

# Learning to Generate Image Source-Agnostic Universal Adversarial Perturbations

Pu Zhao<sup>1</sup>, Parikshit Ram<sup>2</sup>, Songtao Lu<sup>2</sup>, Yuguang Yao<sup>3</sup>,  
 Djallel Bouneffouf<sup>2</sup>, Xue Lin<sup>1</sup> and Sijia Liu<sup>2,3</sup>

<sup>1</sup>Northeastern University, <sup>2</sup>IBM Research, <sup>3</sup> Michigan State University

zhao.pu@northeastern.edu, {parikshit.ram, songtao}@ibm.com, yaoyugua@msu.edu,  
 djallel.bouneffouf@ibm.com, xue.lin@northeastern.edu, liusiji5@msu.edu

## Abstract

Adversarial perturbations are critical for certifying the robustness of deep learning models. A “universal adversarial perturbation” (UAP) can simultaneously attack multiple images, and thus offers a more unified threat model, obviating an image-wise attack algorithm. However, the existing UAP generator is underdeveloped when images are drawn from different image sources (e.g., with different image resolutions). Towards an *authentic universality across image sources*, we take a novel view of UAP generation as a customized instance of “few-shot learning”, which leverages bilevel optimization and learning-to-optimize (L2O) techniques for UAP generation with improved attack success rate (ASR). We begin by considering the popular model agnostic meta-learning (MAML) framework to meta-learn a UAP generator. However, we see that the MAML framework does not directly offer the universal attack across image sources, requiring us to integrate it with another meta-learning framework of L2O. The resulting scheme for meta-learning a UAP generator (i) has better performance (50% higher ASR) than baselines such as Projected Gradient Descent, (ii) has better performance (37% faster) than the vanilla L2O and MAML frameworks (when applicable), and (iii) is able to simultaneously handle UAP generation for different victim models and image data sources.

## 1 Introduction

Adversarial perturbations are imperceptible changes to input examples (such as images) aimed at manipulating the predictions of a “victim model” [Madry *et al.*, 2017; Carlini and Wagner, 2017]. These are essential for evaluating the worst-case robustness of deep learning (DL) models, which is critical when deploying such models to real world scenarios. The usual focus in adversarial attack generation is on the perturbation of an individual data sample for a single victim model [Madry *et al.*, 2017]. A more powerful threat model is that of *universal adversarial perturbation* (UAP) [Moosavi-Dezfooli *et al.*, 2017] to simultaneously perturb multiple examples. This is often accomplished with standard optimizations such as Projected Gradient Descent (PGD) [Madry *et al.*, 2017], optimizing for a *single* perturbation that simultaneously attacks a set of provided examples.

**Motivation.** Although UAP has been widely studied in the literature [Li *et al.*, 2020; Khrulkov and Oseledets, 2018; Liu *et al.*, 2019b], three fundamental **challenges** (C1-C3) remain.

(C1) *Lack of generalization ability to unseen images*: UAP can attack other previously unseen examples from the same distribution on the same victim model to some degree [Moosavi-Dezfooli *et al.*, 2017]. However, existing UAP methods are usually unable to attack unseen images with high success rate when generated with limited seen examples.

(C2) *Less effective on new images*: For a new set of examples, the UAP usually needs to be regenerated by rerunning the generation algorithm, which is less effective. It is desirable to develop more effective UAP optimization algorithms on new data examples within one- or few-step updates.

(C3) *Non-applicability to diverse image sources*: Beyond existing work, a more advantageous and authentic UAP generation scheme should be broadly applicable to any image source and corresponding victim model – the same attack generator can simultaneously generate perturbations for images from different data sources (with different resolutions).

**Research idea and rationale.** The challenges C1-C3 bring us to the central question we aim to answer in this paper:

(Q) “*Can we develop a powerful threat model in the form of a UAP generator that can improve the attack performance on unseen examples with a small set of seen examples (C1), and generate UAPs more effectively within a few steps (C2) for images from different sources (C3)?*”

To address (C1) and (C2), we formulate the UAP generation process as an instance of *few-shot learning*, and explore the use of *model-agnostic meta-learning* (MAML) techniques [Finn *et al.*, 2017] to warm-start the learning when only few examples are available. We show that MAML enables us to learn a good meta-model of UAP with an *explicit goal of fast adaptation* – the ability to quickly learn generalizable UAP with just a few examples. We highlight that this is different from existing computationally-intensive UAP generators, requiring to perturb a large volume of images for improved attack generalizability.

To address (C3), namely, accomplish a source-agnostic UAP model, we present an extension of MAML to ‘*incongruous tasks*’ – tasks drawn from diverse image sources – by leveraging a different meta-learning framework, *learning-to-optimize* (L2O) [Li and Malik, 2017; Andrychowicz *et al.*, 2016]. The conventional MAML only applies to ‘*congruous tasks*’ where the meta-learning and task-specific learning (fine-tuning) occur on the same set of learnt parameters. In UAP generation, the task-specific parameters are the image perturbations, whose size depend on the image sizes. To

remain agnostic to the image source and size, it is not possible to share the learnt parameters between different image sources. Thus, MAML is not directly applicable. To tackle this, we employ L2O to meta-learn the UAP generator over incongruous tasks. Our rationale is that L2O provides us with a *learnt optimizer*, which use gradients or zeroth-order gradient estimates [Liu *et al.*, 2020], and can operate on objectives with *different set of optimizee variables* (task-specific parameters), allowing meta-learning across incongruous tasks.

We highlight that MAML and L2O make complementary contributions to UAP: MAML focuses on meta-learning “for better generalization”; L2O meta-learns “how to learn”.

**Contributions.** We outline our contributions as follows: ( $\S$  refers to section number):

- ▶ ( $\S 3$ ) We propose a novel interpretation of UAP threat model as a few-shot learning problem.
- ▶ ( $\S 4$ ) Algorithm-wise, we show how meta-learning across incongruous tasks can be developed using an extended bilevel optimization by integrating L2O with MAML and applied to UAP generation. Theory-wise, we quantify how our meta-learned fine-tuner (LFT) differs from L2O.
- ▶ ( $\S 5$ ) We demonstrate the improved attack performance of our meta-learning based UAP generator against standard optimization based generators and other meta-learning based UAP generators (when applicable in limited scope).

**Notation.** In the few-shot learning setup, we denote the  $i^{\text{th}}$  individual task as  $\mathcal{T}_i \sim P(\mathcal{T})$  sampled from a task distribution  $P(\mathcal{T})$  with the task specific (i) learning parameters (or *optimizee variables*)  $\theta_i \in \Theta_i$ , (ii) data domain  $D_i$ , (iii) support (training/fine-tuning/seen) set  $\mathcal{D}_i^{\text{tr}} \in D_i$ , (iv) query (validation/test/unseen) set  $\mathcal{D}_i^{\text{val}} \in D_i$ , (v) objective function  $f_i : \Theta_i \times D_i \rightarrow \mathbb{R}$ . In the MAML framework, the tasks are congruous, and share the optimizee variables, hence the optimizee domains are the same, that is,  $\Theta_i = \Theta_j = \Theta \quad \forall i, j$ , and a single  $\theta \in \Theta$  is meta-learned and fine-tuned to  $\theta_i \in \Theta$  for any task  $\mathcal{T}_i$ . In the L2O framework, the tasks can be incongruous, and only share the optimizer parameters  $\phi$  while maintaining their own separate *optimizee variables*  $\theta_i$ .

## 2 Related work

**Adversarial attacks and UAP.** There exist an extensive amount of work on the design of adversarial attacks, ranging from white-box attacks to black-box attacks [Carlini and Wagner, 2017; Ruan *et al.*, 2020; Goodfellow *et al.*, 2015; Papernot *et al.*, 2016; Zhao *et al.*, 2019; Chen *et al.*, 2017; Xu *et al.*, 2019]. In the context of UAP, various attack generation methods were proposed [Li *et al.*, 2020; Khrulkov and Oseledets, 2018; Liu *et al.*, 2019b; Hashemi *et al.*, 2020; Matachana *et al.*, 2020]. For example, Li *et al.* proposed regionally homogeneous perturbations. Khrulkov and Oseledets leveraged the singular vectors of the Jacobian matrices of deep features to construct UAP. Liu *et al.* designed a robust UAP generation by fully exploiting the model uncertainty. To boost the attack generalizability, most UAP methods require to perturb a large volume of images simultaneously. However, this makes UAP generation computationally-intensive at test time. Moreover, existing UAP generation is restricted to a single image source, leaving image source-agnostic UAP an open question.

**MAML.** MAML [Finn *et al.*, 2017] has been extremely useful in supervised and reinforcement learning (RL), and

widely extended and studied both theoretically [Liu *et al.*, 2019a; Balcan *et al.*, 2019; Khodak *et al.*, 2019] and empirically [Nichol *et al.*, 2018]. It has been extended to model uncertainty [Finn *et al.*, 2018] and handle the online setting [Finn *et al.*, 2019]. The second order derivatives in MAML have been handled in multiple ways [Rajeswaran *et al.*, 2019; Fallah *et al.*, 2019]. Specific to RL, various enhancements obviate the second order derivatives of the RL reward function, such as variance reduced policy gradients [Liu *et al.*, 2019a] and Monte Carlo zeroth-order Evolution Strategies gradients [Song *et al.*, 2020]. However, all extensions and applications of MAML focus on congruous tasks where different few-shot tasks share the same parameters and optimizee domain.

**L2O.** Learnt optimizers have long been considered in training neural networks [Bengio *et al.*, 1990; Thrun and Pratt, 2012]. More recent work has posed optimization *with gradients* as a RL problem [Li and Malik, 2017] or as learning a recurrent neural network (RNN) [Andrychowicz *et al.*, 2016] instead of leveraging the usual hand-crafted optimizers (such as SGD or Adam [Kingma and Ba, 2015]). The RNN based optimizers have been improved [Wichrowska *et al.*, 2017; Chen *et al.*, 2020] by – (i) using hierarchical RNNs to better capture parameter structure of DL models, (ii) using hand-crafted-optimizer-inspired inputs to RNN (such as momentum), (iii) using a diverse set of optimization objectives (with different hardness levels) to meta-learn the RNN, and (iv) leveraging different training techniques such as curriculum learning and imitation learning. The learnt optimizers have also been successful with particle swarm optimization [Cao *et al.*, 2019] and zeroth-order gradient estimates [Ruan *et al.*, 2020]. Learnt optimizers have been used for meta-learning with congruous few-shot tasks [Ravi and Larochelle, 2016], but MAML has been shown to outperform it. In the context of adversarial robustness, L2O has also been leveraged to design instance-wise attack generation [Ruan *et al.*, 2020; Jiang *et al.*, 2018] and adversarial defense (such as adversarial training [Xiong and Hsieh, 2020]).

## 3 UAP Design as a Few-Shot Problem

Many existing attack generators execute in an instance-wise manner, namely, requesting repeated invocations of an iterative optimizer to acquire adversarial perturbations with respect to an *individual* input example [Chen *et al.*, 2017; Croce and Hein, 2020]. To circumvent the limitation of instance-wise attack generator, the problem of UAP arises, which seeks a *single* perturbation pattern to manipulate the DNN outputs over *multiple examples simultaneously* [Moosavi-Dezfooli *et al.*, 2017; Matachana *et al.*, 2020].

**Problem setup.** Let  $\mathcal{T}_i$  denote an attack generation task, which constitutes a support set  $\mathcal{D}_i^{\text{tr}}$  used to generate the UAP  $\theta_i$  and a query set  $\mathcal{D}_i^{\text{val}}$  (from the same classes as in  $\mathcal{D}_i^{\text{tr}}$ ) on which we evaluate the ASR (attack success rate) of the generated UAP  $\theta_i$ . Our *goal* is to learn a UAP with a few examples in  $\mathcal{D}_i^{\text{tr}}$  that can successfully attack *unseen* examples  $\mathcal{D}_i^{\text{val}}$  with the same victim model. This motivates us to view UAP generation as a **few-shot learning** problem. Mathematically, given a few-shot UAP generation task  $\mathcal{T}_i$  with training set  $\mathcal{D}_i^{\text{tr}}$ , the task-specific UAP  $\theta_i$  is obtained by solving:

$$\underset{\theta_i}{\text{minimize}} \quad f_i(\theta_i, \mathcal{D}_i^{\text{tr}}) := \sum_{(\mathbf{x}, y) \in \mathcal{D}_i^{\text{tr}}} \ell_{\text{atk}}(\theta_i, \mathbf{x}, y) + \lambda \|\theta_i\|_1, \quad (1)$$

where  $y$  is the true label of  $\mathbf{x}$ ,  $\lambda > 0$  is a regularization parameter, and  $\ell_{\text{atk}}$  is the C&W attack loss [Carlini and Wagner, 2017], which is 0 (indicating a successful attack) with an incorrect predicted class. The second term is an  $\ell_1$  norm regularizer, penalizing the perturbation strength.

In our few-shot notation, the attack loss (1) is considered as the task-specific loss  $f_i$  with  $\theta_i$  as the task-specific parameters. With multiple few-shots tasks  $\mathcal{T}_i$  and corresponding  $(\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{val}})$ , we wish to meta-learn a UAP generator that can solve new few-shot UAP generation tasks.

To achieve this, we leverage MAML to meta-learn an initialization of *optimizee* variables  $\theta$  (i.e., UAP variables) that enables fast adaptation to new tasks when fine-tuning the optimizee from this learned initialization with a few new examples. Formally, with  $N$  few-shot learning tasks  $\{\mathcal{T}_i\}_{i=1}^N$ , when meta-learning with  $\mathcal{T}_i$ , (i) the support set  $\mathcal{D}_i^{\text{tr}}$  is used for the task-specific *inner level* in MAML to fine-tune the initial optimizee  $\theta$ , and (ii) the query set  $\mathcal{D}_i^{\text{val}}$  is used in the *outer level* for evaluating the fine-tuned optimizee  $\theta_i^*(\theta)$  to meta-update  $\theta$ . Thus, **MAML-oriented UAP generation** solves the following **bilevel optimization problem**:

$$\begin{aligned} & \underset{\theta}{\text{minimize}} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{(\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{val}}) \sim \mathcal{T}_i} [f_i(\theta_i^*(\theta); \mathcal{D}_i^{\text{val}})], \\ & \text{subject to } \theta_i^*(\theta) = \arg \min_{\theta_i} f_i(\theta_i(\theta); \mathcal{D}_i^{\text{tr}}) \end{aligned} \quad (2)$$

where  $\theta_i(\theta)$  is the task-specific optimizee and  $f_i(\theta_i(\theta); \mathcal{D})$  is the task-specific loss evaluated on data  $\mathcal{D}$  using variable  $\theta_i(\theta)$  obtained from fine-tuning the meta-learned initialization  $\theta$ . We will use  $\theta_i := \theta_i(\theta)$  from hereon.

**Why (2)? Formulation (1) is not able to generate transferable UAP.** We provide a warm-up example showing why we choose to rely on the bilevel few-shot formulation (2), rather than (1). We thus examine if solving the task-specific problem (1) is sufficient to generate UAP with “attack generalizability” when applied to unseen test example. We consider a UAP generation on CIFAR-10 and use PGD to generate the task-specific UAP  $\theta_i$  over 1000 tasks each with 4 images. We compare the average ASR of  $\theta_i$  on  $\mathcal{D}_i^{\text{tr}}$  (with which the UAP was generated) and  $\mathcal{D}_i^{\text{val}}$  (unseen images from same classes). Since the UAP  $\theta_i$  is generated with  $\mathcal{D}_i^{\text{tr}}$ , the ASR achieves 100% on  $\mathcal{D}_i^{\text{tr}}$ . However, the same  $\theta_i$  when applied to  $\mathcal{D}_i^{\text{val}}$  achieves *less than 30% ASR* – this highlights how the PGD-generated UAP that only solves problem (1) has low ASR on unseen images from the same distribution on the same victim model. Formulation (2) explicitly minimizes the loss on unseen examples in the outer-loop, thereby explicitly promoting attack generalizability, allowing us to address C1 in §1.

## 4 Learning Optimizers for Fast Adaptation

While (2) improves attack generalizability and addresses challenges C1 and C2 stated in Sec. 1, it is only applicable to ‘congruous tasks’ where all UAP generation tasks are from the same image source. In this section, we generalize (2) to ‘**incongruous tasks**’ with learned fine-tuners (LFTs) by integrating with L2O, thereby addressing C3 to yield image source-agnostic UAP generators.

**MAML and beyond.** To solve problem (2), the conventional approach [Finn *et al.*, 2017; Yin *et al.*, 2020] relies on the approximation of the inner problem with a  $K$ -step gradient

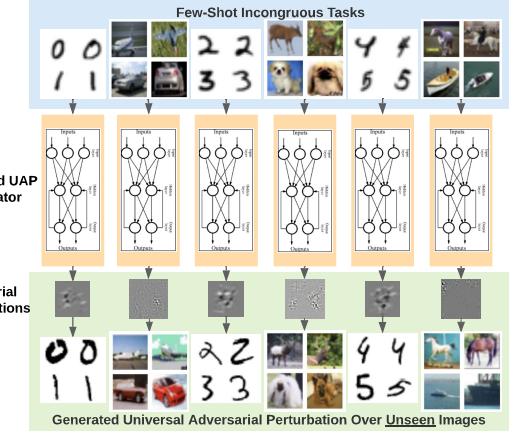


Figure 1: Desired UAP generation setup with *incongruous few-shot* tasks across  $28 \times 28$  MNIST and  $3 \times 32 \times 32$  CIFAR-10. The meta-learned optimizer parameters are shared by all incongruous tasks to find task-specific UAP patterns even with various image sizes.

descent (GD) with the initial iterate  $\theta_i^{(0)} \leftarrow \theta$ , the final iterate  $\theta_i^* \leftarrow \theta_i^{(K)}$  and  

$$\theta_i^{(k)} = \theta_i^{(k-1)} - \alpha \nabla_{\theta_i} f_i(\theta_i^{(k-1)}; \mathcal{D}_i^{\text{tr}}) \text{ for } k \in [K], \quad (3)$$

where  $\theta_i^{(k)}$  is the  $k^{\text{th}}$ -step optimizee fine-tuned with  $\mathcal{D}_i^{\text{tr}}$  from  $\theta_i^{(0)} = \theta$ ,  $\alpha > 0$  is a learning rate, and  $[K] = \{1, 2, \dots, K\}$ .

It is clear from (3) that both levels of the optimization (with respect to  $\theta$  and  $\theta_i$ ) in (2) must operate on the same-type optimizee variables, and accordingly, few-shot tasks  $\{\mathcal{T}_i\}_{i=1}^N$  are restricted to problems which *share the same optimizee domain* (i.e., *congruous* tasks drawn from the same image source). However, in the general meta-learning setting, similar tasks could be from related yet *incongruous* domains corresponding to *different objectives*  $\{f_i\}$  with optimizee variables of *different domains* (such as different image sizes and number of channels) that *cannot be shared between tasks* as in our UAP generation tasks illustrated in Figure 1 – UAP parameters cannot be shared between images from different data sources with different resolutions. In such cases, meta-learning the *initial iterate is not possible*. Next, we will leverage L2O to meta-learn an optimizer – the fine-tuner – for fast adaptation of the task-specific optimizee  $\theta_i$  in a few-shot setting even when meta-learning across *incongruous* tasks.

**L2O-enabled learned fine-tuner for MAML.** L2O allows us to replace the hand-designed GD (3) with a learnable RNN parameterized by  $\phi$ . For any task  $\mathcal{T}_i$ ,  $\text{RNN}_\phi(\cdot)$  mimics a hand-crafted gradient based optimizer to output a descent direction  $\Delta\theta_i$  to update task-specific optimizee variable  $\theta_i$  given the function gradients as input. Thus, we replace (3) with

$$\begin{aligned} \Delta\theta_i^{(k)}, \mathbf{h}_i^{(k)} &= \text{RNN}_\phi \left( g_i(\theta_i^{(k-1)}; \mathcal{D}_i^{\text{tr}}), \mathbf{h}_i^{(k-1)} \right), \\ \theta_i^{(k)} &= \theta_i^{(k-1)} - \Delta\theta_i^{(k)}, \quad \forall k \in [K], \end{aligned} \quad (4)$$

where  $\mathbf{h}_i^{(k)}$  denotes the state of  $\text{RNN}_\phi$  at the  $k^{\text{th}}$  RNN unrolling step,  $g_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{tr}})$  is the gradient  $\nabla_{\theta_i} f_i(\theta_i; \mathcal{D}_i^{\text{tr}})$  or gradient estimate [Liu *et al.*, 2018]. Each task-specific  $\theta_i$  is initialized with a random  $\theta_i^{(0)}$ . Note that  $\Delta\theta_i^{(k)} := \Delta\theta_i^{(k)}(\phi)$  is a function of  $\phi$  in (4) and hence  $\theta_i^{(k)} := \theta_i^{(k)}(\phi)$  depends on  $\phi$ .

We term  $\text{RNN}_\phi$  in (4) *learned fine-tuner* (LFT) for incongruous few-shot learning. Combining (4) with (2), we can

**Algorithm 1** Meta-learning LFT with problem (5)

---

```

1: Input: UAP generation tasks  $\{\mathcal{T}_i\}_{i \in [N]}$ , # meta-learning steps
    $T$ , # fine-tuning steps per task  $K$ , initial  $\phi$ ,  $\{\theta_i^{(0)}, \mathbf{h}_i^{(0)}\}_{i \in [N]}$ ,
   meta-learning rate  $\beta > 0$ 
2: for  $t \leftarrow 1, 2, \dots, T$  do
3:   for tasks  $i \in$  sampled task batch  $\mathcal{B}_t \subseteq N$  do
4:     sample data  $\mathcal{D}_i^{\text{tr}} \sim \mathcal{T}_i$  for fine-tuning
5:     generate task-specific optimizee (the UAP) with  $\mathcal{D}_i^{\text{tr}}$  via
       (4) for  $K$  steps to obtain  $\{\theta_i^{(k)}, \mathbf{h}_i^{(k)}\}_{k=1}^K$ 
6:     sample data  $\mathcal{D}_i^{\text{val}} \sim \mathcal{T}_i$  for the meta-update
7:     obtain task-specific fast-adaptation gradient w.r.t. LFT
       parameters  $\phi$  with  $\mathcal{D}_i^{\text{val}}$ 
        $\mathbf{g}^{(i)} \leftarrow \nabla_{\phi} \sum_{k=1}^K [w_k f_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{val}})] \quad (6)$ 
8:   end for
9:   update LFT parameters  $\phi$ :  $\phi \leftarrow \phi - \beta \sum_{i \in \mathcal{B}_t} \mathbf{g}^{(i)}$  { // Adam
      can also be used}
10: end for
11: Output:  $\text{RNN}_{\phi}$ 

```

---

cast the meta-learning of a LFT as

$$\underset{\phi}{\text{minimize}} \quad \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{(\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{val}}) \sim \mathcal{T}_i} \sum_{k=1}^K w_k f_i(\theta_i^{(k)}(\phi); \mathcal{D}_i^{\text{val}}) \quad (5)$$

:=  $\hat{F}(\phi)$

subject to  $\theta_i^{(k)}(\phi)$  is given by (4),

where  $w_k$  is an importance weight for the  $k^{\text{th}}$  unrolled step in (4). We can set (i)  $w_k = 1$  [Andrychowicz *et al.*, 2016], (ii)  $w_k = k$  [Ruan *et al.*, 2020], or (iii)  $w_k = \mathbb{I}[k = K]$  [Lv *et al.*, 2017]. Choice (iii) matches the objective (2), focusing on the final fine-tuned solution. However, unlike MAML, problem (5) meta-learns the *fine-tuner*  $\phi$  instead of an initialization  $\theta$ . We elide the  $\phi$  argument in the sequel for brevity.

The fine-tuner  $\phi$  acquired from (5) yields a *UAP generator* with generalized universality across hybrid image sources. First, it can generate a UAP  $\theta_i^{(k)}(\phi)$  with a few “seen” images in  $\mathcal{D}_i^{\text{tr}}$  that can successfully attack “unseen” images in  $\mathcal{D}_i^{\text{val}}$  (C1) within a few updating steps (C2). Second, as will be evident later, its form given by a RNN-based optimizer enables us to handle incongruous few-shot tasks (that is, design attacks) with images for different data sources (C3).

## 4.1 Methodologies

The meta-learning problem (5) is still a bilevel optimization, similar to MAML (2). However, both inner and outer levels are distinct from MAML: In the inner level, we update a task-specific optimizee  $\theta_i$  by unrolling  $\text{RNN}_{\phi}$  for  $K$  steps from a random initial state  $\theta_i^{(0)}$ ; by contrast, MAML uses GD to update  $\theta_i$  from the meta-learned initialization. In the outer level, we minimize the objective (5) w.r.t. the optimizer  $\phi$  instead of the optimizee initialization  $\theta$ . We present our proposed scheme in Algorithm 1. In each outer iteration  $t \in [T]$  of the meta-learning, we sample a batch  $\mathcal{B}_t$  of few-shot UAP generation tasks (line 3), and for each sampled task  $\mathcal{T}_i, i \in \mathcal{B}_t$ , we obtain a set of seen images  $\mathcal{D}_i^{\text{tr}}$  (line 4) and use it with the LFT to generate a (sequence of) UAP  $\theta_i^{(k)}$  using the update rule in (4). Then we sample a set of unseen images  $\mathcal{D}_i^{\text{val}}$  (line 6) and generate the *fast adaptation gradient* w.r.t.  $\phi$  (line 7)

that explicitly takes into account the utility of the generated UAP  $\theta_i^{(k)}$  on unseen images. These gradients w.r.t.  $\phi$  are aggregated across all tasks in the batch and used to update the LFT parameters  $\phi$  (line 9). In what follows, we discuss our proposed meta-learning (Alg. 1), showcasing its (i) ability to meta-learn across incongruous tasks, (ii) applicability to zeroth-order (ZO) optimization, (iii) theoretical advantage from the fast adaptation gradient (6).

**Incongruous meta-learning.** When fine-tuning the task-specific optimizee  $\theta_i$  by  $\text{RNN}_{\phi}$  (Algorithm 1, Step 5), we use an *invariant* RNN architecture to tolerate the task-specific variations in the domains  $\{\Theta_i\}_{i=1}^N$  (e.g., dimensions) of optimizee variables  $\{\theta_i\}_{i=1}^N$ . Recall from (4) that  $\text{RNN}_{\phi}$  uses the gradient or gradient estimate  $g_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{tr}})$  as an input, which has the same dimension as  $\theta_i$ . At first glance, a single  $\text{RNN}_{\phi}$  seems incapable of handling incongruous  $\{\mathcal{T}_i\}_{i=1}^N$  defined over optimizee variables of different dimensionalities. However, a  $\text{RNN}_{\phi}$  configured as a *coordinate-wise* Long Short Term Memory (LSTM) network [Andrychowicz *et al.*, 2016], is *invariant* to the dimensionality of optimizee variables  $\{\theta_i\}_{i=1}^N$  by using a separate LSTM to independently operate on each coordinate of  $\theta_i$  for any task  $i \in [N]$ , but requiring all LSTMs to *share their weights*. This weight-sharing allows the RNN to operate on tasks with optimizee variables of different dimensionalities. In contrast to MAML, the invariant  $\text{RNN}_{\phi}$  expands meta-learning for fast adaptation beyond congruous tasks to incongruous ones such as designing UAPs across incongruous attack tasks.

**Derivative-free meta-learning.** L2O in (4) allows us to update the task-specific optimizee variable  $\theta_i$  using not only first-order (FO) information (gradients) but also zeroth-order (ZO) information (function values) if the loss function  $f_i$  is a black-box objective function. We can estimate the gradient  $\nabla f_i(\theta_i; \mathcal{D}_i)$  with finite-differences of function values  $g_i(\theta_i; \mathcal{D}_i)$  [Liu *et al.*, 2020; Ruan *et al.*, 2020]:

$$g_i(\theta_i; \mathcal{D}_i) = \frac{\sum_{j=1}^n [\mathbf{u}_j (f_i(\theta_i + \mu \mathbf{u}_j; \mathcal{D}_i) - f_i(\theta_i; \mathcal{D}_i))] }{\mu n}, \quad (7)$$

where  $\mu > 0$  is a small smoothing parameter,  $\mathbf{u}_j, j \in [n]$  are  $n$  random directions with entries from  $\mathcal{N}(0, 1)$ . The function  $g_i$  can also be sophisticated quantities derived from gradients or gradient estimates [Wichrowska *et al.*, 2017; Lv *et al.*, 2017; Cao *et al.*, 2019]. The support for ZO optimization is crucial when explicit gradient are computationally difficult or infeasible in the black-box attack setting.

Since Alg. 1 meta-learns the optimizer variable  $\phi$  rather than the initialization of optimizee variable  $\theta$ , it requires a different meta-learning gradient  $\nabla_{\phi} \sum_k w_k f_i(\theta_i^{(k)}, \mathcal{D}_i^{\text{val}})$ . Focusing only on  $\nabla_{\phi} f_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{val}})$  in (6):

$$\nabla_{\phi} f_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{val}}) = \underbrace{\frac{\partial f_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{val}})}{\partial \theta_i}}_{g_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{val}})} \bullet \underbrace{\frac{\partial \theta_i^{(k)}}{\partial \phi}}_{\mathbf{G}_i^{(k)}} \quad (8)$$

where  $\bullet$  denotes a matrix product that the chain rule obeys [Petersen *et al.*, 2008]. The computation of  $\mathbf{G}_i^{(k)}$  calls for the first-order (or second-order) derivative of  $f_i$  w.r.t.  $\theta_i$  if  $g_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{tr}})$  denotes the gradient estimate of  $f_i$  in (4). See Appendix A for details<sup>1</sup>.

<sup>1</sup>A full version is available at <https://arxiv.org/abs/2009.13714>.

**Theoretical advantage of fast adaptation gradient.** We make use of the fast adaptation gradient (4) explicitly evaluating the performance of the UAP (generated by fine tuning with  $\mathcal{D}_i^{\text{tr}}$ ) on unseen images in  $\mathcal{D}_i^{\text{val}}$ . We want to highlight the advantage of this explicit choice – we could have instead used the gradient of the fine-tuning objective  $f_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{tr}})$  w.r.t.  $\phi$  by solving the following problem:

$$\begin{aligned} \underset{\phi}{\text{minimize}} \quad & F(\phi) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \sim \mathcal{T}_i} \sum_{k=1}^K w_k f_i(\theta_i^{(k)}; \mathcal{D}_i^{\text{tr}}), \\ \text{s.t. } & \theta_i^{(k)} \text{ defined as (4).} \end{aligned} \quad (9)$$

Different from (5), this is the standard L2O in the context of few-shot tasks. Under mild assumptions, we show the following result, quantifying the difference between L2O and our proposed meta-learning in terms of the meta-learning gradient w.r.t.  $\phi$ , highlighting the explicit effect of the fast adaptation gradient (see Appendix B for details):

**Proposition 1.** Consider the meta-learning objectives defined in (5) and (9). Suppose that the gradient size is bounded by  $G$  and gradient estimate per sample has uniformly bounded variance  $\sigma^2$ . Then, for any  $\phi$ , we have

$$\left\| \nabla_{\phi} F(\phi) - \nabla_{\phi} \hat{F}(\phi) \right\| \leq \sqrt{2G\sigma} \sqrt{\frac{1}{D_{\text{tr}}} + \frac{1}{D_{\text{val}}}} \quad (10)$$

where  $D_{\text{tr}} := \min_{i \in [N]} |\mathcal{D}_i^{\text{tr}}|$  and  $D_{\text{val}} := \min_{i \in [N]} |\mathcal{D}_i^{\text{val}}|$  denote the minimum batch size of per-task datasets  $\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{val}} \sim \mathcal{T}_i, i \in [N]$ .

**Remark 1.** When the data size is small – the few-shot regime – there could be a significant difference in the gradients w.r.t.  $\phi$ , resulting in a significantly different solutions, especially for the case where  $\sigma$  or  $G$  is large. This explains the significant difference between empirical performance of L2O and our LFT in evaluation over few-shot learning problems.

## 5 Experiments

We demonstrate the effectiveness of LFT through extensive experiments on the UAP design in the black-box setup [Chen *et al.*, 2017], where the internal configurations and parameters of the DNN are not revealed to the attacker. Thus, the only interaction of the adversary with the system is via submission of inputs and receiving the corresponding predicted outputs. LFT is implemented using ZO gradient estimates.

**Experimental setup.** We wish to evaluate the ability of UAP generated with a small set of seen images to successfully attack unseen images from the same image source (for a fixed victim model). For this, we generate 100 UAP generation tasks  $\mathcal{T}_i, i = 1, \dots, 100$ , each with a set  $\mathcal{D}_i^{\text{tr}}$  of seen images (to be used to generate the UAP) and a set  $\mathcal{D}_i^{\text{val}}$  of unseen images (on which the generated UAP is evaluated). In both  $\mathcal{D}_i^{\text{tr}}$  &  $\mathcal{D}_i^{\text{val}}$ , 2 image classes with 2 samples per class are randomly selected. As image sources, we utilize MNIST and CIFAR-10, and as victim architectures, we utilize LeNet [Lecun *et al.*, 1998] and VGG-11 [Simonyan and Zisserman, 2014]. Note that the same architecture applied to different image sources will lead to different victim models. For meta-learning, we separately generate an additional 1000 UAP generation tasks from each image source, ensuring that there is no overlap between images present in the evaluation (meta-test) and the meta-learning tasks.

Table 1: Attack success rate (ASR) of UAPs generated by different schemes using 200 steps aggregated over 100 meta-test tasks from different image sources. We highlight the superior ASR of LFT.

Meta-learning	Meta-testing	MNIST (LeNet-5)	CIFAR-10 (LeNet-5)	CIFAR-10 (VGG-11)
	S-UAP	63 ± 5	42 ± 4	38 ± 4
MNIST (LeNet-5)	PGD	50 ± 7	25 ± 0	25 ± 0
	MAML	100 ± 0	-	-
	L2O	100 ± 0	40 ± 10	36 ± 2
	LFT (Ours)	100 ± 0	50 ± 0	48 ± 3

**Baselines.** We consider two non-meta-learning based standard UAP generators – PGD-based (termed PGD) and singular vector based (termed S-UAP) [Khrulkov and Oseledets, 2018]. In addition to our proposed meta-learning based LFT scheme, we also consider standard L2O (given by (9)) to ablate the effect of the fast adaptation gradient<sup>2</sup>. We also consider standard MAML (given by (2)) to ablate the effect of using a meta-learnt optimizer (with random initial iterate) instead of a meta-learnt initialization. Note that MAML is only applicable if all the meta-learning and meta-testing UAP generation tasks come from a single image source. When the meta-learning tasks are generated from multiple image sources, we meta-learn a MAML per image source and term it *ensemble MAML*. Note that this requires us to solve the bilevel problem (2) separately for each image source.

**Learned optimizer details.** We use a one-layer LSTM with 10 hidden units, and one additional linear layer to project the RNN hidden state to the output. We use Adam with an initial learning rate of 0.001 to meta-learn the RNN with truncated backpropagation through time (BPTT) by unrolling the RNN for 20 steps and running each optimization for 200 steps.

**Evaluation metrics.** We report (a) averaged attack success rate (ASR), (b)  $\ell_1$  perturbation strength, (c) convergence in terms of optimization steps needed to reach 100% ASR. We also provide visualization of generated UAPs in Appendix C.

### Experimental questions and results.

Can LFT outperform standard non-meta-learning schemes for few-shot UAP generation tasks, even for tasks from images sources & victim models not encountered at meta-learning?

To answer this, we meta-learn only with tasks generated from MNIST (with LeNet-5 as the victim model), and perform the meta-testing on tasks generated from (i) MNIST with LeNet-5, (ii) CIFAR-10 with LeNet-7, and (iii) CIFAR-10 with VGG-11, where (ii) and (iii) involve image sources and victim models not encountered during the meta-learning. We summarize the results in Table 1. Note that MAML is only applicable in case (i). For case (i) where meta-testing tasks are from the image source used for meta-learning, all meta-learning schemes achieve 100% ASR compared to the 63% for S-UAP (PGD is much lower). This demonstrates the improved attack generalizability from meta-learning. For image sources not encountered during meta-learning (cases (ii) & (iii)), LFT still continues to outperform S-UAP by 8–10%. LFT outperforms L2O by over 10%, highlighting the gain from the fast adaptation gradient.

Can LFT effectively meta-learn with multiple image sources and improve ASR & convergence with lower UAP  $\ell_1$  norms?

Here, we consider meta-learning with two meta-learning setups (L1) only tasks from MNIST and (L2) meta-

<sup>2</sup>Unseen images  $\mathcal{D}_i^{\text{val}}$  are not utilized in L2O (9), so we add it to  $\mathcal{D}_i^{\text{tr}}$  to ensure that L2O gets the same images for meta-learning.

Table 2: Aggregated ASR and  $\ell_1$ -norm of UAPs by various schemes. We highlight the best performance at each (meta-learning and -testing) scenario, measured by (i) highest ASR within 100 steps, (ii) distortion at 100 steps, and (iii) steps required to first reach 100% ASR.

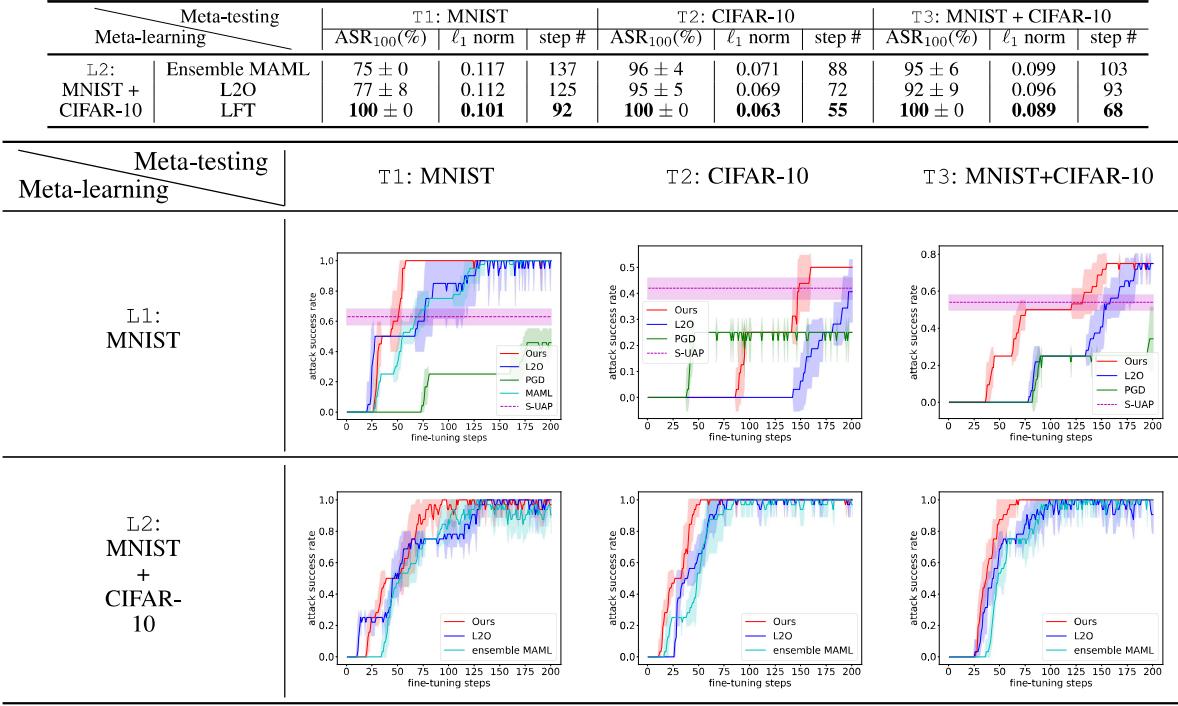


Figure 2: Aggregated ASR of UAPs generated by different schemes in various meta-learning & meta-testing settings (*higher is better*). Each cell presents results for schemes meta-learnt with tasks from row-specific image source(s) and meta-tested with tasks from column-specific image source(s) with increasing UAP generation steps. S-UAP does not have an iterative process and thus we show the final ASR (dashed).

learning with tasks from both MNIST and CIFAR-10 (MNIST+CIFAR-10), and three meta-testing setups with tasks from (T1) MNIST, (T2) CIFAR-10, and (T3) both (MNIST+CIFAR-10). In these cases, the victim architecture is LeNet-5 for MNIST and LeNet-7 (LeNet with more CONV layers) for CIFAR-10 to demonstrate the ability to work with various victim architectures. **Figure 2** presents ASR (on unseen images  $D_i^{\text{val}}$ ), aggregated over 100 UAP generation tasks, with increasing number of steps in iterative UAP generation schemes (using the seen images  $D_i^{\text{tr}}$ ); S-UAP is a non-iterative scheme. Note that MAML is only applicable in the T1 case when meta-learning with MNIST tasks in L1 case; when meta-learning with tasks from MNIST+CIFAR-10 in L2 case, we utilize ensemble MAML.

**Figure 2** indicates that our proposed LFT converges first to the best ASR, with L2O and (ensemble) MAML performing comparably to each other when applicable. As seen in Table 1, the non-meta-learning schemes (PGD & S-UAP) are unable to match the ASR of the meta-learning schemes in all meta-testing cases T1, T2, T3 – PGD converges to less than half the ASR achieved by LFT. Note that, when meta-learning with MNIST+CIFAR-10 in case L2, the meta-test performance of LFT on MNIST in case T1 is not significantly affected, indicating that meta-learning from multiple image sources *does not hurt* the performance of LFT. Furthermore, meta-learning from MNIST+CIFAR-10 leads to improved ASR when meta-testing with CIFAR-10 in T2 or MNIST+CIFAR-10 in T3, *highlighting the improvement in LFT from meta-learning with multiple image sources*.

The relative performances of the meta-learning based

schemes are further summarized in **Table 2**. These results indicate that, compared to other meta-learning baselines, LFT achieves better ASR with limited number of steps (5-18% improved ASR), converges to 100% ASR faster (24-37% reduction in number of steps), and does so while producing perturbations with smaller  $\ell_1$  norms (up to 10% smaller). We present further results in Appendix C which show that LFT produces smaller  $\ell_1$  norms than PGD as well. We also study the effect of the number of tasks available per image source for meta-learning in Appendix D.

**Summary of evaluation.** We mitigate challenge C1 by demonstrating improved attack generalizability (improved ASR on unseen images) with meta-learning. Improved convergence (to 100% ASR) of LFT with smaller perturbation strength mitigates challenge C2. By demonstrating improved ASR for LFT (i) by meta-learning with tasks from different image sources, and (ii) when handling UAP generation tasks from image sources not encountered during meta-learning, we mitigate challenge C3.

## 6 Conclusion

In this paper, we focus on the various challenges in UAP generation, and present a meta-learning scheme LFT that extends MAML with the learning-to-optimize framework. This LFT improves the attack universality or generalizability of UAPs generated from a small set of images, and can be meta-learned from multiple image sources and be applied to image sources and architectures not seen during meta-learning, thereby extending the universality of UAP generators.

## 7 Acknowledgments

This work is partly supported by the National Science Foundation Award CNS-1932351.

## References

- [Andrychowicz *et al.*, 2016] Marcin Andrychowicz, Misha Denil, et al. Learning to learn by gradient descent by gradient descent. In *NeurIPS*, 2016.
- [Balcan *et al.*, 2019] Maria-Florina Balcan, Mikhail Khodak, and Ameet Talwalkar. Provable guarantees for gradient-based meta-learning. In *ICML*, 2019.
- [Bengio *et al.*, 1990] Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. *Learning a synaptic learning rule*. Citeseer, 1990.
- [Cao *et al.*, 2019] Yue Cao, Tianlong Chen, et al. Learning to optimize in swarms. In *NeurIPS*, 2019.
- [Carlini and Wagner, 2017] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *SP*, 2017.
- [Chen *et al.*, 2017] Pin-Yu Chen, Huan Zhang, et al. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *ACM Workshop on Artificial Intelligence and Security*, 2017.
- [Chen *et al.*, 2020] Tianlong Chen, Weiyi Zhang, et al. Training stronger baselines for learning to optimize. *NeurIPS*, 2020.
- [Croce and Hein, 2020] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *ICML*, 2020.
- [Fallah *et al.*, 2019] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the convergence theory of gradient-based model-agnostic meta-learning algorithms. *arXiv preprint*, 2019.
- [Finn *et al.*, 2017] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017.
- [Finn *et al.*, 2018] Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. In *NeurIPS*, 2018.
- [Finn *et al.*, 2019] Chelsea Finn, Aravind Rajeswaran, et al. Online meta-learning. In *ICML*, 2019.
- [Goodfellow *et al.*, 2015] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. 2015 *ICLR*, arXiv preprint arXiv:1412.6572, 2015.
- [Hashemi *et al.*, 2020] Atiye Sadat Hashemi, Andreas Bär, et al. Transferable universal adversarial perturbations using generative models. *arXiv preprint*, 2020.
- [Jiang *et al.*, 2018] Haoming Jiang, Zhehui Chen, et al. Learning to defense by learning to attack. *arXiv preprint*, 2018.
- [Khodak *et al.*, 2019] M Khodak, M Balcan, et al. Adaptive gradient-based meta-learning methods. In *NeurIPS*, 2019.
- [Khrulkov and Oseledets, 2018] Valentin Khrulkov and Ivan Oseledets. Art of singular vectors and universal adversarial perturbations. In *CVPR*, 2018.
- [Kingma and Ba, 2015] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [Lecun *et al.*, 1998] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, Nov 1998.
- [Li and Malik, 2017] Ke Li and Jitendra Malik. Learning to optimize. *ICLR*, 2017.
- [Li *et al.*, 2020] Yingwei Li, Song Bai, et al. Regional homogeneity: Towards learning transferable universal adversarial perturbations against defenses. In *ECCV*, 2020.
- [Liu *et al.*, 2018] S. Liu, J. Chen, et al. Zeroth-order online admm: Convergence analysis and applications. In *AISTATS*, 2018.
- [Liu *et al.*, 2019a] Hao Liu, Richard Socher, and Caiming Xiong. Taming maml: Efficient unbiased meta-reinforcement learning. In *ICML*, 2019.
- [Liu *et al.*, 2019b] Hong Liu, Rongrong Ji, et al. Universal adversarial perturbation via prior driven uncertainty approximation. In *ICCV*, 2019.
- [Liu *et al.*, 2020] Sijia Liu, Pin-Yu Chen, et al. A primer on zeroth-order optimization in signal processing and machine learning. *IEEE Signal Processing Magazine*, 2020.
- [Lv *et al.*, 2017] Kaifeng Lv, Shunhua Jiang, and Jian Li. Learning gradient descent: Better generalization and longer horizons. In *ICML*, 2017.
- [Madry *et al.*, 2017] Aleksander Madry, Aleksandar Makelov, et al. Towards deep learning models resistant to adversarial attacks. *arXiv preprint*, 2017.
- [Matachana *et al.*, 2020] Alberto G Matachana, Kenneth T Co, et al. Robustness and transferability of universal attacks on compressed models. *arXiv preprint*, 2020.
- [Moosavi-Dezfooli *et al.*, 2017] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, et al. Universal adversarial perturbations. In *CVPR*, 2017.
- [Nichol *et al.*, 2018] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- [Papernot *et al.*, 2016] Nicolas Papernot, Ian Goodfellow, et al. cleverhans v1.0.0: an adversarial machine learning library. *arXiv preprint*, 2016.
- [Petersen *et al.*, 2008] KB Petersen, MS Pedersen, et al. The matrix cookbook, vol. 7. *Technical University of Denmark*, 15, 2008.
- [Rajeswaran *et al.*, 2019] Aravind Rajeswaran, Chelsea Finn, et al. Meta-learning with implicit gradients. In *NeurIPS*, 2019.
- [Ravi and Larochelle, 2016] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2016.
- [Ruan *et al.*, 2020] Yangjun Ruan, Yuanhao Xiong, et al. Learning to learn by zeroth-order oracle. In *ICLR*, 2020.
- [Simonyan and Zisserman, 2014] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint*, 2014.
- [Song *et al.*, 2020] Xingyou Song, Wenbo Gao, et al. Es-maml: Simple hessian-free meta learning. *ICLR*, 2020.
- [Thrun and Pratt, 2012] Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 2012.
- [Wichrowska *et al.*, 2017] Olga Wichrowska, Niru Maheswaranathan, et al. Learned optimizers that scale and generalize. In *ICML*, 2017.
- [Xiong and Hsieh, 2020] Yuanhao Xiong and Cho-Jui Hsieh. Improved adversarial training via learned optimizer. In *European Conference on Computer Vision*, pages 85–100. Springer, 2020.
- [Xu *et al.*, 2019] Kaidi Xu, Sijia Liu, et al. Structured adversarial attack: Towards general implementation and better interpretability. In *ICLR*, 2019.
- [Yin *et al.*, 2020] Mingzhang Yin, George Tucker, et al. Meta-learning without memorization. In *ICLR*, 2020.
- [Zhao *et al.*, 2019] Pu Zhao, Sijia Liu, et al. On the design of black-box adversarial examples by leveraging gradient-free optimization and operator splitting method. In *ICCV*, 2019.

## Appendix

### A Gradients of MAML loss with respect to RNN parameters

Based on  $\mathbf{G}^{(k)} = \frac{\partial \boldsymbol{\theta}^{(k)}}{\partial \phi} \in \mathbb{R}^{d \times |\phi|}$  and (4), we obtain

$$\mathbf{G}^{(k)} = \mathbf{G}^{(k-1)} - \frac{\partial \Delta \boldsymbol{\theta}^{(k)}}{\partial \phi}. \quad (\text{S1})$$

For ease of presentation, we use  $\text{RNN}_\phi(g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})$  to represent  $\Delta \boldsymbol{\theta}^{(k)} \in \mathbb{R}^d$ , and the RNN output  $\mathbf{h}^{(k)}$  of  $\text{RNN}_\phi$  is omitted when its meaning can clearly be inferred from the context. We then have

$$\frac{\partial \Delta \boldsymbol{\theta}^{(k)}}{\partial \phi} = \frac{\partial \text{RNN}_\phi(g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial \phi} \quad (\text{S2})$$

$$= \frac{\partial \text{RNN}_\phi(g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial g} \cdot \frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \phi} \quad (\text{S3})$$

$$+ \frac{\partial \text{RNN}_\phi(g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}_{k-1}}{\partial \phi} \quad (\text{S4})$$

$$+ \left. \frac{\partial \text{RNN}_\phi(g, h)}{\partial \phi} \right|_{g=g(\boldsymbol{\theta}^{(k-1)}), h=\mathbf{h}_{k-1}}, \quad (\text{S5})$$

where the equality holds by chain rule [Petersen *et al.*, 2008],  $\cdot$  denotes a matrix product that the chain rule obeys, and the term (S5) denotes the derivative w.r.t.  $\phi$  by fixing  $g(\boldsymbol{\theta}^{(k-1)})$  and  $\mathbf{h}_{k-1}$  as constants.

$$\frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \phi} = \frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \boldsymbol{\theta}} \cdot \frac{\partial \boldsymbol{\theta}^{(k-1)}}{\partial \phi} = \frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \boldsymbol{\theta}} \cdot \mathbf{G}^{(k-1)}. \quad (\text{S6})$$

Next, we simplify the term (S4). Let  $\mathbf{H}^{(k)} = \frac{\partial \mathbf{h}_k}{\partial \phi} \in \mathbb{R}^{|\mathbf{h}_k| \times |\phi|}$ . Note that  $\mathbf{h}_k$  depends on  $\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1}$ . So we write  $\mathbf{h}_k = \Pi(\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})$

$$\begin{aligned} \mathbf{H}^{(k)} &= \frac{\partial \mathbf{h}_k}{\partial \phi} = \frac{\partial}{\partial \phi} \Pi(\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1}) \\ &= \frac{\partial \Pi(\phi, g, h)}{\partial \phi} \Big|_{g=g(\boldsymbol{\theta}^{(k-1)}), h=\mathbf{h}_{k-1}} \\ &\quad + \frac{\partial \Pi(\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial g} \cdot \frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \phi} \\ &\quad + \frac{\partial \Pi(\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}_{k-1}}{\partial \phi} \\ &\stackrel{(\text{S6})}{=} \frac{\partial \Pi(\phi, g, h)}{\partial \phi} \Big|_{g=g(\boldsymbol{\theta}^{(k-1)}), h=\mathbf{h}_{k-1}} \\ &\quad + \frac{\partial \Pi(\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial g} \cdot \frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \boldsymbol{\theta}} \mathbf{G}^{(k-1)} \\ &\quad + \frac{\partial \Pi(\phi, g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial \mathbf{h}} \cdot \mathbf{H}^{(k-1)}. \end{aligned} \quad (\text{S7})$$

Substituting (S6) and (S7) into (S2), we can then express (S1) as:

$$\begin{aligned} \mathbf{G}^{(k)} &= \mathbf{G}^{(k-1)} - \frac{\partial \text{RNN}_\phi(g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial g} \cdot \frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \boldsymbol{\theta}} \cdot \mathbf{G}^{(k-1)} \\ &\quad - \frac{\partial \text{RNN}_\phi(g(\boldsymbol{\theta}^{(k-1)}), \mathbf{h}_{k-1})}{\partial \mathbf{h}} \cdot \mathbf{H}^{(k-1)} \\ &\quad - \left. \frac{\partial \text{RNN}_\phi(g, h)}{\partial \phi} \right|_{g=g(\boldsymbol{\theta}^{(k-1)}), h=\mathbf{h}_{k-1}}, \end{aligned} \quad (\text{S8})$$

where  $\mathbf{H}^{(k-1)}$  is determined by the recursion (S7).

It is clear from (S7) and (S8) that the second order derivative would at most be involved due to the presence of  $\frac{\partial g(\boldsymbol{\theta}^{(k-1)})}{\partial \boldsymbol{\theta}}$  if  $g(\boldsymbol{\theta}^{(k-1)})$  is specified by the first-order derivative w.r.t.  $\boldsymbol{\theta}$ . By contrast, if it is specified by the ZO gradient estimate, then there will only be first-order derivatives involved in (S7) and (S8). Lastly, we remark that the recursive forms of (S7) and (S8) facilitate our computation, and  $\mathbf{G}^{(0)} = \mathbf{0}$  and  $\mathbf{H}^{(0)} = 0$ .

## B Proof of Proposition 1

Before showing the theoretical results, we first give the following a blanket of assumptions.

### B.1 Assumptions

In practice, the size of data and variables are limited and the function is also bounded. To proceed, we have the following standard assumptions for quantifying the gradient difference between L2O and LFT.

A1. We assume that gradient estimate is unbiased, i.e.,

$$\left[ \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{D}_i} \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^\cdot) \right] = \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i^{(k)}) := \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}), \forall k, \quad (\text{S9})$$

where  $\mathcal{D}_i$  denotes the training/validation data sample of the  $i$ th task, and  $\boldsymbol{\theta}$  stands for  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N$ .

A2. We assume that the gradient estimate has bounded variance for both  $\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{val}}, \forall i, k$ , i.e.,

$$\mathbb{E}_{\mathcal{D}_i^\cdot} \left[ \left\| \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^\cdot) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \right] \leq \sigma^2, \forall i, k. \quad (\text{S10})$$

The same assumption is also applied for  $\partial \boldsymbol{\theta}_i^{(k)} / \partial \boldsymbol{\phi}$ :

$$\mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \frac{\partial(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^\cdot)}{\partial \boldsymbol{\phi}} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\partial \boldsymbol{\theta}_i^{(k)}}{\partial \boldsymbol{\phi}} := \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\phi}}, \forall i, k, \quad (\text{S11})$$

$$\mathbb{E} \left[ \left\| \frac{\partial(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^\cdot)}{\partial \boldsymbol{\phi}} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\phi}} \right\|^2 \right] \leq \sigma^2, \forall i, k. \quad (\text{S12})$$

A3. We assume that the size of gradient is uniformly upper bounded, i.e.,  $\mathbb{E} \|\nabla f_i(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^\cdot)\| \leq G$ ,  $\mathbb{E} \|\frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^\cdot)}{\partial \boldsymbol{\phi}}\| \leq G, \forall i, k$ .

*Proof.* Assume that A1–A3 hold. Let  $w_k = \frac{1}{K}$ . From the definitions of  $F(\phi)$  and  $\widehat{F}(\phi)$ , we have

$$\begin{aligned} & \|\nabla_{\phi} \widehat{F}(\phi) - \nabla_{\phi} F(\phi)\|^2 \\ &= \left\| \frac{1}{N} \sum_{i=1}^N \frac{1}{K} \sum_{k=1}^K \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \nabla_{\phi} f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}}) - \nabla_{\phi} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \end{aligned} \quad (\text{S13})$$

$$\stackrel{(a)}{\leq} \frac{1}{K} \sum_{k=1}^K \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \frac{1}{N} \sum_{i=1}^N \nabla_{\phi} f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}}) - \nabla_{\phi} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \quad (\text{S14})$$

$$\stackrel{(b)}{\leq} \frac{1}{K} \sum_{k=1}^K \left\| \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \nabla_{\phi} f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}}) - \nabla_{\phi} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \quad (\text{S15})$$

$$\stackrel{(c)}{\leq} \frac{1}{K} \sum_{k=1}^K \left\| \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \frac{\partial f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \nabla_{\phi} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \quad (\text{S16})$$

$$\begin{aligned} & \leq \frac{1}{K} \sum_{k=1}^K \frac{1}{N} \sum_{i=1}^N \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \frac{\partial f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{tr}})}{\partial \phi} \right. \\ & \quad \left. + \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \nabla_{\phi} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \end{aligned} \quad (\text{S17})$$

$$\stackrel{(d)}{\leq} \frac{2G^2}{K} \sum_{k=1}^K \left( \frac{\sigma^2}{D_{\text{val}}} + \frac{\sigma^2}{D_{\text{tr}}} \right) \quad (\text{S18})$$

$$\leq 2G^2 \sigma^2 \left( \frac{1}{D_{\text{val}}} + \frac{1}{D_{\text{tr}}} \right) \quad (\text{S19})$$

where in (a) we use Jensen's inequality, in (b) we apply the triangle inequality, in (c) we use the chain rule, (d) is true because

$$\begin{aligned} & \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \frac{\partial f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} \right\|^2 \\ & \stackrel{(i)}{\leq} \left\| \mathbb{E}_{\mathcal{D}_i^{\text{val}}} \mathbb{E}_{\mathcal{D}_i^{\text{tr}} | \mathcal{D}_i^{\text{val}}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} \right\|^2 \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}}} \mathbb{E}_{\mathcal{D}_i^{\text{val}} | \mathcal{D}_i^{\text{tr}}} \frac{\partial f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}})}{\partial \boldsymbol{\theta}_i^{(k)}} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \right\|^2 \end{aligned} \quad (\text{S20})$$

$$\stackrel{(ii)}{\leq} G^2 \mathbb{E}_{\mathcal{D}_i^{\text{tr}}} \mathbb{E}_{\mathcal{D}_i^{\text{val}} | \mathcal{D}_i^{\text{tr}}} \left\| \frac{\partial f(\boldsymbol{\theta}_i^{(k)}; \mathcal{D}_i^{\text{val}})}{\partial \boldsymbol{\theta}_i^{(k)}} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \right\|^2 \quad (\text{S21})$$

$$\leq \frac{G^2 \sigma^2}{D_{\text{val}}} \quad (\text{S22})$$

where in (i) we use Cauchy-Schwarz inequality, in (ii) we use Jensen's inequality; and similarly we have

$$\begin{aligned} & \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}} \mathcal{D}_i^{\text{val}}} \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \nabla_{\phi} f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)}) \right\|^2 \\ & \leq \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}}} \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \left( \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \phi} \right) \right\|^2 \end{aligned} \quad (\text{S23})$$

$$\leq \left\| \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \boldsymbol{\theta}_i^{(k)}} \right\|^2 \left\| \mathbb{E}_{\mathcal{D}_i^{\text{tr}}} \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \phi} \right\|^2 \quad (\text{S24})$$

$$\leq G^2 \mathbb{E}_{\mathcal{D}_i^{\text{tr}}} \left\| \frac{\partial(\boldsymbol{\theta}_i^{(k)}, \mathcal{D}_i^{\text{tr}})}{\partial \phi} - \frac{\partial f(\boldsymbol{\theta}_1^{(k)}, \dots, \boldsymbol{\theta}_N^{(k)})}{\partial \phi} \right\|^2 \quad (\text{S25})$$

$$\leq \frac{G^2 \sigma^2}{D_{\text{tr}}} \quad (\text{S26})$$

Then, we can have

$$\|\nabla_{\phi} F(\phi) - \nabla_{\phi} \widehat{F}(\phi)\| \leq \sqrt{2}G\sigma \sqrt{\frac{1}{D_{\text{tr}}} + \frac{1}{D_{\text{val}}}} \sim \mathcal{O}\left(\sqrt{\frac{1}{D_{\text{tr}}} + \frac{1}{D_{\text{val}}}}_{\epsilon(D_{\text{tr}}, D_{\text{val}})}\right). \quad (\text{S27})$$

Therefore, it is shown that when  $D_{\text{tr}}$  and  $D_{\text{val}}$  are both large, then the difference between  $\nabla_{\phi} F(\phi)$  and  $\nabla_{\phi} \widehat{F}(\phi)$  are quite small. If we assume that L2O converges to the stationary points in the sense the L2O algorithm is able to find a good solution for equation (9) in the main paper by optimizing  $\phi$ , then it is implied that our approach will also converge to the stationary point up some error, i.e.,  $\epsilon(D_{\text{tr}}, D_{\text{val}})$  which is small. ■

### C $\ell_1$ norm and loss comparison

Figure A1 shows  $\ell_1$  norm of UAP learnt by LFT and other baselines in various training and evaluation scenarios. In general, the  $\ell_1$  norm of the UAP learnt by LFT is smaller than that of using L2O or other baselines, except the case (MNIST, MNIST+CIFAR10). However, we recall from Figure 2 that the ASR obtained by L2O is much poorer than LFT in the case (MNIST, MNIST+CIFAR10). Figure A2 further demonstrates the fine-tuning loss of the UAP learnt by LFT and other baselines. As expected, LFT yields a fast adaptation of UAP to attack unseen test images, corresponding to the lowest fine-tuning loss. If the ASR can hardly reach 100%, the loss does not decrease. We visualize the UAP patterns by LFT, L2O and PGD on MNIST for the same set of  $\mathcal{D}_{\text{val}}^{\text{UAP}}$  images in Figure A3.

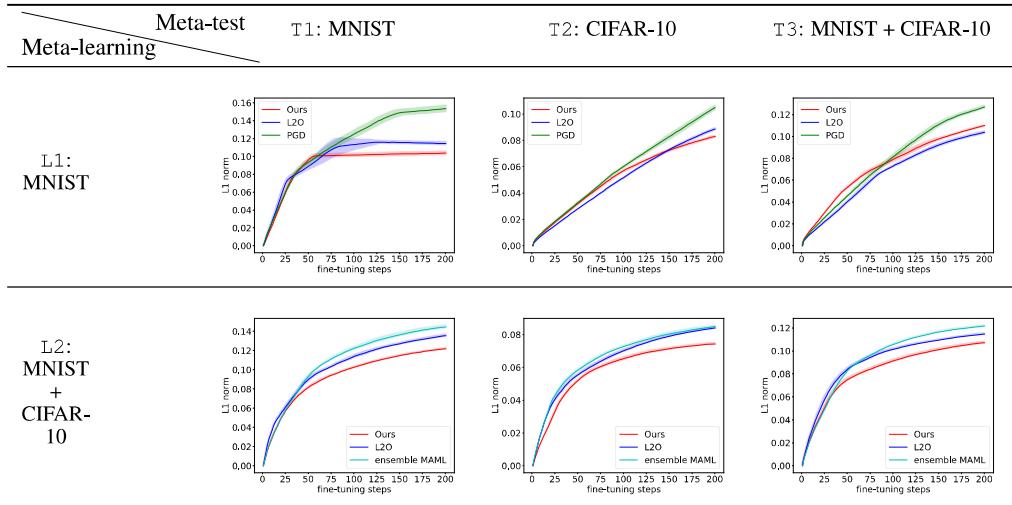


Figure A1: UAP Perturbation strength ( $\ell_1$  norm) by LFT and L2O in various scenarios (*lower is better*). The settings are consistent with Figure 2

### D Effect of the number of tasks in meta-learning

Here we study the effect of the number of tasks  $N$  available for meta-learning on the meta-learning based schemes. For the experiments in §5, we utilize  $N = 1000$  for each meta-learning. We consider LFT as well as ensemble MAML as a baseline (as described in §5) and meta-learn using tasks from both MNIST and CIFAR-10 (case L2 in §5). We consider  $n = 100, 500, 1000$  tasks from each of the image sources and hence for meta-learning each of the MAMLS (by solving (2)) for each image source. For meta-learning the LFT (with Algorithm 1) we will have  $N = 2n = 200, 1000, 2000$  tasks since the same LFT can be meta-learned with tasks from multiple image sources.

We meta-test each of the schemes with 100 UAP generation tasks (obtained as described in §5) from MNIST (case T1 in §5), CIFAR-10 (case T2 in §5) and MNIST+CIFAR (case T3 in §5). We present the ASR at 50 steps and 100 steps, and the number of steps needed to converge to 100% ASR for each value of  $N$  in Table A1. As a reference point for non-meta-learning based schemes, the best average ASR achieved (i) by S-UAP is 63% for T1, 42% for T2 and 53% for T3, and (ii) by PGD is 50% for T1, 25% for T2 and 38% for T3. The results indicate that, for all values of  $N$ , our proposed LFT outperforms our ensemble MAML in terms of ASR (18-50% higher ASR at 50 steps, 5-25% higher ASR at 100 steps) as well as requires significantly smaller number of steps to converge to 100% ASR. Overall, both methods improve the ASR with increasing  $N$ . However, even with just 100 tasks per image source (corresponding to 2 classes  $\times$  2 images per class  $\times$  2 sets for  $\mathcal{D}_i^{\text{tr}} \& \mathcal{D}_i^{\text{val}} \times 100 = 800$  total images per source), LFT is able to achieve 100% in almost all meta-test scenarios, demonstrating that LFT shows benefit

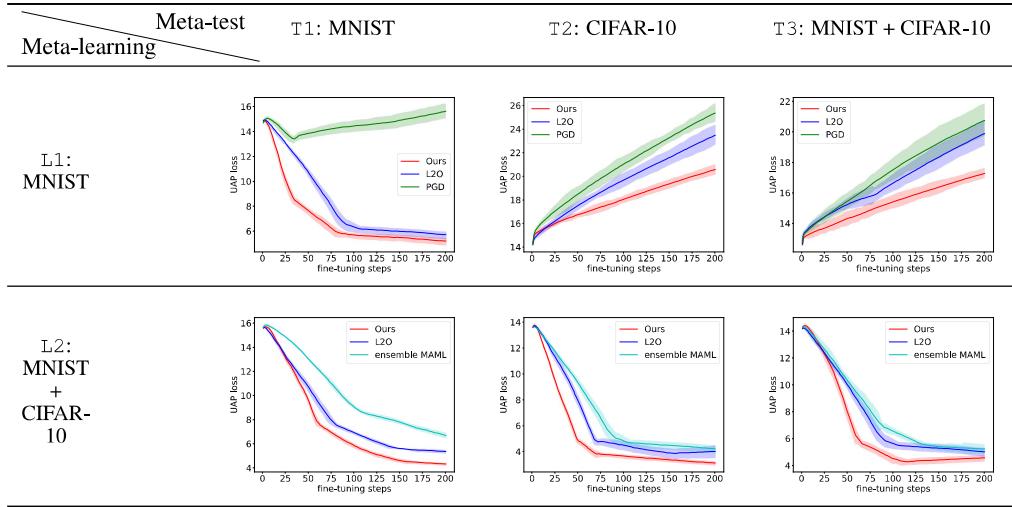


Figure A2: Fine-tuning loss of UAP by LFT and L2O in various scenarios (*lower is better*). The settings are consistent with Figure 2.

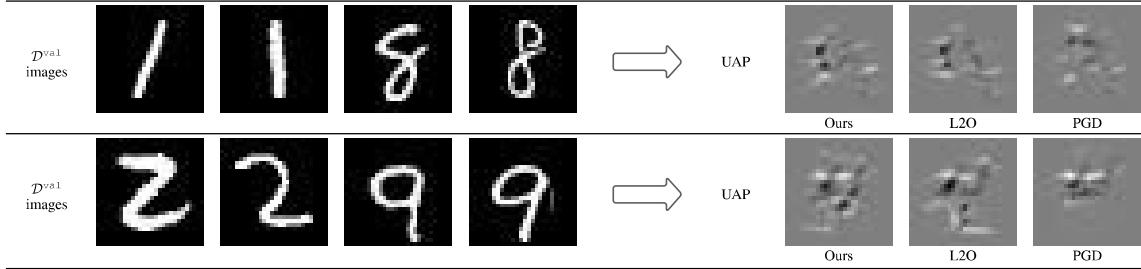


Figure A3: Visualization of UAP patterns by LFT, L2O and PGD on MNIST for the same set of  $\mathcal{D}_{\text{val}}^1$  images.

even when the number of tasks available for meta-learning is relatively small. These numbers also show that, even with a small number of meta-learning tasks, LFT continues to outperform PGD and S-UAP baselines.

Table A1: ASR with the standard deviation of UAPs generated by LFT and ensemble MAML, following the setting of Table 2.  $N$  denotes the number of few-shot tasks for meta-learning. The performance is measured by (i) highest ASR within 50 steps (ASR<sub>50</sub>), (ii) highest ASR within 100 steps (ASR<sub>100</sub>), (iii) step # when first reaching 100% ASR. As a point of comparison, the best average ASR achieved (i) by S-UAP is 63% for T1, 42% for T2 and 53% for T3, and (ii) by PGD is 50% for T1, 25% for T2 and 38% for T3.

Method	Testing	T1: MNIST			T2: CIFAR-10			T3: MNIST + CIFAR-10		
		ASR <sub>50</sub>	ASR <sub>100</sub>	step #	ASR <sub>50</sub>	ASR <sub>100</sub>	step #	ASR <sub>50</sub>	ASR <sub>100</sub>	step #
Ours	$N = 200$	$50 \pm 9$	$100 \pm 2$	98	$68 \pm 21$	$100 \pm 0$	79	$31 \pm 12$	$100 \pm 0$	92
	$N = 1000$	$60 \pm 17$	$100 \pm 0$	92	$96 \pm 7$	$100 \pm 0$	55	$86 \pm 17$	$100 \pm 0$	68
	$N = 2000$	$82 \pm 8$	$100 \pm 0$	75	$92 \pm 11$	$100 \pm 0$	52	$93 \pm 5$	$100 \pm 0$	58
Ensemble MAML	$N = 200$	$0 \pm 0$	$50 \pm 10$	142	$50 \pm 0$	$92 \pm 7$	106	$5 \pm 9$	$82 \pm 17$	120
	$N = 1000$	$32 \pm 15$	$75 \pm 0$	137	$57 \pm 18$	$96 \pm 4$	88	$54 \pm 11$	$95 \pm 6$	103
	$N = 2000$	$50 \pm 13$	$100 \pm 0$	78	$75 \pm 16$	$100 \pm 0$	71	$63 \pm 7$	$100 \pm 0$	73