

---

# Guided Learning of Nonconvex Models through Successive Functional Gradient Optimization

---

Rie Johnson<sup>1</sup> Tong Zhang<sup>2</sup>

## Abstract

This paper presents a framework of successive functional gradient optimization for training nonconvex models such as neural networks, where training is driven by mirror descent in a function space. We provide a theoretical analysis and empirical study of the training method derived from this framework. It is shown that the method leads to better performance than that of standard training techniques.

## 1. Introduction

This paper presents a new framework to train nonconvex models such as neural networks. The goal is to learn a vector-valued function  $f(\theta; x)$  that predicts an output  $y$  from input  $x$ , where  $\theta$  is the model parameter. For example, for  $K$ -class classification where  $y \in \{1, 2, \dots, K\}$ ,  $f(\theta; x)$  is  $K$ -dimensional, and it can be linked to conditional probabilities via the soft-max logistic function. Given a set of training data  $S$ , the standard method for solving this problem is to use stochastic gradient descent (SGD) for finding a parameter that minimizes on  $S$  a loss function  $L(f(\theta; x), y)$  with a regularization term  $R(\theta)$ : 
$$\min_{\theta} \left[ \frac{1}{|S|} \sum_{(x,y) \in S} L(f(\theta; x), y) + R(\theta) \right].$$

In this paper, we consider a new framework that *guides training through successive functional gradient descent* so that training proceeds with alternating the following:

- Generate a guide function so that it is ahead (but not too far ahead) of the current model with respect to the minimization of the loss. This is done by functional gradient descent.
- ‘Push’ the model towards the guide function.

---

<sup>1</sup>RJ Research Consulting, Tarrytown, New York, USA <sup>2</sup>Hong Kong University of Science and Technology, Hong Kong. Correspondence to: Rie Johnson <riejohnson@gmail.com>, Tong Zhang <tongzhang@tongzhang-ml.org>.

Our original motivation was functional gradient learning of additive models in *gradient boosting* (Friedman, 2001). In our framework, essentially, training proceeds with repeating a local search, which limits the searched parameter space to the functional neighborhood of the current parameter at each iteration, instead of searching the entire space at once as the standard method does. This is analogous to  $\epsilon$ -boosting where the use of a very small step-size (for successively expanding the ensemble of weak functions) is known to achieve better generalization (Friedman, 2001).

For measuring the distances between models, we use the Bregman divergence (see e.g., (Bubeck, 2015)) by applying it to the model output. Given a convex function  $h$ , the Bregman divergence  $D_h$  is defined by

$$D_h(u, v) = h(u) - h(v) - \nabla h(v)^\top (u - v). \quad (1)$$

This is the difference between  $h(u)$  and the approximation of  $h(u)$  based on the first-order Taylor expansion around  $v$ . This means that when  $u - v$  is small,

$$D_h(u, v) \approx \frac{1}{2} (u - v)^\top (\mathbf{H}(h(v))) (u - v), \quad (2)$$

where  $\mathbf{H}(h(v))$  denotes the Hessian matrix of  $h$  with respect to  $v$ . Therefore, use of the Bregman divergence has the beneficial effect of utilizing the second-order information.

We show that the parameter update rule of an induced method generalizes that of *distillation* (Hinton et al., 2014). That is, our framework subsumes iterative *self-distillation* as a special case.

Distillation was originally proposed to *transfer knowledge* from a high-performance but cumbersome model to a more manageable model. Various forms of self-distillation, which applies distillation to the models of the same architecture, has been empirically studied (Xu & Liu, 2019; Yang et al., 2019a; Lan et al., 2018; Furlanello et al., 2018; Anil et al., 2018; Zhang et al., 2018; Tarvainen & Valpola, 2017; Yim et al., 2017). One trend is to add to the original scheme, e.g., adding a term to the update rule, data distortion/division, more models for mutual learning, and so forth. However, we are not aware of any work on theoretical understanding such as a convergence analysis of the basic self-learning scheme.

Our theoretical analysis of the proposed framework provides a new functional gradient view of self-distillation, and we show a version of the generalized self-distillation procedure converges to a stationary point of a regularized loss function. Our empirical study shows that the iterative training of the derived method goes through a ‘smooth path’ in a restricted region with good generalization performance. This is in contrast to standard training, where the entire (and therefore much larger) parameter space is directly searched, and thus complexity may not be well controlled.

**Notation**  $\nabla h(v)$  denotes the gradient of a scalar function  $h$  with respect to  $v$ . We omit the subscript of  $\nabla$  when the gradient is with respect to the first argument, e.g., we write  $\nabla f(\theta; x)$  for  $\nabla_{\theta} f(\theta; x)$ .  $H(h(v))$  denotes the Hessian matrix of a scalar function  $h$  with respect to  $v$ . We use  $x$  and  $y$  for input data and output data, respectively. We use  $\langle \cdot \rangle$  to indicate the mean, e.g.,  $\langle F(x, y) \rangle_{(x, y) \in S} = \frac{1}{|S|} \sum_{(x, y) \in S} F(x, y)$ .  $L(u, y)$  is a loss function with  $y$  being the true output. We also let  $L_y(u) = L(u, y)$  when convenient.

## 2. Guided Learning through Successive Functional Gradient Optimization

In this section, after presenting the framework in general terms, we develop concrete algorithms and analyze them.

### 2.1. Framework

We first describe the framework in general terms so that the models to be trained are not limited to parameterized ones. Let  $f$  be the model we are training. Starting from some initial  $f$ , training proceeds by repeating the following:

1. Generate a *guide function*  $f^*$  by applying functional gradient descent for reducing the loss to the current model  $f$ , so that  $f^*$  is an improvement over  $f$  in terms of loss but not too far from  $f$ .
2. Move the model  $f$  in the direction of the guide function  $f^*$  according to some distance measure.

We use the Bregman divergence  $D_h$ , defined in (1), for representing the distances between models.

**Step 1: Guide going ahead** We formulate Step 1 as

$$f^*(x, y) := \arg \min_q [D_h(q, f(x)) + \alpha \nabla L_y(f(x))^\top q], \quad (3)$$

where  $\alpha$  is a meta-parameter. The second term pushes the guide function towards the direction of reducing loss, and the first term pulls back the guide function towards the current model  $f$ . Thus,  $f^*$  is ahead of  $f$  but not too far ahead. Note that we use the knowledge of the true output  $y$  here; therefore,  $f^*$  takes  $y$  as the second argument. The function value for each data point  $(x, y)$  can be found approximately

by solving the optimization problem by SGD if there is no analytical solution. Also, this formulation is equivalent to finding  $f^*$  such that

$$\nabla h(f^*(x, y)) = \nabla h(f(x)) - \alpha \nabla L_y(f(x)). \quad (4)$$

This is mirror descent (see e.g., (Bubeck, 2015)) performed in a function space.

Due to the relation of the Bregman divergence to the Hessian matrix stated in (2), (3) implies that

$$f^*(x, y) \approx f(x) - \alpha (H(h(f(x))))^{-1} \nabla L_y(f(x)). \quad (5)$$

Therefore, if we set  $h(f) = L_y(f)$ , (5) becomes

$$f^*(x, y) \approx f(x) - \alpha (H(L_y(f(x))))^{-1} \nabla L_y(f(x)), \quad (6)$$

which is approximately a second-order functional gradient step (one step of the relaxed Newton method) with step-size  $\alpha$  for minimizing the loss.

If we set  $h(f) = \frac{1}{2} \|f\|^2$ , then the optimization problem (3) has an analytical solution

$$f^*(x, y) = f(x) - \alpha \nabla L_y(f(x)),$$

which is a first-order functional gradient step with step-size  $\alpha$  for minimizing the loss.

**Taking  $m$  steps in Step 1** For further generality, let us also consider  $m$  steps of functional gradient descent by extending  $f^*$  in (3) to  $f_m^*$  recursively defined as follows.

$$\begin{aligned} f_0^*(x, y) &:= f(x) \\ f_{i+1}^*(x, y) &:= \arg \min_q [D_h(q, f_i^*(x)) + \alpha \nabla L_y(f_i^*(x))^\top q]. \end{aligned}$$

Then, in parallel to (5), we have

$$f_{i+1}^*(x, y) \approx f_i^*(x) - \alpha (H(h(f_i^*(x))))^{-1} \nabla L_y(f_i^*(x)).$$

**Step 2: Following the guide** Using the Bregman divergence  $D_h$ , we formulate Step 2 above as an update of the model  $f$  to reduce

$$\langle D_h(f(x), f^*(x, y)) \rangle_{(x, y) \in S} + R(f) \quad (7)$$

so that the model  $f$  approaches the guide function  $f^*$  in terms of the Bregman divergence.  $R(f)$  is a regularization term.

**Parameterization** Although there can be many variations of this scheme, in this work, we parameterize the model  $f$  so that we can train neural networks. Thus, we replace  $f(x)$  by  $f(\theta; x)$  with parameter  $\theta$ . This does not affect Step 1, and to reduce (7) in Step 2, we repeatedly update the model parameter  $\theta$  by descending the stochastic gradient

$$\nabla_{\theta} \left[ \langle D_h(f(\theta; x), f^*(x, y)) \rangle_{(x, y) \in B} + R(\theta) \right], \quad (8)$$

where  $B$  is a mini-batch sampled from a training set  $S$ .

## 2.2. Algorithms

Putting everything together, we obtain Algorithm 1, which performs mirror descent in a function space in Line 3. We call it (and its derivatives) a method of *GUIded Learning through successive Functional gradient optimization (GULF)*. We now instantiate function  $h$  used by the Bregman divergence  $D_h$  to derive concrete algorithms. In general we allow  $h$  to vary for each data point. That is, it may depend on  $(x, y)$ . Here we use two functions discussed above, which correspond to the first-order and the second-order methods, respectively; however, note that choice of  $h$  is not limited to these two.

---

**Algorithm 1** GULF in the most general form. **Input:**  $\theta_0$ , training set  $S$ . Meta-parameters:  $m, \alpha, T$ . **Output:**  $\theta_T$ .

---

```

1:  $\theta \leftarrow \theta_0$ 
2: for  $t = 0$  to  $T - 1$  do
3:   Define  $f_m^*$  by:  $f_0^*(x, y) := f(\theta_t; x)$ ,  $f_{i+1}^*(x, y) :=$ 
      $\arg \min_q [D_h(q, f_i^*(x, y)) + \alpha \nabla L_y(f_i^*(x, y))^\top q]$ 
4:   repeat
5:     Sample a mini-batch  $B$  from  $S$ .
6:     Update  $\theta$  by descending the stochastic gradient
      $\nabla_\theta [\langle D_h(f(\theta; x), f_m^*(x, y)) \rangle_{(x,y) \in B} + R(\theta)]$ 
     for optimizing
      $Q_t(\theta) := \langle D_h(f(\theta; x), f_m^*(x, y)) \rangle_{(x,y) \in S} + R(\theta)$ .
7:   until some criteria are met
8:    $\theta_{t+1} \leftarrow \theta$ 
9: end for

```

---

**GULF1 (1st-order, Algorithm 2)** With  $h(u) = \frac{1}{2} \|u\|^2$ , we obtain Algorithm 2. Derivation is straightforward. This algorithm performs  $m$  steps of first-order functional gradient descent (Line 3) to push the guide function ahead of the current model and then let the model follow the guide by reducing the 2-norm between them.

---

**Algorithm 2** GULF1 ( $h(u) = \frac{1}{2} \|u\|^2$ ): **Input:**  $\theta_0$ , training set  $S$ . Meta-parameters:  $m, \alpha, T$ . **Output:**  $\theta_T$ .

---

```

1:  $\theta \leftarrow \theta_0$ 
2: for  $t = 0$  to  $T - 1$  do
3:   Define  $f_m^*$  by:  $f_0^*(x, y) = f(\theta_t; x)$ ,
      $f_{i+1}^*(x, y) = f_i^*(x, y) - \alpha \nabla L_y(f_i^*(x, y))$ 
4:   repeat
5:     Sample a mini-batch  $B$  from  $S$ .
6:     Update  $\theta$  by descending the stochastic gradient
      $\nabla_\theta [\langle \frac{1}{2} \|f(\theta; x) - f_m^*(x, y)\|^2 \rangle_{(x,y) \in B} + R(\theta)]$ 
7:   until some criteria are met
8:    $\theta_{t+1} \leftarrow \theta$ 
9: end for

```

---

**GULF2 (2nd order, Algorithm 3)** We consider the case of  $h(p) = L_y(p)$  (i.e.,  $h$  returns loss given prediction  $p$ ). (6) has shown that in this case Step 1 becomes approximately the second-order functional gradient descent. Also, with this choice of  $h$ , Algorithm 1 can be converted to a simpler

---

**Algorithm 3** GULF2 ( $h(p) = L_y(p)$ ): **Input:**  $\theta_0$ , training set  $S$ . Meta-parameters:  $\alpha \in (0, 1), T$ . **Output:**  $\theta_T$ . Notation:  $f_\theta = f(\theta; x)$  and  $f_{\theta_t} = f(\theta_t; x)$ .

---

```

 $\theta \leftarrow \theta_0$ 
for  $t = 0$  to  $T - 1$  do
  repeat
    Sample a mini-batch  $B$  from  $S$ .
    Update  $\theta$  by descending the stochastic gradient
     $\nabla_\theta [\langle D_{L_y}(f_\theta, f_{\theta_t}) + \alpha \nabla L_y(f_{\theta_t})^\top f_\theta \rangle_{(x,y) \in B} + R(\theta)]$ 
  until some criteria are met
   $\theta_{t+1} \leftarrow \theta$ 
end for

```

---

form where we do not have to compute the values of the guide function  $f_m^*$  explicitly, and where we have one fewer meta-parameter. This simpler form is shown in Algorithm 3, which has the following relationship to Algorithm 1.

**Proposition 2.1** When  $h(p) = L_y(p)$  that returns loss given prediction  $p$ , Algorithm 1 with  $\alpha = \gamma$  is equivalent to Algorithm 3 with  $\alpha = 1 - (1 - \gamma)^m$ .

The proofs are all provided in the supplementary material.

To simplify notation, let  $f_\theta = f(\theta; x)$ , which is the model that we are updating, and  $f_{\theta_t} = f(\theta_t; x)$ , which is a model that was frozen when time changed from  $t - 1$  to  $t$ . In the stage associated with time  $t$ , Algorithm 3 minimizes

$$\langle D_{L_y}(f_\theta, f_{\theta_t}) + \alpha \nabla L_y(f_{\theta_t})^\top f_\theta \rangle_{(x,y) \in S} + R(\theta) \quad (9)$$

approximately through mini-batch SGD. The second term  $\alpha \nabla L_y(f_{\theta_t})^\top f_\theta$  pushes the model  $f_\theta$  towards the direction of reducing loss, and the first term  $D_{L_y}(f_\theta, f_{\theta_t})$  pulls it back towards the frozen model  $f_{\theta_t}$ . With a certain family of loss functions, (9) can be further transformed as follows.

**Proposition 2.2** Let  $y$  be a vector representation such as a  $K$ -dim vector representing  $K$  classes. Assume that the gradient of the loss function can be expressed as

$$\nabla L(f, y) = \nabla L_y(f) = p(f) - y \quad (10)$$

with  $p(f)$  not depending on  $y$ . Let

$$J_t(\theta) = \langle D_{L_y}(f_\theta, f_{\theta_t}) + \alpha \nabla L_y(f_{\theta_t})^\top f_\theta \rangle_{(x,y) \in S} \quad (11)$$

$$J'_t(\theta) = \langle (1 - \alpha)L(f_\theta, p(f_{\theta_t})) + \alpha L_y(f_\theta) \rangle_{(x,y) \in S} \quad (12)$$

Then we have

$$J_t(\theta) = J'_t(\theta) + c_t,$$

where  $c_t$  is independent of  $\theta$ . This implies that

$$\arg \min_\theta [J_t(\theta) + R(\theta)] = \arg \min_\theta [J'_t(\theta) + R(\theta)].$$

Both the cross-entropy loss and squared loss satisfy (10). In particular, when  $L_y(f)$  is the cross-entropy loss,  $p(f)$  becomes the soft max function. In this case, (12) is the *distillation* formula with the frozen model  $f_{\theta_t}$  playing the role of a cumbersome source model, and therefore, the parameter update rule of Algorithm 3 involving (11) becomes that of distillation. Thus, Algorithm 3 can be regarded as a generalization of self-distillation for arbitrary loss functions.

### 2.3. Convergence Analysis

Let us define  $\alpha$ -regularized loss

$$\ell_\alpha(\theta) := \langle L(f(\theta; x), y) \rangle_{(x,y) \in S} + \frac{1}{\alpha} R(\theta). \quad (13)$$

The following theorem shows that Algorithm 1 with step-size  $\alpha$  always approximately reduces the  $\alpha$ -regularized loss if  $\alpha$  is appropriately set.

**Theorem 2.1** *In the setting of Algorithm 1 with  $m = 1$ , assume that there exists  $\beta > 0$  such that  $D_h(f, f') \geq \beta D_{L_y}(f, f')$  for any  $f$  and  $f'$ , and assume that  $\alpha \in (0, \beta]$ . Assume also that  $Q_t(\theta)$  defined in Algorithm 1 is  $1/\eta$  smooth in  $\theta$ :  $\|\nabla Q_t(\theta) - \nabla Q_t(\theta')\| \leq (1/\eta)\|\theta - \theta'\|$ .*

*Assume that  $\theta_{t+1}$  is an improvement of  $\theta_t$  with respect to minimizing  $Q_t$  so that  $Q_t(\theta_{t+1}) \leq Q_t(\tilde{\theta})$ , where*

$$\tilde{\theta} = \theta_t - \eta \nabla Q_t(\theta_t). \quad (14)$$

*Then we have*

$$\ell_\alpha(\theta_{t+1}) \leq \ell_\alpha(\theta_t) - \frac{\alpha\eta}{2} \|\nabla \ell_\alpha(\theta_t)\|^2.$$

For Algorithm 3, we have  $h(\cdot) = L_y(\cdot)$  and thus  $\beta = 1$ , leading to  $\alpha \in (0, 1]$ . (14) is the parameter update step of Algorithm 1 except that the algorithm stochastically estimates the mean over  $S$  from a mini-batch  $B$  sampled from  $S$ . Therefore, the theorem indicates that each stage (corresponding to  $t$ ) of the algorithm approximately reduces the  $\alpha$ -regularized loss  $\ell_\alpha(\theta)$ . In other words, while the guide function changes from stage to stage, a quantity that does *not* depend on the guide function goes down throughout training, namely, the  $\alpha$ -regularized loss  $\ell_\alpha$ .

Furthermore, we obtain from Theorem 2.1 that

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla \ell_\alpha(\theta_t)\|^2 \leq \frac{2(\ell_\alpha(\theta_0) - \ell_\alpha(\theta_T))}{\alpha\eta T}.$$

Assuming  $\ell_\alpha(\theta) \geq 0$ , this implies that as  $T$  goes to infinity, the right-hand side goes to zero, and so Algorithm 3 converges with  $\nabla \ell_\alpha(\theta_T) \rightarrow 0$ . Therefore, when  $T$  is sufficiently large,  $\theta_T$  finds a stationary point of  $\ell_\alpha$ .

The convergence result indicates that having a regularization term  $R(\theta)$  in the algorithm effectively causes minimization

of the  $\alpha$ -regularized loss. However, our empirical results (shown later) indicate that GULF models are very different from standard models trained directly to minimize the  $\alpha$ -regularized loss. For example, standard models trained with  $\ell_{0.01}$  suffers from severe underfitting, but GULF model with  $\alpha=0.01$  produces high performance. This is because each step of guided learning tries to find a good solution which is near the previous solution (guidance). The complexity of each iterate is better controlled, and hence this approach leads to better generalization performance. We will come back to this point in the next section.

## 3. Empirical study

While the proposed framework is general, our empirical study places a major focus on GULF2 (Algorithm 3) with the cross-entropy loss, due to its connection to distillation (Proposition 2.2). In particular, we set up our implementation so that one instance of GULF2 coincides with self-distillation to provide empirical insight into it from a functional gradient viewpoint.

First, with the goal of understanding the empirical behavior of the algorithm, we examine obtained models in reference to our theoretical findings. We use relatively small neural networks for this purpose. Next, we study the case of larger networks with consideration of practicality.

### 3.1. Implementation

To implement the algorithms presented above, methods of parameter initialization and optimization need to be considered. To observe the basic behavior, our strategy in this work is to keep it as simple as possible.

**Initial parameter  $\theta_0$**  As the functions of interest are nonconvex, the outcome depends on the initial parameter  $\theta_0$ . The most natural (and simplest) choice is random parameters. This option is called ‘ini:random’ below. We also considered two more options. One is to start from a *base model* obtained by regular training, called ‘ini:base’. This option enables study of self-distillation. The other is to start from a shrunk version of the base model, and details of this option will be provided later.

**Parameter update** To update parameter  $\theta$  by descending the stochastic gradient, standard techniques can be used such as momentum, Rmsprop (Tieleman & Hinton, 2012), Adam (Kingma & Ba, 2015), and so forth. As is the case for regular training, learning rate scheduling is beneficial. Among many possibilities, we chose to repeatedly use for each  $t$ , the same method that works well for regular training. For example, a standard method for CIFAR10 is to use momentum and decay the learning rate only a few times, and therefore we use this scheme for each stage on CIFAR10. That is, the learning rate is reset to the initial rate for each  $t$ ; however

**Algorithm 4** base-loop (simplified SGDR): **Input:**  $\theta_0$ , training set  $S$ . Meta-parameter:  $T$ . **Output:**  $\theta_T$ .

**for**  $t = 0$  **to**  $T - 1$  **do**  
 $\theta_{t+1} \leftarrow \arg \min_{\theta} [\langle L_y(f(\theta; x)) \rangle_{(x,y) \in S} + R(\theta)]$   
 where  $\theta$  is initialized by  $\theta_t$ .  
**end for**

	#class	train	dev.	test
CIFAR10	10	49000	1000	10000
CIFAR100	100	49000	1000	10000
SVHN	10	599388	5000	26032
ImageNet	1000	1271167	10000	50000

Table 1. Data. For each dataset, we randomly split the official training set into a training set and a development set to use the development set for meta-parameter tuning. For ImageNet, following custom, we used the official validation set as our ‘test’ set.

note that  $\theta$  is not reset. Although this is perhaps not the best strategy in terms of computational cost, its advantage is that at the end of each stage, we obtain “clean” intermediate models with  $\theta_t$  that were optimized for intermediate goals. (If instead, we used one decay schedule from the beginning to the end, the convergence theorem still holds, but  $\theta_t$  would be noisy when the learning rate is still high.) This strategy enables to study how a model changes as the guide function gradually goes ahead, and also relates the method to self-distillation.

Since  $\theta$  is not reset when the learning rate is reset, this schedule can be regarded as a simplified fixed-schedule version of *SGD with warm restarts (SGDR)* (Loshchilov & Hutter, 2017b). (SGDR instead does sophisticated scheduling with cosine-shape decay and variable epochs.) For comparison, we test the same schedule with the standard optimization objective (‘base-loop’; Algorithm 4).

**Enabling study of self-distillation** We study classification tasks with the standard cross-entropy loss, which satisfies the condition of Proposition 2.2. Combined with the choice of learning rate scheduling above, GULF2 with the ini:base option (which initializes  $\theta_0$  with a trained model) essentially becomes self-distillation. Thus, one aspect of our experiments is to study self-distillation from the viewpoint of functional gradient learning.

### 3.2. Experimental setup

Table 1 summarizes the data we used. As for network architectures, we mainly used ResNet (He et al., 2016a;b) and wide ResNet (WRN) (Zagoruyko & Komodakis, 2016). Following the original work, the regularization term  $R(\theta)$  was set to be  $R(\theta) = \frac{\lambda}{2} \|\theta\|^2$  where  $\lambda$  is the weight decay. We fixed mini-batch size to 128 and used the same learning rate decay schedule for all but ImageNet. Due to the page limit, details are described in the supplementary material. However, note that the schedule we used for all but ImageNet is

3–4 times longer than those used in the original ResNet or WRN study for CIFAR datasets. This is because we used the “train longer” strategy (Loshchilov & Hutter, 2017a), and accordingly, the base model performance visibly improved from the original work. This, in fact, made it harder to obtain large performance gains over the base models (not only for GULF but also for all other tested methods) as the bar was set higher. We feel that this is more realistic testing than using the original shorter schedule.

We applied the standard mean/std normalization to images and used the standard image augmentation. In particular, for ImageNet, we used the same data augmentation scheme as used for training the pre-trained models provided as part of TorchVision, since we used these models as our base model.

The default value of  $\alpha$  is 0.3.

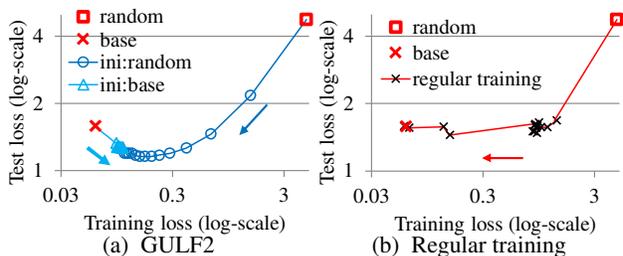


Figure 1. Test loss in relation to training loss. The arrows indicate the direction of time flow. CIFAR100, ResNet-28.

### 3.3. Smooth path

We start with examining training of a relatively small network ResNet-28 (0.4M parameters) on CIFAR100. In this setting, optimization is fast, and so a relatively large  $T$  (the number of stages) is feasible.

We performed GULF2 training with  $T=25$  starting from random parameters (ini:random) as well as starting from a base model obtained by regular training (ini:base). Figure 1a plots test loss of these two runs in relation to training loss. Each point represents a model  $f(\theta_t; x)$  at time  $t = 1, 3, 5, \dots, 25$ , and the arrows indicate the direction of time flow. We observe that training proceeds on a *smooth path*. ini:random( $\circ$ ), which starts from random parameters ( $\square$ ), reduces both training loss and test loss. ini:base( $\triangle$ ) starts from the base model ( $\times$ ) and increases training loss, but reduces test loss. ini:random and ini:base meet and complete one *smooth path* from a random state ( $\square$ ) to the base model ( $\times$ ). ini:random goes forward on this path while ini:base goes backward, and importantly, the path goes through the region where test loss is lower than that of the base model. The test error plotted against training loss also forms a U-shape path. Similar U-shape curves were observed across datasets and network architectures. The supplementary material shows a test error curve and a few more examples of test loss curves including a case of

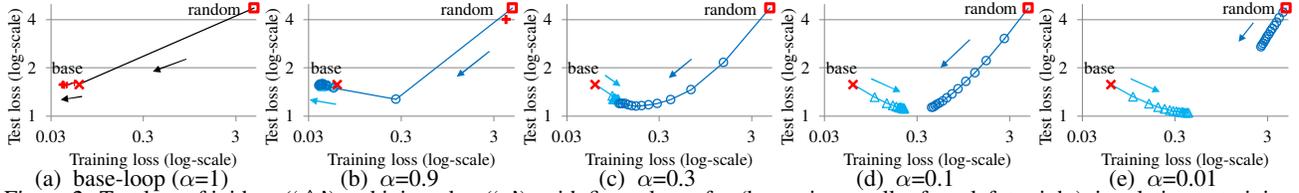


Figure 2. Test loss of ini:base(‘ $\triangle$ ’) and ini:random(‘ $\circ$ ’), with five values of  $\alpha$  (becoming smaller from left to right), in relation to training loss. GULF2.  $T=25$ . CIFAR100. ResNet-28. As  $\alpha$  becomes smaller, the (potential) meeting point shifts further away from the base model. The left-most figure shows base-loop, which is equivalent to  $\alpha=1$ .

DenseNet (Huang et al., 2017).

In the middle of this path, a number of models with good generalization performance lie. One might wonder if regular training also forms such a path. Figure 1b shows that this is not the case. This figure plots the loss of intermediate models in the course of regular training so that the  $i$ -th point represents a model after  $20K \times i$  steps of mini-batch SGD with the learning rate being reduced twice. The path of regular training from random initialization ( $\square$ ) to the final model ( $\times$ ) is rather bumpy and the test loss generally stays as high as the final outcome. The bumpiness is due to the fact that the learning rate is relatively high at the beginning of training. Comparing Figures 1a and 1b, GULF training clearly takes a very different path from regular training.

### 3.4. In relation to the theory

**Going forward, going backward** It might look puzzling why ini:base goes *backward* in the direction of *increasing* the training loss. Theorem 2.1 suggests that this is the effect of the regularization term  $R(\theta)$ , in this case  $R(\theta) = \frac{\lambda}{2} \|\theta\|^2$  with weight decay  $\lambda$ . The theory indicates that for  $\alpha \in (0, 1]$ , the  $\alpha$ -regularized loss

$$\ell_\alpha(\theta) = \langle L_y(f(\theta; x)) \rangle_{(x,y) \in S} + R(\theta)/\alpha$$

goes down and eventually converges as GULF2 proceeds. By contrast, The base model is a result of minimizing

$$\langle L_y(f(\theta; x)) \rangle_{(x,y) \in S} + R(\theta).$$

As we always set  $\alpha < 1$  (0.3 in this case), i.e.,  $1/\alpha > 1$ , GULF2 prefers smaller parameters than the base model does. Consequently, when GULF2 (with small  $\alpha$ ) starts from the base model (which has low training loss and high  $R(\theta)$ ), GULF2 is likely to reduce  $R(\theta)/\alpha$  at the expense of increasing loss (going backward). When GULF2 starts from random parameters, whose training loss is high, GULF2 is likely to reduce loss (going forward) at the expense of increasing  $R(\theta)/\alpha$ .

**Effects of changing  $\alpha$**  With GULF2, the guide function  $f^*$  satisfies

$$f^* \approx f_{\theta_t} - \alpha(\mathcal{H}(L_y(f_{\theta_t})))^{-1} \nabla L_y(f_{\theta_t}),$$

thus,  $\alpha$  serves as a step-size of functional gradient descent for *reducing loss*. The effects of changing  $\alpha$  are shown in Figure 2 with  $T$  fixed to 25. The left-most graph is base-loop, which is equivalent to GULF2 with  $\alpha=1$  in this implementation. There are three things to note. First, with a very small step-size  $\alpha=0.01$  (the right most), ini:random cannot reach far from the random state for  $T=25$ . This is a straightforward effect of a small step size. Second, as step-size  $\alpha$  becomes smaller (from left to right), the (potential) meeting/convergence point shifts further away from the base model; the convergence point of  $\|\theta_t\|^2$  also shifts away from the base model and decreases (supplementary material). This is the effect of larger  $R(\theta)/\alpha$  for smaller  $\alpha$ . Finally, with a large step-size (0.9 and 1), the curve flattens and it no longer goes through the high-performance regions *slowly* or *smoothly*, and the benefit diminishes/vanishes.

**$\alpha$ -regularized loss  $\ell_\alpha(\theta)$**  Figure 3a confirms that, as suggested by the theory,  $\ell_\alpha(\theta)$  goes down and almost converges as training proceeds. This fact motivates examining standard models trained with this  $\ell_\alpha(\theta)$  objective, which we call base- $\lambda/\alpha$  models. We found that base- $\lambda/\alpha$  models do not perform as well as GULF2 at all. In particular, with a very small  $\alpha=0.01$ , which 100 times tightens regularization, test error of base- $\lambda/\alpha$  drastically degrades due to underfitting; in contrast, ini:base with  $\alpha=0.01$  performs well. Moreover, base- $\lambda/\alpha$  models are very different from GULF2 models with corresponding  $\alpha$  even with a moderate  $\alpha$ . For example, Figure 3b plots the parameter size  $\|\theta_t\|^2$  in relation to training loss for  $\alpha=0.3$ . base- $\lambda/\alpha$  is clearly far away from where ini:base and ini:random converge to.

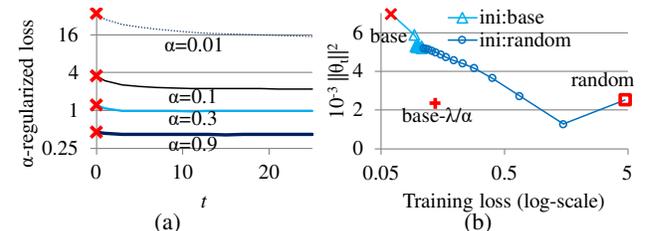


Figure 3. (a)  $\alpha$ -regularized loss  $\ell_\alpha(\theta)$  in relation to time  $t$ . GULF2 ini:base. (b)  $\|\theta_t\|^2$  and training loss of base- $\lambda/\alpha$  in comparison with GULF2.  $\alpha=0.3$ . CIFAR100. ResNet-28.

**Benefit of guiding** This fact illustrates the merit of guided learning (including self-distillation). GULF (indirectly and locally) minimizes the  $\alpha$ -regularized loss  $\ell_\alpha(\theta)$ , but it does

this against the restraining force of *pulling the model back* to the current model. This serves as a form of regularization. Without such a force, training for, say,  $\ell_{0.01}$  would make a big jump to rapidly reduce the parameter size and end up with a radical solution that suffers severely from underfitting. This is what happens with  $\text{base-}\lambda/\alpha$ . By contrast, guided learning finds a more moderate solution with good generalization performance, and this is the benefit of extra regularization (in the form of pulling back) provided through the guide function. The regularization effect of distillation has been mentioned (Hinton et al., 2014), and our framework formalizes the notion through the functional gradient learning viewpoint.

### 3.5. With smaller networks

Now we review the test error results of using relatively small networks in Table 2.  $T$  for GULF2 and base-loop was fixed to 25 on CIFAR10/100 and 15 on SVHN. Step-size  $\alpha$  was fixed to 0.3 for ini:random and chosen from  $\{0.01, 0.03\}$  for ini:base. GULF2 is consistently better than the base model (Row 1) and generally better than the three baseline methods (Row 2–4). The  $\text{base-}\lambda/\alpha$  results (Row 2) were obtained by  $\alpha=0.3$ , and they are generally not much different from the base model. base-loop (Row 3) generally makes small improvement over the base model, but it generally falls short of GULF2. A common technique, label smoothing (Row 4) (Szegedy et al., 2016), ‘softens’ labels by taking a small amount of probability from the correct class and distributing it equally to the incorrect classes. It generally worked well, but the improvements were small. That is, the three baseline methods produced performance gains to some extent, but their gains are relatively small, and they are not as consistent as GULF2 across datasets.

**ini:random** In these experiments, ini:random performed as well as ini:base. This fact cannot be explained from the knowledge-transfer viewpoint of distillation, but it can be explained from our functional gradient learning viewpoint, as in the previous section.

### 3.6. With larger networks

The neural networks and the size of images ( $32 \times 32$ ) used above are relatively small. We now consider computation-

		C10	C100	SVHN		
1	baselines	base model	6.42	30.90	1.86	1.64
2		$\text{base-}\lambda/\alpha$	6.60	30.24	1.78	1.67
3		base-loop	6.20	30.09	1.93	<b>1.53</b>
4		label smooth	6.66	30.52	1.71	1.60
5	GULF2	ini:random	5.91	<b>28.83</b>	1.71	<b>1.53</b>
6		ini:base	<b>5.75</b>	29.12	<b>1.65</b>	1.56

Table 2. Test error (%). Median of 3 runs. Resnet-28 (0.4M parameters) for CIFAR10/100, and WRN-16-4 (2.7M parameters) for SVHN. Two numbers for SVHN are without and with dropout.  $\text{base-}\lambda/\alpha$ : weight decay  $\lambda/\alpha$ . base-loop: Algorithm 4.

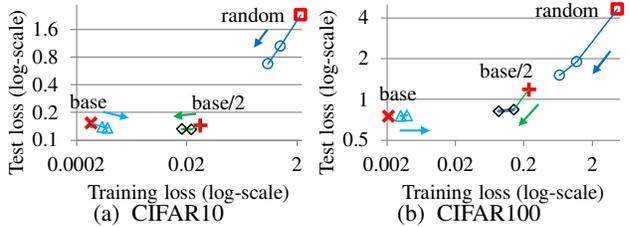


Figure 4. Test loss in relation to training loss. WRN-28-10 on CIFAR10 and CIFAR100. GULF2. ini:base/2 ( $\diamond$ ) fills the gap between ini:random ( $\circ$ ) and ini:base ( $\triangle$ ).

		CIFAR10	CIFAR100
1	base model	3.82	18.55
2	$\text{base-}\lambda/\alpha$	3.70	27.89
3	base-loop	3.70	18.91
4	lab smooth	4.13	19.44
5	GULF1	<b>3.46</b>	18.14
6	GULF2	3.63	<b>17.95</b>

Table 3. Test error (%) results on CIFAR10 and CIFAR100. WRN-28-10 (36.5M parameters) without dropout. Median of 3 runs.

ally more expensive cases.

**Parameter shrinking** ini:random, the most natural option from the functional gradient learning viewpoint, unfortunately, turned out to be too costly in this large-network situation. Moreover, in this setting, it is useful to have an option of starting somewhere between the two end points (‘random’ and ‘base’) since that is where good models tend to lie according to our study with small networks. Therefore, we experimented with ‘rewinding’ a base model by shrinking its weights and bias of the last fully-connected linear layer by dividing them with  $V > 1$  (a meta-parameter). We use these partially-shrunk parameters as the initial parameter  $\theta_0$  for GULF. Since doing so shrinks the model output  $f(\theta_0; x)$  by the factor of  $V$ , this is closely related to *temperature scaling*, for *distillation* (Hinton et al., 2014) and post-training calibration (Guo et al., 2017). Parameter shrinking is, however, simpler than temperature scaling of distillation, which scales logits of both models, and fits well in our framework.

Figure 4 shows training loss (the  $x$ -axis) and test loss (the  $y$ -axis) obtained when parameter shrinking is applied to WRN-28-10 on CIFAR10 and CIFAR100. By shrinking with  $V=2$ , the loss values of the base model change from ‘base’ ( $\times$ ) to ‘base/2’ ( $+$ ). The location of base/2 is roughly the midpoint of two end points ‘base’ and ‘random’. ini:base/2 ( $\diamond$ ), which starts from the shrunk model, explores the space neither ini:random nor ini:base can reach in a few stages.

**Larger ResNets on CIFAR10 and CIFAR100** Table 3 shows test error of ini:base/2 using WRN-28-10 on CIFAR10/100.  $T$  was fixed to 1. Compared with the base model, both GULF1 and 2 consistently improved performance, while the baseline methods mostly failed to make

improvements. GULF1 and 2 produced similar performances. This is the best WRN non-ensemble results on CIFAR10/100 among the self-distillation related studies that we are aware of.

**ImageNet** To test further scale-up, we experimented with ResNet-50 (25.6M parameters) and WRN50-2 (68.9M parameters) on the ILSVRC-2012 ImageNet dataset. As ImageNet training is resource-consuming, we only tested selected configurations, which are GULF2 with ini:base and ini:base/2 options. In these experiments,  $\alpha$  was set to 0.5, but partial results suggested that 0.3 works well too. We used models pre-trained on ImageNet provided as part of TorchVision<sup>1</sup> as the base models. Table 4 shows that GULF2 consistently improves error rates over the base model. The best-performing ini:base/2 achieved lower error rates than a twice deeper counterpart of each network, ResNet-101 for ResNet-50 and WRN-101-2 for WRN-50-2 trained in a standard way (Rows 11–12). Thus, we confirmed that GULF2 scales up and brings performance gains on ImageNet. To our knowledge, this is one of the largest-scale ImageNet experiments among the self-distillation related studies.

methods		Resnet-50		WRN50-2		
1	base model	23.87	7.14	21.53	5.91	
2	base-loop	$t=1$	23.73	6.95	21.99	6.11
3		$t=2$	23.50	6.93	–	–
4		$t=3$	23.36	6.78	–	–
5	ini:base	$t=1$	22.79	6.43	21.17	5.65
6		$t=2$	22.49	6.27	–	–
7		$t=3$	22.31	6.28	–	–
8	ini:base/2	$t=1$	22.50	6.25	<b>20.69</b>	<b>5.35</b>
9		$t=2$	22.31	6.18	–	–
10		$t=3$	<b>22.08</b>	<b>6.10</b>	–	–
11	Resnet-101†	22.63	6.44	–	–	
12	WRN-101-2†	–	–	21.16	5.72	

Table 4. ImageNet 224×224 single-crop results on the validation set. GULF2. top-1 and top-5 errors (%).

† The ResNet-101 and WRN-101-2 performances are from the description of the pre-trained torchvision models.

**Additional experiments on text** Finally, the experiments in this section used image data. Additional experiments using text data are presented in the supplementary material.

## 4. Discussion

**Guided exploration of landscape** GULF is an informed/guided exploration of the loss landscape, where the guidance is successively given as interim goals set in the neighborhood of the model at the time, and such guidance is provided by gradient descent in a function space. Another view of this process is an accumulation of successive greedy optimization. Instead of searching the entire space for the ultimate goal of loss minimization at once, guided learning

proceeds with repeating a local search, which limits the space to be searched and leads to better generalization. Its benefit is analogous to that of  $\epsilon$ -boosting.

**GULF1** GULF2 uses the second-order information of loss in the functional gradient step for generating the guide function, and GULF1 does not. GULF2’s update rule is equivalent to that of distillation, and GULF1’s is not. GULF1 also differs from the logit least square fitting version of distillation. In our experiments (though limited due to our focus on self-distillation study), GULF1 performed as well as GULF2. If this is a general trend, this indicates that inclusion of the second-order information is not particularly helpful. If so, this could be because the second-order information is useful for accelerating optimization, but we would like to proceed slowly to obtain better generalization performance. This motivates further investigation of GULF1 as well as other instantiations of the framework.

**Computational cost** From a practical viewpoint, a shortcoming of the particular setup tested here (but not the general framework of GULF) is computational cost. Since we used the same learning rate scheduling as regular training in each stage, GULF training with  $T$  stages took more than  $T$  times longer than regular training. It is conceivable that training in each stage can be shortened without hurting performance since optimization should be easier as a results of aiming at a nearby goal. Schemes that decay the learning rate throughout the training without restarts or hybrid approaches might also be beneficial for reducing computation. Note that Theorem 2.1 does not require each stage to be performed to the optimum. On the other hand, testing (i.e., making predictions) of the models trained with GULF only requires the same cost as regular models. As shown in the ImageNet experiments, a model trained with GULF could perform better than a much larger (and so slow-to-predict) model; in that case, GULF can save the overall computational cost since the cost for making predictions can be significant for practical purposes.

**Relation to other methods** The proposed method seeks to improve generalization performances in a principled way that limits the searched parameter space. The relation to existing methods for similar purposes is at least two-fold. First, we view that this work gives theoretical insight into related methods such as self-distillation and label smoothing, which we hope can be used to improve them. Second, methods derived from this framework can be used *with* existing techniques that are based on different principles (e.g., weight decay and dropout) for *further* improvements.

**Distillation** Due to the connection discussed above, our theoretical and empirical analyses of GULF2 provide a new functional gradient view of distillation. Here we discuss a few self-distillation studies from this new viewpoint. (Furlanello et al., 2018) showed that iterative self-distillation

<sup>1</sup><https://pytorch.org/docs/stable/torchvision/models.html>

improves performance over the base model. They set  $\alpha$  to 0 (in our terminology) and reported that there were no performance gains on CIFAR10. According to our theory, when  $\alpha$  goes to 0, the quantity reduced throughout the process is not the  $\alpha$ -regularized loss but merely  $R(\theta)$ . Such an extreme setting might be risky. In *deep mutual learning* (Zhang et al., 2018), multiple models are simultaneously trained by reducing loss and aligning each other’s model output. They were surprised by the fact that ‘no prior powerful teacher’ was necessary. This fact can be explained by our functional gradient view by relating their approach to our ini:random. Finally, the regularization effect of distillation has been noticed (Hinton et al., 2014). Our framework formalized the notion through the functional gradient learning viewpoint.

## 5. Conclusion

This paper introduces a new framework for guided learning of nonconvex models through successive functional gradient optimization. A convergence analysis is established for the proposed approach, and it is shown that our framework generalizes the popular self-distillation method. Since the guided learning approach learns nonconvex models in restricted search spaces, we obtain better generalization performance than standard training techniques.

## Acknowledgements

We thank Professor Cun-Hui Zhang for his support of this research.

## References

- Anil, R., Perea, G., Passos, A., Ormandi, R., Dahl, G. E., and Hinton, G. E. Large scale distributed neural network training through online distillation. In *Proceedings of International Conference on Machine Learning (ICML)*, 2018.
- Bubeck, S. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8: 231–358, 2015.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2019.
- Friedman, J. H. Greedy function approximation: a gradient boosting machine. *Ann. Statist.*, 29(5):1189–1232, 2001. ISSN 0090-5364.
- Furlanello, T., Lipton, Z. C., Tsehannen, M., Itti, L., and Anandkumar, A. Born-again neural networks. In *Proceedings of International Conference on Machine Learning (ICML)*, 2018.
- Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. In *Proceedings of International Conference on Machine Learning (ICML)*, 2017.
- He, K., Zhang, X., Ren, S., and Sun, J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016a.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2016b.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. In *Proceedings of Deep Learning and Representation Learning Workshop: NIPS 2014*, 2014.
- Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2017.
- Johnson, R. and Zhang, T. Deep pyramid convolutional neural networks for text categorization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2017.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.
- Lan, X., Zhu, X., and Gong, S. Knowledge distillation by on-the-fly native ensemble. In *Advances in Neural Information Processing Systems 31 (NeurIPS 2018)*, 2018.
- Loshchilov, I. and Hutter, F. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, pp. 1731–1741, 2017a.
- Loshchilov, I. and Hutter, F. SGDR: Stochastic gradient descent with warm restarts. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2017b.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pp. 2818–2826, 2016.

- Tarvainen, A. and Valpola, H. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 2017.
- Tieleman, T. and Hinton, G. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 4, 2012.
- Xie, Q., Dai, Z., Hovy, E., Luong, M.-T., and Le, Q. V. Unsupervised data augmentation for consistency training. *arXiv:1904.12848*, 2019.
- Xu, T.-B. and Liu, C.-L. Data-distortion guided self-distillation for deep neural networks. In *Proceedings of The 33rd AAAI Conference on Artificial Intelligence*, 2019.
- Yang, C., Xie, L., Qiao, S., and Yuille, A. Training deep neural networks in generations: A more tolerant teacher educates better students. In *Proceedings of AAAI 2019*, 2019a.
- Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., and Le, Q. V. XLNet: Generalized autoregressive pre-training for language understanding. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019b.
- Yim, J., Joo, D., Bae, J., and Kim, J. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2016.
- Zhang, X., Zhao, J., and LeCun, Y. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems 28 (NIPS 2015)*, 2015.
- Zhang, Y., Xiang, T., Hospedales, T. M., and Lu, H. Deep mutual learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2018.