

---

# Does MAML Only Work via Feature Re-use?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 It has been recently observed that a good embedding is all we need to solve many  
2 few-shot learning benchmarks. In addition, other work has strongly suggested  
3 that Model Agnostic Meta-Learning (MAML) mostly works via this same method  
4 – by learning a good embedding. This highlights our lack of understanding of  
5 what meta-learning algorithms are doing and when they work. In this work we  
6 provide empirical results that shed some light towards understanding meta-learning  
7 algorithms better. In particular we identify three interesting properties: 1) In  
8 contrast to previous work, we show that it is possible to define a family of synthetic  
9 benchmarks that result in a low degree of feature re-use – suggesting that current  
10 few-shot learning benchmarks *might not have the properties* needed for the success  
11 of meta-learning algorithms; 2) meta-overfitting occurs when the number of classes  
12 (or concepts) are finite, and this issue disappears once the task has an *unbounded*  
13 number of concepts (e.g. online learning); 3) more adaptation at meta-test time  
14 with MAML does not necessarily result in a significant representation change  
15 or even an improvement in meta-test performance – even when training on our  
16 proposed synthetic benchmarks. Finally, we suggest that, to understand meta-  
17 learning algorithms better, it is imperative that we go beyond tracking only absolute  
18 performance and in addition formally quantify the degree of meta-learning and  
19 track both metrics together. Reporting results in future work this way will help  
20 us identify the sources of meta-overfitting more accurately, and hopefully design  
21 more flexible meta-learning algorithms that learn beyond fixed feature re-use.

## 22 1 Introduction

23 Few-shot learning is a research challenge that assesses an artificial intelligence (AI) model’s capacity  
24 to quickly adapt to new tasks or new environments. This has been the leading area where AI  
25 researchers apply meta-learning algorithms - where a strategy that learns to learn quickly is likely to  
26 be the most promising. However, it was recently shown by Tian et al. [20] that a model that only has  
27 a good embedding is able to match and beat many modern sophisticated meta-learning algorithms. In  
28 addition, there seems to be growing evidence that this is a real phenomena [1, 3, 7, 13]. Furthermore,  
29 analysis of the representations learned by Model Agnostic Meta-Learning (MAML) [9] (on few-shot  
30 learning tasks) revealed that MAML mainly works by learning a feature that is re-usable for many  
31 tasks [19] – what we are calling a good embedding in this paper.

32 These discoveries reveal a lack of understanding on when and why meta-learning algorithms work  
33 and are the main motivations for this work. In particular our contributions are:

- 34 1. It is possible to define a synthetic task that results in lower degree of feature re-use, thus  
35 suggesting that current few-shot learning benchmarks might not have the properties needed  
36 for the success of meta-learning algorithms;

- 37 2. Meta-overfitting occurs when the number of classes (or concepts) are finite, and the issue  
38 disappears once the tasks have an unbounded number of concepts;  
39 3. More adaptation for MAML does not necessarily result in representations that change  
40 significantly or even perform better at meta-test time.

## 41 2 Unified Framework for Studying Meta-Learning and Absolute 42 Performance

43 We propose that future work on meta-learning should not only report absolute performance, but also  
44 quantify and report the degree of meta-learning. In addition, for us to be able to understand and trust  
45 such a system, we need metrics that can diagnose basic issues, e.g. if the system is meta-overfitting  
46 – defined as when the system has a high degree of meta-learning *coupled* with a high gap between  
47 meta-train and meta-test errors.

48 In this work, we make an important first step, inspired by [19], that defines the degree of meta-learning  
49 by measuring the normalized degree of change in the representation of a neural network  $nn_\theta$  after  
50 using meta-adaptation  $A$ :

$$ML(nn_\theta) = \text{Diff}(nn_\theta, A(nn_\theta)). \quad (1)$$

51 In this work we set  $ML(nn_\theta)$  to be distance based Canonical Correlation Analysis (dCCA) [17].  
52 Note that dCCA is simply 1 minus CCA to switch the similarity based metric to a difference based  
53 metric is between 0 and 1.

## 54 3 Benchmarks that Require Meta-Learning

### 55 3.1 Background

56 **Model-Agnostic Meta-Learning (MAML).** The MAML algorithm [9] attempts to meta-learn an  
57 initialization of parameters for a neural network that is primed for quick gradient descent adaptation.  
58 It consists of two main optimization loops: 1) an outer loop used to prime the parameters for fast  
59 adaptation, and 2) an inner loop that does the fast adaptation. During meta-testing, only the inner  
60 loop is used to adapt the representation learned by the outer loop.

61 **Feature re-use.** In the context of MAML, this term usually means that the inner loop provides little  
62 adaptation during meta-testing, when solving an unseen task. In particular, Raghu et al. [19] showed  
63 that MAML has little representation change as measured with CCA and CKA after adaptation, during  
64 meta-testing on the mini-ImageNet few-shot learning benchmark.

### 65 3.2 Motivation for Our Work

66 The analysis by Raghu et al. [19] showing that MAML works mainly by feature re-use is the  
67 main motivation for our work. However, we argue that their conclusion is highly dependent on  
68 the benchmark used. This motivates us to construct a different benchmark and show that by *only*  
69 constructing a different benchmark, we can exhibit lower degrees of feature re-use in a statistically  
70 significant way. Therefore, our goal will be to show a lower degree of feature re-use than them.  
71 In particular, their work [19] showed that the representation layer of a neural network trained with  
72 MAML had a dCCA of  $0.1 \pm 0.02$  [19]. *Therefore, our concrete goal will be to show that the dCCA*  
73 *on our task is greater than 0.12.* If this is achieved, it is good evidence that this new benchmark  
74 benefits from meta-learning and can be detected at a higher degree than previous work [19] in a  
75 statistically significant way. This is our main result of this section and is discussed in detail in Section  
76 3.3.3.

### 77 3.3 Synthetic Task that Requires Meta-learning

#### 78 3.3.1 Overview and Goal

79 The main idea is to sample functions to be approximated, such that the final layer needs little or no  
80 adaptation but the feature layers require a large amount of adaptation. This type of task would forcibly  
81 require that the meta-learner at least learns a representation that needs the feature layers to change to

82 achieve good meta-test performance (i.e. it cannot rely solely on feature re-use). Therefore, to perform  
 83 well, not only would it be good to adapt the representation layers, but additionally performance is  
 84 likely to be obtained from a (meta-learned) initialization that is primed to changed flexibly. In other  
 85 words, tasks must not all have the same shared representation to be solved for meta-learning to be  
 86 most useful and detectable.

### 87 3.3.2 Definition

88 In this section we describe a family of benchmarks that exhibits detectable meta-learning and requires  
 89 more than a re-usable representation layer to be solved. We propose a set of regression functions  
 90 specified as a fully connected neural network (FCNN), such that the magnitude of parameters of the  
 91 representation are larger than the head. In particular we sample the parameters of the representation  
 92 layer from a Gaussian with a larger standard deviation, compared to the parameter sampling of the  
 93 head. We define the representation layer to be the first  $L - 1$  layers, and the head to be the final layer.

94 Next we describe the process to sample one function (regression task) from a Gaussian distribution.  
 95 We have two pairs of benchmark parameters  $[(\mu^{(1)}, \sigma^{(1)}), (\mu^{(2)}, \sigma^{(2)})]$ :  $(\mu^{(1)}, \sigma^{(1)})$  to sample the  
 96 parameters for the representation layer, and  $(\mu^{(2)}, \sigma^{(2)})$  to sample the parameters for the final layer.  
 97 Then each regression task  $f^{(t)}$  (with index  $t$ ) is sampled as follows:

- 98 • Sample the representation parameters  $w^{(l)} \sim N(\mu^{(1)}, \sigma^{(1)})$  for each layer  $l \in [L - 1]$  in  
 99 the representation layers
- 100 • Sample the final layer parameters  $w^{(L)} \sim N(\mu^{(2)}, \sigma^{(2)})$

101 The idea is that for some  $c \in \mathbb{R}$  we have  $\sigma^{(1)} > c \cdot \sigma^{(2)}$  such that the variance in tasks is due to  
 102 the representation layers, and therefore adapting the representation layers is necessary. For all our  
 103 experiments  $\sigma^{(2)} = 1.0$ . An example task can be seen in Figure 1. During meta-training, points  
 104 are uniformly sampled from  $[-1, 1]$ , and the standard support set and query set are constructed by  
 105 computing  $f_w^{(t)}(x)$ .

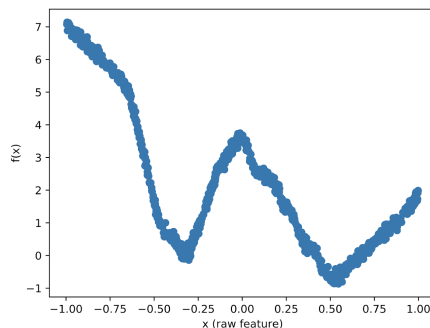


Figure 1: An example regression task constructed as described in Section 3.3.2. Addressing such tasks requires high degree of meta-learning.

### 106 3.3.3 Results on Benchmarks that Require Meta-Learning

107 In this section we show a higher degree of meta-learning and a lower degree of feature re-use from an  
 108 initialization trained with MAML on the benchmarks described in Section 3.3.2. In particular we  
 109 show this in Figure 2 because the dCCA value exhibited is much larger than 0.12 of previous work  
 110 [19]. *Most importantly, the results are statistically significant, because the error bars do not intersect*  
 111 *with the red dotted line with (worst case) dCCA value of 0.12.* The red dotted line is the top error  
 112 band of previous work - i.e. the mean plus the standard deviation.

113 Note that a dCCA higher than 0.12 was observed across all of our experiments in over sixteen different  
 114 benchmarks. In particular this happened even in models that had meta-overfitted e.g. see Figure 3.  
 115 This is strongly suggestive that the benchmarks we defined in Section 3.3.2 require meta-learning,  
 116 since they do not solely rely on feature re-use to be solved.

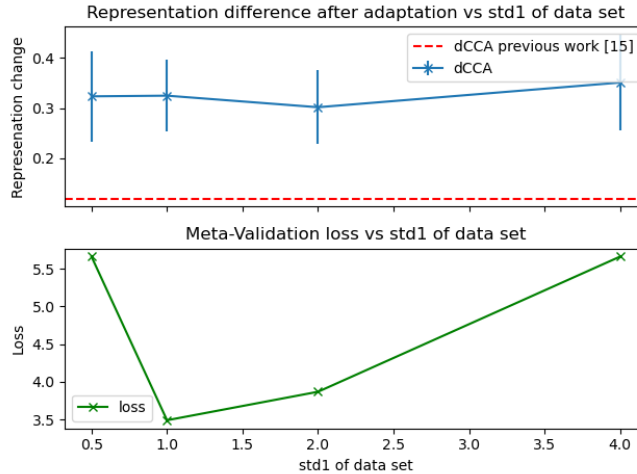


Figure 2: Shows the of lack of feature re-use and a higher degree of meta-learning, as the standard deviation of the representation layer  $\sigma^{(1)}$  for generating regression. The x-axis is the standard deviation (std) of the parameter  $\sigma^{(1)}$  for generating the tasks for the data sets. The models used for each point in the plot are models selected from early stopping (using the meta-validation MSE loss) when meta-trained with MAML. The models are the same architecture as the target function (4 layers fully connected neural network) with ReLU activation function. We also show the meta-validation loss vs the standard deviation of the task. The dCCA was computed by from the average and standard deviation over the representation layers, in this case the first three layers. The average is across different runs using the same meta-learned initialization. The red dotted line shows the value of 0.12 that our models have to be statistically significant. The only difference of this figure with respect to figure 3 is that we selected a model with the best validation here and in the figure 3 we selected the model in last step.

## 117 4 Meta-Overfitting

118 In this section, we show how being armed with the additional metric discussed in Section 2, we are  
 119 able to identify an increasing gap between the meta-test and meta-train losses/accuracy – a term we  
 120 refer to as *meta-overfitting*. In particular, this phenomena is observed when we meta-train models  
 121 with MAML, and becomes more pronounced as the number of iterations increases. We attribute  
 122 this to the adaptation, because this increase in the meta-generalization gap is observed in conjunction  
 123 to the low degree of feature re-use (as discussed in Section 3.3.3), and is most noticeable in our  
 124 synthetic benchmarks compared to in mini-ImageNet [19]. Note that the dCCA of the models was  
 125 much larger in our synthetic benchmarks than in mini-ImageNet. In addition, we show that if the  
 126 number of regression tasks (in this case functions) is not fixed, then the meta-overfitting issue is no  
 127 longer observed.

### 128 4.1 Finite Number of Tasks

129 When the number of regression tasks (functions) is finite (200 in our experiments), we consistently  
 130 observe meta-overfitting. We show this in Figure 4 by increasing the meta-generalization gap (i.e. an  
 131 increase in the difference between the meta-train and the meta-validation losses). This is consistently  
 132 observed in over 30 experiments with a finite number of regression tasks.

133 Furthermore, meta-overfitting is also observed in a few-shot image recognition benchmark. This is  
 134 shown in Figure 5 with mini-ImageNet. With a Pytorch ResNet-18 model, one can observe a meta-  
 135 generalization gap of about 30%. With a state-of-the-art ResNet-12 [20], the meta-generalization gap  
 136 is instead about 20%.

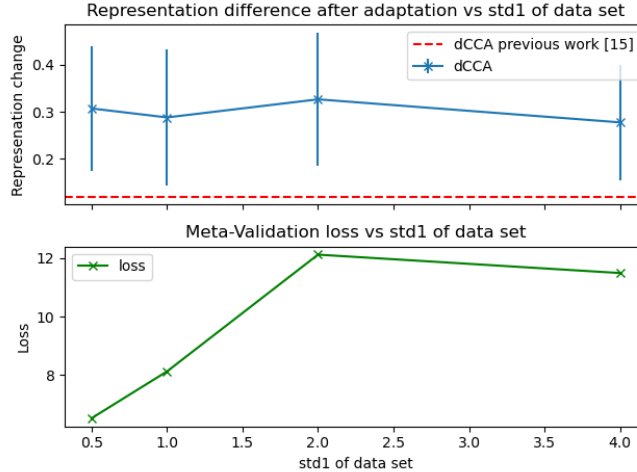


Figure 3: This figure supports the main result of the paper because a higher degree of meta-learning and a lack of feature re-use are present – even in models that are meta-overfitted. A meta-overfitted model can be easily obtained in our experiments by selecting a model at the final iteration. The x-axis is the standard deviation (std) of the parameter  $\sigma^{(1)}$  for generating the tasks for the data sets. The red dotted line shows the value of 0.12 that our models have to be above for statistically significant results that support our claims. The only difference of this figure with respect to figure 2 is that we selected a model in last step (after trough and it had meta-overfitted) while in in figure 2 we select the model with lowest meta-validation loss.

#### 137 4.2 Infinite Number of Tasks

138 We believe it is important to highlight that meta-overfitting was not observed when the number of  
 139 regression tasks is unbounded, as shown in Figure 6. This suggests that, when the number of tasks  
 140 is unbounded but sampled from a related set of tasks, meta-learning algorithms can leverage their  
 141 power to adapt without meta-overfitting.

142 To measure the amount of meta-learning and the lack of feature re-use, we compute the dCCA value  
 143 of the model as in Section 3.3.3 and observe a value of  $0.31 \pm 0.11$ . This also implies that the degree  
 144 of meta-learning is higher when the number of tasks is unbounded.

### 145 5 Effects of More Meta-Adaptation

146 In this section, we show that increasing the number of inner steps for MAML during adaptation does  
 147 not necessarily change the representation further as measured with dCCA (as in Equation 1). In  
 148 addition, the meta-validation performance also does not change.

149 To show this, we obtain a single neural network meta-trained with MAML using a dataset as described  
 150 in Section 3.3.2. Then we plot how the representation changes and how the meta-validation error  
 151 changes as a function of the inner steps. We show this in Figures 7 and 8. We observe that the MAML  
 152 neural networks are robust to meta-overfitting with respect to the inner steps of its inner adaptation  
 153 rule.

154 Note that this is different from what was observed in Section 4.1, because that section shows it as a  
 155 function of the meta iterations (what is sometimes called outer iterations). In addition, it is important  
 156 to emphasize that the representation change in the plots is above the 0.12 compared to previous work  
 157 [19], supporting the main results of section 3.3.3.

### 158 6 Related Work

159 Oh et al. [18] shows that one can encourage models to use less feature re-use purely algorithmically  
 160 by setting the inner learning rate to zero for the final layer. They showed BOIL outperforms MAML

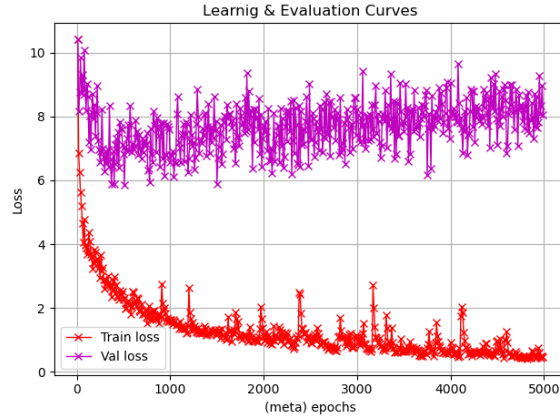


Figure 4: Shows meta-overfitting when the number of tasks (functions) is finite at 200 regression tasks because the meta-validation loss increases as the meta-train loss decreases. In particular the dCCA for this models was  $0.36 \pm 0.12$  corresponding to  $\sigma^{(1)} = 1.0$ . The plot is the learning curve for a 4-layered fully connected neural network trained with MAML [9] using episodic meta-learning. Note that we use a (large) meta-batch size of 75 to decrease the noise during training in the figure. The main difference of this figure with figure 6 is that in this one has a finite set of tasks using our synthetic benchmark while the other has an infinite set of tasks using the sinusoid tasks suggested in [9].

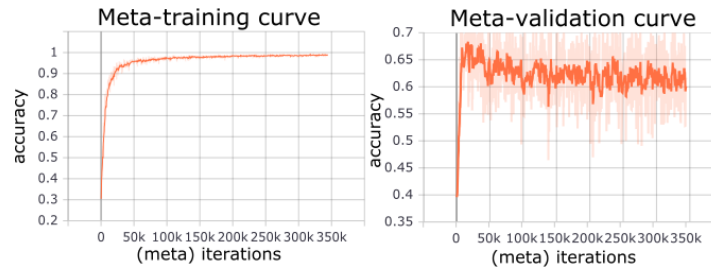


Figure 5: Shows that meta-overfitting is a real phenomenon in mini-Imagenet. We interpret this due to the peak in the meta-validation accuracy followed by a decline as the number of iterations increases. Importantly, the meta-train accuracy continues to increase as it converges. The model trained is an out-of-the-box PyTorch ResNet-18. Note that the higher noise of the meta-validation accuracy is due to having a meta-batch size of 2 to speed up experiments. We smoothed the meta-validation curve with a tensorboard smoothing weight of 0.8. We consistently saw that increases in meta-batch size lead to decreases in noise in the learning curves but we didn't re-run these experiments since it can take up to a week to reproduce a episodic meta-learning run - even on a Quadro RTX 6000.

161 in both traditional few-shot learning (e.g. meta-trained on mini-ImageNet then meta-tested on mini-  
 162 ImageNet) and cross-domain few-shot learning (meta-trained on mini-ImageNet then meta-tested  
 163 on tiered ImageNet). In particular, their cross-domain few-shot learning is similar in spirit to the  
 164 synthetic task we propose in section 3.3.2. However, note that we show that even MAML (and  
 165 algorithm that has been shown to work by feature re-use [19, 18]) can exhibit large representation  
 166 changes if it is trained solely on a task that requires large feature changes. Concisely, we encourage  
 167 rapid learning by only changing the task while Oh et al [18] do encourage it by changing the algorithm  
 168 itself.

169 Guo et al.'s [11] work is similar to ours in that they focus on defining a benchmark more appropriate  
 170 for meta-learning and transfer learning. They propose that meta-learning should be done in a fashion  
 171 where the distribution of tasks sampled changes considerably when moving from meta-training to  
 172 meta-evaluation. Our work is different in that we emphasize more that the meta-training tasks  
 173 themselves need to have diversity to be able to encourage meta-learning. Although Guo et al.'s [11]  
 174 meta-evaluation procedure is excellent, we believe - based on our results - that their benchmark won't

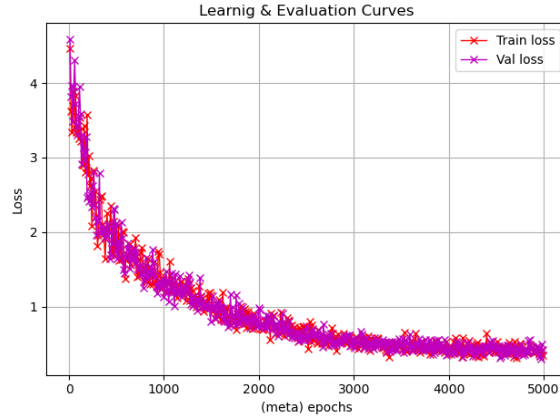


Figure 6: Shows that meta-overfitting does not occur and perfect meta-generalization occurs when the number of tasks (functions) is unbounded when training with MAML. In other words the meta-train and meta-validation error are indistinguishable and decrease together as the meta-iterations increases. The main difference of this figure with figure 4 is that in this one has a finite set of tasks using our synthetic benchmark while the other has an infinite set of tasks using the sinusoid tasks suggested in [9].

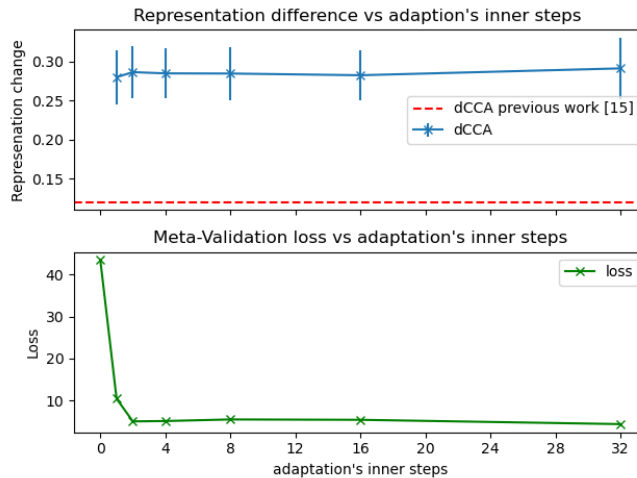


Figure 7: Shows 1) the lack of representation change and b) meta-validation change as the number of inner steps increases. 1 is shown by the relative flatness of the blue and orange lines in the upper plot. Similarly 2 is shown by the flatness of the green line in the lower plot. In particular notice that we exponentially increase the inner steps from 1 to 2 to 32. The models used are 4 layered FCNN tained with MAML with 1 inner step and 0.1 inner learning rate, selected using early stopping using the meta-validation set with the Sigmoid activation function. The only difference of this figure with figure 8 is that this figure uses a sigmoid activation and the other one uses a ReLU. Note that this is the model used for figure 4. Note the dCCA value remains above 0.12 suggesting lower degree of feature re-use.

175 have enough diversity to encourage large representation changes during meta-training. However, we  
 176 believe that combining our ideas and their's to design one benchmark is a promising step for creating  
 177 a better benchmark for meta-learning.

178 Similar work by Triantafillou et al. [21] attempt to improve benchmarks by merging more data sets  
 179 but we believe their data sets are not diverse enough to achieve this. In terms of methods our work  
 180 is most similar to Raghu et al. [19], but they lack an analysis of the role of the tasks in explaining  
 181 their observations. There is also other work [1, 3, 7, 13] that shows that a good representation is

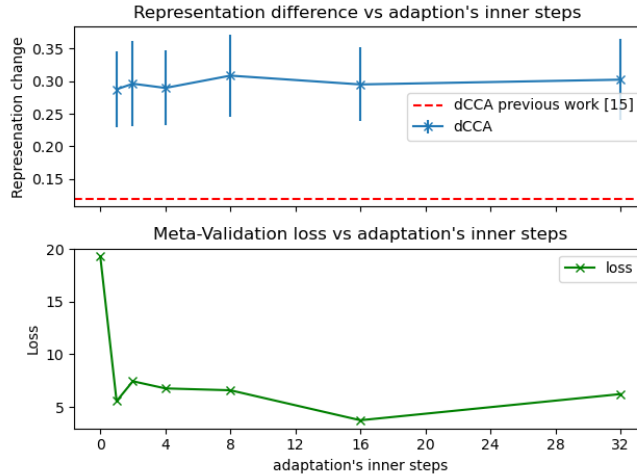


Figure 8: Shows 1) the lack of representation change and b) meta-validation change as the number of inner steps increases. 1 is shown by the relative flatness of the blue and orange lines in the upper plot. Similarly, 2 is shown by the flatness of the green line in lower figure. We want to emphasize that we exponentially increase the inner steps from 1 to 2 to 32. The models used are 4 layered FCNN trained with MAML with 1 inner step and 0.1 inner learning rate, selected using early stopping using the meta-validation set with the ReLU activation function. The only difference of this figure with figure 7 is that this figure uses a ReLU activation and the other one uses a sigmoid. Note the dCCA value remains above 0.12 suggesting lower degree of feature re-use.

182 sufficient to achieve a high meta-accuracy on modern few-shot learning tasks e.g. mini-Imagenet,  
 183 tiered-Imagenet, Cifar FS, FC100, Omniglot, [20], which we hope to analyze in the future. We also  
 184 believe in is imperative that a definition of meta-learning is developed and conned to the general  
 185 intelligence. Chollet [4] takes this direction but to our understanding the proposed definition is mostly  
 186 focused for program synthesis. We also hope that in the future a metric for AI safety is ubiquitously  
 187 reported as suggested in Miranda et al. [15].

## 188 7 Discussion

189 We believe it's exciting that by only changing the few shot learning benchmark one can consistently  
 190 showed higher degrees of representation changes as measured by two different metrics. We believe  
 191 this is the case because the meta-learning system has to be trained explicitly with a task that demands  
 192 it to learn to adapt.

193 An important discussion point is the lack of an authoritative definition for measuring meta-learning in  
 194 our work and in the general literature. In particular in our work we decided to not report any results  
 195 with CKA. We decided this because Ding et al [8] showed that it's possible remove up to 97% of the  
 196 principal components of the weights of a layer until CKA starts to detect it. Thus, we used dCCA  
 197 which doesn't have the problem. It instead has a higher variance but it's easier to address this with  
 198 experiment repetition sand error bars (which we did). However, we believe it would be interesting to  
 199 use and extend Orthogonal Procrustes as suggest by [8] in future work.

200 The most obvious gap in our work is a thorough analysis with a real world vision data set. We hope  
 201 to repeat our work with the hinted extension in section 6 benchmarks as suggested in [11, 21].

202 In addition, Figures 7, 8 shows that as the number of inner steps increases, the dCCA does not  
 203 increase. This is somewhat surprising given the meta-overfitting results observed in section 4 and  
 204 further experiments would be valuable.



## 205 8 Broader Impact

### 206 8.1 Quantifying general intelligence through meta-learning

207 There is valuable efforts that try to make benchmarks which require higher level cognition e.g. [22].  
208 An example of work that tries to quantify AGI and proposes a benchmark is [4]. We believe the  
209 second approach is likely to have more impact in the long run because it also deliberately quantifies  
210 general intelligence. We believe that suggesting benchmarks without clearly specifying the long term  
211 goal or measuring the metric we are trying to optimize is a suboptimal approach. However, we do  
212 believe grounding benchmarks on tasks that humans are able to perform is a good idea but suggest to  
213 augment these proposals with metrics and explicit discussions of general intelligence.

214 Another approach we believe has high potential is program synthesis [2] and theorem proving [16, 5]  
215 because humans create higher abstractions that are composed and re-use thus suggesting to meta-  
216 learning might be taking place. We believe that higher level cognition tasks are a challenging to  
217 assess meta-learning algorithms.

### 218 8.2 Quantifying AI safety

219 We also believe quantifying and tracking metrics for AI safety as early as possible is crucial. Few-shot  
220 learning is likely one of simplest - and arguably the atomic building blocks for general intelligence.  
221 We believe AI safety could be enriched if research community deliberately tracks, discusses and  
222 report it in all it's research - especially in meta-learning research. For a brief discussion see [15].

### 223 8.3 Summary of Broader Impact

224 We hope that this discussion inspires the AI community - but especially the meta-learning research  
225 community - to always report their progress using, what Miranda et al. [15] call the "the big three":  
226 1) the score for absolute performance (to ensure usefulness) 2) the score for general skill acquisition  
227 (to ensure flexibility and general intelligence) and 3) the AI safety score (to ensure positive outcome).

## 228 References

- 229 [1] Wei-Yu Chen et al. "A Closer Look at Few-shot Classification". In: *7th International Confer-*  
230 *ence on Learning Representations, ICLR 2019* (2019). URL: [http://arxiv.org/abs/1904.](http://arxiv.org/abs/1904.04232)  
231 [04232](http://arxiv.org/abs/1904.04232).
- 232 [2] Xinyun Chen et al. "Compositional generalization via neural-symbolic stack machines". In:  
233 *NeurIPS* (2020).
- 234 [3] Yinbo Chen et al. *A New Meta-Baseline for Few-Shot Learning*. Tech. rep. URL: [https :](https://github.com/)  
235 [//github.com/](https://github.com/).
- 236 [4] François Chollet. "On the Measure of Intelligence". In: (2019). eprint: 1911.01547. URL:  
237 <http://arxiv.org/abs/1911.01547>.
- 238 [5] Maxwell Crouse et al. *A Deep Reinforcement Learning Approach to First-Order Logic Theorem*  
239 *Proving*. Tech. rep. URL: <https://arxiv.org/abs/1911.02065>.
- 240 [6] Tristan Deleu et al. *Torchmeta: A Meta-Learning library for PyTorch*. Available at:  
241 <https://github.com/tristandeleu/pytorch-meta>. 2019. URL: [https://arxiv.org/abs/1909.](https://arxiv.org/abs/1909.06576)  
242 [06576](https://arxiv.org/abs/1909.06576).
- 243 [7] Guneet S. Dhillon et al. "A Baseline for Few-Shot Image Classification". In: (2019). URL:  
244 <http://arxiv.org/abs/1909.02729>.
- 245 [8] Frances Ding, Jean-Stanislas Denain, and Jacob Steinhardt. "Grounding Representation Simi-  
246 larity with Statistical Testing". In: (). arXiv: 2108.01661v1. URL: [https://github.com/](https://github.com/js-d/sim_metric..)  
247 [js-d/sim\\_metric..](https://github.com/js-d/sim_metric..)
- 248 [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. "Model-Agnostic Meta-Learning for Fast  
249 Adaptation of Deep Networks". In: (2017). URL: <http://arxiv.org/abs/1703.03400>.
- 250 [10] Edward Grefenstette et al. "Generalized Inner Loop Meta-Learning". In: *arXiv preprint*  
251 *arXiv:1910.01727* (2019).
- 252 [11] Yunhui Guo et al. "A Broader Study of Cross-Domain Few-Shot Learning". In: (Dec. 2019).  
253 arXiv: 1912.07200. URL: <http://arxiv.org/abs/1912.07200>.

- 254 [12] Ryuichiro Hataya. *anatome, a PyTorch library to analyze internal representation of neural*  
 255 *networks*. 2020. URL: <https://github.com/moskomule/anatome>.
- 256 [13] Shaoli Huang and Dacheng Tao. “All you need is a good representation: A multi-level and  
 257 classifier-centric representation for few-shot learning”. In: (2019). URL: <http://arxiv.org/abs/1911.12476>.
- 258
- 259 [14] Brando Miranda. “DiMO : Differentiable Model Optimization and metaDiMO”. In: *Illinois*  
 260 *Digital Environment for Access to Learning and Scholarship (IDEALS)* (2019).
- 261 [15] Brando Miranda. “Establishing the foundations of Meta-learning - a Proposal”. In: *Illinois*  
 262 *Digital Environment for Access to Learning and Scholarship (IDEALS)* (2020). URL: <https://doi.org/10.1007/s11633-017-1054-2>.
- 263
- 264 [16] Brando Miranda. “Sketching: a Cognitively inspired Compositional Theorem Prover that  
 265 Learns to Prove - a Proposal”. In: *Illinois Digital Environment for Access to Learning and*  
 266 *Scholarship (IDEALS)* (2020). URL: [https://www.ideals.illinois.edu/handle/](https://www.ideals.illinois.edu/handle/2142/109134)  
 267 [2142/109134](https://www.ideals.illinois.edu/handle/2142/109134).
- 268 [17] Ari S Morcos et al. *Insights on representational similarity in neural networks with canonical*  
 269 *correlation*. Tech. rep.
- 270 [18] Jaehoon Oh et al. “Does MAML really want feature reuse only?” In: (Aug. 2020). URL:  
 271 <http://arxiv.org/abs/2008.08882>.
- 272 [19] Aniruddh Raghu et al. *Rapid Learning or Feature Reuse? Towards Understanding the Effec-*  
 273 *tiveness of MAML*. Tech. rep. URL: <https://arxiv.org/abs/1909.09157>.
- 274 [20] Yonglong Tian et al. *Rethinking Few-Shot Image Classification: a Good Embedding Is All You*  
 275 *Need?* 2020.
- 276 [21] Eleni Triantafillou et al. “Meta-Dataset: A Dataset of Datasets for Learning to Learn from Few  
 277 Examples”. In: (2019). URL: <http://arxiv.org/abs/1903.03096>.
- 278 [22] Ramakrishna Vedantam et al. *CURI: A Benchmark for Productive Concept Learning Under*  
 279 *Uncertainty*. Tech. rep. URL: <https://arxiv.org/abs/2010.02855v1>.

## 280 Checklist

281 The checklist follows the references. Please read the checklist guidelines carefully for information on  
 282 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
 283 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
 284 the appropriate section of your paper or providing a brief inline description. For example:

- 285
- 286 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
  - 287 • Did you include the license to the code and datasets? **[No]** The code and the data are  
 288 proprietary.
  - 289 • Did you include the license to the code and datasets? **[N/A]**

289 Please do not modify the questions and only use the provided macros for your answers. Note that the  
 290 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
 291 block and only keep the Checklist section heading above along with the questions/answers below.

- 292 1. For all authors...
- 293 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 294 contributions and scope? **[Yes]**
- 295 (b) Did you describe the limitations of your work? **[Yes]**
- 296 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
- 297 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 298 them? **[No]**
- 299 2. If you are including theoretical results...
- 300 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 301 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 302 3. If you ran experiments...

- 303 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
 304 imental results (either in the supplemental material or as a URL)? [No] We include  
 305 hyperparams in the supp section. But are happy better instructions for reproducibility  
 306 and the our weight and biases (wandb) repo with recordings of our runs etc and we are  
 307 happy to make them available upon acceptance
- 308 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 309 were chosen)? [No] Not every single detail but every details has been recorded and we  
 310 are happy to make them available upon acceptance
- 311 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 312 ments multiple times)? [Yes]
- 313 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 314 of GPUs, internal cluster, or cloud provider)? [Yes] More details in supp
- 315 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 316 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 317 (b) Did you mention the license of the assets? [N/A]
- 318 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 319
- 320 (d) Did you discuss whether and how consent was obtained from people whose data you're  
 321 using/curating? [N/A]
- 322 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 323 information or offensive content? [N/A]
- 324 5. If you used crowdsourcing or conducted research with human subjects...
- 325 (a) Did you include the full text of instructions given to participants and screenshots, if  
 326 applicable? [N/A]
- 327 (b) Did you describe any potential participant risks, with links to Institutional Review  
 328 Board (IRB) approvals, if applicable? [N/A]
- 329 (c) Did you include the estimated hourly wage paid to participants and the total amount  
 330 spent on participant compensation? [N/A]

## 331 A Supplementary Material

### 332 A.1 Experimental and hyper-parameter details

#### 333 A.1.1 Details for experiments on our benchmark that requires more than feature re-use

334 All models were trained on a CPU cluster with intel CPUs. All models were 4 layered FCNN. All  
 335 models had batch-normalization and collected running statistics during meta-training but used batch  
 336 statistics during training. MAML models for figure 2 and 3 had inner learning rate of 0.1 and 1  
 337 inner step. One inner super was chosen to further emphasizes the feature re-use since it is the lowest  
 338 inner step we can choose (nothing lower exists except 0 which doesn't exhibit adaptation). Adam  
 339 outer optimizer was used with learning rate 0.001 and default parameters. No learning schedulers  
 340 were used but would be interesting to experiment with. Since there were 200 regression tasks, we  
 341 trained the models with meta-epochs instead of meta-iterations. This means that we reported errors,  
 342 losses etc. after all tasks were seen. Note that the input and target values were guaranteed to be  
 343 novel because during meta-train and meta-testing we sample a function and generate data on the fly -  
 344 similar to online learning. Note that this is very similar to how classification tasks for few-shot work  
 345 (e.g. mini-Imagenet) because those tasks have a very limited number of image classes and thus results  
 346 in highly correlated tasks. In addition, we showed how both exhibited similar meta-overfitting. For  
 347 CCA value computation we used a query set size of 100 due to numerical issues with the Anatomy  
 348 library [12]. We did not use first order MAML. All models were trained until convergence with about  
 349 200,000 meta-epochs.

350 For models with 8 and 7, we meta-trained with MAML but used 2 inner steps. The remaining  
 351 hyperparameters remain the same. We did not use first order MAML.

352 **A.1.2 Details for experiments on ResNet-18 meta-overfitting on mini-Imagenet**

353 The ResNet-18 is a standard out of the box ResNet-18 downloaded from PyTorch. We trained the  
354 models for 5,000,000 meta-iterations. We used 1 inner step with 0.1 inner learning rate. We used  
355 meta-batch size of 4 and 2 for meta-training and meta-testing respectively. We used an outer learning  
356 rate of 0.001 and Adam with default parameters. We did not use first order MAML.

357 **A.2 Role of Backbone on meta-accuracy**

358 In this section we describe the relation of the depth of a Pytorch ResNet model with the meta-test  
359 accuracy. The motivation for these experiments is that if we can close the gap on mini-Imagenet  
360 to over 90% by only increasing the back bone depth then this would provide strong evidence that  
361 such benchmarks really only need a good embedding. However, we discovered that for the ResNets  
362 used in [20] it seems that accuracy saturates at 80% (results not shown in paper) but when using the  
363 Pytorch models we see meta-overfitting and decreasing meta-test error 9. This suggests that even this  
364 simple scenario of few-shot learning still has still space for meta-learning to be a solution.

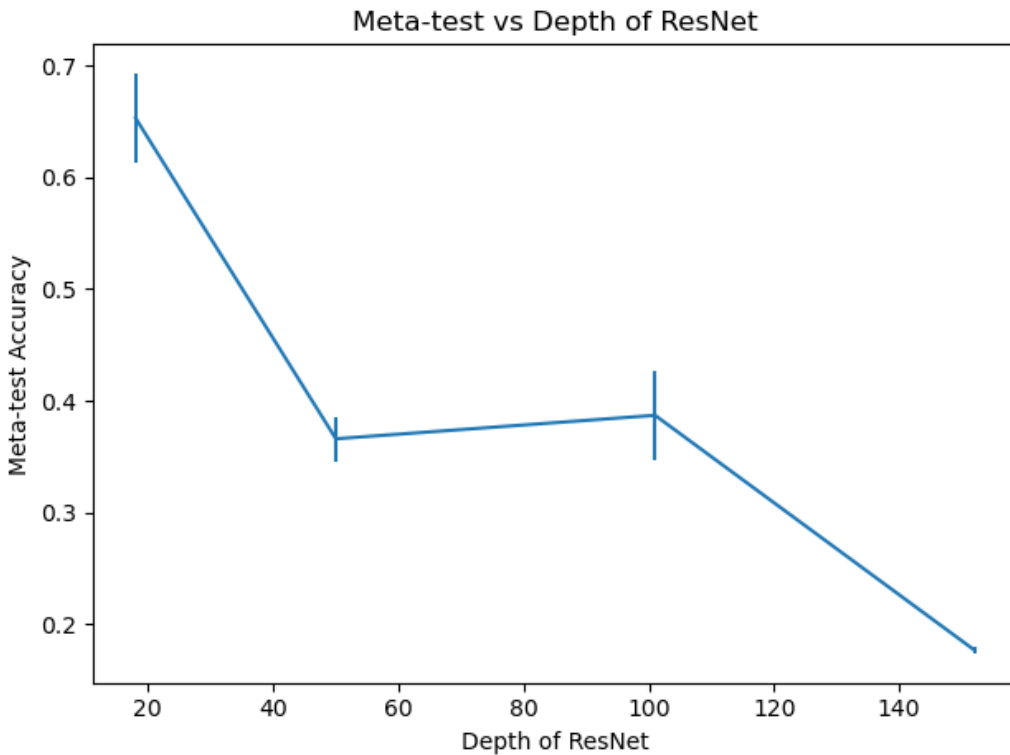


Figure 9: Shows that as the backbone of the Pytorch ResNets increases to 152 the meta-accuracy on mini-Imagenet decreases. These models were trained with supervised union training as in [20]. The meta-adaption algorithm used logistic regression and was adapted to convergence on the final layer as in [20]. When using the Pytorch ResNet models instead of the special ResNets designed for mini-Imagenet [20] we observe that the meta-accuracy decreases 9.

365 **A.3 Analysis of meta-learned initialization**

366 In this section we do an experiment where we fix the ResNet18 meta-learned initialization and use the  
367 adaption that only adapts the final layer as in [20]. The results in figure 10 are mixed but it is interesting  
368 to note MAML with no inner steps performs worse than a random neural network. This result is  
369 interesting because this is very similar to supervised pre-training in that no meta-learner is present  
370 during training but instead of seeing all 64 images it sees 5 randomly (but uses no meta-learner).

371 We would have expected that the initialization obtained would have been equivalent to one with  
 372 supervised pre-training. Since they are not it shows a MAML is at the very least capable of learning a  
 373 representation that is invariant to concept permutation.

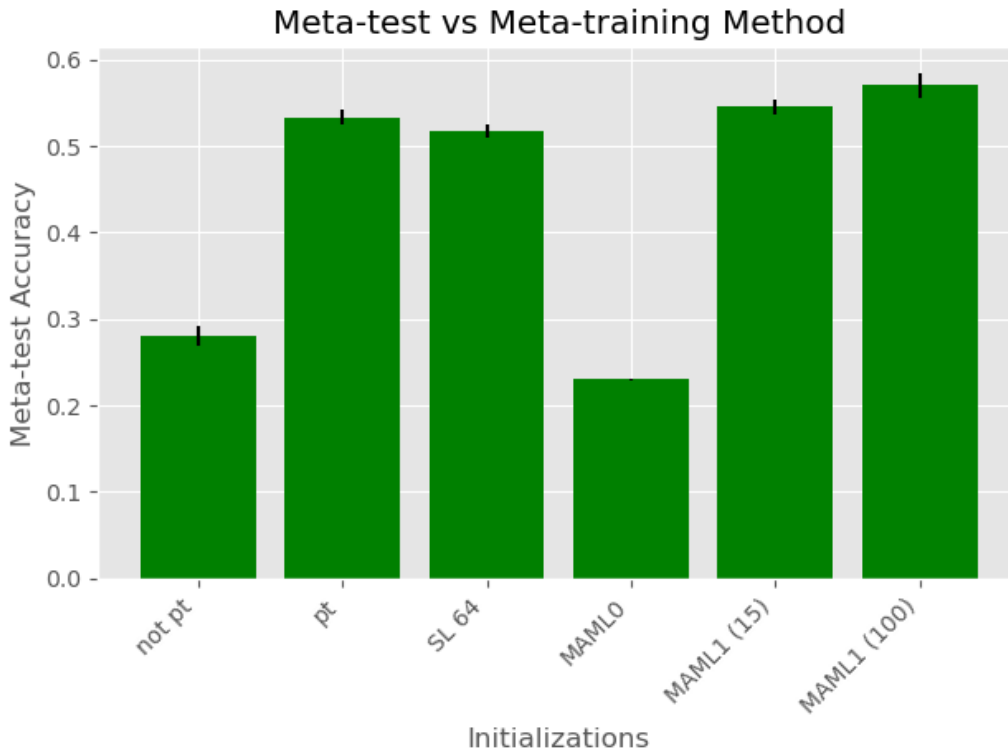


Figure 10: Shows relation of meta-test accuracy with models with a different meta-learned initialization. PT stands for Pre-trained on Imagenet. Not Pt stands for a random model. All models are ResNet18s from PyTorch. SL 64 stands for supervised union pre-training on mini-imagenet using all 64 labels during meta-training. MAMLO stands for only using episodic meta-training (i.e. MAML with zero inner steps). MAML1 (15) and MAML1 (100) stand for training using MAML with a query set of size 15 to 100. The meta-adaptation is the same as in [20] (training logistic regression in the final layer to convergence).

#### 374 A.4 Training with zero number of inner steps

375 We believe it is an interesting observation that MAML with 0 inner steps (MAMLO) (i.e. only using  
 376 episodic meta-training) resulted in very different meta-learned initialization compared to MAML with  
 377 1 inner step (MAML1) on mini-Imagenet. Previous work observed that supervised pre-training [20]  
 378 with all 64 images during meta-training results in a strong baseline. With this in mind it is natural to  
 379 ask: what is the difference between seeing all 64 images during supervised pre-training or seeing only  
 380 5 using episodic training? With this in mind we trained MAMLO and obtained a model that performs  
 381 at chance. Figure 11 compares MAMLO with MAML1 to show that MAMLO obtains a model that  
 382 has a very high meta-training loss. Additionally, figure 10 shows such an initialization performed  
 383 even worse than random. This is surprising but it seems that meta-learned initialization with MAML1  
 384 learn at least a model that is invariant to permutation of the order of the classes. Unfortunately, this  
 385 result seems to only be reproducible in classification since training MAMLO in a synthetic regression  
 386 task did converge to have model with low meta-train loss 12. This suggests future studies would be  
 387 interesting to disentangle the casual factors.

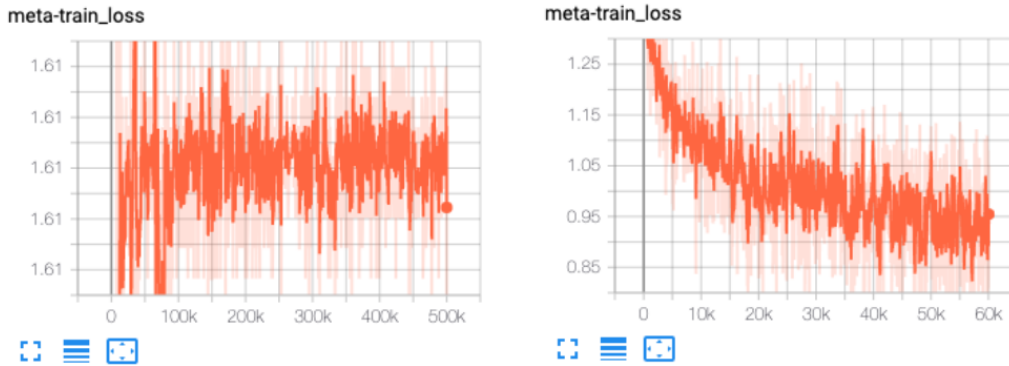


Figure 11: Compares MAML0 (only episodic training) vs MAML1 (MAML with 1 inner step). MAML0 remains close to chance with a high loss while MAML1 converges. This suggests MAML0 is not equivalent to supervised pre-training and that MAML1 does learn a representation that is invariant to class order permutation.

### 388 A.5 Tips and tricks for episodic meta-training

389 From our experiments we suggest the following (when episodic meta-training [9]):

- 390 1. Use a large number of query examples e.g. greater than popular 15 (since they often speeds up convergence of the meta-learning algorithm).
- 391
- 392 2. A large meta-batch size (since it's important to be able to have a low level of noise when tracking the meta-validation error/loss for doing early stopping). We found empirically for
- 393 75 – 100 tasks to be a good meta-batch size.
- 394
- 395 3. Episodic training as suggested in [9] is expensive and takes at least a week to train on mini-Imagenet on a Quadro RTX 6000 using torchmeta and higher [6, 10], so these suggestions
- 396 are important.
- 397
- 398 4. Experiments with synthetic data were run with CPU only.

### 399 A.6 Future work

#### 400 A.6.1 Summary

- 401 1. Defining a synthetic benchmark that is a classification problem that also requires meta-learning (or rapid learning with MAML).
- 402
- 403 2. We also hope to construct a (real) benchmark from images that requires meta-learning. Formally, we propose a good start would be a benchmark where the probability of two task
- 404 having the same class be small, otherwise we are more likely to see overfitting. Alternatively, a benchmark that requires the tasks to be different by at least requiring a different
- 405 representation. We believe compositionality is an ideal benchmark since this would allow sophisticated re-use of lower level representations and simultaneously have an unbounded
- 406 number of tasks. Humans are able to richly and flexibly cope with both. Additionally, it would be interesting to be able to quantify the distance between two different N-way, K-shot
- 407 tasks to make these ideas more rigorous.
- 408
- 409 3. An interesting benchmark with a large number of classes with real images is taking the union of many vision classification tasks and re-scaling all images to be of the size of
- 410 mini-Imagenet.
- 411
- 412 4. Plotting the meta-generalization gap (with a synthetic classification task) and demonstrate it decreases as the number of tasks increases would be interesting (note however we already
- 413 have the limiting case when the number of tasks is unbounded and the meta-generalization
- 414 gap is zero).
- 415
- 416
- 417
- 418

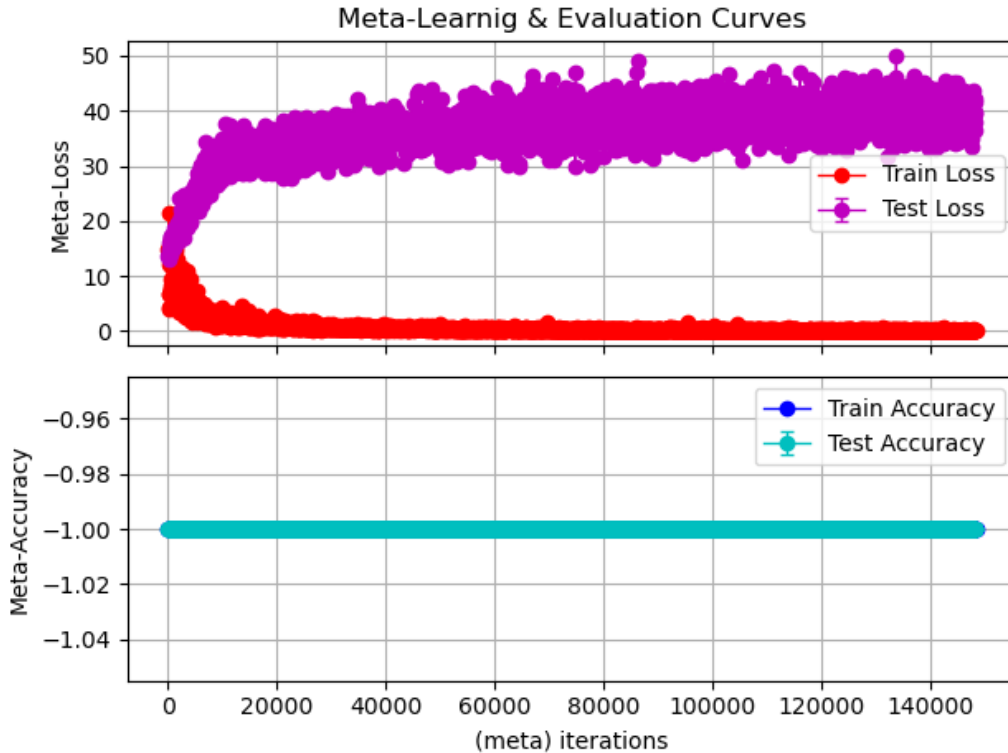


Figure 12: Shows MAML0 (only episodic training) getting zero meta-train loss (red curve) for a synthetic regression task. This suggests that meta-learning in regression and classification might not be entirely equivalent. Note that meta-overfitting is still observed (purple curve). This is a regression task so the blue curves can be ignored.

- 419 5. An interesting experiment would be to train a deep neural network with the episodic training  
 420 (but without the MAML inner loop) but have an unbounded number of tasks and see if the  
 421 test error keeps increases (or stays at chance as observed when this is done with mini-Imagent  
 422 11).
- 423 6. An interesting hypothesis to investigate is if meta-learning algorithms get representation  
 424 that are optimal for their respective meta-learner (or adaptation rule). If this is true it means  
 425 methods like [20] can be improved by making the entire pipeline differential and learning it  
 426 end-to-end [14].
- 427 7. Test meta-learning algorithms in domains where higher level cognition is required and thus  
 428 compositionality is essential e.g. program synthesis [2] and theorem proving [16, 5].

#### 429 A.6.2 Proposal on Synthetic classification task that possibly require meta-learning

430 Synthetic tasks that use classification instead of regression are not hard to define. Two possible  
 431 alternatives are: 1) a mixture of Gaussians but the standard deviation controls the radius of limit  
 432 of how far the classes can be from each other 2) another option is the similar as with a mixture of  
 433 Gaussians but have the (vector) samples be weights of a Neural Networks (so that the goal is to  
 434 identify from which Neural Network data is coming from)