

519 A Further Discussions

520 We'd like to emphasize that our synthetic experiments are promising because we can systematically
521 differentiate meta-learning algorithms from transfer learning algorithms – which supports our action-
522 able suggestion to: 1. use the diversity coefficient to effectively study meta-learning and transfer
523 learning algorithms, and 2. to use the diversity coefficient to design better benchmarks. In addition,
524 this problematizes the observations that fo-*proto*-MAML in meta-data set [44] is better than transfer
525 learning solutions – since our synthetic experiments show MAML is not better than USL in the high
526 diversity regime. To further problematize, we want to point out that meta-learning methods are not
527 better than transfer learning as observed by [20] – as observed in our synthetic experiments. We
528 hypothesize however that the two scenarios in [44] are different [20]. The first one focuses on the
529 same meta-training and meta-testing conditions, while the latter focuses on a cross-domain. We
530 hypothesize that the cross-domain scenario might benefit from a meta-learning which lower variance
531 (e.g., a fixed embedding [43]) – which might explain why sophisticated meta-learning solutions might
532 perform worse on the cross-domain setting as observed in [20]. Further research is needed in both
533 benchmarks – especially from a problem centric perspective with quantitative methods like the ones
534 we suggest.

535 In addition, we hypothesize that diversity might be a good proxy to predict the difference between
536 meta-learning and transfer learning methods. More precisely, we conjecture that in a low diversity
537 setting meta-learning methods are equivalent at meta-test time to transfer learning methods but their
538 difference increases as the diversity of tasks in a benchmark increases. In the high diversity regime
539 we conjecture that the difference between meta-learning and transfer learning methods increases as
540 the diversity increases. We are optimistic that meta-learning algorithms might outperform transfer
541 learning methods, once we start comparing them in more thoughtfully designed benchmarks. It
542 is possible that despite our efforts, meta-learning algorithms – as currently designed – are too
543 sophisticated and in fact lead to meta-overfitting, as shown in previous work [32].

544 We'd like to emphasize, that up until now, the meta-learning community has evaluated meta-learning
545 algorithms in benchmarks that might not be the most appropriate. We conjecture high diversity
546 benchmarks are more appropriate, since they might capture the meta-learning inductive prior: high
547 diversity means that adaptation is required by construction. Thus, we conjecture that previous
548 conclusions should be taken with a grain of salt until a more in depth study can be made in the
549 high diversity regime – especially with benchmarks with real world data that have been analyzed
550 extensively with metrics like the diversity coefficient that we propose. We conjecture that we can
551 finally do meta-learning research effectively – given that a regime where meta-learning and transfer
552 learning methods can be differentiated has been discovered and previous low diversity benchmarks
553 have been understood.

554 We also conjecture that meta-learning research is different from classical machine learning research.
555 Historically, a seminal paper is the one where AlexNet was proposed [25]. In that time we had low
556 performance on a fixed task e.g., Imagenet and couldn't even interpolate the data (i.e., reach zero
557 train error). We conjecture that meta-learning is different because if we have a diversity so large
558 where all possible tasks are incorporated, then we should reach the no-free lunch theorem regime
559 [46] – where all algorithms should perform the same on average. Therefore, we hypothesize that a
560 deliberate and quantitative efforts to design benchmarks is essential. A great example of such an
561 attempt is the Abstraction and Reasoning Corpus (ARC) benchmark [11] – which was made very
562 thoughtfully with Artificial General Intelligence (AGI) in mind. We conjecture meta-learning is the
563 most promising path in that direction, and hope this work inspires the design of benchmarks that lead
564 to actionable and deliberate attempts to make progress to build such AGI technologies.

565 B Related Work (continued)

566 The meta-learning literature is growing quickly and hope to provide a wider coverage here.

567 In terms of benchmarks, we'd like to start with the ARC benchmark [11]. ARC was designed with
568 AGI in mind – arguably the ultimate meta-learner. Its focus is primarily on visual reasoning using
569 program synthesis techniques. We hypothesize it's a very promising path but our work inspires
570 extension that go beyond program synthesis approaches. The meta-data set benchmark is an attempt
571 to make the data set for few-shot learning at a larger scale and more diverse [44]. The main difference

572 of their work and ours is that we propose a quantitative metric to measure the intrinsic diversity in
573 the data and go beyond data set size or number of classes. They also showed that a meta-learning
574 algorithms – fo-Proto-MAML – is capable of beating transfer learning. However, they also showed
575 transfer learning baselines are in fact quite difficult to beat. The IBM Cross-Domain few-shot
576 learning benchmark [20] is a fascinating benchmark to evaluate meta-learning algorithms. Their
577 central premise however is transfer from a source domain to a different target domain – instead of our
578 setting where tasks are created from the same meta-distribution. This is why their paper is considered,
579 in addition to few-shot learning, a *cross-domain benchmark*. We believe this is an essential scenario
580 to think about, but consider it different from our setting or the setting of meta-data set. We’d also like
581 to emphasize that they do not employ a metric like our diversity coefficient that quantitatively asses
582 the diversity of their benchmarks. These two last benchmarks, although fascinating, are missing the
583 essential quantitative analysis of the data itself we are trying to propose.

584 The work by [8] give to the best of our knowledge – the first non-vacuous generalization bounds for
585 the (supervised) meta-learning setting. Their statements apply to non-convex loss function and use
586 stability theory at the task level. The bound depends on the mutual information on the input data vs
587 the output data of the meta-learner. The results, although fascinating, are not built to separate classes
588 of meta-learning – like our work attempts to do empirically.

589 The work by [45] proposes the idea of global labels as a way to indirectly optimize for the met-
590 learning objective for a fixed feature extractor. Global labels is equivalent to the concept
591 we call USL in this paper. They show that pre-training (i.e. using USL/global labels) provides
592 excellent meta-test results – including with their method (named MeLa) that can infer global labels
593 given only local labels provided at in episodic meta-training. Their theoretical analysis depends
594 on a fixed feature extractor, instead of considering the whole end-to-end meta-learner as a whole –
595 meaning two different deep learning models cannot be used in their analysis. Therefore, their analysis
596 fails to separate how different feature extractors might be trained, e.g. comparing USL vs MAML
597 directly in an end-to-end fashion. In contrast, we instead tackle this question head on theoretically [L
598 (with limited results) but instead show that the feature extractors indeed are different empirically [E

599 The work by [13] proposes a theoretical treatment of meta-learning using meta-learners with closed-
600 form equations derived from ridged regularization using fixed features. They formulate the conditional
601 and unconditional formulation using side information for the task (e.g. the support set) and show the
602 conditional method is superior. In relation to our work, they do not provide characterizations of the
603 role of a neural network doing end-to-end meta-learning (in their empirical or theoretical analysis).
604 In contrast, our findings make an explicit effort in understanding the role of the neural network in
605 meta-learning in an end-to-end fashion through emperical analysis. Another contrast is that their
606 results are highly theoretical, while ours focus on empirical results. In addition, their results are on
607 synthetic experiments and do not explore their findings in the context of modern few-shot learning
608 benchmarks like MiniImagenet or Cifar-fs.

609 The work by [19] provide strong evidence that adaptation at test time is best done when the meta-
610 trained model matches the adaptation it was meta-trained with. This is shown because their classically
611 pre-trained nets cannot perform better than the MetaOpt models with any fine-tuning method. How-
612 ever, their results cannot beat [43] and thus does not help separate the role of meta-training and union
613 supervised learning (USL). Their Resnet12 results do provide further support to our hypothesis that
614 large enough neural networks all perform the same, since 78.63 (Goldblum) vs 79.74 (RFS) have
615 very close errors, in line with our findings. However, we hypothesize it is not due to the model size in
616 accordance with our experiments [3]

617 The work by [18] provides theoretical bounds of when the expected risk of MAML and DRS (Domain
618 Randomized Search) by bounding the gradient norm. DRS attempts to model USL but fails to do
619 so completely, because USL is capable of modeling adaptation because the final layer is capable of
620 adaption. Thus, it does not address the capabilities of the feature extractor being able to learn all the
621 information needed to meta-learn. Concisely, their analysis is not capable of separating performance
622 of MAML and USL. Even if it hypothetically could, their analysis remains an upper bounds (with
623 assumptions). This raising the question if their method truly explain the observations that transfer
624 learning methods – like USL – beat meta-learning methods. In addition, they do not provide in
625 depth empirical analysis with respect to any real few-shot learning benchmarks like MiniImagenet or
626 Cifar-fs.

627 The work by [26] provides an exploration of the effects of diversity in meta-learning. The main
 628 difference with our work is that they focus mostly on sampling strategies, and it's effect on diversity,
 629 while we focused on the intrinsic diversity in the benchmarks themselves.

630 The work by [40] provides a theoretical analysis on the difference between interpolation and extrapo-
 631 lation in transfer learning (and domain generalization). We believe this type of theory may be helpful
 632 as an inspiration to explore why in the high diversity regime there seems to be a difference between
 633 the performance of meta-learning and transfer learning methods.

634 C Convergence of Learning Curves for Fair Comparison

635 C.1 Convergence of Learning Curves for MiniImagenet and Cifar-fs

636 In this section, we have the plots showing the learning curves achieving convergence for the models
 637 used in figure 1 and 2. Note, the learning curves for the models trained with MAML look noisier
 638 because the distributed training reduces the size (meta) batch size for logging purposes. In addition,
 639 due to episodic meta-training, (meta) batch sizes have to be smaller compared to batch sizes used in
 640 USL.

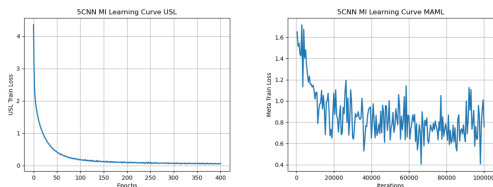


Figure 5: Plot showing convergence of 5CNN on MiniImagenet.

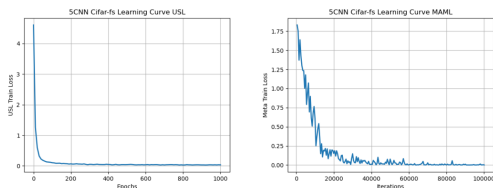


Figure 6: Plot showing convergence of 5CNN on Cifar-fs.

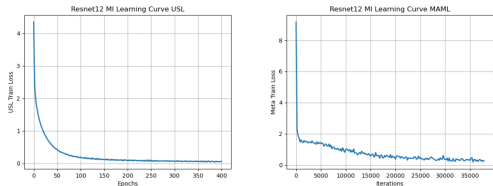


Figure 7: Plot showing convergence of Resnet12 on MiniImagenet.

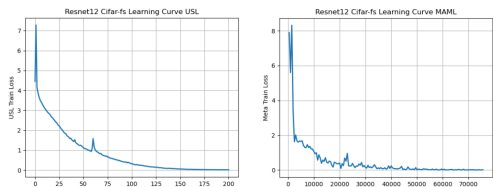


Figure 8: Plot showing convergence of Resnet12 on Cifar-fs.

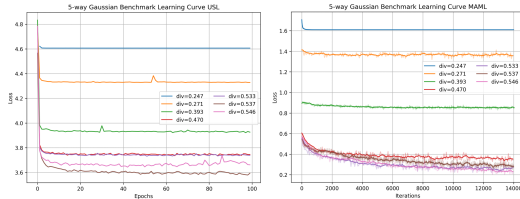


Figure 9: Plot showing convergence of a custom 3-layer fully connected network, for all synthetic Gaussian benchmarks tested. Each curve represents a different synthetic Gaussian benchmark tested on either MAML (left plot) or USL (right plot). The curves are then color-coded by the value of the Task2Vec-based task diversity coefficient of the Gaussian benchmark tested in that particular run.

641 D Experimental Details

642 D.1 Experimental Details on MiniImagenet and Cifar-fs

643 We will explain the details for the four models we trained on figure 1 and 2.

644 **Summary:** We trained a five layer CNN (5CNN) and Resnet12 on both MiniImagenet and Cifar-
 645 fs to convergence. We used the Adam optimizer with learning rate $1e-3$. We used the standard
 646 MiniImagenet and Cifar-fs data augmentations as provided in [5] matching [43].

647 **Experimental Details for 5CNN on MiniImagenet:** We used the five layer CNN from [17, 39]. We
 648 used 32 filters as used in previous work. We used the Adam optimizer with learning rate $1e-3$ for both
 649 MAML and USL. We used no scheduler. We trained the USL model for 1000 epochs. We trained the
 650 MAML model for 100,000 episodic iterations (outer loop iterations). We used a batch size of 128 for
 651 USL and a (meta) batch size of 8 for MAML. For MAML we used an inner learning rate of $1e-1$ and
 652 5 inner learning steps. We did not use first order MAML. It took 3 hours 5 minutes 5 seconds to train
 653 USL to convergence with a single GPU. It took 1 day 6 hours 21 minutes 8 seconds to train MAML
 654 to convergence with 4 NVIDIA GeForce GTX TITAN X GPUs.

655 **Experimental Details for Resnet12 for MiniImagenet:** We used the Resnet12 provided by [43]. We
 656 used the Adam optimizer with learning rate $1e-3$ for both MAML and USL. We used the same cosine
 657 scheduler as in [43] for USL and no cosine scheduler for MAML. We trained the USL model for 186
 658 epochs. We trained the MAML model for 37,800 episodic iterations (outer loop iterations). We used
 659 a batch size of 512 for USL and a (meta) batch size of 4 for MAML. For MAML we used an inner
 660 learning rate of $1e-1$ and 4 inner learning steps. We did not use first order MAML. It took 1 day 17
 661 hours 2 minutes 41 seconds to train USL to convergence with a single dgx A100-SXM4-40GB GPU.
 662 The MAML model was trained with Torchmeta [12] which didn't support multi gpu training when
 663 we ran this experiment, so we estimate it took 1-2 weeks to train on a single GPU. In addition, it
 664 was ran with an earlier version of our code, so we unfortunately did not record the type of GPU but
 665 suspect it was either an A100, A40 or Quadro RTX 6000.

666 **Experimental Details for 5CNN for Cifar-fs:** We used the five layer CNN from [17, 39] provided by
 667 [5]. But we used 1024 filters instead of 32 (to speed up convergence). We used the Adam optimizer
 668 with learning rate $1e-3$ for both MAML and USL. We used the same cosine scheduler as in [43] for
 669 MAML and no cosine scheduler for USL. We trained the USL model for 1000 epochs. We trained
 670 the MAML model for 100,000 episodic iterations (outer loop iterations). We used a batch size of 256
 671 for USL and a (meta) batch size of 8 for MAML. For MAML we used an inner learning rate of $1e-1$
 672 and 5 inner learning steps. We did not use first order MAML. It took 10 hours 43 minutes 31 seconds
 673 to train USL to convergence with a single GPU dgx A100-SXM4-40GB. It took 2 days 9 hours 26
 674 minutes 27 seconds to train MAML to convergence with 4 Quadro RTX 6000 GPUs.

675 **Experimental Details for Resnet12 for Cifar-fs:** We used the Resnet12 provided by [43]. We used
 676 the Adam optimizer with learning rate $1e-3$ for both MAML and USL. We used the cosine scheduler
 677 used in [43] for both USL and MAML. We trained the USL model for 200 epochs. We trained the
 678 MAML model for 75,500 episodic iterations (outer loop iterations). We used a batch size of 1024 for
 679 USL and a (meta) batch size of 8 for MAML. For MAML we used an inner learning rate of $1e-1$ and
 680 5 inner learning steps. We did not use first order MAML. It took 45 minutes 54 seconds to train USL

681 to convergence with a single GPU. It took 1 day 19 hours 29 minutes 31 seconds to train MAML to
682 convergence with 4 dgx A100-SXM4-40GB GPUs.

683 **Why the Adam optimizer?** We hypothesize that the Adam optimizer is the most appropriate
684 optimizer for a fair comparison for various reasons. First, the Adam optimizer is widely used –
685 making our results most relevant and broadly applicable. Adam is generally a stable optimizer
686 – especially for sophisticated meta-learning algorithms like MAML. It is not uncommon to have
687 SGD result in exploding gradients or end up diverging – especially for MAML. Most importantly
688 however, we hope to stay faithful whenever possible to how modern transformer models are trained
689 – because they have been shown to be good meta-learners, e.g., gpt-3 is often cited as a zero-shot
690 learner [7]. These type of models do use more complicated learning schemes besides only Adam
691 (e.g., warm-ups, decay rates etc.) but we hypothesize using Adam is a good first step. We conjecture
692 that the small benefits that SGD might provide are negligible compared to the stability that Adam
693 provides, especially as the scale of the data sets starts to increase. Without Adam we conjecture
694 it would be hard to even perfectly fit the data for large scale data sets as it’s usually done in Deep
695 Learning. This was definitively true in our own experiments. Therefore, we decided to use Adam
696 for our experiments, since it would be too hard to use SGD reliably at scale or with sophisticated
697 meta-learning algorithms.

698 D.2 Experimental Details on N-way Gaussian Tasks

699 We used a custom 3-layer fully connected network derived from Learn2Learn’s OmniglotFC model
700 [5], with parameters $input_size = 1$, $output_size = 5$, and hidden layer sizes $sizes = [128, 128]$.
701 We used the Adam optimizer with learning rate 1e-3 for both MAML and USL, but did not use a
702 cosine scheduler for either USL or MAML. We trained the USL model for 100 epochs. We trained
703 the MAML model for 14,000 episodic iterations (outer loop iterations). We used a batch size of 100
704 for USL and a (meta) batch size of 100 for MAML. For MAML we used an inner learning rate of
705 1e-1 and 5 inner learning steps. We did not use first order MAML. It took 19 minutes 24 seconds to
706 train USL to convergence with a single Titan X. It took 2 days and 13 hours 6 minutes 27 seconds to
707 train MAML to convergence with a single Titan X.

708 E Feature Extractor Analysis of USL and MAML

709 Figure 10 significant difference between the feature extractor layers of a MAML trained model vs. a
710 union supervised learned model.

711 F Background of Few-shot Learning Basics

712 The goal of few-shot learning is to learn to classify from a limited set of training samples. A few-shot
713 benchmark is utilized to evaluate few-shot learning algorithms and typically contains many classes
714 and a smaller number of samples per class. Typically, few-shot learning algorithms learn in episodes,
715 where in each episode, a task consisting of a train (or support) set and a held-out validation (or query)
716 set is sampled. In particular, a task is a n -way k -shot classification problem, means that the support
717 and query sets each consist of n classes sampled from the benchmark, and each of the n classes are
718 represented by k shots or examples. The learner uses the support set to adapt to the task, and the
719 query set to evaluate the performance on the given task.

720 G Synthetic Gaussian Benchmark and N-way Gaussian Tasks

We create a series of synthetic few-shot benchmarks, where each Gaussian benchmark B is defined by
four parameters $B = (\mu_m, \sigma_m, \mu_s, \sigma_s)$. To form the dataset of our benchmark, we first sample 100
meta-train, 100 meta-test, and 100 meta-validation classes, where class $1 \leq i \leq 300$ is a Gaussian
parameterized by $(\mu_{class_i}, \sigma_{class_i})$ where

$$\mu_{class_i} \sim N(\mu_m, \sigma_m), \sigma_{class_i} \sim |N(\mu_s, \sigma_s)|$$

Then, for each class i , we sample 1000 data points $(x_{i,1}, i) \dots (x_{i,1000}, i)$ where each datapoint (x, y)
is composed of a input value $x \in \mathbb{R}$ and class label $1 \leq y \leq 300$. The input values $x_{i,1} \dots x_{i,1000}$

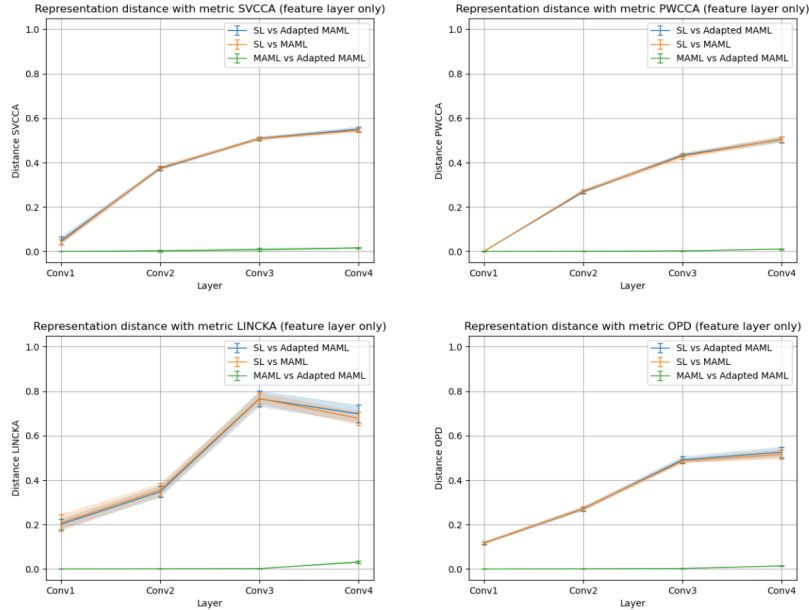


Figure 10: Shows the significant difference between the feature extractor layers of a MAML trained model vs. a union supervised learned model – especially in contrast to the small change in the adapted MAML model (green line). This figure suggests that although benchmark diversity is small, a meta-learned representation still learns through a different mechanism than a supervised learned representation. Note that the green line is our reproduction of previous work [37] that showed that a MAML trained model does not change after using the MAML adaptation. They term this observation as “feature re-use”.

are each sampled from class i 's class distribution:

$$x_{i,1} \dots x_{i,1000} \sim N(\mu_{class_i}, \sigma_{class_i})$$

721 Having defined our dataset underlying our benchmark, we may now sample individual tasks from
 722 our benchmark. Each task in our benchmark is 5-way, 10-shot - that is, each task is formed by
 723 first sampling 5 ways from the benchmark dataset, then sampling 10 shots from each of the 5 ways.
 724 The goal of each task is to correctly predict which of the 5 ways an input value $x \in \mathbb{R}$ falls into.
 725 We conducted experiments using 7 different benchmarks, with each benchmark defined by four
 726 parameters and its corresponding Hellinger distribution diversity coefficient and Task2Vec task
 727 diversity coefficient, as listed in Table 3.

728 H Distribution-based Diversity Metrics

729 In addition to the task diversity methods (such as Task2Vec) that we chose as a measure of diversity
 730 across our experiments in our main paper, we would like to introduce an additional class of diversity
 731 metrics that we call *distribution diversity*. Unlike task diversity, which quantifies diversity through
 732 the expected distance between any two distinct tasks sampled from the benchmark, distribution
 733 diversity quantifies diversity through the expected distance between any two distinct *distributions* that
 734 underlie the benchmark. In our synthetic Gaussian experiments, we define the distribution diversity
 735 of our Gaussian benchmark as the expected Hellinger distance between two distinct Gaussian class
 736 distributions sampled from the benchmark - we describe the calculation of the distribution diversity
 737 of our synthetic Gaussian benchmark in more detail in Section 4.

738 I Hellinger Diversity Coefficient and Hellinger Distance

An alternative metric to the Task2Vec-based task diversity metric is the Hellinger-based distribution diversity metric. The Hellinger-based distribution diversity of our Gaussian benchmark is obtained

Benchmark Parameters ($\mu_m, \sigma_m, \mu_s, \sigma_s$)	Hellinger-based Distribution Diversity	Task2Vec-based Task Diversity
(0, 0.01, 1, 0.01)	$7.475e-05 \pm 4.891e-07$	$0.247 \pm 1.04e-3$
(0, 1, 1, 0.01)	$0.183 \pm 1.24e-3$	$0.271 \pm 1.15e-3$
(0, 3, 1, 0.01)	$0.574 \pm 2.28e-3$	$0.393 \pm 1.79e-3$
(0, 10, 1, 0.01)	$0.860 \pm 1.75e-3$	$0.470 \pm 2.35e-3$
(0, 20, 1, 0.01)	$0.929 \pm 1.31e-3$	$0.533 \pm 2.47e-3$
(0, 30, 1, 0.01)	$0.952 \pm 1.10e-3$	$0.537 \pm 2.57e-3$
(0, 1000, 1, 0.01)	$0.998 \pm 2.07e-4$	$0.546 \pm 2.74e-3$

Table 3: **Benchmarks of increasing diversity are created by increasing σ_m , or the standard deviation of the class mean** A larger σ_m increases the variance of the class means, making their respective class distributions farther apart on average and causing both the Hellinger-based distribution diversity and Task2Vec-based task diversity coefficients to increase. We varied σ_m from 0.01 to 1000 and fixed all remaining benchmark parameters to obtain 7 different Gaussian benchmarks. The corresponding Hellinger-based distribution diversity coefficients were obtained by numerically approximating the expected Hellinger distance between two classes sampled from the benchmark and computing the 95% confidence interval of the approximation. We also computed Task2Vec-based task diversity coefficients as an alternative measure to diversity using a random 3-layer fully connected probe network described in Section D. Figure 12 visualizes the Task2Vec task diversities among the synthetic benchmarks via a heatmap showing the relative pairwise distance between sampled tasks.

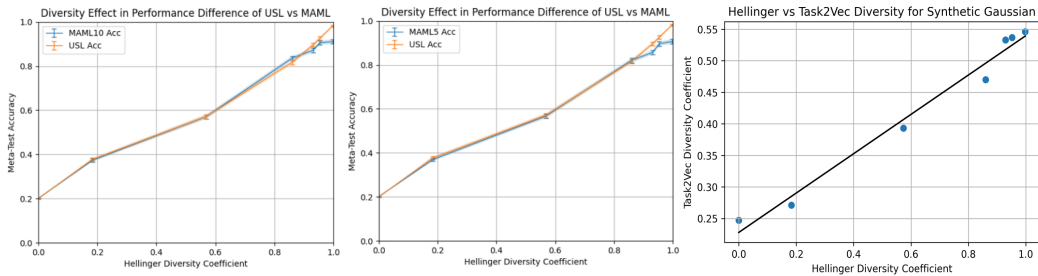


Figure 11: **Shows the strong relation between Hellinger distribution diversity and the Task2Vec task diversity coefficients, as both coefficients may be used interchangeably as a measure for the diversity of a given synthetic Gaussian benchmark** The first two plots show the relation between the Hellinger-based distribution diversity of a synthetic Gaussian benchmark and the benchmark’s performance on the MAML5, MAML10, and USL methods. These first two plots are noticeably similar to Figure 4 (where Task2Vec-based task diversity was used as a measure of diversity instead of Hellinger-based distribution diversity), which indicates that our Hellinger-based distribution diversity metric also serves as a good proxy for task diversity. The rightmost plot shows a strong positive correlation between Hellinger-based distribution diversity and Task2Vec task diversity (Pearson $r = 0.990$).

by computing the expected Hellinger distance between any two classes sampled from the benchmark. That is, for some benchmark parameterized by $B = (\mu_m, \sigma_m, \mu_s, \sigma_s)$, the diversity coefficient is given by

$$\text{div}(B) = \mathbb{E}_{\mu_1, \mu_2 \sim N(\mu_m, \sigma_m)} \mathbb{E}_{\sigma_1, \sigma_2 \sim |N(\mu_s, \sigma_s)|} [H^2(N(\mu_1, \sigma_1), N(\mu_2, \sigma_2))]$$

739 where H^2 denotes the squared Hellinger distance metric and $N(\mu_1, \sigma_1), N(\mu_2, \sigma_2)$ denote the
740 distributions of the two classes sampled from the benchmark. The Hellinger-based distribution
741 diversity metric provides an intuitive, model-agnostic characterization of the diversity of a benchmark
742 - the larger the diversity, the less similar any two classes within the benchmark are, and the easier it is
743 to distinguish between two classes. Conversely, the lower the diversity, the more similar any two
744 classes within the benchmark are, and the harder it is to distinguish between two classes due to a
745 larger overlap between the two classes’ distributions.

746

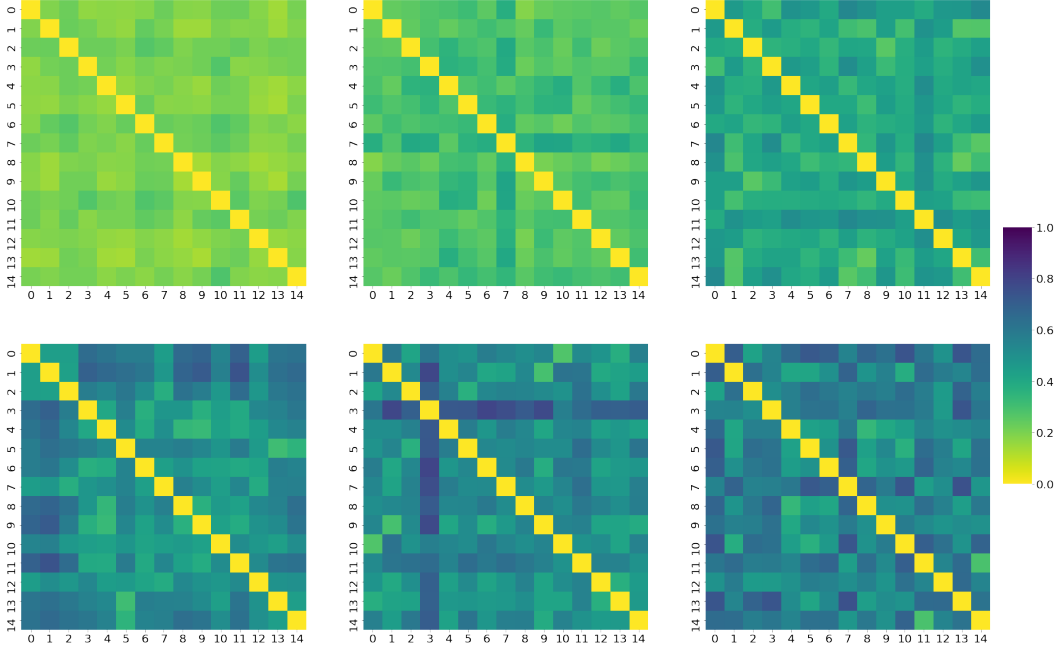


Figure 12: **Heatmaps show how benchmarks with larger Task2Vec-based task diversity coefficient show more heterogeneity between sampled tasks** Each heatmap below shows the pairwise distance between fifteen 5-way, 10-shot few-shot learning tasks sampled from the various synthetic Gaussian benchmarks described in Table 3. Note that as σ_m (the standard deviation of the class mean) increases, the distance between two tasks becomes larger on average and more varied, which can be seen as the heatmaps become more heterogeneous. This increase in expected distance among different tasks in turn increases the Task2Vec-based task diversity coefficient, which summarizes the average distance between any two tasks. From left to right, top to bottom, the benchmarks tested have parameters $\sigma_m = 0.01, 1, 3, 20, 30, 1000$ and Task2Vec-based task diversity coefficient parameters $div = 0.247, 0.271, 0.393, 0.533, 0.537, 0.546$.

Note that the closed-form equation for the Hellinger distance between the two class distributions $N(\mu_1, \sigma_1), N(\mu_2, \sigma_2)$ is given by

$$H^2(N(\mu_1, \sigma_1), N(\mu_2, \sigma_2)) = 1 - \sqrt{\frac{2\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2}} e^{-\frac{1}{4} \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}}$$

747 However, there is no simple closed-form equation for computing the diversity $div(B)$ itself. As a
 748 result, we computed the diversity coefficient as a numerical approximation by repeatedly sampling
 749 two classes from the benchmark distribution and calculating the Hellinger distance between the two
 750 classes. These samples ultimately provide a 95% confidence interval that represents the expected
 751 Hellinger distance between two classes sampled from the benchmark.
 752

753 We also compared our Hellinger-based distribution diversity coefficient with the Task2Vec-based task
 754 diversity coefficient for each of the synthetic Gaussian benchmarks tested in Table 3. We observe a
 755 strong positive correlation between the Hellinger-based distribution diversity and Task2Vec-based
 756 task diversity coefficients according to Figure 11 indicating that the Hellinger-based distribution
 757 diversity serves as an effective proxy for task diversity when the number of ways and shots of all tasks
 758 are fixed.

759 J Analysis of distribution of task distances in few-shot learning benchmarks

760 J.1 Heat Maps show Low Diversity and Homogeneity of tasks from MiniImagenet and 761 Cifar-fs

762 In this section, we show the heat maps showing the distances between 5-way, 20-shot few-shot
763 learning tasks from MiniImagenet and Cifar-fs in figure 13 and 14. We used 20-shots because we do
764 not need to separate the data into support and query set to compute the diversity coefficient. We show
765 that tasks sampled from these benchmarks create not only a low diversity coefficient on average, but
766 also at the level of individual distances between pairs of tasks. In addition, the heat map’s uniform
767 coloring reveals that it is also justifiable to call the tasks from these benchmarks *homogeneous*. Low
768 diversity is shown because the distances are between 0.07-0.12 given that max is 1.0 and minimum is
769 0.0.

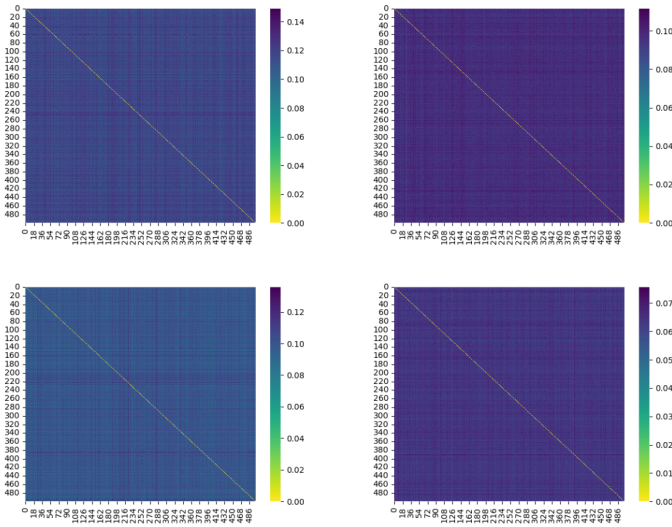


Figure 13: Shows homogeneity and low diversity of 5-way, 20-shot tasks from MiniImagenet using the Task2Vec distance [4]. The top left heat map uses a Resnet18 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The top right heat map uses a Resnet18 with random weights to compute the Task2Vec distance between tasks. The bottom left heat map uses a Resnet34 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The bottom right heat map uses a Resnet34 with random weights on Imagenet to compute the Task2Vec distance between tasks. Homogeneity is shown because of the uniform color shown in the heat map. Low diversity is shown because the distance is between 0.07-0.12 given that max is 1.0 and minimum is 0.0. Note the diagonal is exactly zero because it is comparing the same tasks. The axis indices indicate the arbitrary name for the tasks. Indices between heat maps do not indicate the same task. We used the cosine distance between task Task2Vec embeddings.

770 J.2 Histograms of distances of tasks in the synthetic Gaussian Benchmark, MiniImagenet and 771 Cifar-fs

772 In this section, we show the histograms of the cosine distances between pairs of tasks for the Gaussian
773 Benchmark, MiniImagenet and Cifar-fs. The main purpose of this is to argue that a (relatively small)
774 sample of the tasks is sufficient to estimate population statistics – like the expected distance between
775 tasks i.e. diversity coefficient.

776 For ease exposition of the argument, consider the case where we have 500 distance from a large
777 population of size $\binom{64}{5} = 7,624,512$. The goal is to argue that 500 samples are enough to make
778 strong statistical inferences about the population – even if it’s as large as 7,624,512. If we assume
779 the distribution of the data is Gaussian, then we expect to see a single mode with an approximate
780 bell curve. Therefore, if we plot the histogram of task pair distances of the 500 tasks and see this

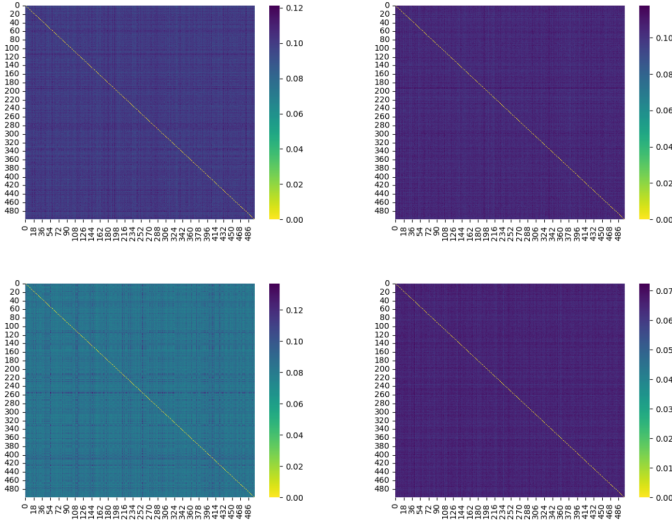


Figure 14: Shows homogeneity and low diversity of 5-way, 20-shot tasks from Cifar-fs using the Task2Vec distance. The top left heat map uses a Resnet18 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The top right heat map uses a Resnet18 with random weights to compute the Task2Vec distance between tasks. The bottom left heat map uses a Resnet34 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The bottom right heat map uses a Resnet34 with random weights on Imagenet to compute the Task2Vec distance between tasks. Homogeneity is shown because of the uniform color shown in the heat map. Low diversity is shown because the distance is between 0.07-0.12 given that max is 1.0 and minimum is 0.0. The axis indices indicate the arbitrary name for the tasks. Indices between heat maps do not indicate the same task. We used the cosine distance between task Task2Vec embeddings.

781 then we can infer our Gaussian assumption is approximately correct. Given that we do see that in
 782 figures 15, 16, 17 then we can infer our assumption is approximately correct. This implies we can
 783 make strong statistical assumptions about the population – in particular, that we have a good estimate
 784 of the diversity coefficient using 500 samples. Additionally, in histograms also discard the presence
 785 of outlier tasks.

786 As a remark, we want to emphasize that if we use T tasks to estimate the diversity coefficient, we
 787 in fact will use $\frac{T^2-T}{2} = O(T^2)$ distances to compute the diversity coefficient. This increases the
 788 computation cost, but does make the number of sample to compute the mean larger. For a subsample
 789 size of, $T = 500$ we in fact use $\frac{T^2-T}{2} = 124,750$ distances to compute the diversity coefficient.

790 K Background on distance metrics

791 K.1 Neuron Vectors

792 The representation of a neuron d in layer l is the vector $z_d^{(l)}(X) \in \mathbb{R}^N$ of activations for a set of N
 793 examples, where $X \in \mathbb{R}^{N,D}$ is the data matrix with N examples.

794 K.2 Layer Matrix

795 A layer matrix L for layer l is a matrix of neuron vectors $z_d^{(l)}(X) \in \mathbb{R}^N$ with shape, $[N, D_i]$ i.e.
 796 $L \in \mathbb{R}^{N, D_i}$. In other words, the layer matrix L is the subspace of \mathbb{R}^N spanned by its neuron vectors
 797 $z_d^{(l)}(X)$. In short, L is the layer matrix $[z_{d_1}^l; \dots; z_{d_i}^l] \in \mathbb{R}^{N, D_i}$ with neuron vector $z_{d_i}^l$.

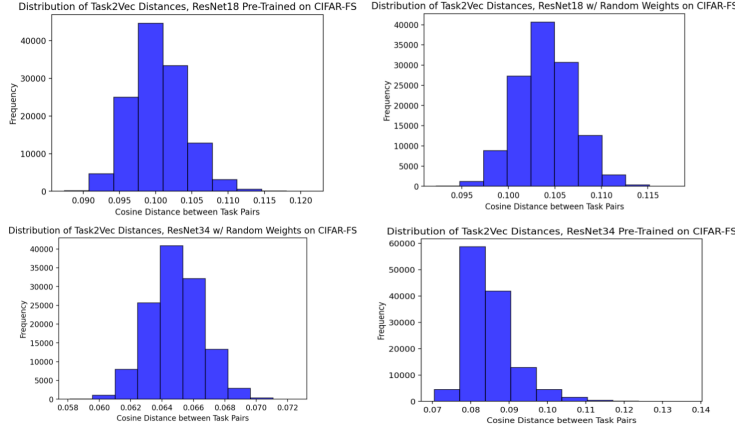


Figure 15: **Histogram of distances of 5-way, 20-shot tasks from Cifar-fs using the Task2Vec distance. This plot justifies the use of a subsample of the population to estimate the diversity coefficient because of its approximate Gaussian distribution.** For the full argument, see the main text, section J.2. The top left histogram uses a Resnet18 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The top right histogram uses a Resnet18 with random weights to compute the Task2Vec distance between tasks. The bottom left histogram uses a Resnet34 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The bottom right histogram uses a Resnet34 with random weights on Imagenet to compute the Task2Vec distance between tasks.

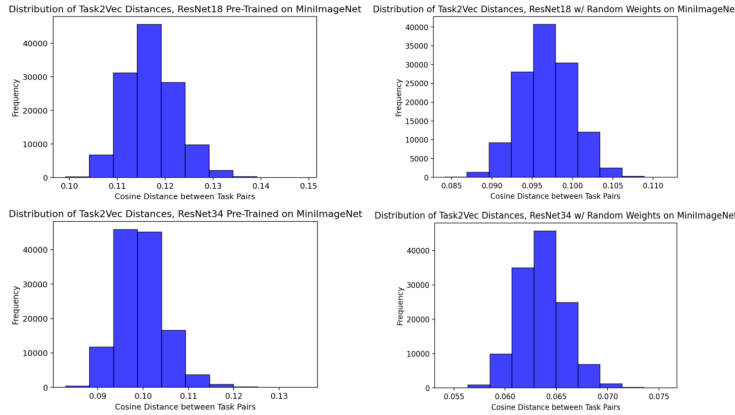


Figure 16: **Histogram of distances of 5-way, 20-shot tasks from MiniImageNet using the Task2Vec distance. This plot justifies the use of a subsample of the population to estimate the diversity coefficient because of its approximate Gaussian distribution.** For the full argument, see the main text, section J.2. The top left histogram uses a Resnet18 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The top right histogram uses a Resnet18 with random weights to compute the Task2Vec distance between tasks. The bottom left histogram uses a Resnet34 pre-trained on Imagenet to compute the Task2Vec distance between tasks. The bottom right histogram uses a Resnet34 with random weights on Imagenet to compute the Task2Vec distance between tasks.

798 **K.3 CCA**

799 Canonical Correlation Analysis (CCA) is a well established statistical technique for comparing the
 800 (linear) correlation of two sets of random variables (or vectors of random variables). In the empirical
 801 case, however, one computes the correlations between two sets of data sets (e.g. two matrices
 802 $X \in \mathbb{R}^{N, D_1}$ and $Y \in \mathbb{R}^{N, D_2}$ with N examples and D_1, D_2 features or layer matrices).

True distribution based Canonical Correlation Analysis (CCA): What we call true distribution based CCA is the standard CCA measure using the true but known distribution of the data $p^*(x)$ and $p^*(y)$. In this case, CCA searches for a pair of linear combinations a^*, b^* of two set of random

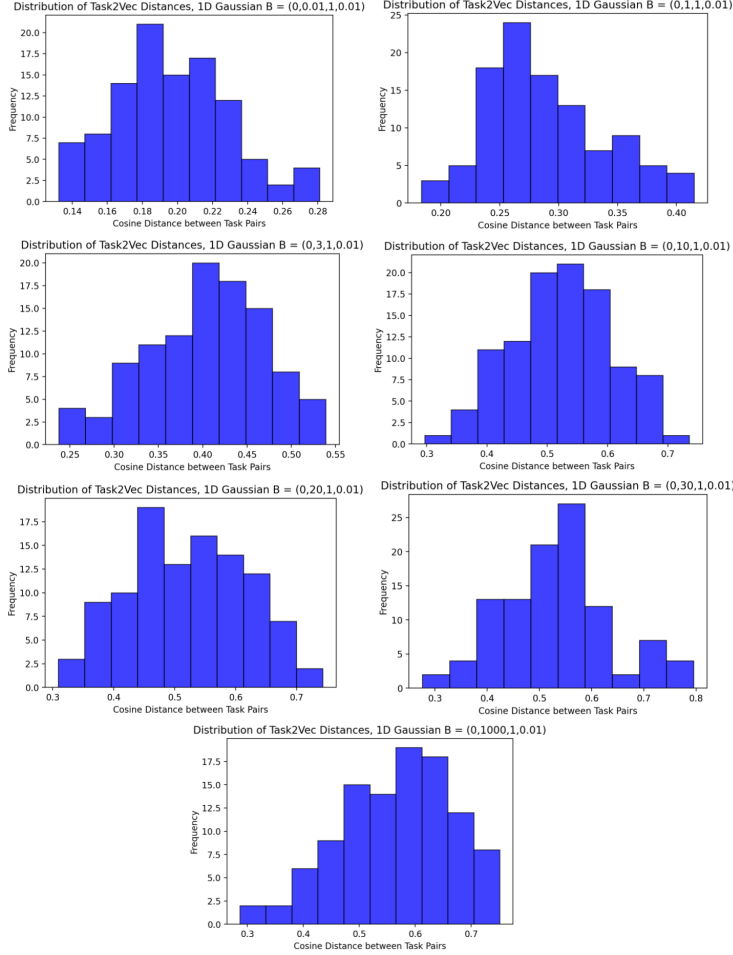


Figure 17: **Histogram of distances of 5-way, 20-shot tasks from the Synthetic Gaussian benchmark using the Task2Vec distance. This plot justifies the use of a subsample of the population to estimate the diversity coefficient because of its approximate Gaussian distribution.** For the full argument, see the main text, section J.2. The meta parameters generating tasks for each benchmark are denoted by $B = (0, x, 1, 0.01)$ where x is in the list $[0.01, 1, 3, 10, 20, 30, 1000]$ indicating the mean to generate the mean of the Gaussian tasks. For full details of the synthetic Gaussian benchmark, see section I.

variables (or vectors of random variables) $\mathbf{x} = [X_1, \dots, X_{D_1}]$ and $\mathbf{y} = [Y_1, \dots, Y_{D_2}]$ that maximizes the Pearson correlation coefficient:

$$a^*, b^* = \arg \max_{a, b} \frac{\mathbb{E}_{X, Y}[(a^\top \mathbf{x})((b^\top \mathbf{y}))]}{\sqrt{\mathbb{E}_X[(a^\top \mathbf{x})^2]} \sqrt{\mathbb{E}_Y[(b^\top \mathbf{y})^2]}} = \arg \max_{w_1, w_2} \frac{a^\top \Sigma_{X, Y} b}{\sqrt{a^\top \Sigma_{X, X} a} \sqrt{b^\top \Sigma_{Y, Y} b}}$$

803 where $\Sigma_{X, Y}, \Sigma_{X, X}, \Sigma_{Y, Y}$ are the (true) covariance and variance matrices respectively (e.g.
 804 $\Sigma_{X, Y}[i, j] = \text{Cov}[X_i, X_j] = [X_i Y_j]$ for centered random variables). All of these can be replaced by
 805 empirical data matrices in the obvious way.

806 K.4 SVCCA

807 At a high level, SVCCA is a similarity measure of two matrices that aims in removing redundant
 808 neurons (i.e. redundant features) with the truncated SVD by keeping 0.99 of the variance and then
 809 measure the overall similarity by averaging the top C CCA values.

810 **SV:** Given two matrices $L_1 \in \mathbb{R}^{N, D_1}, L_2 \in \mathbb{R}^{N, D_2}$ (e.g. layer matrices) first reduce the effective
 811 dimensionality of the matrix via a low rank approximation $L'_1 \in \mathbb{R}^{N, D'_1}, L'_2 \in \mathbb{R}^{N, D'_2}$ by choosing

812 the top k singular values that keeps 0.99 of the variance. In particular, for each layer matrix, L_i keep
 813 the top D'_i singular values (and vectors) such that $\sum_{j=1}^{D'_i} |\sigma_j| \geq 0.99 \sum_{j=1}^{\text{rank}(L_i)} |\sigma_j|$.

814 **SVCCA:** SVCCA is a statistical technique for the measuring the (linear) similarity of two sets
 815 of data sets $L_1 \in \mathbb{R}^{N, D_1}, L_2 \in \mathbb{R}^{N, D_2}$ (e.g. data matrices, layer matrices) by first reducing the
 816 effective dimensionality of the matrix via a low rank approximation $L'_1 \in \mathbb{R}^{N, D'_1}, L'_2 \in \mathbb{R}^{N, D'_2}$ (e.g.
 817 by choosing the top k singular values that keeps 0.99 of the variance) and then applying the standard
 818 empirical CCA to the resulting matrices. This is repeated $C = \min(D'_1, D'_2)$ times and the overall
 819 similarity of the two matrices is computed as the average CCA: $svcca = \text{sim}(L'_1, L'_2) = \frac{1}{C} \sum_{c=1}^C \rho_c$

820 Concretely:

- 821 1. Get the D'_i components that keep 0.99 of the variance (i.e. D'_i such that $\sum_{j=1}^{D'_i} |\sigma_j| \geq$
 822 $0.99 \sum_{j=1}^{\text{rank}(L_i)} |\sigma_j|$).
- 823 2. Get the SVD: $U_1, \Sigma_1, V_1^\top = \text{SVD}(L_1)$ and $U_2, \Sigma_2, V_2^\top = \text{SVD}(L_2)$
- 824 3. Then produce the SVD dimensionality reduction by $L'_1 = L_1 V_1 [1 : k_i] \in \mathbb{R}^{N, D_1}$ and
 825 $L'_2 = L_2 V_2 [1 : k] \in \mathbb{R}^{N, D_2}$ where $V_i [1 : D_i]$ gets the top D_i columns of a layer matrix i .
- 826 4. Get the CCA of the reduced layer matrix: $[\rho_c]_{c=1}^C = \text{CCA}(L'_1, L'_2)$ where, $C =$
 827 $\min(D'_1, D'_2)$
- 828 5. Finally return the mean CCA: $svcca = \frac{1}{C} \sum_{c=1}^C \rho_c$, where is the k -th CCA value of the
 829 reduced layer matrix.

830 K.5 PWCCA

831 At a high level, PWCCA was developed to increase the robustness (to noise) of SVCCA in the context
 832 of deep neural networks. In particular, Maithra et al. [35] noticed that when the performance of the
 833 neural networks stabilized, so did the set of CCA vectors (or principle neuron vectors) related to the
 834 network stabilized on the data set in question. Thus, they suggest to give higher weighting to the
 835 canonical correlation ρ_c of these stable CCA vectors – in particular to the ones that are similar to the
 836 final output layer matrix, e.g. L_1 . Note this is simpler than trying to track the stability of these CCA
 837 vectors during training and then give those higher weighting.

838 **PWCCA:** Formally let L_1 be the layer matrix $[z_d^l; \dots; z_{D_1}^l] \in \mathbb{R}^{N, D_1}$ with neuron vectors z_d^l
 839 for some layer l . Recall that the k -th left CCA vector for layer matrix L_1 is defined as follows,
 840 $\tilde{x}_c = L_1 a_c = L_1 (\Sigma^{-\frac{1}{2}} u_c)$ where a_c is the c th CCA direction and u_c is the c -th left singular value
 841 from the matrix $M = \Sigma_{L_1}^{-\frac{1}{2}} \Sigma_{L_1, L_2} \Sigma_{L_2}^{-\frac{1}{2}} = U \Lambda V^\top$. Then, PWCCA can be computed as follows:

- 842 1. Calculate the CCA vectors $\tilde{x}_c = L_1 a_c = L_1 (\Sigma_{L_1}^{-\frac{1}{2}} u_c)$ and explicitly orthonormalize with
 843 Gram-Schmidt for numerical stability.
- 844 2. Compute the weight $\tilde{\alpha}_c$ of how much the layer matrix L_1 is account for by each CCA vector
 845 \tilde{x}_k with equation $\tilde{\alpha}_c(h_c, L_1) = \sum_{c=1}^C |\langle \tilde{x}_c, z_c^l \rangle_{\mathbb{R}^N}|$ where z_c^l is the c -th column of the layer
 846 matrix L_1
- 847 3. Normalize the weight indicating how much each CCA vector h_c accounts for L_1 and denote
 848 it with, $\alpha_c(\tilde{x}_k, L_1) = \frac{\tilde{\alpha}_c(\tilde{x}_k, L_1)}{\sum_{c=1}^C \tilde{\alpha}_c(\tilde{x}_k, L_1)}$
- 849 4. Finally return the mean CCA weighted by $\alpha_c(\tilde{x}_k, L_1)$: $pwcca = \sum_{c=1}^C \alpha_c(\tilde{x}_k, L_1) \rho_c$ where
 850 $C = \min(D'_1, D'_2)$.

851 The original authors could have used the right CCA vectors, i.e. $\tilde{y}_c = L_2 b_k = L_2 (\Sigma^{-\frac{1}{2}} v_k)$ and in
 852 fact the details of their code suggest they choose the one that would have lead to less values removed
 853 by SVD. This choice seems to already be robust to noise, as shown in [35]. Note that the CCA vectors
 854 \tilde{x}_k, \tilde{y}_k are of size \mathbb{R}^N and thus could be viewed as the principle neuron vectors that correlate two
 855 layers L_1, L_2 . With this view, PWCCA computes the mean CCA normalized by of the c principle
 856 neuron vectors are account for the output layer matrix the most.

857 K.6 CCA for CNNs

858 The input to CCA are two data matrices, but CNNs have intermediate representations that are
859 4D tensors. Therefore, some justification is needed in how to create the data matrices needed for
860 computing CCA for CNNs. Note that it's the same reasoning for both SVCCA and PWCCA.

861 **Each channel as the dimensionality of the data matrix:** One option is to get the intermediate
862 representation of size $[M, C, H, W]$ and get a layer matrix of size $[MHW, C]$. Thus, MHW is the
863 effective number of data points and the channels (or number of filters) is the effective dimensionality
864 of the (layer) matrix. In this view, each patch of an image processed by the CNN is effectively
865 considered a data point. This view is very natural because it also considers each filter as its own
866 "neuron" – which seems reasonable considering that each filter uniquely responds to each stimulus
867 (e.g., data patch). This view results in HW images for every sample in the data set (or batch) of size
868 M and C effective neurons.

869 Although the original authors suggest this metric as a good metric mainly for comparing two layers
870 that are the same – we hypothesize it is also good for comparing different layers (as long as the
871 effective number of data points match for the two layer matrices). The reason is that CCA tries to
872 compute the maximum correlation of two data sets (or sets of random variables) and assumes no
873 meaning in the ordering of the data points and assumes no process for generating each individual
874 sample for the set of random variables, thus meaning that this metric (CCA) can be used for any two
875 layers in a matrix. Overall, in this view, we are comparing the representation learned in each channel.

876 **Each activation as the dimensionality of the data matrix:** One option is to get the intermediate
877 representation of size $[M, C, H, W]$ and get a layer matrix of size $[M, CHW]$. Thus, M is the
878 effective number of data points (which matches the number of samples in the data set or batch) and
879 therefore each activation value is the effective dimensionality of the (layer) matrix. In this view,
880 each activation is viewed as a neuron of size M and we have CHW effective neurons for each
881 activation. The authors suggest this metric for comparing different layers (potentially at different
882 depths). However, because CCA assumes no correspondence between the data points nor the same
883 dimensionality in the data matrix – we hypothesize this way to define the data matrix is as valid as
884 the previous definition for comparisons between any models at any layer. One disadvantage however
885 is that it will often result in data matrices that are very large due to CHW being very large – which
886 results in artificially high CCA similarity values. Potential ways to deal with it are noticing that there
887 is no correspondence between the data matrices, so a cross comparison of every data point with every
888 other data point in CCA is possible (resulting in $O(M^2)$ comparisons for the empirical covariance
889 matrix). Alternatively one can pool in the spatial dimensions $[H, W]$ resulting in potential smaller
890 layer matrices e.g. of shape $[M, C]$ with a pool over the entire spatial dimension. For these reasons
891 and the fact that we hypothesize an image patch being its own image – we prefer to interpret the
892 number of channels as the natural way to compare CNNs so that the layer matrices results of size
893 $[MHW, C]$.

894 **Subsampling of representations for channels as dimensionality:** In this section, we review the
895 subsampling we did when comparing the representations learned in each channel, i.e. the layer
896 matrix has size $[MHW, C]$. The effective number of data points MHW will often be much larger
897 than needed (e.g. for 16 data samples $M = 16$ and $H = W = 84$ results in $MHW = 112, 896$),
898 especially compared to the number of filters/channels (e.g. $C = 64$). Previous work [35, 38] suggest
899 using the number of effective data points to be from 5-10 times the size of the dimensionality in a
900 layer matrix of size $[N', D']$ that means $N' = 10D'$. Based on our reproductions of that number, we
901 choose $N' = 20D'$ which results in $NHW = 20C$

902 K.7 Centered Kernel Alignment (CKA)

903 At a high level, CKA is based on the insight that one can first measure the similarity between every
904 pair of examples in each representation separately and then use the similarity structure to compute an
905 overall similarity metric. In our case, we can treat the examples as the neuron vectors and compare
906 all neuron vectors using some kernel function. Usually this will end in a kernel matrix of size M, M'
907 where M and M' are the number of examples. In our case, they would be D, D' for the number of
908 neurons of each layer matrix. Note, the layers matrices can correspond to neurons of different layers
909 in a neural network.

910 **Linear CKA:** We use the linear kernel function as used in previous work [16, 24]. Given two layer
 911 matrices $X_1 \in \mathbb{R}^{N, D_l}$ and $X_2 \in \mathbb{R}^{N, D_{l'}}$ for layers l, l' , we compute the linear kernel $X_1^\top X_2$ to get
 912 the D_l by $D_{l'}$ kernel matrix indicating the (linear) similarity per neuron vector for the two layers.
 913 Then to obtain a single distance value we compute the Frobenius norm of the kernel matrix and
 914 subtract by one after normalization:

$$d_{linearCKA}(A, B) = 1 - \frac{\|A^\top B\|_F}{\|A^\top A\|_F \|B^\top B\|_F} \quad (3)$$

915 Note that depending on how the examples in matrices A, B are organized the cross-product could
 916 be computed with AB^\top instead. Other kernel functions have been tested (e.g., the RBF kernel) for
 917 CKA but similar results are obtained, resulting in linear CKA being the most popular CKA method
 918 [16, 24] to the best of our knowledge.

919 K.8 Orthogonal Procrustes Distance (OPD)

920 At a high level, the orthogonal Procrustes distance computes the distances between two matrices after
 921 using for the best orthogonal matrix that tries to match the two. Usually this is done after centering
 922 and dividing by the Frobenius norm of the matrices, i.e. normalizing the matrices. In addition, previous work
 923 [16] finds that OPD is a better metric at detecting changes that matter functionally and robust against
 924 changes that do not matter.

925 **OPD:** Formally, the Orthogonal Procrustes Distance is the smallest distance between two matrices X
 926 and Y (with columns as the vectors in question) found by finding the orthogonal matrix Q which most
 927 closely maps A to B . Therefore, the OPD distance is the distance value from solving the orthogonal
 928 Procrustes problem:

$$d_{OPD}(X, Y) = \min_Q \|X - YQ\|_F \quad (4)$$

929 where $\|\cdot\|_F$ is the Frobenius norm. When matrices are normalized (centered and divided by their
 930 Frobenius norm) this is called the general Procrustes problem. However, the closed form equation we
 931 used is the following:

$$d'_{OPD}(X, Y) = \frac{1}{2} (\|X\|_F + \|Y\|_F - 2\|X^\top Y\|_*) \quad (5)$$

932 where $\|\cdot\|_*$ is the nuclear norm, i.e. the sum of singular values $\sum_i \sigma(A)_i = \|A\|_*$. The division
 933 by 2 is to guarantee that the OPD distance is between $[0, 1]$ instead of $[0, 2]$. We do the standard
 934 normalization of the matrices before computing the OPD distance – by centering and dividing by
 935 the Frobenius norm of the matrix. This is done because the orthogonal matrix in the orthogonal
 936 Procrustes problem does not allow for translation or rescaling of the matrices. Therefore, this
 937 normalization enforces invariance to this type of transformations – i.e. we don't want large OPD
 938 values due to rescaling or translation (and even if present, the orthogonal matrix wouldn't be able to
 939 reflect it).

940 Therefore, the final equation for OPD we use is:

$$d_{OPD}(X, Y) = 1 - \|X^\top Y\|_* \quad (6)$$

941 **Why OPD?** We use OPD due to the findings of [16]. They find that OPD is a more robust metric
 942 (compared to SVCCA, PWCCA, and CKA) because it is *sensitive* to changes that affect real functional
 943 behavior (so it detects changes to behavior that “matter”) and it's *specific* against changes that do
 944 not. As a summary, some of the evidence that they provide for this is that OPD is able to detect
 945 when 0.75 of the principal components are removed, while CKA cannot detect removal of principal
 946 components until 0.97 are removed. CCA like metrics on the other hand are not specific – even
 947 random initialization noise overwhelms the distances it reports, while OPD is more robust to this
 948 random noise. For the last point, this means that even if we compare two different layers with CCA,
 949 the noise will dominate the distance reported instead of the difference caused by comparing different
 950 layer.

951 K.9 Correctly using Feature Based Distances

952 When comparing two layers of a neural network using two layer matrices, one needs to be careful with
 953 the number of data points (or batch size) being used. This is because metrics like CCA intrinsically

954 are formulated as an optimization and if the number of examples is not larger than the number
 955 of dimensionality of the examples – then the similarity can be pathologically be perfect (e.g., the
 956 distance is zero when it’s actually not zero). Therefore, we follow the suggestions by the original
 957 authors of SVCCA [38] and always use at least 10 times more examples than there are features for our
 958 feature based comparisons. We call this value the *safety margin*. To illustrate this idea, we produce
 959 two random matrices and compare how the similarity (SVCCA) values varies as a function of the
 960 dimensionality of the data and the number of points. Since the two matrices are completely random,
 961 we know they should not be very similar and thus SVCCA should report a high similarity value (or
 962 low distance value). Therefore, we can see in figure 18 how as the dimensionality increases, the
 963 similarity value approaches a perfect similarity of 1.0. In figure 19 we can see how as the number of
 964 points increases, we approach a smaller similarity – closer to the true similarity for random matrices.

965 In general, given two matrices $X, Y \in \mathbb{R}^{M', D'}$ with the number of (effective) data points M' and
 966 (effective) dimensionality (number of features) D' – we want the number of points to be larger by
 967 a safety factor s . Formally, it must satisfy this inequality to avoid the pathological case for feature
 968 based distances:

$$D' \leq sM' \quad (7)$$

969 where we suggest to use $s \leq 10$ (as used in previous work [38]). Note the effective number of data
 970 points used and dimensionality can be different depending on how one reshapes the CNN tensors
 971 to produce layer matrices as explained in section K.6. For example, if one uses the channel as
 972 the dimensionality (i.e., use the image patches as an effective data point) then one has to obey the
 973 following inequality:

$$C \leq sMHW \quad (8)$$

974 where M is the batch size, H, W is the height and width of the images, and C is the number of
 975 filters/channels for the current layer. This means that for a given architecture processing images of a
 976 given size that the only parameter we can change to make the above inequality true is the batch size
 977 M .

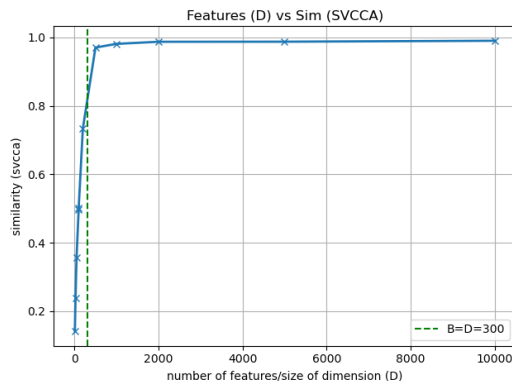


Figure 18: Shows that as the dimensionality of a random data matrix increases – the SVCCA similarity approaches the pathological case by falsely reports the similarity is perfect. The green line indicates when the number of examples and dimensionality are equal (and equal to 300). D denotes the dimensionality of the simulated data and B the size of the batch size/number of points.

978 L A Statistical Decision view of the differences between Supervised Learning 979 and Meta-learning

980 Recent work in meta-learning implies that feature-reuse might be all we need to solve modern
 981 few-shot learning benchmarks [43]. However, what it also reveals is our poor understanding of
 982 meta-learning algorithms. Therefore, in this section, we take the most foundational perspective
 983 to formulate and analyze meta-learning algorithms by analyzing them from an optimal statistical
 984 decision theory perspective?

985 We hope that this can help clarify the results from [43] and therefore help meta-learning researchers
 986 design better meta-learning benchmarks and meta-learning algorithms.

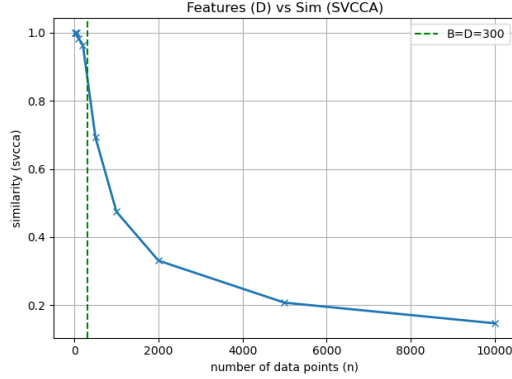


Figure 19: **Show how to avoid the pathological case when using feature based similarities by increasing the number of data points (or batch size).** In particular, as the number of data points in two random data matrix increases – the true similarity approaches the true low similarity value. The green line indicates when the number of examples and dimensionality are equal (and equal to 300). D denotes the dimensionality of the simulated data and B the size of the batch size/number of points.

987 **L.1 Supervised Meta-Learning problem set-up**

988 In this section, we introduce the notation for supervised meta-learning. Intuitively, we seek to find a
 989 function that minimize the expected risk over tasks and the data in the tasks. To formalize it, we will
 990 use three formulations:

991 **Monolithic meta-learner:** for a monolithic decision rule g (or meta-learner), we want to find the
 992 optimal g by minimizing the *supervised meta-learning expected risk*:

$$R_{Mono}(g) = \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x, y \sim p(x, y | \tau)} [l(g(x, \tau), y)] \quad (9)$$

993 where g is a single monolithic function, $p(\tau)$ is the true but unknown distribution of tasks, $p(x, y | \tau)$
 994 is the true, but unknown distribution of data pair given a task τ and (x, y) is the data pair of input and
 995 target value sampled from a task.

996 **Meta-learned meta-learner:** for a meta-learned decision rule we usually have an adaptation rule A
 997 (e.g. SGD in MAML) and a function approximator h (e.g. a neural network) and minimize the follow
 998 over both:

$$R_{ML}(A, h) = \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x, y \sim p(x, y | \tau)} [l(A(h, \tau)(x), y)] \quad (10)$$

999 $p(\tau)$ is the true but unknown distribution of tasks, $p(x, y | \tau)$ is the true, but unknown distribution of
 1000 data pair given a task τ and (x, y) is the data pair of input and target value sampled from a task.

1001 **Fixed representation meta-learner without adaptation:** one can also solve [9](#) using a single decision
 1002 rule f that does not take the task τ as input as follows:

$$R_{SL}(f) = \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x, y \sim p(x, y | \tau)} [l(f(x), y)] \quad (11)$$

1003 where f is a function to be adapted (e.g. a neural network), $p(\tau)$ is the true but unknown distribution
 1004 of tasks, $p(x, y | \tau)$ is the true, but unknown distribution of data pair given a task τ and (x, y) is the
 1005 data pair of input and target value sampled from a task.

1006 **Fixed representation meta-learner with a final adaptation layer:** one can also solve [9](#) using a
 1007 single feature extractor g that does not take the task τ as input with a feature extractor g :

$$R_{SLA}(f, g) = \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x, y \sim p(x, y | \tau)} [l((f(\tau) \circ g)(x), y)] \quad (12)$$

1008 where g is the feature extractor from the raw inputs (e.g. a neural network), f the final layer adapted
 1009 (e.g. a linear layer), $p(\tau)$ is the true but unknown distribution of tasks, $p(x, y | \tau)$ is the true, but
 1010 unknown distribution of data pair given a task τ and (x, y) is the data pair of input and target value
 1011 sampled from a task.

1012 *Remark L.1.* Note that in practice, the meta-learner does not usually take the full task τ as input, but
 1013 instead a train and test set (often referred to as support set and query set) sampled from the task τ .

1014 The goal of this work is to clarify the difference between [11](#) and [10](#) under the framework of statistical
 1015 decision theory. Arguably the most important comparison between [10](#) and [12](#) is left for future work.

1016 **L.2 Main Result: Difference between the Supervised Learned and Meta-learned decision rule**

1017 The proof sketch is as follows: we first show the optimal decision rules for both supervised learning
 1018 and meta-learning when minimizing the expected meta-risk from equations [11](#) and [10](#) and then
 1019 highlight that the main difference between them is that the meta-learned solution can act optimally if
 1020 it identifies the task τ while the supervised learned solution has no capabilities of this since it learns
 1021 an average based on tasks instead.

1022 **Theorem L.2.** *The minimizer to equation [10](#) is:*

$$A(h, \tau)(x) = \bar{y}_{y|x, \tau}^* = \mathbb{E}_{y \sim p(y|x, \tau)} [y] \quad (13)$$

1023 where $\bar{y}_{y|x, \tau}^* = \mathbb{E}_{y \sim p(y|x, \tau)} [y]$ and l is the squared loss $l(\hat{y}, y) = (\hat{y} - y)^2$.

Proof. The proof is the same as the standard decision rule textbook proof but instead of minimizing it point-wise w.r.t. x we minimize it point-wise w.r.t. (x, τ) . In particular, we have:

$$\begin{aligned} R_{ML}(A, h) &= \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x, y \sim p(x, y|\tau)} [l(A(h, \tau)(x), y)] \\ \min_{A, h} \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x \sim p(x|\tau)} \mathbb{E}_{y \sim p(y|x, \tau)} [(A(h, \tau)(x) - y)^2] \end{aligned}$$

without loss of generality (WLOG) and for clarity of exposition consider the special case for discrete variables:

$$\min_{A, h} \sum_{\tau} p(\tau) \sum_x p(x | \tau) \mathbb{E}_{y \sim p(y|x, \tau)} [(A(h, \tau)(x) - y)^2]$$

At this point we notice we can minimize the above point-wise w.r.t (x, τ) and ignore h . To do that, take the derivative of $R(A, h)$ with respect to $A(h, \tau)(x)$ because that $A(h, \tau)(x) \in \mathbb{R}$ and set it to zero:

$$\begin{aligned} \frac{d}{dA(h, \tau)(x)} \mathbb{E}_{y \sim p(y|x, \tau)} [(A(h, \tau)(x) - y)^2] &= 0 \\ \mathbb{E}_{y \sim p(y|x, \tau)} [(A(h, \tau)(x) - y)] &= 0 \\ \mathbb{E}_{y \sim p(y|x, \tau)} [(A(h, \tau)(x))] &= \mathbb{E}_{y \sim p(y|x, \tau)} [y] \\ A(h, \tau)(x) &= \mathbb{E}_{y \sim p(y|x, \tau)} [y] = \bar{y}_{y|x, \tau}^* \end{aligned}$$

1024 as desired. □

1025 **Corollary L.3.** *For a monolithic meta-learner defined in section [L.1](#) the solution to the meta
 1026 supervised learning problem is the same as in equation [13](#) for the squared loss $l(\hat{y}, y) = (\hat{y} - y)^2$ i.e.
 1027 $g(\tau, x) = \bar{y}_{y|x, \tau}^* = \mathbb{E}_{y \sim p(y|x, \tau)} [y]$.*

1028 *Proof.* Proof is trivial, replace $A(h, \tau)(x)$ with $g(\tau, x)$ since h is not used. In this case, there is
 1029 no difference with having an adaptation rule A equipped with another function h and a monolithic
 1030 meta-learner g . □

1031 **Theorem L.4.** *The minimizer to equation [11](#):*

$$f(x) = \mathbb{E}_{\tau \sim p(\tau|x)} [\bar{y}_{y|x, \tau}^*] \quad (14)$$

1032 where $\bar{y}_{y|x, \tau}^* = \mathbb{E}_{y \sim p(y|x, \tau)} [y]$ and l is the squared loss $l(\hat{y}, y) = (\hat{y} - y)^2$.

Proof. WLOG, consider the minimizer of equation [11](#) in the discrete case. In particular, we have:

$$\begin{aligned} R_{SL}(f) &= \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x, y \sim p(x, y|\tau)} [l(f(x), y)] \\ \min_f \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x \sim p(x|\tau)} \mathbb{E}_{y \sim p(y|x, \tau)} [(f(x) - y)^2] \\ \min_f \sum_x \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x, \tau)} [(f(x) - y)^2] \end{aligned}$$

Note we can minimize the above point-wise w.r.t. x only (and not also w.r.t. τ as we did in proof [L.2](#)). Thus, we have want:

$$f(x) = \min_{f(x) \in \mathbb{R}} \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x, \tau)} [(f(x) - y)^2]$$

at this point it is interesting to observe the disadvantage of supervised learning methods with fixed functions without dependence on the task is that they are forced to consider all task τ at once. We proceed to take derivatives as in proof [L.2](#) but with this objective:

$$\begin{aligned} & \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x,\tau)} [(f(x) - y)^2] \\ & \frac{d}{df(x)} \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x,\tau)} [(f(x) - y)^2] = 0 \\ & \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x,\tau)} [f(x)] = \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x,\tau)} [y] \\ & f(x) \mathbb{E}_{\tau \sim p(\tau)} [p(x | \tau)] = \mathbb{E}_{\tau \sim p(\tau)} p(x | \tau) \mathbb{E}_{y \sim p(y|x,\tau)} [y] \end{aligned}$$

1033

$$f(x) = \mathbb{E}_{\tau \sim p(\tau)} \left[\frac{p(x | \tau)}{\mathbb{E}_{\tau \sim p(\tau)} [p(x | \tau)]} \mathbb{E}_{y \sim p(y|x,\tau)} [y] \right] \quad (15)$$

We proceed by noticing that $\mathbb{E}_{\tau \sim p(\tau)} [p(x | \tau)] = p(x)$, thus:

$$\begin{aligned} f(x) &= \mathbb{E}_{\tau \sim p(\tau)} \left[\frac{p(x | \tau)}{p(x)} \mathbb{E}_{y \sim p(y|x,\tau)} [y] \right] \\ f(x) &= \sum_{\tau} p(\tau) \frac{p(x | \tau)}{p(x)} [\mathbb{E}_{y \sim p(y|x,\tau)} [y]] \\ f(x) &= \sum_{\tau} \frac{p(\tau)}{p(x)} \frac{p(x, \tau)}{p(\tau)} [\mathbb{E}_{y \sim p(y|x,\tau)} [y]] \\ f(x) &= \sum_{\tau} \frac{p(x, \tau)}{p(x)} [\mathbb{E}_{y \sim p(y|x,\tau)} [y]] \\ f(x) &= \sum_{\tau} p(x | \tau) [\mathbb{E}_{y \sim p(y|x,\tau)} [y]] \\ f(x) &= \mathbb{E}_{\tau \sim p(x|\tau)} [\mathbb{E}_{y \sim p(y|x,\tau)} y] \\ f(x) &= \mathbb{E}_{\tau \sim p(x|\tau)} [\bar{y}_{y|x,\tau}^*] \end{aligned}$$

1034 as required by the rightmost RHS of equation [14](#). □

1035 **Theorem L.5.** The minimizer in equation [14](#) reduces to an expectation only over w.r.t. $p(\tau)$ of $\bar{y}_{y|x,\tau}^*$
1036 under benchmarks that are balanced. Formally

$$f(x) = \mathbb{E}_{\tau \sim p(\tau)} [\bar{y}_{y|x,\tau}^*] = \mathbb{E}_{\tau \sim p(\tau)} [\mathbb{E}_{y \sim p(y|x,\tau)} [y]] \quad (16)$$

1037 where $\bar{y}_{y|x,\tau}^* = \mathbb{E}_{y \sim p(y|x,\tau)} [y]$ and under assumption A1: $p(x | \tau)$ is a constant, i.e. $p(x | \tau) =$
1038 $k_{XT} \in \mathbb{R}, \forall x \in X, \forall \tau \in T$ and l is the squared loss $l(\hat{y}, y) = (\hat{y} - y)^2$.

Proof. Recall equation [14](#):

$$f(x) = \mathbb{E}_{\tau \sim p(\tau|x)} [\bar{y}_{y|x,\tau}^*]$$

due to Bayes's rule we have $p(\tau | x) = \frac{p(\tau)p(x|\tau)}{p(x)}$ and equation [14](#) can be re-written as follows:

$$f(x) = \mathbb{E}_{\tau \sim p(\tau)} \left[\frac{p(x | \tau)}{p(x)} \bar{y}_{y|x,\tau}^* \right]$$

under assumption A1 we have that $p(x | \tau)$ does not depend on as a function of x or τ . Thus, we have:

$$p(x) = \sum_{\tau} p(\tau) p(x | \tau) = p(x | \tau) \sum_{\tau} p(\tau) = p(x | \tau)$$

Thus we have:

$$\begin{aligned} f(x) &= \mathbb{E}_{\tau \sim p(\tau)} \left[\frac{p(x)}{p(x)} \bar{y}_{y|x,\tau}^* \right] \\ f(x) &= \mathbb{E}_{\tau \sim p(\tau)} [\bar{y}_{y|x,\tau}^*] \end{aligned}$$

1039 as required. □

1040 *Remark L.6.* Note that assumption A1 holds for the common MiniImagenet few-shot learning data
 1041 set, where $p(x | \tau) = \frac{1}{600}$.

1042 *Remark L.7.* In addition, because all classes are equally likely (e.g. $p(class) = \frac{1}{64}$ for the meta-train
 1043 set) we have $p(\tau)$ is the same constant independent of the task τ . Proof in the appendix, lemma [L.8](#).

1044 **Theorem L.8.** *If the tasks are equally likely, then equation [16](#) becomes an average over conditional
 1045 predictions over all tasks. Formally, if $p(\tau) = \frac{1}{T}$ then equation [16](#) becomes:*

$$f(x) = \frac{1}{T} \sum_{\tau} \bar{y}_{y|x,\tau}^* \quad (17)$$

1046 *under the squared loss $l(\hat{y}, y) = (\hat{y} - y)^2$.*

1047 *Proof.* Since $f(x) = \mathbb{E}_{\tau \sim p(\tau)} [\bar{y}_{y|x,\tau}^*]$ then, plugging $p(\tau) = \frac{1}{T}$ completes proof. \square

1048 *Remark L.9.* It is interesting to note that without adaptation or dependence on the task τ being
 1049 solved, the supervised learned meta-learner is suboptimal compared to the meta-learned solution.
 1050 The proof is simple, and it follows because the meta-learned decision rule was chosen to minimize
 1051 each term individually, but the supervised learned decision is not of that form. Proof in appendix
 1052 [L.11](#). Unfortunately, note that this does not necessarily apply to previous work [\[43\]](#).

1053 *Remark L.10.* Note that remark [L.9](#) does not apply to work [\[43\]](#) because that work does depend on a
 1054 task τ during meta-test time by adapting the final layer even if the representation is fixed.

1055 *Remark L.11.* The supervised learning decision rule is suboptimal compared to the meta-learned
 1056 decision rule.

1057 **L.3 The supervised Learning Solution is equivalent to the Meta-Learning solution when there** 1058 **is low task diversity**

1059 **Sketch argument:** The main idea is that because all tasks are very similar (task diversity is low)
 1060 – it essentially means that τ is not truly an input to the adaptation rule or monolithic meta-learner).
 1061 Equivalently, the problem is essentially a single task problem, so the task is implicitly an input to any
 1062 method used. Therefore, since the task conditioning does not exist, then the optimization problem is
 1063 the same for the meta-learned solution and when there is a fixed supervised learning feature extractor.

1064 **Theorem L.12.** *Assume $\tau_1 = \tau_2$ for any tasks in T and the data sets are balanced (i.e. same number
 1065 of images x for each task). Then we have the meta-learned solution is the same as the supervised
 1066 learning solution with shared embeddings: $f_{sl}(x) = A(f_{ml}, \tau)(x)$.*

Proof. Consider the optimization problem, for supervised learning:

$$\min_{A,h} \mathbb{E}_{\tau \sim p(\tau)} \mathbb{E}_{x \sim p(x|\tau)} \mathbb{E}_{y \sim p(y|x,\tau)} [(A(h, \tau)(x) - y)^2]$$

1067 If every pair of tasks is equal, it means their distributions are equal $p(x, y | \tau) = p(x, y)$ (mean-
 1068 ing τ can be ignored). Thus, the solution to the supervised learning problem is: $f_{sl}(x) =$
 1069 $\mathbb{E}_{\tau} \mathbb{E}_{p(x,y)} [y] = \mathbb{E}_{p(x,y)} [y] = y_x^*$. Now for the meta-learning problem we have: $A(f_{ml}, \tau)(x) =$
 1070 $y_{x,\tau}^* = \mathbb{E}_{y \sim p(y|x,\tau)} [y]$ but due to every pair of tasks being equal means $p(x, y | \tau) = p(x, y)$ (i.e.
 1071 all task share the same distributions) we have: $A(f_{ml}, \tau)(x) = \mathbb{E}_{y \sim p(y|x,\tau)} [y] = \mathbb{E}_{y \sim p(y|x)} [y] = y_x^*$
 1072 which is the same as the solution as in f_{sl} . Thus $f_{sl}(x) = A(f_{ml}, \tau)(x)$. \square

1073 *Remark L.13.* Proofs were presented in the discrete case clarity, but it is trivial to expand them to the
 1074 continuous case – e.g., using integrals instead of summations.

1075 **M Summary of Compute Required**

1076 We used an internal compute cluster with wide varied of GPUs. We used Titan X GPUs for most
 1077 five layer CNN experiments. We used A40 and dgx-A100 GPUs for Resnet12 experiments, with 48
 1078 GB and 40 GB GPU memory respectively. We did notice that the Resnet12 architecture we used
 1079 from previous work [\[43\]](#) required more memory than Resnet18 and Resnet34 used in Task2Vec
 1080 [\[4\]](#). By requiring more memory, we mean we did not have many memory out of bound issues with

1081 Resnet18/Resnet34 but did have memory issues with Resnet12. In addition, our episodic meta-
1082 learning training for MAML used Learn2Learn’s [5] distributed training to speed up experiments.
1083 Experiments took 1-2 weeks with MAML in a single GPU to potentially 2-3 days with multiple
1084 GPUs (we used 2, 4 to 8 GPUs depending on availability). For synthetic experiments we used Titan
1085 X GPUs with 16GB of GPU memory. Experiments took around 1-2 days on average with a single
1086 GPU. For more precision check the experimental details section [D].