

---

# Generalization Bounds for Meta-Learning via PAC-Bayes and Uniform Stability

---

**Alec Farid      Anirudha Majumdar**

Department of Mechanical and Aerospace Engineering, Princeton University  
`{afarid, ani.majumdar}@princeton.edu`

## Abstract

We are motivated by the problem of providing strong generalization guarantees in the context of meta-learning. Existing generalization bounds are either challenging to evaluate or provide vacuous guarantees in even relatively simple settings. We derive a probably approximately correct (PAC) bound for gradient-based meta-learning using two different generalization frameworks in order to deal with the qualitatively different challenges of generalization at the “base” and “meta” levels. We employ bounds for uniformly stable algorithms at the base level and bounds from the PAC-Bayes framework at the meta level. The result of this approach is a novel PAC bound that is tighter when the base learner adapts quickly, which is precisely the goal of meta-learning. We show that our bound provides a tighter guarantee than other bounds on a toy non-convex problem on the unit sphere and a text-based classification example. We also present a practical regularization scheme motivated by the bound in settings where the bound is loose and demonstrate improved performance over baseline techniques.

## 1 Introduction

A major challenge with current machine learning systems is the need to acquire large amounts of training data in order to learn a new task. Over the past few decades, meta-learning [62, 70] has emerged as a promising avenue for addressing this challenge. Meta-learning relies on the intuition that a new task often bears significant similarity to previous tasks; hence, a learner can learn to perform a new task very quickly by exploiting data from previously-encountered related tasks. The meta-learning problem formulation thus assumes access to datasets from a variety of tasks during meta-training. The goal of the meta learner is then to learn inductive biases from these tasks in order to train a base learner to achieve few-shot generalization on a new task.

Over the past few years, there has been tremendous progress in practical algorithms for meta-learning (see, e.g., [61, 55, 25, 32]). Techniques such as model-agnostic meta-learning (MAML) [25] have demonstrated the ability to perform few-shot learning in a variety of supervised learning and reinforcement learning domains. However, our theoretical understanding of these techniques lags significantly behind successes on the empirical front. In particular, the problem of deriving *generalization bounds* for meta-learning techniques remains an outstanding challenge. Current methods for obtaining generalization guarantees for meta-learning [5, 34, 77] either (i) produce bounds that are extremely challenging to compute or (ii) produce vacuous or near-vacuous bounds in even highly simplified settings (see Section 5 for numerical examples). Indeed, we note that existing work on generalization theory for meta-learning techniques do not explicitly report numerical values for generalization bounds. This is in contrast to the state of generalization theory in the supervised learning setting, where recent techniques demonstrate the ability to obtain non-vacuous generalization guarantees on benchmark problems (e.g. visual classification problems [24, 78, 54]).

The generalization challenge in meta-learning is similar to, but distinct from, the supervised learning case. In particular, any generalization bound for meta-learning must account for *two levels* of

generalization. First, one must account for generalization at the base level, i.e., the ability of the base learner to perform well on new data from a given task. This is particularly important in the few-shot learning setting. Second, one must account for generalization at the meta level, i.e., the ability of the meta learner to generalize to new tasks not encountered during meta-training. Moreover, the generalization performance at the two levels is coupled since the meta learner is responsible for learning inductive biases that the base learner can exploit for future tasks.

The key technical insight of this work is to bound the generalization error at the two levels (base and meta) using two *different* generalization theory frameworks that each are particularly well-suited for addressing the specific challenges of generalization. At the base level, we utilize the fact that a learning algorithm that exhibits uniform stability [14, 15] also generalizes well in expectation (see Section 4.1 for a formal statement). Intuitively, uniform stability quantifies the sensitivity of the output of a learning algorithm to changes in the training dataset. As demonstrated by [29], limiting the number of training epochs of a gradient-based learning algorithm leads to uniform stability. In other words, a gradient-based algorithm that *learns quickly* is stable. Since the goal of meta-learning is *precisely* to train the base learner to learn quickly, we posit that generalization bounds based on stability are particularly well-suited to bounding the generalization error at the base level. At the meta level, we employ a generalization bound based on *Probably Approximately Correct (PAC)-Bayes* theory. Originally developed two decades ago [43, 38], there has been a recent resurgence of interest in PAC-Bayes due to its ability to provide strong generalization guarantees for neural networks [24, 8, 54]. Intuitively, the challenge of generalization at the meta level (i.e., generalizing to new tasks) is similar to the challenge of generalizing to new data in the standard supervised learning setting. In both cases, one must prevent over-fitting to the particular tasks/data that have been seen during meta-training/training. Thus, the strong empirical performance of PAC-Bayes theory in supervised learning problems makes it a promising candidate for bounding the generalization error at the meta level.

**Contributions.** The primary contributions of this work are the following. First, we leverage the insights above in order to develop a novel generalization bound for gradient-based meta-learning using uniform stability and PAC-Bayes theory (Theorem 3). Second, we develop a regularization scheme for MAML [25] that explicitly minimizes the derived bound (Algorithm 1). We refer to the resulting approach as *PAC-BUS* since it combines PAC-Bayes and Uniform Stability to derive generalization guarantees for meta-learning. Third, we demonstrate our approach on two meta-learning problems: (i) a toy non-convex classification problem on the unit-ball (Section 5.1), and (ii) the *Mini-Wiki* benchmark introduced in [34] (Section 5.2). Even in these relatively small-scale settings, we demonstrate that recently-developed generalization frameworks for meta-learning provide either near-vacuous or loose bounds, while PAC-BUS provides significantly stronger bounds. Fourth, we demonstrate our approach in larger-scale settings where it remains challenging to obtain non-vacuous bounds (for our approach as well as others). Here, we propose a practical regularization scheme which re-weights the terms in the rigorously-derived PAC-BUS upper bound (*PAC-BUS(H)*; Algorithm 3 in the appendix). Recent work [7] introduces a challenging variant of the *Omniglot* benchmark [35] which highlights and tackles challenges with *memorization* in meta-learning. We show that *PAC-BUS(H)* is able to prevent memorization on this variant (Section 5.3).

## 2 Problem formulation

**Samples, tasks, and datasets.** Formally, consider the setting where we have an unknown meta distribution  $P_t$  over tasks (roughly, “tasks” correspond to different, but potentially related, learning problems). A sampled task  $t \sim P_t$  induces an (unknown) distribution  $P_{z|t}$  over sample space  $\mathcal{Z}$ . We assume that all sampling is independent and identically distributed (i.i.d.). Note that the sample space  $\mathcal{Z}$  is shared between tasks, but the distribution  $P_{z|t}$  may be different. We then sample within-task samples  $z \sim P_{z|t}$  and within-task datasets  $S = \{z_1, z_2, \dots, z_m\} \sim P_{z|t}^m$ . We assume that each sample  $z$  has a single corresponding label  $o(z)$ , where the function  $o$  is an oracle which outputs the correct label of  $z$ . At meta-training time, we assume access to  $l$  datasets, which we call  $\mathbf{S} = \{S_1, S_2, \dots, S_l\}$ . Each dataset  $S_i$  in  $\mathbf{S}$  is drawn by first selecting a task  $t_i$  from  $P_t$ , and then drawing  $S_i \sim P_{z|t_i}^m$ .

**Hypotheses and losses.** Let  $h$  denote a hypothesis and  $L(h, z)$  be the loss incurred by hypothesis  $h$  on sample  $z$ . The loss is computed by comparing  $h(z)$  with the true label  $o(z)$ . For simplicity, we assume that there is no noise on the labels; we can thus assume that all loss functions have access to

the label oracle function  $o$  and thus the loss depends only on hypothesis  $h$  and sample  $z$ . We note that this assumption is not required for our analysis and is made for the ease of exposition. Overloading the notation, we let  $L(h, P_{z|t}) := \mathbb{E}_{z \sim P_{z|t}} L(h, z)$  and  $\widehat{L}(h, S) := \frac{1}{|S|} \sum_{i=1}^{|S|} L(h, z_i)$ .

**Meta-learning.** As with model-agnostic meta-learning (MAML) [25], we let meta parameters  $\theta \in \mathbb{R}^{n_\theta}$  correspond to an initialization of the base learner’s hypothesis. Let  $h_\theta$  be the  $\theta$ -initialized hypothesis. Generally, the initialization  $\theta$  is learned from the multiple datasets we have access to at meta-training time. In this work, we will learn a *distribution*  $P_\theta$  over initializations so that we can use bounds from the PAC-Bayes framework. At test time, a new task  $t \sim P_t$  is sampled and we are provided with a new dataset  $\tilde{S} \sim P_{z|t}^m$ . The base learner uses an algorithm  $A$  (e.g., gradient descent), the dataset  $S$ , and the initialization  $\theta \sim P_\theta$  in order to fine-tune the hypothesis and perform well on future samples drawn from  $P_{z|t}$ . We denote the base learner’s updated hypothesis by  $h_{A(\theta, S)}$ . More formally, our goal is to learn a distribution  $P_\theta$  with the following objective:

$$\min_{P_\theta} \mathcal{L}(P_\theta, P_t) := \min_{P_\theta} \mathbb{E}_{t \sim P_t} \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, P_{z|t}). \quad (1)$$

We are particularly interested in the few-shot learning case, where the number of samples which the base learner can use to adapt is small. A common technique to improve test performance in the few-shot learning case is to allow for validation data at meta-training time. Thus, in addition to a generalization guarantee on meta-learning without validation data, we will derive a bound when allowing for the use of validation data  $S_{\text{va}} \sim P_{z|t}^n$  during meta-training.

### 3 Related work

**Meta-learning.** Meta-learning is a well-studied technique for exploiting similarities between learning tasks [62, 70]. Often used to reduce the need for large amounts of training data, a number of approaches for meta-learning have been explored over decades [11, 13, 16, 31, 72, 61, 55, 32]. Recently, methods based on model-agnostic meta-learning (MAML) [25] have demonstrated strong performance across different application domains and benchmarks such as *Omniglot* [35] and *Mini-ImageNet* [74]. These methods operate by optimizing a set of initial parameters that can be quickly fine-tuned via gradient descent on a new task. The approaches mentioned above typically do not provide any generalization guarantees, and none of them compute explicit numerical bounds on generalization performance. Our approach has the structure of gradient-based meta-learning while providing guarantees on generalization.

**Generalization bounds for supervised learning.** Multiple frameworks have been developed for providing generalization guarantees in the classical supervised learning setting. Early breakthroughs include Vapnik-Chervonenkis (VC) theory [71, 6], Rademacher complexity [65], and the minimum description length principle [12, 56, 36]. More recent frameworks include algorithmic stability bounds [14, 19, 29, 57, 1] and PAC-Bayes theory [67, 43, 64]. The connection between stability and learnability has been established in [66, 73, 29], and suggests that algorithmic stability bounds are a strong choice of generalization framework. PAC-Bayes theory in particular provides some of the tightest known generalization bounds for classical supervised learning approaches such as support vector machines [64, 38, 26, 57, 4, 48]. Since its development, researchers have continued to tighten [38, 42, 54] and generalize the framework [17, 18, 59]. Exciting recent results [24, 45, 46, 10, 8, 54] have demonstrated the promise of PAC-Bayes to provide strong generalization bounds for neural networks on supervised learning problems (see [33] for a recent review of generalization bounds for neural networks). It is also possible to combine frameworks such as PAC-Bayes and uniform stability to derive bounds for supervised learning [39]. We will use these two frameworks to bound generalization in the two levels of meta-learning. In contrast to the standard supervised learning setting, generalization bounds for meta-learning are less common and remain loose.

**Generalization bounds for meta-learning.** As described in Section 1, meta-learning bounds must account for two “levels” of generalization (base level and meta level). The approach presented in [41] utilizes algorithmic stability bounds at both levels. However, this requires both meta and base learners to be uniformly stable. This is a strong requirement that is challenging to ensure at the meta level. Another recent method, known as follow-the-meta-regularized-leader (FMRL) [34], provides guarantees for a regularized meta-learning version of the follow-the-leader (FTL) method for online learning, see e.g. [30]. The generalization bounds provided are derived from the application of online-to-batch techniques [3, 22]. A regret bound for meta-learning using an aggregation technique at the meta-level and an algorithm with a uniform generalization bound at the base level is provided

in [3]. The techniques mentioned do not present an algorithm which makes use of validation data (in contrast to our approach). Using validation data (i.e., held-out data) is a common technique for improving performance in meta-learning and is particularly important for the few-shot learning case.

Another method for deriving a generalization bound on meta-learning is to use PAC-Bayes bounds at both the base and meta levels [52, 53]. In [5], generalization bounds based on such a framework are provided along with practical optimization techniques. However, the method requires one to maintain distributions over distributions of initializations, which can result in large computation times during training and makes it extremely challenging to numerically compute the bound. Moreover, the approach also does not allow one to incorporate validation data to improve the bound. Recent work has made progress on some of these challenges. In [60], the computational efficiency of training is improved but the challenges associated with numerically computing the generalization bound or incorporating validation data are not addressed. State-of-the-art work tightens the two-level PAC-Bayes guarantee, addresses computation times for training and evaluation of the bound, and allows for validation data [77]. However, all of the two-level PAC-Bayes bounds require a separate PAC-Bayes bound for each task, and thus a potentially loose union bound.

We present a framework which, to our knowledge, is the first to combine algorithmic stability and PAC-Bayes bounds (at the base- and meta- levels respectively) in order to derive a meta-learning algorithm with associated generalization guarantees. As outlined in Section I, we believe that the algorithmic stability and PAC-Bayes frameworks are particularly well-suited to tackling the specific challenges of generalization at the different levels. We also highlight that *none* of the approaches mentioned above report numerical values for generalization bounds, even for relatively simple problems. Here, we empirically demonstrate that prior approaches tend to provide either near-vacuous or loose bounds even in relatively small-scale settings while our proposed method provides significantly stronger bounds.

## 4 Generalization bound on meta-learning

We use two different frameworks for the two levels of generalization required in a meta-learning bound. We utilize the PAC-Bayes framework to bound the expected training loss on future tasks, and uniform stability bounds to argue that if we have a low training loss when using a uniformly stable algorithm, then we achieve a low test loss. The following section will introduce these frameworks independently. We then present the overall meta-learning bound and associated algorithm to find a distribution over initialization parameters (i.e., meta parameters) that minimizes the upper bound.

### 4.1 Preliminaries: two generalization frameworks

#### 4.1.1 Uniform stability

Let  $S = \{z_1, z_2, \dots, z_m\} \in \mathcal{Z}^m$  be a set of  $m$  elements of  $\mathcal{Z}$ . Let  $S^i = \{z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m\}$  be identical to dataset  $S$  except that the  $i^{th}$  sample  $z_i$  is replaced by some  $z'_i \in \mathcal{Z}$ . Note that our analysis can be extended to allow for losses bounded by some finite  $M$ , but we work with losses bounded within  $[0, 1]$  for the sake of simplicity. With these precursors, we define an analogous notion of *uniform stability* to [29, Definition 2.1] for deterministic algorithms  $A$  and distributions  $P_\theta$  over initializations<sup>1</sup>

**Definition 1** (Uniform Stability) *A deterministic algorithm  $A$  has  $\beta > 0$  uniform stability with respect to loss  $L$  if  $\forall z \in \mathcal{Z}$ ,  $\forall S \in \mathcal{Z}^m$ ,  $\forall i \in \{1, \dots, m\}$ , and all distributions  $P_\theta$  over initializations, the following holds:*

$$\mathbb{E}_{\theta \sim P_\theta} |L(h_{A(\theta, S)}, z) - L(h_{A(\theta, S^i)}, z)| \leq \beta. \quad (2)$$

We define  $\beta_{\text{US}}$  as the minimal such  $\beta$ .

In this work, we will bound  $\beta_{\text{US}}$  as a function of the algorithm, form of the loss, and number of samples that the algorithm uses (See Appendix A.4 for further details on the bounds on  $\beta_{\text{US}}$  for our setup). We then establish a relationship between uniform stability and generalization in expectation. The following is adapted from [29, Theorem 2.2] for the notion of uniform stability presented in Definition I.

<sup>1</sup>We use deterministic algorithms to avoid excess computation when calculating the provided meta-learning upper bounds. See Appendix A.4 for further details.

**Theorem 1** (Algorithmic Stability Generalization in Expectation) *Fix a task  $t \sim P_t$ . The following inequality holds for hypothesis  $h_{A(\theta, S)}$  learned using  $\beta_{\text{US}}$  uniformly stable algorithm  $A$  with respect to loss  $L$ :*

$$\mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, P_{z|t}) \leq \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S)}, S) + \beta_{\text{US}}. \quad (3)$$

*Proof.* The proof is similar to the one presented for [29, Theorem 2.2] and is presented in Appendix A.1.  $\square$

#### 4.1.2 PAC-Bayes theory

For the meta-level bound, we make use of the PAC-Bayes generalization bound introduced in [43]. Note that other PAC-Bayes bounds such as the quadratic variant [58] and PAC-Bayes- $\lambda$  variant [69] may be used and substituted in the following analysis. We first present a general version of the PAC-Bayes bound and then specialize it to our meta-learning setting in Section 4.2. Let  $f(\theta, s)$  be an arbitrary loss function which only depends on parameters  $\theta$  and the sample  $s$  which has been drawn from an arbitrary distribution  $P_s$ . The following bound is a tightened version of the bound presented in [43] for when  $l \geq 8$ .

**Theorem 2** (PAC-Bayes Generalization Bound [40]) *For any data-independent prior distribution  $P_{\theta,0}$  over  $\theta$ , some loss function  $f$  where  $0 \leq f(\theta, s) \leq 1, \forall s, \forall \theta, l \geq 8$ , and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over a sampling of  $\{s_1, s_2, \dots, s_l\} \sim P_s^l$ , the following holds simultaneously for all distributions  $P_\theta$  over  $\theta$ :*

$$\mathbb{E}_{s \sim P_s} \mathbb{E}_{\theta \sim P_\theta} f(\theta, s) \leq \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\theta \sim P_\theta} f(\theta, s_i) + R_{\text{PAC-B}}(P_\theta, P_{\theta,0}, \delta, l), \quad (4)$$

where the PAC-Bayes “regularizer” term is defined as follows

$$R_{\text{PAC-B}}(P_\theta, P_{\theta,0}, \delta, l) := \sqrt{\frac{D_{\text{KL}}(P_\theta \| P_{\theta,0}) + \ln \frac{2\sqrt{l}}{\delta}}{2l}}, \quad (5)$$

and  $D_{\text{KL}}$  is the Kullback-Leibler (KL) divergence.

#### 4.2 Meta-learning bound

In order to obtain a generalization guarantee for meta-learning, we utilize the two frameworks above. We first specialize the PAC-Bayes bound in Theorem 2 to bound the expected training loss on future tasks. We then utilize Theorem 1 to demonstrate that if we have a low expected training loss when using a uniformly stable algorithm, then we achieve a low expected test loss. These two steps allow us to combine the generalization frameworks above to derive an upper bound on (1) which can be computed with known quantities. With the following assumption, the resulting generalization bound is presented in Theorem 3.

**Assumption 1** (Bounded loss.) *The loss function  $L$  is bounded:  $0 \leq L(h, z) \leq 1$  for any  $h$  in the hypothesis space for the given problem, and any  $z$  in the sample space.*

**Theorem 3** (Meta-Learning Generalization Guarantee) *For hypotheses  $h_{A(\theta, S)}$  learned with  $\beta_{\text{US}}$  uniformly stable algorithm  $A$ , data-independent prior  $P_{\theta,0}$  over initializations  $\theta$ , loss  $L$  which satisfies Assumption 1,  $l \geq 8$ , and  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over a sampling of the meta-training dataset  $S \sim P_S^l$ , the following holds simultaneously for all distributions  $P_\theta$  over  $\theta$ :*

$$\mathcal{L}(P_\theta, P_t) \leq \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S_i)}, S_i) + R_{\text{PAC-B}}(P_\theta, P_{\theta,0}, \delta, l) + \beta_{\text{US}}. \quad (6)$$

*Proof.* The proof can be split into three steps:

**Step 1.**

Let  $P_s$  in Theorem 2 be the marginal distribution  $P_S$  over datasets of size  $m$  (see Appendix A.2 for details) and note that sampling  $S \sim P_S$  is equivalent to first sampling  $t \sim P_t$  and then sampling  $S \sim P_{z|t}^m$ . Additionally let  $f(\theta, S) := \widehat{L}(h_{A(\theta, S)}, S)$  where  $A(\theta, S)$  is any deterministic algorithm.

Plugging in these definitions into Inequality (4) results in the following inequality which holds under the same assumptions as Theorem 2, and with probability at least  $1 - \delta$  over the sampling of  $\mathbf{S} \sim P_S^l$ :

$$\begin{aligned} \mathbb{E}_{S \sim P_S} \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S)}, S) &= \mathbb{E}_{t \sim P_t} \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S)}, S) \\ &\leq \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S_i)}, S_i) + R_{\text{PAC-B}}(P_\theta, P_{\theta,0}, \delta, l). \end{aligned} \quad (7)$$

**Step 2.**

Now assume that algorithm  $A$  is  $\beta_{\text{US}}$  uniformly stable. For a fixed task  $t \sim P_t$  we have the following by Theorem 1:

$$\mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, P_{z|t}) \leq \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S)}, S) + \beta_{\text{US}}.$$

Take the expectation over  $t \sim P_t$ . We then have:

$$\mathbb{E}_{t \sim P_t} \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, P_{z|t}) \leq \mathbb{E}_{t \sim P_t} \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S)}, S) + \beta_{\text{US}}, \quad (8)$$

since  $\mathbb{E}_{t \sim P_t} \beta_{\text{US}} = \beta_{\text{US}}$ . This establishes a bound on the true expected loss for a new task after running algorithm  $A$  on a training dataset corresponding to the new task.

**Step 3.**

Note that (7) provides an upper bound on the first term of the RHS of (8) when algorithm  $A$  is  $\beta_{\text{US}}$  uniformly stable. Thus we have the following by plugging (7) in the RHS of (8):

Under the same assumptions as both Theorems 1 and 2, and with probability at least  $1 - \delta$  over the sampling of  $\mathbf{S} \sim P_S^l$ :

$$\mathbb{E}_{t \sim P_t} \mathbb{E}_{S \sim P_{z|t}^m} \mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, P_{z|t}) \leq \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S_i)}, S_i) + R_{\text{PAC-B}}(P_\theta, P_{\theta,0}, \delta, l) + \beta_{\text{US}},$$

completing the proof.  $\square$

Theorem 3 is presented for any distributions  $P_\theta$  and  $P_{\theta,0}$  over initializations. However, in practice we will use multivariate Gaussian distributions for both. The specialization of Theorem 3 to Gaussian distributions is provided in Appendix A.3.1. Next, we allow for validation data  $S_{\text{va}} \sim P_{z|t}^n$  at meta-training time so that the bound is more suited to the few-shot learning case. We compute the upper bound using the evaluation data  $S_{\text{ev}} = \{S, S_{\text{va}}\}$  sampled from the marginal distribution  $P_{S_{\text{ev}}}$  over datasets of size  $m+n$ . However, we still only require  $m$  samples at meta-test time; see Appendix A.3.2 for the derivation. Note that the training data  $S$  is often excluded from the data used to update the meta-learner. However, this is necessary for our approach to obtain a guarantee on few-shot learning performance. The result is a guarantee with high probability over a sampling of  $\mathbf{S}_{\text{ev}} \sim P_{S_{\text{ev}}}^l$ :

$$\mathcal{L}(P_\theta, P_t) \leq \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\theta \sim P_\theta} \widehat{L}(h_{A(\theta, S_i)}, S_{\text{ev},i}) + R_{\text{PAC-B}}(P_\theta, P_{\theta,0}, \delta, l) + \frac{m\beta_{\text{US}}}{m+n}. \quad (9)$$

### 4.3 PAC-BUS algorithm

Recall that we aim to find a distribution  $P_\theta$  over initializations that minimizes  $\mathcal{L}(P_\theta, P_t)$  as stated in Equation (1). We cannot minimize  $\mathcal{L}(P_\theta, P_t)$  directly due to the expectations taken over unknown distributions  $P_t$  and  $P_{z|t}$  for sampled task  $t$ , but we may indirectly minimize it by minimizing the upper bounds in Inequalities (6) or (9).

Computing the upper bound requires evaluating an expectation taken over  $\theta \sim P_\theta$ . In general, this is intractable. However, we aim to minimize this upper bound to provide the tightest guarantee possible. Similar to the method in [24], we use an unbiased estimator of  $\mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, \cdot)$ . Let  $P_\theta$  be a multivariate Gaussian distribution over initializations  $\theta$  with mean  $\mu$  and covariance  $\text{diag}(s)$ ; thus  $P_\theta = \mathcal{N}(\mu, \text{diag}(s))$  and  $P_{\theta,0} = \mathcal{N}(\mu_0, \text{diag}(s_0))$ . Further, let  $\psi := (\mu, \log(s))$ , and use the shorthand  $\mathcal{N}_{\psi_0}$  for the prior and  $\mathcal{N}_\psi$  for the posterior distribution over initializations. We use the following estimator of  $\mathbb{E}_{\theta \sim P_\theta} L(h_{A(\theta, S)}, \cdot)$ :

$$L(h_{A(\theta, S)}, \cdot), \quad \theta \sim \mathcal{N}_\psi. \quad (10)$$

---

**Algorithm 1** PAC-BUS: meta-learning via PAC-Bayes and Uniform Stability

---

**Input:** Fixed prior distribution  $\mathcal{N}_{\psi_0}$  over initializations  
**Input:**  $\beta_{\text{US}}$  uniformly stable Algorithm  $A$   
**Input:** Meta-training dataset  $\mathbf{S}$ , learning rate  $\gamma$   
**Initialize:**  $\psi \leftarrow \psi_0$   
**Output:** Optimized  $\psi^*$

$$B(\psi, \theta'_1, \theta'_2, \dots, \theta'_l) := \frac{1}{l} \sum_{i=1}^l \hat{L}(h_{\theta'_i}, S_i) + R_{\text{PAC-B}}(\mathcal{N}_{\psi}, \mathcal{N}_{\psi_0}, \delta, l) + \beta_{\text{US}}$$

**while** not converged **do**

- Sample  $\theta \sim \mathcal{N}_{\psi}$
- for**  $i = 1$  **to**  $l$  **do**
- $\theta'_i \leftarrow A(\theta, S_i)$
- end for**
- $\psi \leftarrow \psi - \gamma \nabla_{\psi} B(\psi, \theta'_1, \theta'_2, \dots, \theta'_l)$

**end while**

---

We present the resulting training technique in Algorithm 1. This algorithm can be used to learn a distribution over initializations that minimizes the upper bound presented in Theorem 3 and its specializations. This is presented for the case when  $A$  is  $\beta_{\text{US}}$  uniformly stable for some  $\beta_{\text{US}}$ . For gradient-based algorithms, the learning rate  $\alpha$  often appears directly in the bound for  $\beta_{\text{US}}$  [29]. Thus it is potentially beneficial to update  $\alpha$  as well. We present Algorithm 1 without learning the learning rate. To meta-learn the learning rate, we can augment  $\psi_0$  to include a parameterization of a prior distribution over learning rates and update it using the same gradient step presented in 1 for  $\psi$ .

Determining the gradient of  $B(\psi, \theta'_1, \theta'_2, \dots, \theta'_l)$  with respect to  $\psi$  requires computing the Hessian of the loss function if algorithm  $A(\theta, S)$  uses a gradient update to compute  $\theta'_i$ . First order approximations often perform similarly to the second-order meta-learning techniques [27, 25, 47], and can be used to speed up the training. Additionally, Algorithm 1 can be modified to use mini-batches of tasks instead of all tasks in the meta update to improve training times; we present an algorithm which uses mini-batches of tasks in Appendix A.5.1.

In practice, we are interested in algorithms such as stochastic gradient descent (SGD) and gradient descent (GD) for the base learner. We can obtain bounds on the uniform stability constant  $\beta_{\text{US}}$  when using gradient methods with the results from [29]. See Appendix A.4 for details on the  $\beta_{\text{US}}$  bounds we use in this work. With a bound on  $\beta_{\text{US}}$ , we can calculate all the terms in  $B(\psi, \theta'_1, \theta'_2, \dots, \theta'_l)$  and use Algorithm 1 to minimize the meta-learning upper bound. When evaluating the upper bound, we use the sample convergence bound [37, 24] to upper bound the expectation taken over  $\theta \sim P_{\theta}$ . See Appendix A.6 for details.

## 5 Examples

We demonstrate our approach on three examples below. All examples we provide are few-shot meta-learning problems. To adapt at the base level,  $m$  examples from each class are given for an “ $m$ -shot” learning problem. If applicable,  $n$  samples can be given as validation data for each task during the meta-training step. In the first two examples, our primary goal is to demonstrate the tightness of our generalization bounds compared to other meta-learning bounds. We also present empirical test performance on held-out data; however, we emphasize that the focus of our work is to obtain improved generalization guarantees (and not necessarily to improve empirical test performance). In the third example, we present an algorithm that is motivated by our theoretical framework and demonstrate its ability to improve empirical performance on a challenging task. All the code required to run the following examples is available at <https://github.com/irom-lab/PAC-BUS>.

### 5.1 Example: classification on the unit ball

We evaluate the tightness of the generalization bound in Equation (9) on a toy two-class classification problem where the sample space  $\mathcal{Z}$  is the unit ball  $B^2(0, 1)$  in two dimensions with radius 1 and centered at the origin. Data points for each task are sampled from  $P_{z|t}$ , where a task corresponds to a particular concept which labels the data as (+) if within  $B^2(c_t, r_t)$  and (−) otherwise. Center  $c_t$  is sampled uniformly from the  $y \geq 0$  semi-ball  $B^2_{y \geq 0}(0, 0.4)$  of radius 0.4. The radius  $r_t$  is then sampled uniformly from  $[0.1, 1 - \|c_t\|]$ . Notably, the decision boundary between classes is nonlinear. Thus, generalization bounds which rely on convex losses (such as [34]) will have difficulty with

Table 1: We present the generalization bounds (for  $\delta = 0.01$ ) provided by each method if applicable, and use the sample convergence bound [37] for MR-MAML, and PAC-BUS, but not MLAP-M.<sup>2</sup> Note that for these methods, we specifically minimize their respective meta-learning bounds. We also report the meta-test loss (the softmax activated cross-entropy loss –  $\text{CEL}_s$ ) for all methods. We present the mean and standard deviation after 5 trials. We highlight that our approach provides the strongest generalization guarantee.

Classification on Ball	MAML [25]	MLAP-M [5]	MR-MAML [77]	PAC-BUS (ours)
Bound ↓	None	$1.0538 \pm 0.0012^2$	$0.3422 \pm 0.0006$	<b><math>0.2213 \pm 0.0012</math></b>
Test Loss ↓	$0.1701 \pm 0.0070$	$0.1645 \pm 0.0045$	<b><math>0.1584 \pm 0.0012</math></b>	$0.1657 \pm 0.0014$

providing guarantees for networks that perform well. We choose the softmax-activated cross-entropy loss,  $\text{CEL}_s$ , as the loss function. Before running Algorithm 1, we address a few technical challenges that arise from Assumption 1 as well as computing  $c_L$  and  $c_S$ . We address these in Appendix A.5.

We then apply Algorithm 1 using the few-shot learning bound in Inequality (9). We present the guarantee on the meta-test loss associated with each training method in Table 1. In addition, we present the average meta-test loss after training with 10 samples. We compare our bounds and empirical performance with the meta-learning by adjusting priors (MLAP) technique [5] and the meta-regularized MAML (MR-MAML) technique [77]. All methods are given held-out data to learn a prior before minimizing their respective upper bounds (see Appendix A.11.1 for further details on the prior training step). Additionally, since all bounds require the loss to be within  $[0, 1]$ , networks  $N$  are constrained such that the Frobenius norm of the output is bounded by  $r$ , i.e.,  $\|N(z)\|_F \leq r$ . We compare the aforementioned methods' meta-test loss to MAML with weights constrained in the same manner (note that MAML does not provide a guarantee). Upper bounds which use the PAC-Bayes framework are computed with many evaluations from the posterior distribution. This allows us to apply the sample convergence bound [37] (as in Equation (35) for our bound) unless otherwise noted.

We find that PAC-BUS provides a significantly stronger guarantee compared with the other methods. Note that the guarantee provided by MLAP-M [5] is vacuous because the meta-test loss is bounded between 0 and 1, while the guarantee is above 1.

## 5.2 Example: Mini-Wiki

Next, we present results on the *Mini-Wiki* benchmark introduced in [34]. This is derived from the Wiki3029 dataset presented in [9]. The dataset is comprised of 4-class,  $m$ -shot learning tasks with sample space  $\mathcal{Z} = \{z \in \mathbb{R}^d \mid \|z\|_2 = 1\}$ . Sentences from various Wikipedia articles are passed through the continuous-bag-of-words GloVe embedding [51] into dimension  $d = 50$  to generate samples. For this learning task, we use a  $k$ -class version of  $\text{CEL}_s$  and logistic regression. Since this example is convex, we can use GD and bound  $\beta_{\text{US}}$  with Theorem 4 in the appendix [29]. We keep the loss bounded by constraining the network  $\|N(z)\|_F \leq r$  and scale the loss as in the previous example. The tightness of the bounds on  $c_L$  and  $c_S$  affected the upper bound in Inequality (9) more than in the previous example, so we bound them as tightly as possible. See Appendix A.9 for the calculations.

We apply Algorithm 1 using the bound which allows for validation data, Inequality (9), to learn on 4-way *Mini-Wiki*  $m = \{1, 3, 5\}$ -shot. The results are presented in Table 2. We compare our results with the FMRL variant which provides a guarantee [34], follow-the-last-iterate (FLI)-Batch, and with MR-MAML [77]. FLI-Batch does not require bounded losses explicitly, but requires that the parameters of the network lie within a ball of radius  $r$ . For the logistic regression used in the example, this is equivalent to  $\|N(z)\|_F \leq r$ . Thus, we scale the loss and use the same  $r$  for each method to provide a fair comparison. We also show the results of training with MAML constrained in the same way for reference. Each method is given the same amount of held-out data for training a prior (see Appendix A.11.2 for further details on training the prior).

<sup>2</sup>Due to high computation times associated with estimating the MLAP upper bound, this value is not computed with the sample convergence bound as the other upper bounds are. Thus, the value presented does not carry a guarantee, but would be similar if computed with the sample convergence bound. The value is shown to give a qualitative sense of the guarantee.

Table 2: We compare the generalization bounds (for  $\delta = 0.01$ ) provided by each method where applicable and use the sample convergence bound for MR-MAML and PAC-BUS. Since we specifically minimize these methods’ upper bounds, we can fairly compare the relative tightness of each bound. We also report the meta-test loss ( $\text{CEL}_s$ ) for each method for exposition. We report the mean and standard deviation after 5 trials. We highlight that our approach provides the strongest guarantee.

4-Way Mini-Wiki	1-shot $\downarrow$	3-shot $\downarrow$	5-shot $\downarrow$
FLI-Batch Bound [34]	$0.6638 \pm 0.0011$	$0.6366 \pm 0.0006$	$0.6343 \pm 0.0014$
MR-MAML Bound [77]	$0.7400 \pm 0.0003$	$0.7312 \pm 0.0003$	$0.7283 \pm 0.0005$
PAC-BUS Bound (ours)	<b><math>0.4999 \pm 0.0003</math></b>	<b><math>0.5058 \pm 0.0002</math></b>	<b><math>0.5101 \pm 0.0002</math></b>
MAML [25]	<b><math>0.3916 \pm 0.0009</math></b>	<b><math>0.3868 \pm 0.0005</math></b>	<b><math>0.3883 \pm 0.0005</math></b>
FLI-Batch [34]	$0.4091 \pm 0.0008$	$0.4078 \pm 0.0005$	$0.4097 \pm 0.0012$
MR-MAML [77]	$0.3922 \pm 0.0009$	$0.3869 \pm 0.0003$	$0.3884 \pm 0.0005$
PAC-BUS (ours)	$0.3922 \pm 0.0009$	$0.3878 \pm 0.0003$	$0.3895 \pm 0.0005$

As in the previous example, PAC-BUS provides a significantly tighter guarantee than the other methods (Table 2). We see similar empirical meta-test loss for MAML [25], MR-MAML [77], and PAC-BUS with slightly higher loss for FLI-Batch [34]. In addition, we computed the meta-test accuracy as the percentage of correctly classified sentences. See Table 4 in Section A.11.2 for these results along with other experimental details.

### 5.3 Example: memorizable Omniglot

We have demonstrated the ability of our approach to provide strong generalization guarantees for meta-learning in the settings above. We now consider a more complex setting where we are unable to obtain strong guarantees. In this example, we employ a learning heuristic based on the PAC-BUS upper bound,  $PAC\text{-}BUS(H)$ ; see Appendix A.5.2 for the details and the Algorithm. We relax Assumption 1 and no longer constrain the network as in previous sections. Instead, we maintain and update estimates of the Lipschitz and smoothness constants of the network, using [68], and incorporate them into the uniform stability regularizer term,  $\beta_{\text{US}}$ . We then scale each regularizer term (i.e.,  $R_{\text{PAC-BUS}}(P_\theta, P_{\theta,0}, \delta, l)$  and  $\beta_{\text{US}}$ ) by hyper-parameters  $\lambda_1$  and  $\lambda_2$  respectively. Analogous to the technique described in [77], we aim to incorporate the form of the theoretically-derived regularizer into the loss, without requiring it to be as restrictive during learning. The result is a regularizer that punishes large deviation from the prior  $P_{\theta,0}$  and too much adaptation at the base-learning level.

We test our method on *Omniglot* [35] for 20-way,  $m = \{1, 5\}$ -shot classification in the non-mutually exclusive (NME) case [77]. In [77], the problem of memorization in meta-learning is explored and demonstrated with non-mutually exclusive learning problems. *NME Omniglot* corresponds to randomization of class labels for a task at test time only. This worsens the performance of any network that memorized class labels; see [77] for more details.<sup>3</sup> We compare our method to an analogous heuristic presented in [77], which also has a  $D_{\text{KL}}(P_\theta \| P_{\theta,0})$  term in the loss. Thus, this heuristic (referred to as MR-MAML(W) [77]) regularizes the change in weights of the network. Additionally, we compare to the heuristic described in [34] (FLI-Online) which performs better in practice than the FLI-Batch method. We do not provide data for training a prior in this case since we do not aim to compute a bound in this example. We use standard MAML as a reference. See Table 3 for the results.

We see that MAML [25] and FLI-Online [34] do not prevent memorization on *NME Omniglot* [77]. This is especially apparent in the 1-shot learning case, where their performance suffers significantly due to this memorization. Both MR-MAML(W) [77] and PAC-BUS(H) prevent memorization, with PAC-BUS(H) outperforming MR-MAML(W). Note that PAC-BUS(H) outperforms MR-MAML(W) by a wider margin in the 1-shot case as compared with the 5-shot case. We believe this is due to the effectiveness of the uniform stability regularizer at the base level. MR-MAML(W) suffers more in the 1-shot case because over-adaptation is more likely with fewer within-task examples. □

## 6 Conclusion and discussion

We presented a novel generalization bound for gradient-based meta-learning: PAC-BUS. We use different generalization frameworks for tackling the distinct challenges of generalization at the two

<sup>3</sup>We use a slightly different task setup as the one in [77]; see Appendix A.11.3 for the details of our setup.

Table 3: We present the meta-test accuracy as a percentage on non-mutually-exclusive *Omniglot* [77]. In contrast to the previous examples, here we aim to achieve the best empirical performance for each method. In particular, this task compares each methods’ ability to prevent memorization. We report the mean and standard deviation after 5 trials.

20-WAY <i>Omniglot</i>	NME 1-SHOT $\uparrow$	NME 5-SHOT $\uparrow$
MAML [25]	$23.4 \pm 2.2$	$75.1 \pm 4.8$
FLI-ONLINE [34]	$22.4 \pm 0.5$	$39.1 \pm 0.5$
MR-MAML(W) [77]	$84.2 \pm 2.2$	$94.3 \pm 0.3$
PAC-BUS(H) (OURS)	<b><math>87.9 \pm 0.5</math></b>	<b><math>95.0 \pm 0.9</math></b>

levels of meta-learning. In particular, we employ uniform stability bounds and PAC-Bayes bounds at the base- and meta-learning levels respectively. On a toy non-convex problem and the *Mini-Wiki* meta-learning task [34], we provide significantly tighter generalization guarantees as compared to state-of-the-art meta-learning bounds while maintaining comparable empirical performance. To our knowledge, this work presents the first numerically-evaluated generalization guarantees associated with a proposed meta-learning bound. On memorizable *Omniglot* [35, 77], we show that a heuristic based on the PAC-BUS bound prevents memorization of class labels in contrast to MAML [25], and better performance than meta-regularized MAML [77]. We believe our framework is well suited to the few-shot learning problems for which we present empirical results, but our framework is potentially applicable to a broad range of different settings (e.g., reinforcement learning).

We note a few challenges with our method as motivation for future work. Our bound is vacuous on larger scale learning problems such as *Omniglot*. This is partially caused by a larger KL-divergence term in the PAC-Bayes bound when using deep convolutional networks (due to the increased dimensionality of the weight vector). In addition, we do not have a theoretical analysis on the convergence properties of the algorithms presented, so we must experimentally determine the number of samples required for tight bounds. In the results of Section 5.1 and 5.2, despite an improved bound over other methods, our method does not necessarily improve empirical test performance. We emphasize that our focus in this work was on deriving stronger generalization guarantees rather than improving empirical performance. However, obtaining approaches that provide both stronger guarantees and empirical performance is an important direction for future work.

Future work can also explore ways in which to incorporate tighter PAC-Bayes bounds or those with less restrictive assumptions. One interesting avenue is to extend PAC-BUS by using a PAC-Bayes bound for unbounded loss functions for the meta-generalization step (e.g. as presented in [28]). Another promising direction is to incorporate regularization on the weights of the network directly (e.g.,  $L_2$  regularization or gradient clipping) to create networks with smaller Lipschitz and smoothness constants. Additionally, it would be interesting to explore learning of the base-learner’s algorithm while maintaining uniform stability. For example, one could parameterize a set of uniformly stable algorithms and learn a posterior distribution over the parameters.

**Broader impact.** The approach we present in this work aims to strengthen performance guarantees for gradient-based meta-learning. We believe that strong generalization guarantees in meta-learning, especially in the few-shot learning case, could lead to broader application of machine learning in real-world applications. One such example is for medical diagnosis, where abundant training data for certain diseases may be difficult to obtain. Another example on which poor performance is not an option is any safety critical robotic system, such as ones which involve human interaction.

Meta-learning methods typically require a lot of data and training time, and ours is not an exception. In our case, it took multiple weeks of computation time on Amazon Web Services (AWS) instances to train and compute all networks and results we present in this paper. This creates challenges with accessibility and energy usage.

## Acknowledgments

The authors are grateful to the anonymous reviewers for their valuable feedback and suggestions, and to Thomas Griffiths for helpful feedback on this work. The authors were supported by the Office of Naval Research [N00014-21-1-2803, N00014-18-1-2873], the NSF CAREER award [2044149], the Google Faculty Research Award, and the Amazon Research Award.

## References

- [1] Karim Abou-Moustafa and Csaba Szepesvari. An Exponential Efron-Stein Inequality for  $L_q$  Stable Learning Rules. *Proceedings of the 30th International Conference on Algorithmic Learning Theory*, 2019.
- [2] Akshay Agrawal, Robin Verschueren, Steven Diamond, and Stephen Boyd. A rewriting system for convex optimization problems. *Journal of Control and Decision*, 5(1):42–60, 2018.
- [3] Pierre Alquier, The Tien Mai, and Massimiliano Pontil. Regret Bounds for Lifelong Learning. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 2017.
- [4] Amiran Ambroladze, Emilio Parrado-hernandez, and John Shawe-taylor. Tighter PAC-Bayes Bounds. *Advances in Neural Information Processing Systems 19*, 2007.
- [5] Ron Amit and Ron Meir. Meta-Learning by Adjusting Priors Based on Extended PAC-Bayes Theory. *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [6] Martin Anthony and Peter L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999.
- [7] Sebastien Arnold, Praateek Mahajan, Debajyoti Datta, Ian Bunner, and Konstantinos Saitas Zarkias. learn2learn: A library for meta-learning research. arxiv preprint. *preprint arXiv:2008.12284*, 2020.
- [8] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger Generalization Bounds for Deep Nets via a Compression Approach. *arXiv preprint arXiv:1802.05296*, 2018.
- [9] Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A Theoretical Analysis of Contrastive Unsupervised Representation Learning. *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- [10] Peter L. Bartlett, Dylan J. Foster, and Matus J Telgarsky. Spectrally-Normalized Margin Bounds for Neural Networks. In *Advances in Neural Information Processing Systems 30*, pages 6240–6249, 2017.
- [11] Samy Bengio, Yoshua Bengio, Jocelyn Cloutier, and Jan Gecsei. On the Optimization of a Synaptic Learning Rule. *Proceedings of the Conference on Optimality in Artificial and Biological Neural Networks*, pages 6–8, 1992.
- [12] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Occam’s Razor. *Information Processing Letters*, 24(6):377–380, 1987.
- [13] Leon Bottou and Vladimir Vapnik. Local Learning Algorithms. *Neural Computation*, 4: 888–900, 1992.
- [14] Olivier Bousquet and Andre Elisseeff. Stability and Generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [15] Olivier Bousquet, Yegor Klochkov, and Nikita Zhivotovskiy. Sharper Bounds for Uniformly Stable Algorithms. *Proceedings of the 33rd Conference on Learning Theory*, 2020.
- [16] Rich Caruana. Multitask Learning. *Machine Learning*, 28:41–75, 1997.
- [17] Olivier Catoni. *Statistical Learning Theory and Stochastic Optimization*. Ecole d’Et  de Probabilit s de Saint-Flour 2001. Springer, 2004.
- [18] Olivier Catoni. *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*, volume 56 of *Lecture notes - Monograph Series*. Institute of Mathematical Statistics, 2007.
- [19] Alain Celisse and Benjamin Guedj. Stability Revisited: New Generalisation Bounds for the Leave-one-Out. *arXiv preprint arXiv:1608.06412*, 2016.
- [20] Andrew Collette. *Python and HDF5*. O’Reilly, 2013.

[21] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.

[22] Giulia Denevi, Dimitris Stamos, Carlo Ciliberto, and Massimiliano Pontil. Online-Within-Online Meta-Learning. *Advances in Neural Information Processing Systems 32*, 2019.

[23] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.

[24] Gintare Karolina Dziugaite and Daniel M. Roy. Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data. *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence*, 2017.

[25] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. *Proceedings of the 34th International Conference on Machine Learning*, 2017.

[26] Pascal Germain, Alexandre Lacasse, François Laviolette, and Mario Marchand. PAC-Bayesian Learning of Linear Classifiers. In *Proceedings of the 26th International Conference on Machine Learning*, pages 353–360. ACM, 2009.

[27] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and Harnessing Adversarial Examples. *Proceedings of the 3rd International Conference on Learning Representations*, 2015.

[28] Maxime Haddouche, Benjamin Guedj, Omar Rivasplata, and John Shawe-Taylor. PAC-Bayes Unleashed: Generalisation Bounds With Unbounded Losses. *arXiv preprint arXiv:2006.07279*, 2020.

[29] Moritz Hardt, Benjamin Recht, and Yoram Singer. Train Faster, Generalize Better: Stability of Stochastic Gradient Descent. *Proceedings of the 33rd International Conference on Machine Learning*, 2016.

[30] Elad Hazan. Introduction to Online Convex Optimization. *Foundations and Trends in Optimization*, 2(3-4):157–325, 2016.

[31] Tom Heskes. Solving a Huge Number of Similar Tasks: a Combination of Multi-Task Learning and a Hierarchical Bayesian Approach. *Proceedings of the 15th International Conference on Machine Learning*, 1998.

[32] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-Learning in Neural Networks: A Survey. *arXiv preprint arXiv:2004.05439*, 2020.

[33] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic Generalization Measures and Where to Find Them. *Proceedings of the 8th International Conference on Learning Representations*, 2020.

[34] Mikhail Khodak, Maria-Florina Balcan, and Ameet Talwalkar. Provable Guarantees for Gradient-Based Meta-Learning. *Proceedings of the 36th International Conference on Machine Learning*, 2019.

[35] B. Lake, R. Salakhutdinov, J. Gross, and J. Tenenbaum. One Shot Learning of Simple Visual Concepts. *Cognitive Science*, 33, 2011.

[36] John Langford. Tutorial on Practical Prediction Theory for Classification. *Journal of Machine Learning Research*, 6(Mar):273–306, 2005.

[37] John Langford and Rich Caruana. (Not) Bounding the True Error. *Advances in Neural Information Processing Systems 14*, 2002.

[38] John Langford and John Shawe-Taylor. PAC-Bayes & margins. *Advances in Neural Information Processing Systems 15*, 2003.

[39] Ben London. A PAC-Bayesian Analysis of Randomized Learning with Application to Stochastic Gradient Descent. *Advances in Neural Information Processing Systems 30*, 2017.

[40] Andreas Maurer. A Note on the PAC Bayesian Theorem. *arXiv preprint arXiv:0411099*, 2004.

[41] Andreas Maurer. Algorithmic Stability and Meta-Learning. *Journal of Machine Learning Research*, 6:967–994, 2005.

[42] David McAllester. A PAC-Bayesian Tutorial with A Dropout Bound. *arXiv preprint arXiv:1307.2118*, 2013.

[43] David A McAllester. PAC-Bayesian Model Averaging. *Proceedings of the 12th Conference on Learning Theory*, 1999.

[44] MOSEK ApS. Mosek fusion api for python 9.0.105, 2019. URL <https://docs.mosek.com/9.0/pythonfusion/index.html>.

[45] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A PAC-Bayesian Approach to Spectrally-Normalized Margin Bounds for Neural Networks. *preprint arXiv:1707.09564*, 2017.

[46] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring Generalization in Deep Learning. In *Advances in Neural Information Processing Systems 30*, pages 5949–5958, 2017.

[47] Alex Nichol, Joshua Achiam, and John Schulman. On First-Order Meta-Learning Algorithms. *arXiv preprint arXiv:1803.02999*, 2018.

[48] Emilio Parrado-Hernández, Amiran Ambroladze, John Shawe-Taylor, and Shiliang Sun. PAC-Bayes Bounds with Data Dependent Priors. *Journal of Machine Learning Research*, 13:3507–3531, 2012.

[49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems 32*, pages 8024–8035, 2019.

[50] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[51] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, 2014.

[52] Anastasia Pentina and Christoph H. Lampert. A PAC-Bayesian Bound for Lifelong Learning. *Proceedings of the 31st International Conference on Machine Learning*, 2014.

[53] Anastasia Pentina and Christoph H. Lampert. Lifelong Learning with Non-i.i.d. Tasks. *Advances in Neural Information Processing Systems 28*, 2015.

[54] María Pérez-Ortiz, Omar Rivasplata, John Shawe-Taylor, and Csaba Szepesvári. Tighter Risk Certificates for Neural Networks. *arXiv preprint arXiv:2007.12911*, 2020.

[55] Sachin Ravi and Hugo Larochelle. Optimization as a Model for Few-Shot Learning. *Proceedings of the 5th International Conference on Learning Representations*, 2017.

[56] Jorma Rissanen. *Stochastic Complexity in Statistical Inquiry*. World Scientific, 1989.

[57] Omar Rivasplata, Emilio Parrado-Hernandez, John Shawe-Taylor, Shiliang Sun, and Csaba Szepesvari. PAC-Bayes Bounds for Stable Algorithms with Instance-Dependent Priors. *Advances in Neural Information Processing Systems 31*, 2018.

[58] Omar Rivasplata, Vikram M. Tankasali, and Csaba Szepesvari. PAC-Bayes with Backprop. *arXiv preprint arXiv:1908.07380*, 2019.

[59] Omar Rivasplata, Ilja Kuzborskij, Csaba Szepesvari, and John Shawe-Taylor. PAC-Bayes Analysis Beyond the Usual Bounds. *Advances in Neural Information Processing Systems* 33, 2020.

[60] Jonas Rothfuss, Vincent Fortuin, and Andreas Krause. PACOH: Bayes-Optimal Meta-Learning with PAC-Guarantees. *arXiv preprint arXiv:2002.05551*, 2020.

[61] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-Learning with Memory-Augmented Neural Networks. *Proceedings of the 33rd International Conference on Machine Learning*, 2016.

[62] Jurgen Schmidhuber. Evolutionary Principles in Self-Referential Learning. On Learning how to Learn: The Meta-Meta-Meta...-Hook. Diploma thesis, Technische Universitat Munchen, Germany, 1987.

[63] Rolf Schneider. *Convex Bodies: The Brunn–Minkowski Theory*. Encyclopedia of Mathematics and its Applications. Cambridge University Press, 2 edition, 2013.

[64] Matthias Seeger. PAC-Bayesian Generalisation Error Bounds for Gaussian Process Classification. *Journal of Machine Learning Research*, 3(Oct):233–269, 2002.

[65] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.

[66] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, Stability and Uniform Convergence. *Journal of Machine Learning Research*, 11:2635–2670, 2010.

[67] John Shawe-Taylor and Robert C. Williamson. A PAC Analysis of a Bayesian Estimator. *Proceedings of the 10th Conference on Computational Learning Theory*, 1997.

[68] Aman Sinha, Hongseok Namkoong, Riccardo Volpi, and John Duchi. Certifying Some Distributional Robustness with Principled Adversarial Training. *arXiv preprint arXiv:1710.10571*, 2020.

[69] Niklas Thiemann, Christian Igel, Olivier Wintenberger, and Yevgeny Seldin. A Strongly Quasiconvex PAC-Bayesian Bound. *Machine Learning Research*, 76:1–26, 2017.

[70] Sebastian Thrun and Lorien Pratt. *Learning to Learn*. Springer Science & Business Media, 1998.

[71] Vladimir N. Vapnik and A. Ya Chervonenkis. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Dokl. Akad. Nauk*, 181(4), 1968.

[72] Ricardo Vilalta and Youssef Drissi. A Perspective View and Survey of Meta-Learning. *Artificial Intelligence Review*, 18:77–95, 2002.

[73] Silvia Villa, Lorenzo Rosasco, and Tomaso Poggio. On Learnability, Complexity and Stability. In *Empirical Inference*, pages 59–69. Springer, 2013.

[74] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching Networks for One Shot Learning. *Advances in Neural Information Processing Systems* 29, 2016.

[75] Jeremy Watt, Reza Borhani, and Aggelos K. Katsaggelos. *Machine Learning Refined: Foundations, Algorithms, and Applications*. Cambridge University Press, 2016.

[76] Wolfram Research, Inc. Mathematica, Version 12.0, 2019. URL <https://www.wolfram.com/mathematica/> Champaign, IL.

[77] Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. Meta-Learning without Memorization. *Proceedings of the 8th International Conference on Learning Representations*, 2020.

[78] Wenda Zhou, Victor Veitch, Morgane Austern, Ryan Adams, and Peter Orbanz. Nonvacuous Generalization Bounds at the ImageNet Scale: A PAC-Bayesian Compression Approach. *Proceedings of the 7th International Conference on Learning Representations*, 2019.