

# Mean-Field Limits for Two-Layer Neural Networks Trained with Consensus-Based Optimization

William De Deyn<sup>1,2</sup>, Michael Herty<sup>1,3</sup>, and Giovanni Samaey<sup>2</sup>

<sup>1</sup>Institut für Geometrie und Praktische Mathematik, RWTH Aachen University, Germany

<sup>2</sup>Department of Computer Science, KU Leuven, Belgium

<sup>3</sup>Extraordinary Professor, Department of Mathematics and Applied Mathematics, University of Pretoria, South Africa

November 27, 2025

## Abstract

We study two-layer neural networks and train these with a particle-based method called consensus-based optimization (CBO). We compare the performance of CBO against Adam on two test cases and demonstrate how a hybrid approach, combining CBO with Adam, provides faster convergence than CBO. In the context of multi-task learning, we recast CBO into a formulation that offers less memory overhead. The CBO method allows for a mean-field limit formulation, which we couple with the mean-field limit of the neural network. To this end, we first reformulate CBO within the optimal transport framework. Finally, in the limit of infinitely many particles, we define the corresponding dynamics on the Wasserstein-over-Wasserstein space and show that the variance decreases monotonically.

**Keywords:** optimization, neural networks, mean-field limits, particle methods, optimal transport

## 1 Introduction

Artificial Intelligence has witnessed remarkable progress over the past decades, both in its capabilities and its range of applications. Today, neural networks are present in a variety of fields. One classical application is function approximation, which is supported by the universal approximation theory [34]. In computer vision, convolutional neural networks form the backbone of most modern architectures [39, 38], while the framework of neural ordinary differential equations has contributed significantly to optimal control problems [17, 10]. In natural language processing and speech recognition, recurrent neural networks and the long short-term memory variants have yielded significant performance improvements [33, 51]. More recently, diffusion models have illustrated to be powerful generative models, with applications ranging from image denoising to video generation [56]. Neural networks have even found their way into scientific computing. The most notably example is physics-informed neural networks, which are capable of solving both forward and inverse problems governed by partial differential equations [50].

A neural network can be viewed, in general, as a function parametrized by a set of weights and biases, which we collectively refer to as parameters. A two-layer neural network, for example, is written as

$$\hat{g}(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{c}^\top \sigma(W\mathbf{x} + \mathbf{b}), \quad (1)$$

with  $\mathbf{x} \in \mathbb{R}^d$ ,  $W \in \mathbb{R}^{M \times d}$ ,  $\mathbf{b}, \mathbf{c} \in \mathbb{R}^M$  and  $\sigma(x) : \mathbb{R} \rightarrow \mathbb{R}$  the activation function. Here,  $\boldsymbol{\theta} = (W, \mathbf{b}, \mathbf{c}) \in \mathbb{R}^{d_o}$  represents a vector containing all parameters of the neural network, with  $d_o = M(d+2)$  denoting the optimization dimension. A key factor behind the success of neural networks across diverse applications is their ability to approximate highly complex functions by suitable choice of  $\boldsymbol{\theta}$ . The process of finding the best parameters  $\boldsymbol{\theta}$  is more commonly known as training the neural network. Training is typically

formulated as the minimization of the empirical risk function

$$\hat{R}(\boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \mathcal{L}(y_s, \hat{g}(\mathbf{x}_s; \boldsymbol{\theta})), \quad (2)$$

where  $\mathcal{L}$  represents a loss function and  $S$  the number of data points [54]. The loss function denotes the discrepancy between a true sample  $y_s$  and the model prediction made from the input sample  $\mathbf{x}_s$ . Currently, the standard method to find the minimizer of (2) is stochastic gradient descent (SGD) or adaptive variants, such as Adam [11, 36]. The optimization of the risk function is a difficult task [30, 29], primarily because the objective function is in general nonconvex. Gradient-based methods, such as SGD, are susceptible to becoming trapped in local minima. Other well-known difficulties include the vanishing gradient and exploding gradient phenomena [45, 5]. One alternative approach is to apply particle-based methods, such as Consensus-Based Optimization (CBO). The CBO method is a gradient free, global optimization method designed for high-dimensional, nonconvex objective functions [48, 12, 14, 28, 13, 35, 37, 25, 26, 27, 8, 15, 7, 31, 32]. The CBO method allows for the passage to the mean-field limit as the number of particles tends to infinity.

In this paper, we explore the feasibility of training neural networks using Consensus-Based Optimization and compare its performance against the popular Adam method. Owing to the natural mean-field formulation of CBO, we also study its mean-field limit. We consider the mean-field limit of two-layer neural networks, similar to the works [41, 43, 53, 52, 2, 19]. In the same spirit, we aim to derive a time-discrete mean-field limit governing the CBO dynamics on infinite wide neural networks.

The remainder of this paper is organized as follows. Section 2 reviews the Barron space, which is a function space containing functions that can be efficiently represented by two-layer neural networks. Next, in Section 3, we review the general training problem in machine learning and shortly discuss Adam and Consensus-Based Optimization. We also present two variants of CBO, such as a hybrid method and Multi-Task CBO. In Section 4, we cover numerical experiments comparing the performance of CBO to Adam in training neural networks. Additionally, we provide an experiment illustrating the capabilities of Multi-Task CBO. Section 5 establishes the mean-field limits in two regimes: first, as the number of hidden neurons tends to infinity; and second as the number of optimization particles tends to infinity. Finally, Section 6 summarizes the main findings and presents the conclusion.

## 2 The Barron Space

In this section, we introduce the Barron space, which is a function space suited for approximation by two-layer neural networks. First, we rewrite and scale the neural network formulation in (1) by considering the two-layer neural network as a function  $\hat{g}_M(\mathbf{x}; \boldsymbol{\theta}) : X \rightarrow \mathbb{R}$  of the form

$$\hat{g}_M(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{M} \sum_{m=1}^M c_m \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m), \quad (3)$$

with  $\mathbf{x} \in X \subseteq \mathbb{R}^d$ ,  $\mathbf{w}_m \in \mathbb{R}^d$ ,  $b_m, c_m \in \mathbb{R}$  [24]. The parameter  $M \in \mathbb{N}$  denotes the number of neurons in the hidden layer, i.e., the width of the two-layer neural network. The neural network is parameterized by weights  $w_{m,j}$ ,  $c_m$  and biases  $b_m$ , which we collect into the parameter vector  $\boldsymbol{\theta} := \{w_{m,j}, b_m, c_m\}_{m=1}^M$ ,  $j = 1, \dots, d$ . For an empirical measure  $\hat{\mu} \in \mathcal{P}_2(\mathbb{R}^{d+2})$  defined by

$$\hat{\mu}(\mathbf{w}, b, c) = \frac{1}{M} \sum_{m=1}^M \delta(\mathbf{w} - \mathbf{w}_m) \delta(b - b_m) \delta(c - c_m), \quad (4)$$

the neural network in (3) can be rewritten as a function  $\hat{g}_M(\mathbf{x}) : X \rightarrow \mathbb{R}$

$$\hat{g}_M(\mathbf{x}) = \int_{\Omega} c \sigma(\mathbf{w}^\top \mathbf{x} + b) d\hat{\mu}(\mathbf{w}, b, c), \quad \Omega = \mathbb{R}^d \times \mathbb{R} \times \mathbb{R}, \quad (5)$$

and where  $\Omega$  represents our parameter space. We restrict the representing measures to  $\mathcal{P}_2(\mathbb{R}^{d+2})$ , following [20]. In the mean-field regime of neural networks, we consider the limit  $M \rightarrow \infty$ . For an arbitrary probability measure  $\mu \in \mathcal{P}_2(\mathbb{R}^{d+2})$ , the mean-field limit yields a representation of a neural network as

$$g(\mathbf{x}) = \int_{\Omega} c\sigma(\mathbf{w}^\top \mathbf{x} + b) d\mu(\mathbf{w}, b, c) = \mathbb{E}_{\mu} [c\sigma(\mathbf{w}^\top \mathbf{x} + b)]. \quad (6)$$

For functions of the form (6) with a RELU activation function, the Barron norm is defined as

$$\|g\|_{\mathcal{B}} := \inf_{\mu \in \mathcal{P}_2(\mathbb{R}^{d+2})} \max_{(\mathbf{w}, b, c) \in \text{supp}(\mu)} |c| (\|\mathbf{w}\|_1 + |b|). \quad (7)$$

Functions with a finite Barron norm (7) form the Barron space [3, 24]. The Barron norm can be bounded, with  $G$  denoting the Fourier transformation of  $g$ , as

$$\|g\|_{\mathcal{B}} \leq 2 \inf_G \int_{\mathbb{R}^d} \|\omega\|_1^2 |G(\omega)| d\omega + 2\|\nabla g(0)\|_1 + 2|g(0)|, \quad (8)$$

which shows that sufficiently regular functions are in the Barron space. Two-layer neural networks  $\hat{g}_M$  can approximate functions  $g$  in the Barron space up to an arbitrary precision, where the approximation error satisfies

$$\|g(\cdot) - \hat{g}_M(\cdot; \boldsymbol{\theta})\|^2 \leq \frac{3\|g\|_{\mathcal{B}}^2}{M}. \quad (9)$$

With the measure-based viewpoint of neural networks, introduced in Eq. (6), we can consider continuous optimization schemes and analyse the training of infinite wide neural networks. Section 5 elaborates further on a particle-based method for optimizing infinitely wide neural networks.

### 3 Optimization Methods

In this section, we formulate the general optimization problem underlying neural network training. Next, we present the Adam optimizer and the Consensus-Based Optimization method. Lastly, we introduce two variants of the CBO method, namely a hybrid method that combines Adam and CBO and a Multi-Task CBO formulation.

#### 3.1 General Training Problem

The central goal in machine learning is to approximate an unknown function  $g : X \rightarrow Y$  based on observed data. We assume that  $g$  can be well represented by two-layer neural networks with parameters  $\boldsymbol{\theta}$ . Training the neural network then amounts to finding parameters  $\boldsymbol{\theta}$  such that neural network  $\hat{g}(\mathbf{x}; \boldsymbol{\theta})$  approximates the unknown function as accurately as possible.

The quality of the approximation is typically measured with a loss function  $\mathcal{L}$ , which quantifies the discrepancy between the neural network prediction  $\hat{\mathbf{y}}_s = \hat{g}(\mathbf{x}_s; \boldsymbol{\theta})$  and the observed label  $\mathbf{y}_s$ . From a theoretical perspective, training a neural network can be understood as minimizing the risk function  $R(\boldsymbol{\theta})$ . The risk function is defined as the expected loss with respect to a data distribution  $P(\mathbf{x}, \mathbf{y})$ :

$$R(\boldsymbol{\theta}) = \int_{X \times Y} \mathcal{L}(\mathbf{y}, \hat{g}(\mathbf{x}; \boldsymbol{\theta})) dP(\mathbf{x}, \mathbf{y}). \quad (10)$$

In practice, the true data distribution is not known. Therefore, the risk function  $R(\boldsymbol{\theta})$  is approximated by the empirical risk

$$\hat{R}(\boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \mathcal{L}(\mathbf{y}_s, \hat{g}(\mathbf{x}_s; \boldsymbol{\theta})), \quad (11)$$

where  $\{(\mathbf{x}_s, \mathbf{y}_s)\}_{s=1}^S$  denotes the training dataset [54].

The choice of  $\mathcal{L}$  depends on the particular task. For regression, where  $Y = \mathbb{R}$ , the squared error loss

$$\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2, \quad (12)$$

leads to the Mean Squared Error (MSE) risk [30]. For classification with  $C$  classes, the standard choice is the cross-entropy loss function

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{c=1}^C y_c \log(\hat{y}_c), \quad (13)$$

where  $\mathbf{y}$  is a one-hot encoded label and  $\hat{\mathbf{y}}$  a probability vector obtained after applying the softmax function to the output of the neural network [42]. There exist many other loss functions in machine learning, such as hinge loss, Huber loss and logistic loss.

Training a neural network can thus be formulated as an optimization problem, where we aim to find

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \hat{R}(\boldsymbol{\theta}), \quad (14)$$

with the empirical risk  $\hat{R}(\boldsymbol{\theta})$  defined as in Eq. (11). The empirical risk function is typically nonconvex due to the nonlinear structure of the neural networks [30]. Standard optimization methods in machine learning rely on the gradient of the empirical risk, with Adam being the most popular [36]. However, gradient-based methods tend to converge to local minima in nonconvex landscapes. Particle-based optimization methods, such as Consensus-Based Optimization, are known as global optimization methods, for which convergence to the global minimizer of certain nonconvex functions is theoretically guaranteed.

### 3.2 Gradient-Based Optimization

Nearly all neural networks are trained with gradient-based methods. The most popular are by far Stochastic Gradient Descent (SGD) and adaptive variants such as Adam. In what follows, we provide a brief overview of these two methods, which form the foundation of modern neural network training.

Stochastic Gradient Descent is an adaptation of the classical Gradient Descent method [30]. Gradient Descent is an iterative method that updates the model parameters in the direction of the negative gradient of the objective function [44]. While classical Gradient Descent often employs a line search to determine an optimal step size, it is common in neural network training to use a fixed step size for simplicity. At iteration  $k$ , the step direction  $\mathbf{d}^k$  is given by

$$\mathbf{d}^k = \frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{y}_s, \hat{g}(\mathbf{x}_s; \boldsymbol{\theta}^k)), \quad (15)$$

and the parameters are updated according to

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \Delta t \mathbf{d}^k, \quad (16)$$

where  $\Delta t > 0$  denotes the time step or the learning rate. A major drawback of Gradient Descent is that computing the step direction at each iteration scales linearly with the dataset size,  $\mathcal{O}(S)$ . Datasets in machine learning can easily contain up to one million data points, making Gradient Descent computationally expensive.

SGD alleviates this issue by considering a minibatch of data points  $\{(\mathbf{x}_s, \mathbf{y}_s)\}_{s=1}^{S'}$  drawn uniformly from the dataset. The minibatch size  $S' \ll S$  is chosen before training. The step direction  $\mathbf{d}^k$  in SGD equals

$$\mathbf{d}^k = \frac{1}{S'} \sum_{s=1}^{S'} \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{y}_s, \hat{g}(\mathbf{x}_s; \boldsymbol{\theta}^k)), \quad (17)$$

which is an unbiased estimator of the full gradient (15). At iteration  $k$ , the gradient is computed using the current minibatch, and the parameters are subsequently updated. At the next iteration  $k+1$ , a new minibatch is sampled. One pass through the complete training dataset, where each data point has been used once for updating the parameters, is referred to as an epoch.

The performance of SGD depends heavily on the choice of the learning rate  $\Delta t$ . There is a trade-off: a larger learning rate yields faster progress, but is more unstable; a smaller learning rate improves stability but converges slower. In practice, the learning rate decays during training according to a predefined schedule. However, there also exist adaptive algorithms, such as Adam, that adapt the learning rate individually for each parameter based on the gradient history. Given the stochastic gradient  $\mathbf{d}^k$  from Eq. (17), Adam estimates the first and second moment as

$$\mathbf{s}^{k+1} = \beta_1 \mathbf{s}^k + (1 - \beta_1) \mathbf{d}^k \quad \mathbf{r}^{k+1} = \beta_2 \mathbf{r}^k + (1 - \beta_2) \mathbf{d}^k \odot \mathbf{d}^k, \quad (18)$$

with the decay parameters  $\beta_1, \beta_2 \in [0, 1)$  and where  $\odot$  represents the elementwise multiplication. The moment estimates are normalized to correct for the initial bias:

$$\hat{\mathbf{s}}^{k+1} = \frac{\mathbf{s}^{k+1}}{1 - \beta_1^{k+1}} \quad \hat{\mathbf{r}}^{k+1} = \frac{\mathbf{r}^{k+1}}{1 - \beta_2^{k+1}}. \quad (19)$$

Finally, Adam updates the parameters  $\boldsymbol{\theta}$  as follows

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \Delta t \frac{\hat{\mathbf{s}}^{k+1}}{\sqrt{\hat{\mathbf{r}}^{k+1} + \delta}}, \quad (20)$$

with  $\delta$  a small constant for numerical stability. The initial values for the first and second moment estimates  $\mathbf{s}^0$  and  $\mathbf{r}^0$  are set to zero. By adapting learning rates per parameter, Adam typically achieves faster convergence than plain SGD.

### 3.3 Consensus-Based Optimization

We introduce Consensus-Based Optimization, a global optimization method well suited for nonconvex, nonsmooth objective functions [48]. We aim to minimize the empirical risk (11). To find the minimum, we consider an ensemble of  $N$  particles  $\boldsymbol{\theta}_n^k \in \mathbb{R}^{M(d+2)}$ ,  $n = 1, \dots, N$  with  $N \in \mathbb{N}$  at time step  $k$ . The state of the particles evolves according to the discretized stochastic differential equation

$$\boldsymbol{\theta}_n^{k+1} = \boldsymbol{\theta}_n^k - \lambda \Delta t (\boldsymbol{\theta}_n^k - \mathbf{V}^k) + \tilde{\sigma} \sqrt{\Delta t} (\boldsymbol{\theta}_n^k - \mathbf{V}^k) \odot \boldsymbol{\xi}_n^k, \quad \boldsymbol{\xi}_n^k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (21)$$

with  $\lambda, \tilde{\sigma} > 0$  representing the drift and diffusion parameters respectively. The dynamics combine a drift towards the consensus point with a diffusion term that promotes exploration. The consensus point  $\mathbf{V} \in \mathbb{R}^{M(d+2)}$  is calculated as a weighted average

$$\mathbf{V}^k = \frac{\frac{1}{N} \sum_{n=1}^N \boldsymbol{\theta}_n^k \exp(-\alpha \hat{R}(\boldsymbol{\theta}_n^k))}{\frac{1}{N} \sum_{n=1}^N \exp(-\alpha \hat{R}(\boldsymbol{\theta}_n^k))} = \sum_{n=1}^N \beta(\boldsymbol{\theta}_n^k) \boldsymbol{\theta}_n^k, \quad \beta(\boldsymbol{\theta}_n^k) = \frac{\exp(-\alpha \hat{R}(\boldsymbol{\theta}_n^k))}{\sum_{n=1}^N \exp(-\alpha \hat{R}(\boldsymbol{\theta}_n^k))}, \quad (22)$$

where the weight  $\beta(\boldsymbol{\theta}_n^k)$  depends on the relative performance of the particle in the optimization landscape. In statistical physics, the parameter  $\alpha > 0$  represents the inverse temperature. The initial positions of the particles  $\boldsymbol{\theta}_n^0$  are distributed independently and identically according to a chosen initial distribution  $\rho^0$ .

The computation of the consensus point  $\mathbf{V}^k$  requires evaluating the empirical risk for each particle  $\boldsymbol{\theta}_n^k$ . This process can become computationally demanding when the number of particles or the size of the training set is large. To reduce the cost, a minibatch strategy analogous to that used in SGD can be applied within CBO. Let  $\{(\mathbf{x}_s, \mathbf{y}_s)\}_{s=1}^{S'}$ ,  $S' \ll S$ , denote a minibatch of training data drawn uniformly from the training dataset. The corresponding minibatch empirical risk is then defined as

$$\hat{R}(\boldsymbol{\theta}_n^k) = \frac{1}{S'} \sum_{s=1}^{S'} \mathcal{L}(\mathbf{y}_s, \hat{g}(\mathbf{x}_s; \boldsymbol{\theta}_n^k)). \quad (23)$$

At iteration  $k$ , the consensus point is computed using the current minibatch, and the particle positions are subsequently updated. At the next iteration  $k+1$ , a new minibatch is sampled. One pass through

the complete training dataset, where each data point has been used once for updating the particles, is referred to as an epoch.

The consensus-based optimization scheme converges under certain conditions. In [14], the authors show that for high dimensional problems and anisotropic noise, we need  $2\lambda > \tilde{\sigma}^2$  to form consensus, which is independent of the optimization dimension. Secondly, we note that the consensus point in (22) is a finite particle discretization of the mean of a Gibbs-type measure. By the Laplace principle [23]

$$\lim_{\alpha \rightarrow \infty} \left( -\frac{1}{\alpha} \log \left( \int \exp \left( -\alpha \hat{R}(\boldsymbol{\theta}) \right) \rho(\boldsymbol{\theta}) d\boldsymbol{\theta} \right) \right) = \min_{\boldsymbol{\theta}} \hat{R}(\boldsymbol{\theta}), \quad (24)$$

the Gibbs measure concentrates exponentially around the global minimiser of the objective function. Consequently, the consensus point will converge towards the global minimiser, provided that  $\boldsymbol{\theta}^* \in \text{supp}(\rho^0)$  [12]. The particle-based nature of CBO enables the formulation of a mean-field limit for  $N \rightarrow \infty$ . For analytical results on the mean-field limit, we refer to [12, 28, 35, 37, 25].

### 3.4 Variants

#### 3.4.1 Hybrid Method

A natural approach to develop a new method is to combine two existing methods. We devise a hybrid method where we combine an Adam step and a CBO step together, similar to [55]. Formally, the update is given by

$$\boldsymbol{\theta}_n^{k+1} = \boldsymbol{\theta}_n^k - \gamma \Delta t \frac{\hat{\mathbf{s}}_n^{k+1}}{\sqrt{\hat{\mathbf{r}}_n^{k+1}} + \delta} + (1 - \gamma) \left( -\lambda \Delta t \left( \boldsymbol{\theta}_n^k - \mathbf{V}^k \right) + \tilde{\sigma} \sqrt{\Delta t} \left( \boldsymbol{\theta}_n^k - \mathbf{V}^k \right) \odot \boldsymbol{\xi}_n^k \right), \quad (25)$$

where  $\gamma$  is a parameter between 0 and 1 and controls how much influence CBO or Adam have on the update. The first and second moment estimates are computed as in (18) and (19), with the gradient for each particle given by

$$\mathbf{d}_n^k = \frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\theta}} \mathcal{L} \left( \mathbf{y}_s, \hat{g}(\mathbf{x}_s; \boldsymbol{\theta}_n^k) \right). \quad (26)$$

In the hybrid method, the minibatch strategy involves using the same minibatch to compute the gradient and the consensus point, as explained in Subsections 3.2 and 3.3.

#### 3.4.2 Multi-Task CBO

Multi-task learning is another area of machine learning where the CBO method can be naturally applied to. Multi-task learning is a multi-objective optimization problem with the goal to improve the generalization of a model [16]. Instead of building a separate model for each task, one model shares parts of its internal representation. The knowledge gained from one task can improve the performance of other tasks. In practice, the multi-task risk is often of the form  $\hat{R} = \sum_p \kappa_p \hat{R}(\mathcal{T}_p)$ , with  $\hat{R}(\mathcal{T}_p)$  the empirical risk associated to task or problem  $\mathcal{T}_p$  and a weight  $\kappa_p$ . Currently, multi-task models are trained with a gradient method. However, the gradient methods have several drawbacks. One key aspect is that the gradients have to be balanced by updating the weights  $\kappa_p$  to avoid one specific task from dominating the training of the model [18]. Another issue is gradients from different tasks conflicting with one another [57]. This can be alleviated by projecting one of the gradients of a task onto another gradient.

The CBO method allows for a natural and memory-efficient implementation of the multi-task training [7]. Figure 1 provides a conceptual illustration. Consider the training of a single neural network with CBO to approximate a given function. A requirement for convergence to the global minimizer is that  $\boldsymbol{\theta}^* \in \text{supp}(\rho^0)$ . Accordingly, we initialize particles by sampling from the distribution  $\rho^0$ .

Now consider a second, related approximation problem to the first. If the task is sufficiently similar, we expect the corresponding global minimizer to lie close to the first minimizer, in particular, within the support of  $\rho^0$ , as illustrated in Figure 1. In this case, there is no need to sample additional new

particles for the second task; the particles used for the first task can be reused for the second task. More generally, the particles can be recycled across multiple tasks, provided that the global minimizers of these tasks remain within  $\text{supp}(\rho^0)$ .

A practical strategy for choosing the number of particles is to set  $N$  equal to the number of tasks. In that case, an increasing number of tasks results in improved optimization accuracy, since each empirical risk function  $\hat{R}(\mathcal{T}_p)$  is optimized with a larger population of particles. In the multi-task context, CBO does not need a particle ensemble for each task, reducing the memory overhead. However, a larger particle count than the number of tasks is also feasible in practice.

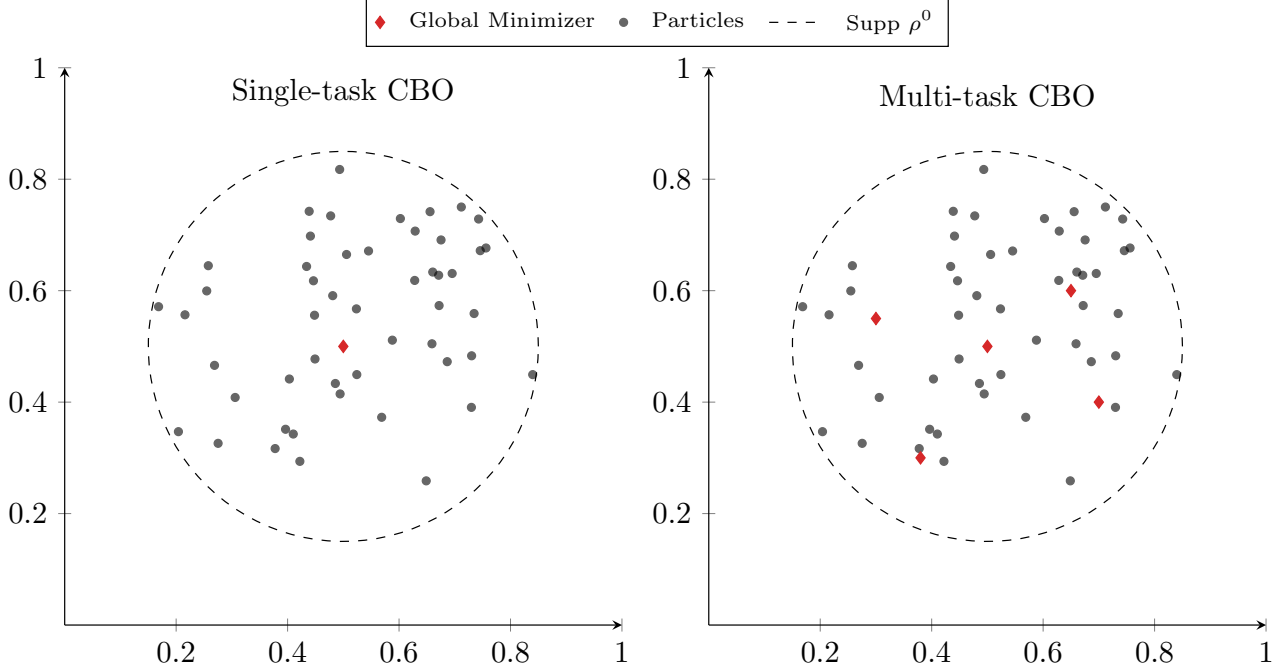


Figure 1: Conceptual illustration of single-task versus Multi-Task CBO. In single-task CBO a single consensus point guides all particles toward the global minimizer of one empirical risk. In Multi-Task CBO the same particle ensemble is recycled across related tasks whose minimizers lie within the support of the common initialization  $\rho^0$ .

We can reformulate the CBO method in the multi-task setting as

$$\begin{cases} \boldsymbol{\theta}_n^{k+1} = \boldsymbol{\theta}_n^k - \lambda \Delta t (\boldsymbol{\theta}_n^k - \mathbf{V}(\mathcal{T}_p^k)) + \tilde{\sigma} \sqrt{\Delta t} (\boldsymbol{\theta}_n^k - \mathbf{V}(\mathcal{T}_p^k)) \odot \boldsymbol{\xi}_n^k \\ \mathcal{T}_p^{k+1} = \mathcal{T}_p^k \end{cases}, \quad n = p = 1, \dots, N, \quad (27)$$

where the number of consensus points equals the number of tasks. The consensus point is given by

$$\mathbf{V}(\mathcal{T}_p^k) = \frac{\frac{1}{N} \sum_{n=1}^N \boldsymbol{\theta}_n^k \exp(-\alpha \hat{R}(\boldsymbol{\theta}_n^k, \mathcal{T}_p^k))}{\frac{1}{N} \sum_{n=1}^N \exp(-\alpha \hat{R}(\boldsymbol{\theta}_n^k, \mathcal{T}_p^k))}. \quad (28)$$

## 4 Numerical Examples

In this section, we present three numerical experiments designed to illustrate and evaluate the optimization methods discussed in Section 3. The implementation is publicly available online [22].

### 4.1 Example 1: Sine Approximation

The goal of the first experiment is to investigate the applicability of CBO to train two-layer neural networks. To this end, we compare CBO with Adam, focusing on both the minimal empirical risk obtained and the stability of the training procedure.

The experiment setup is the following: consider a one-dimensional regression problem, where the goal is to approximate the function  $x \mapsto \sin(2\pi x)$  on the domain  $[0, 1]$  using a finite two-layer neural network, as in Eq. (3). The neural network has a RELU activation function. We take the width of the neural network  $M = 100$ . We sample 8000 data points  $x_s$  uniformly on the interval  $[0, 1]$  and apply the function  $y_s = \sin(2\pi x_s) + 0.01\xi_s$ , with  $\xi_s \sim \mathcal{N}(0, 1)$ , to generate the training dataset  $\{(x_s, y_s)\}_{s=1}^{8000}$ . We apply the minibatch strategy to both Adam and CBO using a minibatch size of  $S'$ . We choose the squared error loss, as in (12), for the loss function. We initialize the particles of the CBO method from the uniform distribution  $\theta_n^0 \sim \mathcal{U}[-1, 1]$ . For Adam, we keep the default initialization of the neural network parameters provided by PyTorch [46].

The training of the two-layer neural network with CBO is carried out with the following parameters:

$$N = 200, \quad \Delta t = 0.1, \quad \alpha = 10^5, \quad \lambda = 1, \quad \tilde{\sigma} = \sqrt{1.6}, \quad S' = 800.$$

Parameters shared between CBO and Adam, such as the minibatch size  $S'$  and the learning rate (or equivalently time step)  $\Delta t$ , are always taken equal. Every 100 epochs, the parameter  $\alpha$  is multiplied by 10 until it reaches  $10^7$ . Lastly, we run the experiment ten times with different seeds and present the median of the results.

Figure 2 shows the empirical risk per epoch, comparing the optimization performance of CBO and Adam. In Figure 3, we present the approximation obtained with a two-layer neural network trained with CBO and Adam.

In Figure 2, we observe that the CBO method converges to a lower empirical risk than Adam. The CBO method also exhibits better stability than Adam. Figure 3 confirms that both CBO and Adam find a good approximation of the sine function. We do note, however, that the computational cost per iteration is higher for CBO, since a forward pass through the neural network is required for each particle.

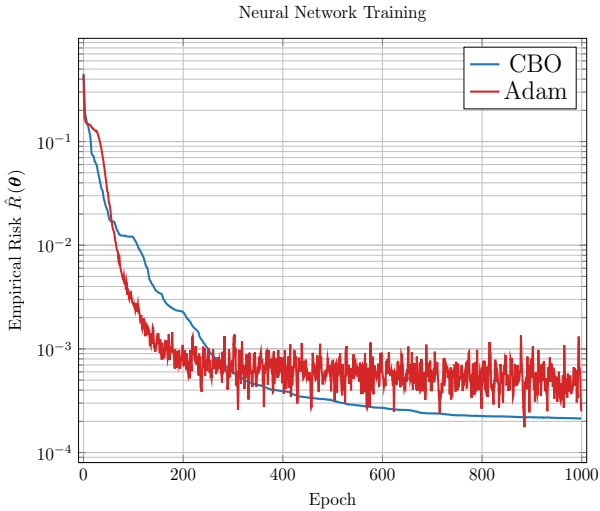


Figure 2: Empirical risk  $\hat{R}(\theta)$  as a function of training epochs for a two-layer neural network trained with Adam and CBO. The figure displays the median empirical risks taken over 10 simulations.

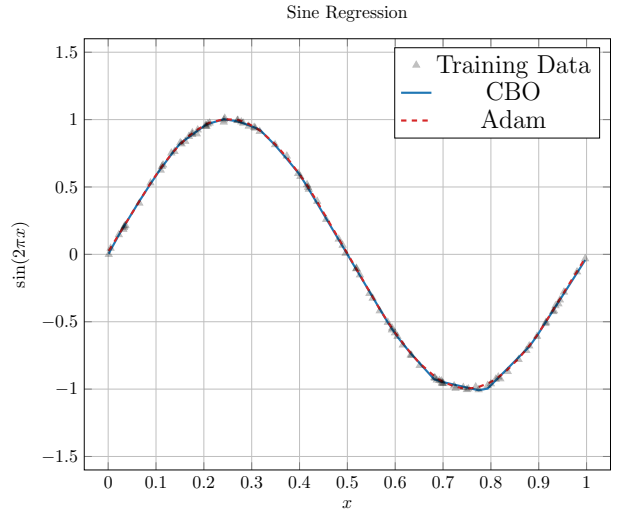


Figure 3: Approximations of the sine function  $\sin(2\pi x)$  obtained with a two-layer neural network with  $M = 100$  trained with CBO and Adam. The plot only contains a subset of the training dataset to improve the clarity.

## 4.2 Example 2: MNIST

In the second experiment, we further investigate the applicability of CBO in training two-layer neural networks to classify the MNIST dataset [40]. This is a standard classification problem in machine learning, and commonly used to benchmark new methods and models. We compare CBO, Adam and the hybrid method on both the minimal empirical risk achieved and the stability of the training.



The experiment setup is the following: consider the multi-class classification of the MNIST dataset with  $C = 10$  classes. We again use a two-layer neural network with a width  $M = 20$  and the RELU activation function. The MNIST dataset contains grayscale images of  $28 \times 28$  pixels. We flatten these images to vectors with 784 elements, which represent the input data points  $\mathbf{x}_s$ . In the MNIST dataset, the true labels are given as class indices, i.e.,  $y_s \in \{1, \dots, C\}$ . We take a subset of 10 000 images to form the training dataset  $\{(\mathbf{x}_s, y_s)\}_{s=1}^{10000}$ . We apply the minibatch strategy to Adam, CBO and the hybrid method using a minibatch size of  $S'$ . The neural network outputs a ten dimensional vector  $\hat{\mathbf{y}}_s$ . We choose the cross-entropy loss function of PyTorch, similar to Eq. (13), which internally applies the softmax function to the predictions [46]. We initialize the particles of the CBO method from the uniform distribution  $\theta_n^0 \sim \mathcal{U}[-1, 1]$ . For Adam, we keep the default initialization of the neural network parameters provided by PyTorch.

The training of the two-layer neural network with CBO is carried out with the following parameters:

$$N = 1000, \quad \Delta t = 0.1, \quad \alpha = 10^5, \quad \lambda = 1, \quad \tilde{\sigma} = \sqrt{1.4}, \quad S' = 1000.$$

Parameters shared between CBO and Adam, such as the minibatch size  $S'$  and the time step  $\Delta t$ , are taken equal. The parameters of the hybrid method are set as:

$$N = 1000, \quad \Delta t = 0.1, \quad \alpha = 10^4, \quad \lambda = 1, \quad \tilde{\sigma} = \sqrt{1.2}, \quad S' = 1000, \quad \gamma = 0.7.$$

Figure 4 illustrates the empirical risk per epoch for CBO, Adam and the hybrid method. We observe that Adam obtains a significantly lower empirical risk than CBO and converges faster. However, we notice that at approximately the same risk values, the Adam method becomes unstable and moves towards sub-optimal areas of the optimization landscape. The hybrid method, which combines CBO and Adam, has a faster convergence than CBO and is more stable than Adam. However, the hybrid method also experiences a decrease in stability near the same risk value, but the effect is considerably less severe.

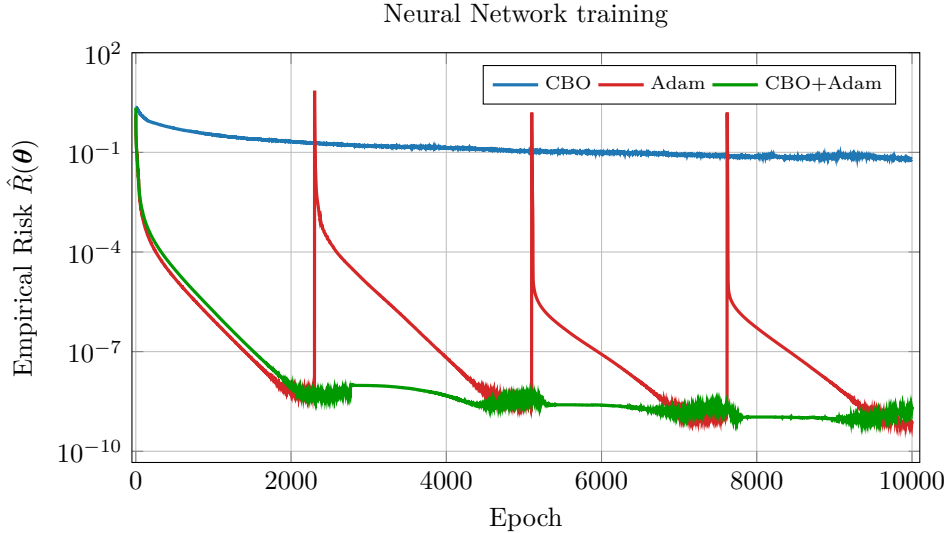


Figure 4: Empirical risk  $\hat{R}(\theta)$  as a function of training epochs for a two-layer neural network trained with Adam, CBO and the hybrid method (Adam + CBO) on the MNIST dataset.

The instability of the Adam method in Figure 4 is likely due to the exploding gradient phenomena [47]. The improved stability of the hybrid method is a consequence of the fact that the minimizer is represented by the consensus point  $\mathbf{V}^k$ . For the consensus point to suddenly jump due to the exploding gradient phenomena, all particles would need to move simultaneously in approximately the same direction. However, before convergence, the particles are distributed across the optimization landscape, so they do not experience the exploding gradients simultaneously; interaction through the consensus point stabilizes the dynamics.

### 4.3 Example 3: Multi-Task CBO

In the third example, we study whether the particle recycling strategy of Multi-Task CBO is able to minimize multiple empirical risk functions. Specifically, we assess whether the empirical risk consistently decreases throughout training for all considered tasks.

We set up the experiment as follows: consider the approximation problem of 100 shifted sine functions, where the goal is to approximate each function using a different two-layer neural network. We take 100 functions of the form  $(x, \Delta y_p) \mapsto \sin(2\pi x) + \Delta y_p$  on the domain  $[0, 1]$ , where the shifts are uniformly spaced in  $[-1, 1]$ ,  $\Delta y_p = -1 + 2(p - 1)/99$ ,  $p = 1, \dots, 100$ . We assume that the shifted sine functions admit similar neural network representations and hence we expect their corresponding global minimizer to be in close proximity to each other.

The neural networks have a RELU activation function and a width of  $M = 100$ . We sample 8000 data points  $x_s$  uniformly on the interval  $[0, 1]$  and apply the function  $y_{s,p} = \sin(2\pi x_s) + \Delta y_p$  to generate 100 training datasets  $\{(x_s, y_{s,p})\}_{s=1}^{8000}$ ,  $p = 1, \dots, 100$ . We apply the minibatch strategy, dividing each dataset into minibatches of size  $S'$ . The loss function for each task  $p$  is the squared error loss, resulting in 100 different MSE risks  $\hat{R}_p$ . We initialize the particles of the CBO method from the uniform distribution  $\theta_n^0 \sim \mathcal{U}[-1, 1]$ .

The parameters of the experiment are:

$$N = 200, \quad \Delta t = 0.2, \quad \alpha = 10^4, \quad \lambda = 1, \quad \tilde{\sigma} = \sqrt{1.8}, \quad S' = 800.$$

Every 100 epochs, the parameter  $\alpha$  is multiplied by 10 until it reaches  $10^7$ . Lastly, we run the experiment ten times to average out as much noise.

Given 100 different problems, the Multi-Task CBO method has 100 different consensus points. When optimizing with 200 particles, it is necessary to determine which particles move towards which consensus point. The particle update strategy of is as follows: the first two particles will move towards the first consensus point, the third and fourth particles will move to the second consensus point and so forth. Multi-Task CBO effectively assigns two particles to each problem.

Figure 5 presents the median and minimum empirical risk per epoch. The median and minimum are taken over the 100 tasks during training. In Figure 6, we display the approximation results for five problems. In Figure 5, we observe that both the minimum and median decrease during training, indicating that the CBO method effectively minimizes the risk for all tasks. In Figure 6, we observe an accurate approximation for these five problems, confirming that CBO successfully trained multiple tasks using the same particle set.

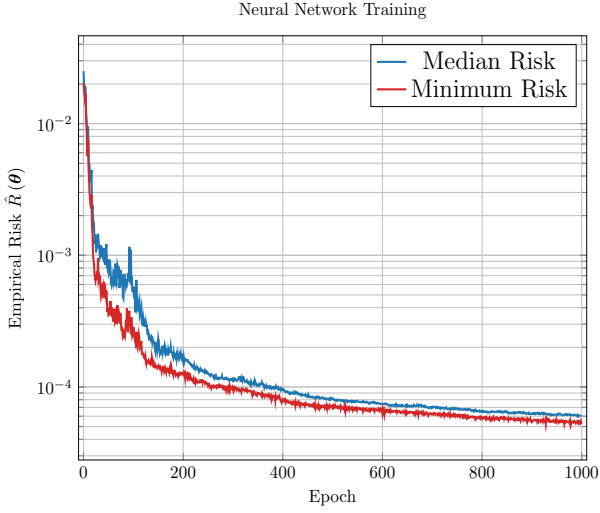


Figure 5: The median and minimum empirical risk  $\hat{R}(\theta)$  as a function of training epochs for two-layer neural networks trained with Multi-Task CBO. The median and minimum are taken over 100 different risk functions.

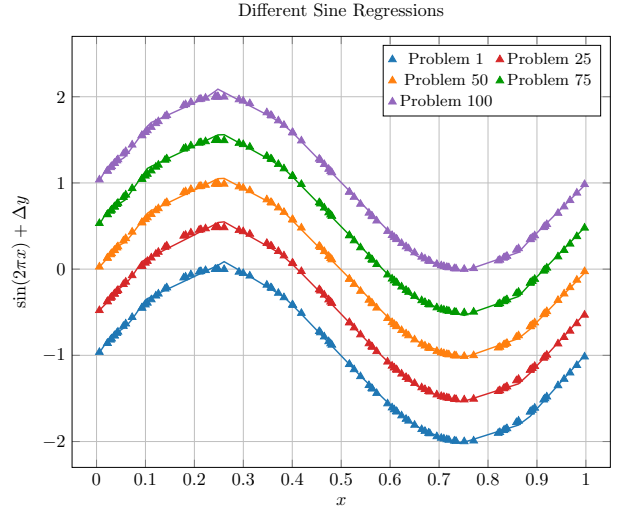


Figure 6: Approximation of five different sine functions obtained with five different two-layer neural networks, each with  $M = 100$ . The neural networks are trained with Multi-Task CBO.

## 5 Mean-field Limits

This section provides a study of mean-field limits, as indicated in Figure 7. We derive an explicit JKO scheme that arises in the limit to infinity of both the width of the neural network  $M$  and the number of particles  $N$ . In Subsection 5.1, we start from the classical CBO formulation, detailed in Subsection 3.3, and take the width of the neural network to infinity ( $M \rightarrow \infty$ ). We obtain an optimal transport (OT) formulation of CBO. Starting from the optimal transport formulation of CBO, Subsection 5.2 derives the time-discrete mean-field limit ( $N \rightarrow \infty$ ). Further, we show that the population variance decreases each iteration.

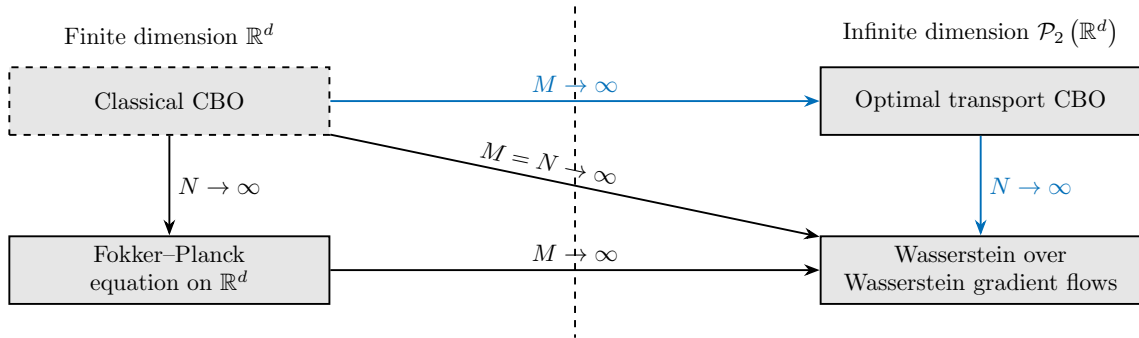


Figure 7: Diagram illustrating possible derivations of the mean-field limits. In this paper, we consider the blue path, starting from the classical CBO formulation (dashed box) of Section 3 and proceeding with  $M$  and then  $N$  to infinity.

### 5.1 The infinite width limit ( $M \rightarrow \infty$ )

The classical CBO formulation described in Section 3 formally breaks down when the number of hidden neurons  $M$  tends to infinity, due to the fact that the optimization dimension equals  $M(d+2)$ . Instead of representing each neural network as a point in  $\mathbb{R}^{M(d+2)}$ , we therefore choose to represent the neural networks by corresponding measures  $\mu_n$ , as introduced in Section 2. We reformulate the dynamics in the Wasserstein space, similar to the work [9]. For an ensemble of  $N$  particles  $\mu_n^k \in \mathcal{P}_2(\mathbb{R}^{d+2})$  and a

time step  $\Delta t \in (0, 1]$  we have

$$\mu_n^{k+1} = ((1 - \Delta t) \text{Id} + \Delta t T_n)_\# \mu_n^k, \quad n = 1, \dots, N, \quad (29)$$

with  $T_n : \mathbb{R}^{d+2} \rightarrow \mathbb{R}^{d+2}$  the optimal transport map defined as

$$\bar{\mu}^k = (T_n)_\# \mu_n^k, \quad n = 1, \dots, N, \quad (30)$$

and  $\bar{\mu}^k$  represents the consensus point, i.e., the barycenter in  $\mathcal{P}_2(\mathbb{R}^{d+2})$ . It is given by

$$\bar{\mu}^k = \arg \min_{\nu} \frac{1}{2} \sum_{n=1}^N \beta(\mu_n^k) W_2^2(\nu, \mu_n^k), \quad \beta(\mu_n^k) = \frac{\exp(-\alpha \hat{R}(\mu_n^k))}{\sum_{n=1}^N \exp(-\alpha \hat{R}(\mu_n^k))}. \quad (31)$$

We note that the dynamics in (29) are fully deterministic. To the best of our knowledge, no formulation of Brownian motion on  $\mathcal{P}_2(\mathbb{R}^{d+2})$  is currently available. As a result, the optimal transport formulation of CBO does not include a diffusion term. However, the noise can be added after discretization with empirical measures  $\hat{\mu}_n^k$ , yielding dynamics that closely resemble the classical CBO formulation.

We consider the optimal transport formulation of CBO in the single-task setting, although it can also be written in the multi-task setting as

$$\begin{cases} \mu_n^{k+1} = ((1 - \Delta t) \text{Id} + \Delta t T_n)_\# \mu_n^k \\ \mathcal{T}_p^{k+1} = \mathcal{T}_p^k \end{cases}, \quad n = p = 1, \dots, N, \quad (32)$$

where now the barycenter depends on the particular task  $\mathcal{T}_p$

$$\bar{\mu}(\mathcal{T}_p^k) = \arg \min_{\nu} \frac{1}{2} \sum_{n=1}^N \beta(\mu_n^k, \mathcal{T}_p^k) W_2^2(\nu, \mu_n^k), \quad \beta(\mu_n^k, \mathcal{T}_p^k) = \frac{\exp(-\alpha \hat{R}(\mu_n^k, \mathcal{T}_p^k))}{\sum_{n=1}^N \exp(-\alpha \hat{R}(\mu_n^k, \mathcal{T}_p^k))}. \quad (33)$$

At any given time, it holds that the consensus point and barycenter are applications of a weighted Fréchet mean in the corresponding metric space:

**Proposition 1.** *The representation of the consensus point in  $\mathbb{R}^d$  (Eq. (22)) and the barycenter in  $\mathcal{P}_2(\mathbb{R}^d)$  (Eq. (31)) are equal in the following sense:*

$$\bar{x} = \arg \min_v \frac{1}{2} \sum_{n=1}^N \beta(x_n) d^2(v, x_n), \quad (34)$$

where  $d(\cdot, \cdot)$  denotes the distance function in  $\mathbb{R}^d$  and  $\mathcal{P}_2(\mathbb{R}^d)$ , respectively.

*Proof.* The consensus point in  $\mathbb{R}^d$  is the result of the minimization of

$$F(\mathbf{v}) = \frac{1}{2} \sum_{n=1}^N \beta(\boldsymbol{\theta}_n) \|\mathbf{v} - \boldsymbol{\theta}_n\|_2^2. \quad (35)$$

Given  $F \in C^1(\mathbb{R}^d)$ , an optimal point must satisfy the first-order necessary condition

$$0 = \nabla F(\mathbf{v}) = - \sum_{n=1}^N \beta(\boldsymbol{\theta}_n) (\boldsymbol{\theta}_n - \mathbf{v}). \quad (36)$$

Solving (36) yields the minimizer

$$\mathbf{v} = \sum_{n=1}^N \beta(\boldsymbol{\theta}_n) \boldsymbol{\theta}_n, \quad (37)$$

which is the classical consensus point, see Equation (22). For the barycenter in  $\mathcal{P}_2(\mathbb{R}^d)$ , Equation (31) fulfills Proposition 1.  $\square$

**Proposition 2.** Assume empirical measures of the form

$$\hat{\mu}_n = \frac{1}{M} \sum_{i=1}^M \delta(x - \mathbf{x}_{n,i}) \in \mathcal{P}_2(\mathbb{R}^d), \quad n = 1, \dots, N. \quad (38)$$

The barycenter in Eq. (31) is given by

$$\bar{\mu} = \frac{1}{M} \sum_{j=1}^M \delta(y - \mathbf{y}_j) \in \mathcal{P}_2(\mathbb{R}^d), \quad (39)$$

with support points

$$\mathbf{y}_j = \sum_{n=1}^N \beta(\hat{\mu}_n) \sum_{i=1}^M (\pi_n^*)_{j,i} \mathbf{x}_{n,i}, \quad j = 1, \dots, M, \quad (40)$$

where  $\pi_n^* \in \mathbb{R}^{d \times d}$  is a permutation matrix.

*Proof.* For any  $M$ , let  $\hat{\mu}_n \in \mathcal{P}_2(\mathbb{R}^d)$  be

$$\hat{\mu}_n^k = \frac{1}{M} \sum_{i=1}^M \delta(x - \mathbf{x}_{n,i}). \quad (41)$$

Then, the barycenter is of the form

$$\hat{\nu} = \frac{1}{M} \sum_{j=1}^M \delta(y - \mathbf{y}_j), \quad (42)$$

and according to Proposition 1, obtained as the minimization of

$$F(\hat{\nu}) = \frac{1}{2} \sum_{n=1}^N \beta(\hat{\mu}_n) W_2^2(\hat{\nu}, \hat{\mu}_n). \quad (43)$$

Substituting Equations (41) and (42) into (43) yields

$$F(\mathbf{y}_1, \dots, \mathbf{y}_M) = \frac{1}{2} \sum_{n=1}^N \beta(\hat{\mu}_n) \min_{\pi_n \in \Pi} \sum_{j=1}^M \sum_{i=1}^M \frac{1}{M} (\pi_n)_{j,i} \|\mathbf{y}_j - \mathbf{x}_{n,i}\|_2^2, \quad (44)$$

where  $\Pi$  denotes the set of permutation matrices of size  $M \times M$ . This set is compact, hence, there exist a minimizer  $\pi^* \in \mathbb{R}^{M \times M}$ . To compute the barycenter, we solve for the optimal coupling  $\pi_n$  and the location of the barycenter points  $(\mathbf{y}_1, \dots, \mathbf{y}_M)$  [21]. Now, for a fixed optimal coupling  $\pi_n^*$ , we obtain

$$F(\mathbf{y}_1, \dots, \mathbf{y}_M) = \frac{1}{2} \sum_{n=1}^N \beta(\hat{\mu}_n) \sum_{j=1}^M \sum_{i=1}^M \frac{1}{M} (\pi_n^*)_{j,i} \|\mathbf{y}_j - \mathbf{x}_{n,i}\|_2^2. \quad (45)$$

The first-order necessary optimality condition reads

$$0 = \frac{\partial F(\mathbf{y}_1, \dots, \mathbf{y}_M)}{\partial \mathbf{y}_j} = \sum_{n=1}^N \beta(\hat{\mu}_n) \sum_{i=1}^M \frac{1}{M} (\pi_n^*)_{j,i} (\mathbf{y}_j - \mathbf{x}_{n,i}), \quad (46)$$

and hence

$$\mathbf{y}_j = \frac{\sum_{n=1}^N \beta(\hat{\mu}_n) \sum_{i=1}^M (\pi_n^*)_{j,i} \mathbf{x}_{n,i}}{\sum_{n=1}^N \beta(\hat{\mu}_n) \sum_{i=1}^M (\pi_n^*)_{j,i}} = \sum_{n=1}^N \beta(\hat{\mu}_n) \sum_{i=1}^M (\pi_n^*)_{j,i} \mathbf{x}_{n,i}. \quad (47)$$

□

### 5.1.1 Example 4: Square Approximation

In the fourth experiment, we investigate whether the optimal transport dynamics proposed in (29) are capable of training arbitrarily wide neural networks. To this end, we focus on the empirical risk during training for neural networks with different values of  $M$ .

In the implementation of the optimal transport CBO formulation, we consider empirical measures of the form:

$$\hat{\mu}_n^k = \frac{1}{M} \sum_{i=1}^M \delta(\mathbf{w} - \mathbf{w}_{n,i}^k) \delta(b - b_{n,i}^k) \delta(c - c_{n,i}^k). \quad (48)$$

The dynamics in (29), for  $\mu_n^k$  equal to the empirical measure (48), takes the form of a noise-free CBO scheme. Written explicitly for  $(\mathbf{w}, b, c)$ , this gives

$$\mathbf{w}_{n,i}^{k+1} = \mathbf{w}_{n,i}^k - \Delta t \left( \mathbf{w}_{n,i}^k - \bar{\mathbf{w}}_{\pi_n^*(i)}^k \right) \quad (49)$$

$$b_{n,i}^{k+1} = b_{n,i}^k - \Delta t \left( b_{n,i}^k - \bar{b}_{\pi_n^*(i)}^k \right) \quad (50)$$

$$c_{n,i}^{k+1} = c_{n,i}^k - \Delta t \left( c_{n,i}^k - \bar{c}_{\pi_n^*(i)}^k \right), \quad (51)$$

where  $\pi_n^*$  denotes the permutation matrix that represents the optimal coupling of the  $i$ -th weight to the barycenter. To facilitate the training, we include a drift parameter  $\lambda$  and add artificial noise to the parameter updates. The complete dynamics are as follows:

$$\mathbf{w}_{n,i}^{k+1} = \mathbf{w}_{n,i}^k - \lambda \Delta t \left( \mathbf{w}_{n,i}^k - \bar{\mathbf{w}}_{\pi_n^*(i)}^k \right) + \tilde{\sigma}^k \sqrt{\Delta t} \boldsymbol{\xi}_{n,i}^k \quad (52)$$

$$b_{n,i}^{k+1} = b_{n,i}^k - \lambda \Delta t \left( b_{n,i}^k - \bar{b}_{\pi_n^*(i)}^k \right) + \tilde{\sigma}^k \sqrt{\Delta t} \xi_{n,i}^k \quad (53)$$

$$c_{n,i}^{k+1} = c_{n,i}^k - \lambda \Delta t \left( c_{n,i}^k - \bar{c}_{\pi_n^*(i)}^k \right) + \tilde{\sigma}^k \sqrt{\Delta t} \xi_{n,i}^k, \quad (54)$$

with  $\lambda, \tilde{\sigma} > 0$  and  $\boldsymbol{\xi}_{n,i}^k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . However, the noise lacks a multiplicative term that decreases as the particles converge to the barycenter, similar to the anisotropic noise in (21). Therefore, we manually reduce  $\tilde{\sigma}^k$  according to a predefined schedule.

We consider a one-dimensional regression problem, where the goal is to approximate the function  $x \mapsto x^2$  on the domain  $[0, 1]$  using a finite two-layer neural network. The neural network has a RELU activation function and we consider various widths  $M$ . We sample 5000 data points  $x_s$  uniformly on the interval  $[0, 1]$  and apply the function  $y_s = x_s^2 + 0.01\xi_s$ , with  $\xi_s \sim \mathcal{N}(0, 1)$ , to generate the training dataset  $\{(x_s, y_s)\}_{s=1}^{5000}$ . We divide the training dataset into minibatches of size  $S'$ . We take the squared error loss as the loss function.

The CBO parameters in the experiment are chosen as follows:

$$N = 100, \quad \Delta t = 0.1, \quad \alpha = 10^4, \quad \lambda = 1, \quad \tilde{\sigma} = \sqrt{1.2}, \quad S' = 2500.$$

For each particle, we choose a uniform distribution as the initial measure and sample the atoms from  $\mathbf{w}_{n,i}, b_{n,i}, c_{n,i} \sim \mathcal{U}[-2, 2]$ . Every 100 epochs, we reduce the noise parameter by a factor of 0.9 every 100 iterations and multiply  $\alpha$  by 10 until it reaches  $10^7$ . Finally, for each value of  $M$ , we perform 10 independent simulations and report the average.

Figure 8 depicts the empirical risk per epoch of neural networks with different widths  $M$  during training. We observe that the empirical risk for each neural network decreases. Additionally, we remark that neural networks with more hidden neurons  $M$  converge to a lower empirical risk and thus obtain a better approximation.

The experiment confirms that the optimal transport CBO formulation is also able to train neural networks. However, computing the barycenter each iteration has a higher computational cost than the classical consensus point, limiting the practical usage of the optimal transport formulation.

## Neural Network Training

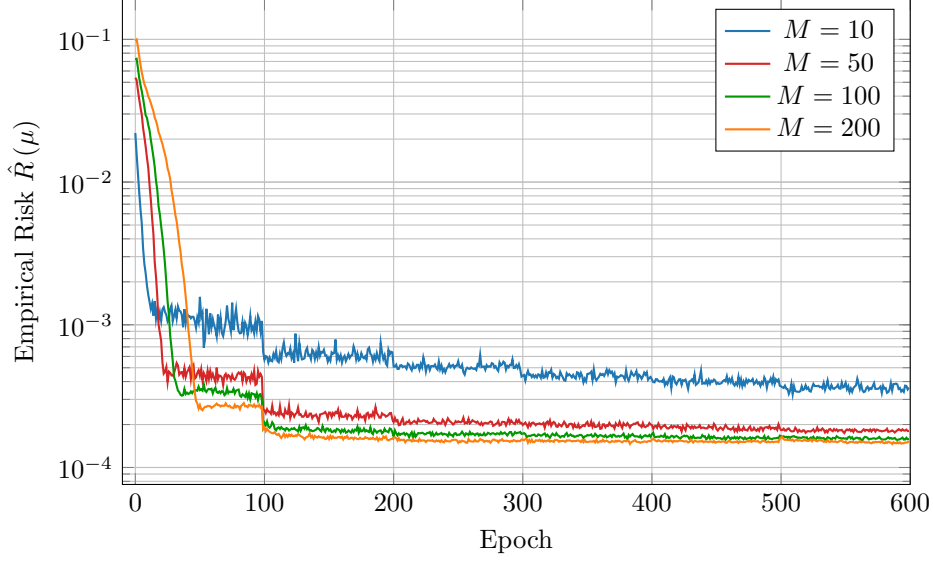


Figure 8: Empirical risk  $\hat{R}(\mu)$  as a function of training epochs for four different neural networks, each with a different width  $M$ . Each neural network is represent by a measure and trained with the optimal transport formulation of CBO. The figure displays the mean empirical risks taken over 10 differet simulations.

### 5.2 The mean-field limit for infinitely many particles ( $N \rightarrow \infty$ )

In Subsection 5.1, we introduced the optimal transport dynamics

$$\mu_n^{k+1} = ((1 - \Delta t) \text{Id} + \Delta t T_n)_\# \mu_n^k, \quad n = 1, \dots, N. \quad (55)$$

with the optimal transport plan  $T_n : \mathbb{R}^{d+2} \rightarrow \mathbb{R}^{d+2}$  given by

$$\bar{\mu}^k = (T_n)_\# \mu_n^k, \quad n = 1, \dots, N. \quad (56)$$

We now consider the mean-field limit as  $N \rightarrow \infty$ . Let  $X := \mathcal{P}_2(\mathbb{R}^{d+2})$  and equip  $X$  with the 2-Wasserstein metric  $W_2$ . Define

$$\mathcal{P}_2(X) := \left\{ \rho \in \mathcal{P}(X) : \int_X W_2^2(\mu, \delta_0) d\rho(\mu) < \infty \right\}. \quad (57)$$

For  $\rho, \phi \in \mathcal{P}_2(X)$  define the Wasserstein-over-Wasserstein distance as [6, 4, 49]

$$\mathbb{W}_2^2(\rho, \phi) := \inf_{\Gamma \in \Pi(\rho, \phi)} \int_{X \times X} W_2^2(\mu, \nu) d\Gamma(\mu, \nu), \quad (58)$$

where

$$\Pi(\rho, \phi) := \{ \Gamma \in \mathcal{P}(X \times X) : (\pi_1)_\# \Gamma = \rho, (\pi_2)_\# \Gamma = \phi \} \quad (59)$$

is the set of couplings between  $\rho$  and  $\phi$ . We denote the law of particles by  $\rho \in \mathcal{P}_2(X)$  and make the following assumptions:

**Assumption 1.** The barycenter  $\bar{\mu} \in X$  of  $\rho$  exists and is the global minimizer of  $F_\rho(\nu) : X \rightarrow \mathbb{R}$ , given by

$$F_\rho(\nu) = \frac{1}{2} \int_X W_2^2(\mu, \nu) d\rho(\mu). \quad (60)$$

**Assumption 2.** Each measure  $\mu, \bar{\mu} \in X$  is absolute continuous with respect to the Lebesgue measure.

For each  $\mu \in X$  we define the measurable map  $\Psi_{\Delta t} : X \rightarrow X$

$$\Psi_{\Delta t}(\mu) := ((1 - \Delta t) \text{Id} + \Delta t T)_\# \mu, \quad (61)$$

with

$$\bar{\mu} = T_{\#}\mu, \quad (62)$$

such that the particle update in Eq. (55) reads  $\mu_n^{k+1} = \Psi_{\Delta t}(\mu_n^k)$ . The law of particles  $\rho$  evolves as

$$\rho^{k+1} = (\Psi_{\Delta t})_{\#}\rho^k. \quad (63)$$

This is the time-discrete mean-field limit, which provides an analytic framework to investigate the optimal transport CBO scheme:

**Proposition 3.** *Let the variance of the measure  $\rho^k$  be given by*

$$\mathcal{V}(\rho^k) = \frac{1}{2} \int_X W_2^2(\mu, \bar{\mu}^k) d\rho^k(\mu). \quad (64)$$

*For a time step  $\Delta t \in (0, 1]$ , it holds that*

$$\mathcal{V}(\rho^{k+1}) \leq (1 - \Delta t)^2 \mathcal{V}(\rho^k). \quad (65)$$

*Proof.* We start by observing that

$$\mathcal{V}(\rho^{k+1}) = F_{\rho^{k+1}}(\bar{\mu}^{k+1}) \leq F_{\rho^{k+1}}(\bar{\mu}^k) = \frac{1}{2} \int_X W_2^2(\mu, \bar{\mu}^k) d\rho^{k+1}(\mu), \quad (66)$$

since  $\bar{\mu}^{k+1}$  is the global minimizer of  $F_{\rho^{k+1}}(\nu)$ . Next, we apply the definition of the pushforward on Eq. (66). This yields

$$\mathcal{V}(\rho^{k+1}) \leq \frac{1}{2} \int_X W_2^2(\Psi_{\Delta t}(\mu), \bar{\mu}^k) d\rho^k(\mu). \quad (67)$$

The map  $\Psi_{\Delta t}$  defines a constant-speed Wasserstein geodesic between  $\mu$  and  $\bar{\mu}$  (Theorem 7.2.2 in [1]), and therefore satisfies

$$W_2^2(\Psi_{\Delta t}(\mu), \bar{\mu}^k) = (1 - \Delta t)^2 W_2^2(\mu, \bar{\mu}^k). \quad (68)$$

Combining (67) and (68), we obtain

$$\mathcal{V}(\rho^{k+1}) \leq \frac{1}{2} (1 - \Delta t)^2 \int_X W_2^2(\mu, \bar{\mu}^k) d\rho^k(\mu) = (1 - \Delta t)^2 \mathcal{V}(\rho^k). \quad (69)$$

□

Proposition 3 demonstrates that the optimal transport formulation of CBO also achieves consensus, as the variance converges to zero.

## 6 Conclusion

This work investigated training two-layer neural networks with Consensus-Based Optimization and analyzed its mean-field limits. On a smooth regression task, CBO achieved competitive final empirical risk values compared to Adam, while the hybrid method on MNIST classification demonstrated greater robustness than Adam and faster convergence than CBO. We hypothesize that the CBO method achieves faster convergence for highly non-convex risk functions, whereas in cases of slow convergence, incorporating local gradient information can be beneficial. In the multi-objective optimization setting, Multi-Task CBO achieves high accuracy with minimal memory overhead by recycling particles across all tasks.

On the theoretical side, we reformulated the CBO dynamics on  $\mathbb{R}^{M(d+2)}$  within the Wasserstein space  $\mathcal{P}_2(\mathbb{R}^{d+2})$ , thereby enabling the training of continuous neural networks with a particle-based method. However, the optimal transport dynamics are currently deterministic, as they do not include Brownian motion. In practice, artificial noise can always be added when considering empirical measures. For the regression of the square function, we demonstrated that the optimal transport scheme with added noise successfully trains neural networks, even as the width increases. Lastly, we presented a time-discrete mean-field formulation over both the neurons ( $M$ ) and the particles ( $N$ ) and proved that the optimal transport scheme also achieves consensus.



## Acknowledgements

The work of W.D.D. is supported by the European Union’s Horizon Europe research and innovation program under the Marie Skłodowska-Curie Doctoral Network Datahyking (Grant No. 101072546). The authors thank the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) for the financial support through 442047500/SFB1481 within the projects B04 (Sparsity fördernde Muster in kinetischen Hierarchien), B05 (Sparsifizierung zeitabhängiger Netzwerkflußprobleme mittels diskreter Optimierung) and B06 (Kinetische Theorie trifft algebraische Systemtheorie). The authors thank the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) for the financial support HE5386/33-1 Control of Interacting Particle Systems, and Their Mean-Field, and Fluid-Dynamic Limits (560288187), and HE5386/34-1 Partikelmethoden für unendlich dimensionale Optimierung (561130572).

## References

- [1] L. Ambrosio, N. Gigli, and G. Savaré. *Gradient Flows: In Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. ETH Zürich. Birkhäuser Verlag, Basel, second edition, 2008.
- [2] F. Bach and L. Chizat. Gradient Descent on Infinitely Wide Neural Networks: Global Convergence and Generalization. In *International Congress of Mathematicians*, Saint-Petersbourg, Russia, July 2022.
- [3] A. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, May 1993.
- [4] M. Beiglböck, G. Pammer, and S. Schrott. A Brenier Theorem on  $(\mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d)), W_2)$  and Applications to Adapted Transport. *arXiv:2509.03506*, Sept. 2025.
- [5] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, Mar. 1994.
- [6] C. Bonet, C. Vauthier, and A. Korba. Flowing datasets with wasserstein over wasserstein gradient flows. In *42nd International Conference on Machine Learning*, 2025.
- [7] G. Borghi, M. Herty, and L. Pareschi. A Consensus-Based Algorithm for Multi-Objective Optimization and Its Mean-Field Description. In *2022 IEEE 61st Conference on Decision and Control (CDC)*, pages 4131–4136, Dec. 2022.
- [8] G. Borghi, M. Herty, and L. Pareschi. Constrained Consensus-Based Optimization. *SIAM Journal on Optimization*, 33(1):211–236, Mar. 2023.
- [9] G. Borghi, M. Herty, and A. Stavitskiy. Dynamics of Measure-Valued Agents in the Space of Probabilities. *SIAM Journal on Mathematical Analysis*, 57(5):5107–5134, Oct. 2025.
- [10] L. Böttcher and T. Asikis. Near-optimal control of dynamical systems with neural ordinary differential equations. *Machine Learning: Science and Technology*, 3(4):045004, Dec. 2022.
- [11] L. Bottou. Online algorithms and stochastic approximations. In D. Saad, editor, *Online Learning and Neural Networks*. Cambridge University Press, Cambridge, UK, 1998.
- [12] J. A. Carrillo, Y.-P. Choi, C. Totzeck, and O. Tse. An analytical framework for consensus-based global optimization method. *Mathematical Models and Methods in Applied Sciences*, 28(06):1037–1066, June 2018.
- [13] J. A. Carrillo, F. Hoffmann, A. M. Stuart, and U. Vaes. Consensus-based sampling. *Studies in Applied Mathematics*, 148(3):1069–1140, Apr. 2022.

- [14] J. A. Carrillo, S. Jin, L. Li, and Y. Zhu. A consensus-based global optimization method for high dimensional machine learning problems. *ESAIM: Control, Optimisation and Calculus of Variations*, 27:S5, 2021.
- [15] J. A. Carrillo, C. Totzeck, and U. Vaes. *Consensus-Based Optimization and Ensemble Kalman Inversion for Global Optimization Problems with Constraints*, volume 40, pages 195–230. WORLD SCIENTIFIC, Feb. 2023.
- [16] R. Caruana. Multitask Learning. *Machine Learning*, 28(1):41–75, July 1997.
- [17] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. Neural ordinary differential equations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, page 6572–6583, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [18] Z. Chen, V. Badrinarayanan, C.-Y. Lee, and A. Rabinovich. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 794–803. PMLR, 2018.
- [19] L. Chizat. Mean-field langevin dynamics: Exponential convergence and annealing. *Transactions on Machine Learning Research*, Aug. 2022.
- [20] L. Chizat and F. Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Proceedings of the Thirty-Third Conference on Learning Theory (COLT)*, volume 125 of *Proceedings of Machine Learning Research*, pages 1305–1338. PMLR, July 2020.
- [21] M. Cuturi and A. Doucet. Fast Computation of Wasserstein Barycenters. *Proceedings of the 31st International Conference on Machine Learning*, 32(2):685–693, June 2014.
- [22] W. De Deyn. Code repository for: Mean-Field Limits for Two-Layer Neural Networks Trained with Consensus-Based Optimization. <https://git.rwth-aachen.de/wdedeyn/mean-field-limits-for-neural-networks>, 2025. Accessed: 05-10-2025.
- [23] A. Dembo and O. Zeitouni. *Large Deviations Techniques and Applications*. Number 38 in Stochastic Modelling and Applied Probability. Springer, Berlin New York, 2nd ed., corr. print edition, 2010.
- [24] W. E, C. Ma, and L. Wu. The Barron Space and the Flow-Induced Function Spaces for Neural Network Models. *Constructive Approximation*, 55(1):369–406, Feb. 2022.
- [25] M. Fornasier, T. Klock, and K. Riedl. Convergence of Anisotropic Consensus-Based Optimization in Mean-Field Law. In *Applications of Evolutionary Computation*, volume 13224, pages 738–754. Springer International Publishing, Cham, 2022.
- [26] M. Fornasier, T. Klock, and K. Riedl. Consensus-Based Optimization Methods Converge Globally. *SIAM Journal on Optimization*, 34(3):2973–3004, Sept. 2024.
- [27] M. Fornasier, L. Pareschi, H. Huang, and P. Sünnen. Consensus-based optimization on the sphere: Convergence to global minimizers and machine learning. *Journal of Machine Learning Research*, 22(237):1–55, 2021.
- [28] N. Gerber, F. Hoffmann, and U. Vaes. Mean-field limits for Consensus-Based Optimization and Sampling. *ESAIM: Control, Optimisation and Calculus of Variations*, July 2025.
- [29] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, Mar. 2010.
- [30] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016.
- [31] S.-Y. Ha, S. Jin, and D. Kim. Convergence of a first-order consensus-based global optimization algorithm. *Mathematical Models and Methods in Applied Sciences*, 30(12):2417–2444, Nov. 2020.

- [32] S.-Y. Ha, S. Jin, and D. Kim. Convergence and error estimates for time-discrete consensus-based optimization algorithms. *Numerische Mathematik*, 147(2):255–282, Feb. 2021.
- [33] S. Hochreiter and J. Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, Nov. 1997.
- [34] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, Jan. 1989.
- [35] H. Huang and J. Qiu. On the mean-field limit for the consensus-based optimization. *Mathematical Methods in the Applied Sciences*, 45(12):7814–7831, Aug. 2022.
- [36] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*, 2015.
- [37] M. Koß, S. Weissmann, and J. Zech. On the mean-field limit of consensus based methods. *arXiv:2409.03518*, Sept. 2024.
- [38] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, May 2017.
- [39] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, May 2015.
- [40] Y. LeCun, C. Cortes, and C. J. Burges. The MNIST Database of handwritten digits. Courant Institute, NYU, Google Labs, New York, Microsoft Research, Redmond, 1998.
- [41] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33), Aug. 2018.
- [42] K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, 2014.
- [43] P.-M. Nguyen and H. T. Pham. A rigorous framework for the mean field limit of multilayer neural networks. *Mathematical Statistics and Learning*, 6(3):201–357, Oct. 2023.
- [44] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York, 2006.
- [45] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1310–1318, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- [46] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [47] G. Philipp, D. Song, and J. G. Carbonell. The exploding gradient problem demystified - definition, prevalence, impact, origin, tradeoffs, and solutions. *arXiv:1712.05577*, 2017.
- [48] R. Pinnau, C. Totzeck, O. Tse, and S. Martin. A consensus-based model for global optimization and its mean-field limit. *Mathematical Models and Methods in Applied Sciences*, 27(01):183–204, Jan. 2017.
- [49] A. Pinzi and G. Savaré. Totally convex functions,  $L^2$ -Optimal transport for laws of random measures, and solution to the Monge problem. *arXiv:2509.01768*, Sept. 2025.
- [50] M. Raissi, P. Perdikaris, and G. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, Feb. 2019.

- [51] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, Oct. 1986.
- [52] J. Sirignano and K. Spiliopoulos. Mean field analysis of neural networks: A central limit theorem. *Stochastic Processes and their Applications*, 130(3):1820–1852, Mar. 2020.
- [53] J. Sirignano and K. Spiliopoulos. Mean Field Analysis of Neural Networks: A Law of Large Numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, Jan. 2020.
- [54] V. Vapnik. An overview of statistical learning theory. *IEEE Transactions on Neural Networks*, 10(5):988–999, 1999.
- [55] J. Wei, F. Wu, and W. Bian. A Consensus-Based Optimization Method for Nonsmooth Nonconvex Programs with Approximated Gradient Descent Scheme. *arXiv.2501.08906*, Jan. 2025.
- [56] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, W. Zhang, B. Cui, and M.-H. Yang. Diffusion Models: A Comprehensive Survey of Methods and Applications. *ACM Computing Surveys*, 56(4):1–39, Apr. 2024.
- [57] T. Yu, S. Kumar, A. Gupta, S. Levine, K. Hausman, and C. Finn. Gradient surgery for multi-task learning. In *Advances in Neural Information Processing Systems*, 2020.