k < S_k$  means that our model performs better than the baseline in the  $k$ -th task, and  
 248  $l_k = 0$  otherwise. We depict our results using different statistics of  $\Delta_k$  such as its mean ( $\text{avg}_k \Delta_k$ ),  
 249 maximum ( $\text{max}_k \Delta_k$ ), and median ( $\text{med}_k \Delta_k$ ) across tasks.

## 250 6.1 Training stability

251 At the end of Section 3.2 we discussed that, by casting  
 252 problem 4 as a Stackelberg game, we have convergence  
 253 guarantees when the rotation optimizer is the slow-learner.  
 254 Next, we empirically show this necessary condition.

255 **Experimental setup.** Similar to [28], we use a multi-task  
 256 version of MNIST [16] where each image is composed  
 257 of a left and right digit, and use as backbone a reduced  
 258 version of LeNet [17] with light-weight heads. Besides  
 259 the left- and right-digit classification proposed in [28], we  
 260 consider three other quantities to predict: i) sum of digits;  
 261 ii) parity of the digit product; and iii) number of active  
 262 pixels. The idea here is to enforce all digit-related tasks  
 263 to cooperate (positive transfer), while the (orthogonal) image-related task should not disrupt these  
 264 learning dynamics. We use negative cross-entropy and accuracy for the left- and right-digit tasks,  
 265 binary cross-entropy and f1-score for the parity task, and mean squared error (MSE) as loss and  
 266 metric for both regression tasks.

267 **Results.** Figure 3 shows the effect averaged over ten independent runs—in terms of test error in the  
 268 sum task, while the rest of tasks are shown in Appendix C.2—of changing the rotations’ learning rate.  
 269 We can observe that, the bigger the learning rate is in comparison to that of the network’s parameters  
 270 ( $1e-3$ ), the higher and more noisy the test error becomes. MSE keeps decreasing as we lower the  
 271 learning rate, reaching a sweet-spot at half the network’s learning rate ( $5e-4$ ). For smaller values,  
 272 the rotations’ learning is too slow and results start to resemble those of the vanilla case, in which no  
 273 rotations are applied (leftmost box in Figure 3).

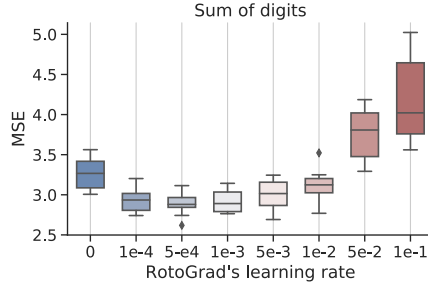


Figure 3: Test error on the sum of digits task for different values of RotoGrad’s learning rate.

## 274 6.2 Rotating a subspace

275 Next, we evaluate the effect of subspace rotations as described at the end of Section 3.4, assessing the  
 276 trade-off between avoiding negative transfer and size of the subspace considered by RotoGrad.

277 **Experimental setup.** We test RotoGrad on a 10-task classification problem on CIFAR10 [15], using  
 278 binary cross-entropy and f1-score as loss and metric, respectively, for all tasks. We use ResNet18 [13]  
 279 without pre-training as backbone ( $d = 512$ ), and linear layers with sigmoids as task-specific heads.

280 **Results** are summarized at the bottom part of Table 1. We can observe that rotating the entire space  
 281 provides the best results, and they worsen as we decrease the size of  $R_k$ . However, rotating only 64  
 282 features (12.5% of the shared feature space) still yields better results than vanilla optimization.

## 283 6.3 Methods comparison

284 We now proceed to compare RotoGrad with the different existing approaches to gradient conflict (for  
 285 both magnitude and direction) in different real-world datasets, showing how RotoGrad outperforms  
 286 existing methods while being on par with existing methods in training time.

287 **Experimental setup.** In order to provide fair comparisons among methods, all experiments use  
 288 identical configurations and random initializations. For all methods we performed a hyper-parameter  
 289 search and chose the best ones based on validation error. Our results are reported using the median and  
 290 standard deviation computed over 5-10 random seeds. Further details can be found in Appendix C.1.

Table 1: Task performance on CIFAR10 for different competing methods (top) and RotoGrad with matrices  $\mathbf{R}_k$  of different sizes (bottom). Table shows median and standard deviation over five runs.

Method	$\text{avg}_k \Delta_k \uparrow$	$\text{med}_k \Delta_k \uparrow$	$\text{max}_k \Delta_k \uparrow$
Vanilla	$2.58 \pm 0.54$	$2.73 \pm 1.37$	$11.14 \pm 3.35$
GradDrop	$3.07 \pm 0.48$	$3.18 \pm 1.07$	$14.03 \pm 2.83$
PCGrad	$2.86 \pm 0.81$	$3.33 \pm 1.68$	$12.01 \pm 3.19$
MGDA	$-1.75 \pm 0.43$	$-4.48 \pm 2.35$	$3.67 \pm 0.98$
GradNorm	$-0.08 \pm 0.95$	$0.09 \pm 2.23$	$8.82 \pm 3.41$
IMTL-G	$2.73 \pm 0.27$	$1.95 \pm 2.21$	$10.20 \pm 2.98$
IMTL-G+ $\mathbf{R}_k$	$3.02 \pm 0.69$	$4.38 \pm 1.11$	$12.76 \pm 1.77$
RotoGrad 64	$2.90 \pm 0.49$	$3.44 \pm 1.51$	$13.16 \pm 2.40$
RotoGrad 128	$2.97 \pm 1.08$	$3.73 \pm 2.14$	$12.64 \pm 3.56$
RotoGrad 256	$3.68 \pm 0.68$	$3.29 \pm 2.18$	$14.01 \pm 3.22$
RotoGrad 512	<b><math>4.48 \pm 0.99</math></b>	<b><math>4.72 \pm 2.84</math></b>	<b><math>15.57 \pm 3.99</math></b>

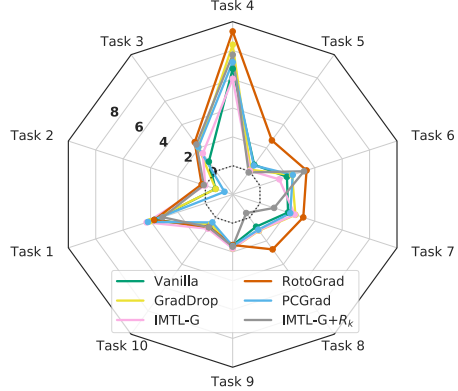


Figure 4: Task improvement (median over five runs) of different methods on CIFAR10. RotoGrad outperforms competing methods on all tasks.

291 **MNIST and SVHN.** We reuse  
 292 the experimental setting from  
 293 Section 6.1—now with multi-  
 294 task versions of MNIST [16] and  
 295 SVHN [26]—in order to evalu-  
 296 ate how disruptive the orthogonal  
 297 image-related task is for differ-  
 298 ent methods. We can observe in  
 299 the results from Table 2 that the  
 300 effect of the image-related task  
 301 is more disruptive in MNIST,  
 302 in which MGDA utterly fails.  
 303 Direction-aware methods (Grad-  
 304 Drop and PCGrad) do not im-  
 305 prove the vanilla results, whereas  
 306 IMTL-G, GradNorm, and RotoGrad obtain the best results.

Table 2: Test performance (median and standard deviation) on two set of unrelated tasks, across ten different runs.

Method	MNIST		SVHN	
	Digits $\text{avg}_k \Delta_k \uparrow$	Act Pix MSE $\downarrow$	Digits $\text{avg}_k \Delta_k \uparrow$	Act Pix MSE $\downarrow$
Single	-	$0.01 \pm 0.01$	-	$0.17 \pm 0.06$
Vanilla	$-2.51 \pm 3.01$	$0.11 \pm 0.01$	$5.14 \pm 0.83$	$2.75 \pm 3.17$
GradDrop	$-2.51 \pm 1.73$	$0.13 \pm 0.02$	$5.68 \pm 1.05$	$1.91 \pm 0.86$
PCGrad	$-3.12 \pm 3.88$	$0.12 \pm 0.02$	$5.50 \pm 0.75$	$2.26 \pm 0.85$
MGDA	$-12.57 \pm 9.97$	$0.06 \pm 0.02$	$5.99 \pm 1.48$	$0.66 \pm 0.75$
GradNorm	$0.13 \pm 2.27$	$0.08 \pm 0.01$	$6.67 \pm 1.02$	$1.41 \pm 0.74$
IMTL-G	$1.17 \pm 2.77$	$0.07 \pm 0.01$	$5.81 \pm 0.85$	$2.47 \pm 1.65$
RotoGrad	$2.12 \pm 2.23$	$0.08 \pm 0.02$	$6.08 \pm 0.48$	$1.61 \pm 2.72$

307 **CIFAR10.** We reuse the setting in Section 6.2 and compare the different MTL methods using five  
 308 different seeds. Results are shown in Table 1 and Figure 4. Unlike the previous setting, scaling  
 309 gradients is not enough to solve the problem. Among existing methods, both direction-aware solutions  
 310 (PCGrad and GradDrop) improve over the vanilla case on all the statistics, whereas most magnitude-  
 311 aware solutions substantially worsen task performance. In stark contrast, RotoGrad improves task  
 312 performance across all ten tasks, as it can be observed both in Table 1 and Figure 4. To further show  
 313 that this is a consequence of gradient homogenization in terms of both magnitudes and directions, we  
 314 introduced an extra-baseline, IMTL-G+ $\mathbf{R}_k$ , which applies IMTL-G to the extended MTL architecture  
 315 (Figure 2), that is, with matrix  $\mathbf{R}_k$  optimizing the  $k$ -th task loss (instead of Equation 3).

316 **NYUv2.** Now, we test all methods using NYUv2 [9] on three different tasks: 13-class semantic  
 317 segmentation; depth estimation; and surface normals. To speed up training, all images were resized  
 318 to  $288 \times 384$  resolution; and data augmentation was applied to alleviate overfitting. As MTL  
 319 architecture, we use SegNet [2] where the decoder is splitted into three convolutional heads. We use  
 320 the same setup as Liu et al. [19]. Like in previous experiments, we observe in Table 3 that RotoGrad  
 321 results in a consistent improvement over all tasks with respect to the vanilla case. MGDA obtains  
 322 the best results in surface normals at the expense of overlooking the other tasks, while GradDrop  
 323 worsens all results and PCGrad obtains minor improvements in all tasks. GradNorm finds a trade-off  
 324 solution instead, improving results in depth estimation and surface normals, yet with worse results in  
 325 semantic segmentation. RotoGrad obtains the best results followed by IMTL-G and, more importantly,  
 326 RotoGrad is the only method resulting in an average positive task improvement—across the three  
 327 tasks—over training three single-task models independently. It is worth mentioning that, with only



Table 3: Results for different methods on the NYUv2 dataset with a SegNet model. RotoGrad obtains the best performance in segmentation and depth tasks on all metrics, while significantly improving the results on normal surfaces with respect to the vanilla case.

Method	Semantic Segmentation $\uparrow$			Depth Estimation $\downarrow$			Surface Normal Within $t^\circ$ $\uparrow$			Hours			
	mIoU	Pix Acc	avg <sub>k</sub> $\Delta_k$ $\uparrow$	Abs Err	Rel Err	avg <sub>k</sub> $\Delta_k$ $\uparrow$	Angle Mean	Distance Median	11.25		22.5	30	avg <sub>k</sub> $\Delta_k$ $\uparrow$
Single Vanilla	0.38	0.63	-	0.59	0.23	-	24.76	18.99	30.11	57.81	69.90	-	11.37
GradDrop	0.37	0.64	-0.62	0.56	0.22	3.68	30.09	26.09	19.74	43.62	57.07	-27.26	<b>3.45</b>
PCGrad	0.37	0.63	-1.55	0.59	0.24	-2.22	30.81	27.19	17.68	41.44	55.15	-31.67	3.55
MGDA	0.39	0.64	1.50	0.54	0.22	4.99	29.85	25.81	19.41	44.02	57.64	-26.68	3.51
GradNorm	0.20	0.51	-32.75	0.73	0.28	-22.33	<b>24.98</b>	<b>19.02</b>	<b>30.57</b>	<b>57.61</b>	<b>69.41</b>	<b>-0.11</b>	3.55
IMTL-G	0.36	0.64	-1.74	0.55	0.23	3.31	25.80	20.30	28.22	54.91	67.21	-5.25	3.54
RotoGrad	0.38	0.65	1.92	0.55	0.23	3.64	26.83	21.96	25.14	51.74	64.76	-11.67	3.60
RotoGrad	<b>0.40</b>	<b>0.66</b>	<b>5.33</b>	<b>0.54</b>	<b>0.20</b>	<b>9.06</b>	26.35	21.27	26.25	53.11	65.99	-8.99	3.85

Table 4: Task f1-score statistics and training hours in CelebA for all competing methods and two different architectures/settings. RotoGrad obtains the best performance in both setups with comparable training time as existing methods.

Method	Convolutional ( $d = 512$ ) task f1-scores (%) $\uparrow$					ResNet18 ( $d = 2048$ ) task f1-scores (%) $\uparrow$				
	min <sub>k</sub>	med <sub>k</sub>	avg <sub>k</sub>	std <sub>k</sub> $\downarrow$	Hours	min <sub>k</sub>	med <sub>k</sub>	avg <sub>k</sub>	std <sub>k</sub> $\downarrow$	Hours
Vanilla	1.62	54.74	58.69	24.18	4.06	15.45	61.52	61.25	22.09	1.49
GradDrop	3.94	55.80	58.62	23.98	4.42	4.46	63.52	63.61	21.79	1.60
PCGrad	2.69	60.30	59.83	23.85	17.03	17.23	61.82	62.74	20.84	5.90
GradNorm	1.83	52.17	54.68	24.94	11.02	14.43	<b>64.10</b>	63.51	21.20	3.59
IMTL-G	3.31	53.05	56.05	26.92	4.90	21.52	62.12	61.98	21.62	1.72
RotoGrad	<b>9.11</b>	<b>62.31</b>	<b>62.45</b>	<b>22.14</b>	11.00	<b>25.72</b>	63.84	<b>65.17</b>	<b>18.99</b>	6.90

three tasks, all methods trained in less than 4 hours; and that this result consolidates RotoGrad’s scalability, as we only rotate the first 1024 dimensions of  $\mathbf{z}$ , out of a total of 7 millions.

**CelebA.** Last, we apply all methods to a 40-class multi-classification problem in CelebA [20] on two different settings: one using a convolutional network as backbone ( $d = 512$ ); and another using ResNet18 [13] as backbone ( $d = 2048$ ). Similar to CIFAR10, we use binary cross-entropy and f1-score as loss and metric for all tasks. Even though we face two completely different architectures, results in Table 4 show that RotoGrad convincingly outperforms all competing methods in all f1-score statistics, independently of the model. Furthermore, since this is a computationally demanding task with 40 tasks—in fact, we omit MGDA as it takes several days to train—we also compare methods in terms of training time. On the one hand, GradDrop and IMTL-G produce little overhead compared with the vanilla case, as expected. On the other hand, GradNorm and PCGrad take, respectively, 2.5 and 4 times longer to train than the vanilla setting. More importantly, RotoGrad outperforms existing methods while staying on par with them in training time, rotating 50% and 75% of the shared feature  $\mathbf{z}$  for the convolutional and residual backbones, respectively, which further demonstrates that RotoGrad can scale-up to real-world settings.

## 7 Conclusions

In this work, we have introduced RotoGrad, an algorithm that tackles negative transfer in MTL by homogenizing task gradients in terms of both magnitudes and directions. RotoGrad enforces a similar convergence rate for all tasks, while at the same time smoothly rotates the shared representation differently for each task in order to avoid conflicting gradients. As a result, RotoGrad leads to stable and accurate MTL. Our empirical results have shown the effectiveness of RotoGrad in many scenarios, staying on top of all competing methods in performance, while being on par in terms of computational complexity with those that better scale to complex networks.

We believe our work opens up interesting venues for future work. For example, it would be interesting to study alternative approaches to further scale up RotoGrad using, for example, diagonal-block or sparse rotation matrices; to rotate the feature space in application domains with structured features (e.g., channel-wise rotations in images); and to combine different methods, for example, by scaling gradients using the direction-awareness of IMTL-G and the “favor slow-learners” policy of RotoGrad.

## References

- 357 [1] Pierre-Antoine Absil, Robert E. Mahony, and Rodolphe Sepulchre. *Optimization Algorithms*  
358 *on Matrix Manifolds*. Princeton University Press, 2008.
- 359 [2] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. “SegNet: A Deep Convolutional  
360 Encoder-Decoder Architecture for Image Segmentation.” In: *IEEE Trans. Pattern Anal. Mach.*  
361 *Intell.* 39.12 (2017), pp. 2481–2495. DOI: [10.1109/TPAMI.2016.2644615](https://doi.org/10.1109/TPAMI.2016.2644615).
- 362 [3] Rich Caruana. “Multitask Learning: A Knowledge-Based Source of Inductive Bias.” In:  
363 *Machine Learning, Proceedings of the Tenth International Conference, University of Mas-*  
364 *sachusetts, Amherst, MA, USA, June 27-29, 1993*. Ed. by Paul E. Utgoff. Morgan Kaufmann,  
365 1993, pp. 41–48. DOI: [10.1016/b978-1-55860-307-3.50012-5](https://doi.org/10.1016/b978-1-55860-307-3.50012-5).
- 366 [4] Mario Lezcano Casado. “Trivializations for Gradient-Based Optimization on Manifolds.”  
367 In: *Advances in Neural Information Processing Systems 32: Annual Conference on Neural*  
368 *Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*  
369 *Canada*. Ed. by Hanna M. Wallach et al. 2019, pp. 9154–9164.
- 370 [5] Mario Lezcano Casado and David Martínez-Rubio. “Cheap Orthogonal Constraints in Neural  
371 Networks: A Simple Parametrization of the Orthogonal and Unitary Group.” In: *Proceedings of*  
372 *Machine Learning Research* 97 (2019). Ed. by Kamalika Chaudhuri and Ruslan Salakhutdinov,  
373 pp. 3794–3803.
- 374 [6] Zhao Chen et al. “GradNorm: Gradient Normalization for Adaptive Loss Balancing in Deep  
375 Multitask Networks.” In: *Proceedings of the 35th International Conference on Machine*  
376 *Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. Ed. by  
377 Jennifer G. Dy and Andreas Krause. Vol. 80. *Proceedings of Machine Learning Research*.  
378 PMLR, 2018, pp. 793–802.
- 379 [7] Zhao Chen et al. “Just Pick a Sign: Optimizing Deep Multitask Models with Gradient Sign  
380 Dropout.” In: *Advances in Neural Information Processing Systems 33: Annual Conference on*  
381 *Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.  
382 Ed. by Hugo Larochelle et al. 2020.
- 383 [8] Jorge Cortés. “Finite-time convergent gradient flows with applications to network consensus.”  
384 In: *Autom.* 42.11 (2006), pp. 1993–2000. DOI: [10.1016/j.automatica.2006.06.015](https://doi.org/10.1016/j.automatica.2006.06.015).
- 385 [9] Camille Couprie et al. “Indoor Semantic Segmentation using depth information.” In: *1st*  
386 *International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA,*  
387 *May 2-4, 2013, Conference Track Proceedings*. Ed. by Yoshua Bengio and Yann LeCun. 2013.
- 388 [10] Tanner Fiez, Benjamin Chasnov, and Lillian J. Ratliff. “Convergence of Learning Dynamics in  
389 Stackelberg Games.” In: *CoRR* abs/1906.01217 (2019).
- 390 [11] Michelle Guo et al. “Dynamic Task Prioritization for Multitask Learning.” In: *Computer*  
391 *Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018,*  
392 *Proceedings, Part XVI*. Ed. by Vittorio Ferrari et al. Vol. 11220. *Lecture Notes in Computer*  
393 *Science*. Springer, 2018, pp. 282–299. DOI: [10.1007/978-3-030-01270-0\\_17](https://doi.org/10.1007/978-3-030-01270-0_17).
- 394 [12] Pengsheng Guo, Chen-Yu Lee, and Daniel Ulbricht. “Learning to Branch for Multi-Task  
395 Learning.” In: *Proceedings of the 37th International Conference on Machine Learning, ICML*  
396 *2020, 13-18 July 2020, Virtual Event*. Vol. 119. *Proceedings of Machine Learning Research*.  
397 PMLR, 2020, pp. 3854–3863.
- 398 [13] Kaiming He et al. “Deep Residual Learning for Image Recognition.” In: *2016 IEEE Conference*  
399 *on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30,*  
400 *2016*. IEEE Computer Society, 2016, pp. 770–778. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- 401 [14] Alex Kendall, Yarin Gal, and Roberto Cipolla. “Multi-task learning using uncertainty to weigh  
402 losses for scene geometry and semantics.” In: *Proceedings of the IEEE conference on computer*  
403 *vision and pattern recognition*. 2018, pp. 7482–7491.
- 404 [15] Alex Krizhevsky, Geoffrey Hinton, et al. “Learning multiple layers of features from tiny  
405 images.” In: (2009).
- 406 [16] Yann LeCun, Corinna Cortes, and CJ Burges. “MNIST handwritten digit database.” In: *ATT*  
407 *Labs [Online]* 2 (2010).
- 408 [17] Yann LeCun et al. “Gradient-based learning applied to document recognition.” In: *Proceedings*  
409 *of the IEEE* 86.11 (1998), pp. 2278–2324.
- 410 [18] Liyang Liu et al. “Towards Impartial Multi-task Learning.” In: *International Conference on*  
411 *Learning Representations*. 2021.

- 412 [19] Shikun Liu, Edward Johns, and Andrew J. Davison. “End-To-End Multi-Task Learning With  
413 Attention.” In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019,*  
414 *Long Beach, CA, USA, June 16-20, 2019.* Computer Vision Foundation / IEEE, 2019, pp. 1871–  
415 1880. DOI: [10.1109/CVPR.2019.00197](https://doi.org/10.1109/CVPR.2019.00197).
- 416 [20] Ziwei Liu et al. “Deep Learning Face Attributes in the Wild.” In: *2015 IEEE International*  
417 *Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015.* IEEE  
418 Computer Society, 2015, pp. 3730–3738. DOI: [10.1109/ICCV.2015.425](https://doi.org/10.1109/ICCV.2015.425).
- 419 [21] David Lopez-Paz and Marc’Aurelio Ranzato. “Gradient Episodic Memory for Continual  
420 Learning.” In: *Advances in Neural Information Processing Systems 30: Annual Conference on*  
421 *Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA.*  
422 Ed. by Isabelle Guyon et al. 2017, pp. 6467–6476.
- 423 [22] Kevis-Kokitsi Maninis, Ilija Radosavovic, and Iasonas Kokkinos. “Attentive Single-Tasking of  
424 Multiple Tasks.” In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*  
425 *2019, Long Beach, CA, USA, June 16-20, 2019.* Computer Vision Foundation / IEEE, 2019,  
426 pp. 1851–1860. DOI: [10.1109/CVPR.2019.00195](https://doi.org/10.1109/CVPR.2019.00195).
- 427 [23] Ishan Misra et al. “Cross-Stitch Networks for Multi-task Learning.” In: *2016 IEEE Conference*  
428 *on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30,*  
429 *2016.* IEEE Computer Society, 2016, pp. 3994–4003. DOI: [10.1109/CVPR.2016.433](https://doi.org/10.1109/CVPR.2016.433).
- 430 [24] Ryan W. Murray, Brian Swenson, and Soumya Kar. “Revisiting Normalized Gradient  
431 Descent: Fast Evasion of Saddle Points.” In: *IEEE Trans. Autom. Control.* 64.11 (2019),  
432 pp. 4818–4824. DOI: [10.1109/TAC.2019.2914998](https://doi.org/10.1109/TAC.2019.2914998).
- 433 [25] Yurii E. Nesterov. *Introductory Lectures on Convex Optimization - A Basic Course.* Vol. 87.  
434 Applied Optimization. Springer, 2004. DOI: [10.1007/978-1-4419-8853-9](https://doi.org/10.1007/978-1-4419-8853-9).
- 435 [26] Yuval Netzer et al. “Reading digits in natural images with unsupervised feature learning.” In:  
436 *NeurIPS Workshop on Deep Learning and Unsupervised Feature Learning* (2011).
- 437 [27] Sebastian Ruder. “An Overview of Multi-Task Learning in Deep Neural Networks.” In: *CoRR*  
438 *abs/1706.05098* (2017).
- 439 [28] Ozan Sener and Vladlen Koltun. “Multi-Task Learning as Multi-Objective Optimization.”  
440 In: *Advances in Neural Information Processing Systems 31: Annual Conference on Neural*  
441 *Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada.*  
442 Ed. by Samy Bengio et al. 2018, pp. 525–536.
- 443 [29] Trevor Standley et al. “Which Tasks Should Be Learned Together in Multi-task Learning?” In:  
444 *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18*  
445 *July 2020, Virtual Event.* Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020,  
446 pp. 9120–9132.
- 447 [30] Ximeng Sun et al. “AdaShare: Learning What To Share For Efficient Deep Multi-Task Learn-  
448 ing.” In: *Advances in Neural Information Processing Systems 33: Annual Conference on Neural*  
449 *Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.* Ed. by  
450 Hugo Larochelle et al. 2020.
- 451 [31] Sebastian Thrun and Joseph O’Sullivan. “Discovering Structure in Multiple Learning Tasks:  
452 The TC Algorithm.” In: *Machine Learning, Proceedings of the Thirteenth International*  
453 *Conference (ICML ’96), Bari, Italy, July 3-6, 1996.* Ed. by Lorenza Saitta. Morgan Kaufmann,  
454 1996, pp. 489–497.
- 455 [32] Simon Vandenhende et al. “Branched Multi-Task Networks: Deciding what layers to share.”  
456 In: *31st British Machine Vision Conference 2020, BMVC 2020, Virtual Event, UK, September*  
457 *7-10, 2020.* BMVA Press, 2020.
- 458 [33] Tianhe Yu et al. “Gradient Surgery for Multi-Task Learning.” In: *Advances in Neural Informa-*  
459 *tion Processing Systems.* Ed. by H. Larochelle et al. Vol. 33. Curran Associates, Inc., 2020,  
460 pp. 5824–5836.
- 461 [34] Tianhe Yu et al. “Meta-World: A Benchmark and Evaluation for Multi-Task and Meta Re-  
462 inforcement Learning.” In: *3rd Annual Conference on Robot Learning, CoRL 2019, Osaka,*  
463 *Japan, October 30 - November 1, 2019, Proceedings.* Ed. by Leslie Pack Kaelbling, Danica  
464 Kragic, and Komei Sugiura. Vol. 100. Proceedings of Machine Learning Research. PMLR,  
465 2019, pp. 1094–1100.

466 [35] Amir Roshan Zamir et al. “Taskonomy: Disentangling Task Transfer Learning.” In: *2018 IEEE*  
467 *Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT,*  
468 *USA, June 18-22, 2018*. IEEE Computer Society, 2018, pp. 3712–3722. DOI: [10.1109/CVPR.](https://doi.org/10.1109/CVPR.2018.00391)  
469 [2018.00391](https://doi.org/10.1109/CVPR.2018.00391).

## 470 Checklist

- 471 1. For all authors...
  - 472 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
473 contributions and scope? [Yes] RotoGrad homogenizes both magnitudes (Equation 1)  
474 and direction (Equation 4), and empirical results in Section 6 and Appendix C demon-  
475 strate our claims.
  - 476 (b) Did you describe the limitations of your work? [Yes] In Section 6.1.
  - 477 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
  - 478 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
479 them? [Yes]
- 480 2. If you are including theoretical results...
  - 481 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Proposi-  
482 tion 3.1’s assumptions are stated in Appendix A, and those regarding RotoGrad’s  
483 stability appear in Appendix B.
  - 484 (b) Did you include complete proofs of all theoretical results? [Yes] Proof of Proposi-  
485 tion 3.1 appears in Appendix A. For the proofs related to Stackelberg games, which are  
486 not a direct contribution of our paper, please refer to [10].
- 487 3. If you ran experiments...
  - 488 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
489 imental results (either in the supplemental material or as a URL)? [Yes] We provide  
490 instructions and code to reproduce our experiments in the supplemental material.
  - 491 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
492 were chosen)? [Yes] In Appendix C.1.
  - 493 (c) Did you report error bars (e.g., with respect to the random seed after running ex-  
494 periments multiple times)? [Yes] We provide statistics for most of our experiments  
495 computed over 5-10 independent runs. Due to time complexity required by larger  
496 datasets on NYUv2 and CelebA, we only report the results for a single random seed,  
497 but still compare the different methods using several performance metrics.
  - 498 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
499 type of GPUs, internal cluster, or cloud provider)? [Yes] All details are provided in  
500 Appendix C.1.
- 501 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - 502 (a) If your work uses existing assets, did you cite the creators? [Yes]
  - 503 (b) Did you mention the license of the assets? [Yes] We only use code from previous  
504 research and licence MIT, which we inherit and acknowledge in our extended version  
505 of the code.
  - 506 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
507 We release the code implementation to reproduce our experiments together with the  
508 supplementary material, and will make it publicly available after the paper acceptance.
  - 509 (d) Did you discuss whether and how consent was obtained from people whose data  
510 you’re using/curating? [N/A] We only use publicly available datasets with no personal  
511 information. Moreover, our experiments only report statistics on the results.
  - 512 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
513 information or offensive content? [N/A] We only use publicly available and broadly  
514 used image datasets.
- 515 5. If you used crowdsourcing or conducted research with human subjects...
  - 516 (a) Did you include the full text of instructions given to participants and screenshots, if  
517 applicable? [N/A]

518  
519  
520  
521

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]