

Orthogonal Over-Parameterized Training

Weiyang Liu^{1,2,*} Rongmei Lin^{3,*} Zhen Liu⁴ James M. Rehg⁵ Liam Paull⁴ Li Xiong³ Le Song⁵ Adrian Weller^{1,6}

¹University of Cambridge ²Max Planck Institute for Intelligent Systems ³Emory University

⁴Mila, Université de Montréal ⁵Georgia Institute of Technology ⁶Alan Turing Institute *Equal Contribution

Abstract

The inductive bias of a neural network is largely determined by the architecture and the training algorithm. To achieve good generalization, how to effectively train a neural network is of great importance. We propose a novel orthogonal over-parameterized training (OPT) framework that can provably minimize the hyperspherical energy which characterizes the diversity of neurons on a hypersphere. By maintaining the minimum hyperspherical energy during training, OPT can greatly improve the empirical generalization. Specifically, OPT fixes the randomly initialized weights of the neurons and learns an orthogonal transformation that applies to these neurons. We consider multiple ways to learn such an orthogonal transformation, including unrolling orthogonalization algorithms, applying orthogonal parameterization, and designing orthogonality-preserving gradient descent. For better scalability, we propose the stochastic OPT which performs orthogonal transformation stochastically for partial dimensions of neurons. Interestingly, OPT reveals that learning a proper coordinate system for neurons is crucial to generalization. We provide some insights on why OPT yields better generalization. Extensive experiments validate the superiority of OPT over the standard training.

1. Introduction

The inductive bias encoded in a neural network is generally determined by two major aspects: how the neural network is structured (*i.e.*, network architecture) and how the neural network is optimized (*i.e.*, training algorithm). For the same network architecture, using different training algorithms could lead to a dramatic difference in generalization performance [36, 60] even if the training loss is close to zero, implying that different training procedures lead to different inductive biases. Therefore, how to effectively train a neural network that generalize well remains an open challenge.

Recent theories [16, 15, 34, 45] suggest the importance of over-parameterization in linear neural networks. For example, [16] shows that optimizing an underdetermined quadratic objective over a matrix M with gradient descent

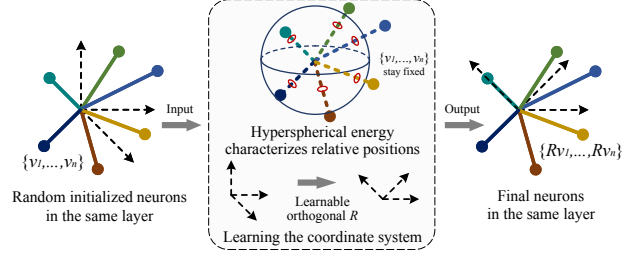


Figure 1: Overview of the orthogonal over-parameterized training framework. OPT learns an orthogonal transformation for each layer in the neural network, while keeping the randomly initialized neuron weights fixed.

on a factorization of M leads to an implicit regularization that may improve generalization. There is also strong empirical evidence [11, 51] that over-parameterizing the convolutional filters under some regularity is beneficial to generalization. Our paper aims to leverage the power of over-parameterization and explore more intrinsic structural priors in order to train a well-performing neural network.

Motivated by this goal, we propose a generic orthogonal over-parameterized training (OPT) framework for neural networks. Different from conventional neural training, OPT over-parameterizes a neuron $w \in \mathbb{R}^d$ with the multiplication of a learnable layer-shared orthogonal matrix $R \in \mathbb{R}^{d \times d}$ and a fixed randomly-initialized weight vector $v \in \mathbb{R}^d$, and it follows that the equivalent weight for the neuron is $w = Rv$. Once each element of the neuron weight v has been randomly initialized by a zero-mean Gaussian distribution [20, 14], we fix them throughout the entire training process. Then OPT learns a layer-shared orthogonal transformation R that is applied to all the neurons (in the same layer). An illustration of OPT is given in Fig. 1. In contrast to standard neural training, OPT decomposes the neuron into an orthogonal transformation R that learns a proper coordinate system, and a weight vector v that controls the specific position of the neuron. Essentially, the weights $\{v_1, \dots, v_n \in \mathbb{R}^d\}$ of different neurons determine the relative positions, while the layer-shared orthogonal matrix R specifies the coordinate system. Such a decoupled parameterization enables strong modeling flexibility.

Another motivation of OPT comes from an empirical observation that neural networks with lower *hyperspherical*

energy generalize better [49]. Hyperspherical energy quantifies the diversity of neurons on a hypersphere, and essentially characterizes the relative positions among neurons via this form of diversity. [49] introduces hyperspherical energy as a regularization in the network but do not guarantee that the hyperspherical energy can be effectively minimized (due to the existence of data fitting loss). To address this issue, we leverage the property of hyperspherical energy that it is independent of the coordinate system in which the neurons live and only depends on their relative positions. Specifically, we prove that, if we randomly initialize the neuron weight \mathbf{v} with certain distributions, these neurons are guaranteed to attain minimum hyperspherical energy in expectation. It follows that OPT maintains the minimum energy during training by learning a coordinate system (*i.e.*, layer-shared orthogonal matrix) for the neurons. Therefore, OPT is able to provably minimize the hyperspherical energy.

We consider several ways to learn the orthogonal transformation. First, we unroll different orthogonalization algorithms such as Gram-Schmidt process, Householder reflection and Löwdin’s symmetric orthogonalization. Different unrolled algorithms yield different implicit regularizations to construct the neuron weights. For example, symmetric orthogonalization guarantees that the new orthogonal basis has the least distance in the Hilbert space from the original non-orthogonal basis. Second, we consider to use a special parameterization (*e.g.*, Cayley parameterization) to construct the orthogonal matrix, which is more efficient in training. Third, we consider an orthogonality-preserving gradient descent to ensure that the matrix \mathbf{R} stays orthogonal after each gradient update. Last, we relax the original optimization problem by making the orthogonality constraint a regularization for the matrix \mathbf{R} . Different ways of learning the orthogonal transformation may encode different inductive biases. We note that OPT aims to utilize orthogonalization as a tool to learn neurons that maintain small hyperspherical energy, rather than to study a specific orthogonalization method. Furthermore, we propose a refinement strategy to reduce the hyperspherical energy for the randomly initialized neuron weights $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$. In specific, we directly minimize the hyperspherical energy of these random weights as a preprocessing step before training them on actual data.

To improve scalability, we further propose the stochastic OPT that randomly samples neuron dimensions to perform orthogonal transformation. The random sampling process is repeated many times such that each dimension of the neuron is sufficiently learned. Finally, we provide some theoretical insights and discussions to justify the effectiveness of OPT. The advantages of OPT are summarized as follows:

- OPT is a generic neural network training framework with strong flexibility. There are many different ways to learn the orthogonal transformations and each one imposes a unique inductive bias. Our paper compares how different

orthogonalizations may affect generalization in OPT.

- OPT is the first training framework where the hyperspherical energy is provably minimized (in contrast to [49]), leading to better empirical generalization. OPT reveals that learning a proper coordinate system is crucial to generalization, and the hyperspherical energy is sufficiently expressive to characterize relative neuron positions.
- There is no extra computational cost for the OPT-trained neural network in inference. In the testing stage, it has the same inference speed and model size as the normally trained network. Our experiments also show that OPT performs well on a diverse class of neural networks and therefore is agnostic to different neural architectures.
- Stochastic OPT can greatly improve the scalability of OPT while enjoying the same guarantee to minimize hyperspherical energy and having comparable performance.

2. Related Work

Orthogonality in Neural Networks. Orthogonality is widely adopted to improve neural networks. [4, 54, 7, 26, 78] use orthogonality as a regularization for neurons. [27, 42, 3, 75, 58, 31] use principled orthogonalization methods to guarantee the neurons are orthogonal to each other. In contrast to these works, OPT does not encourage orthogonality among neurons. Instead, OPT utilizes principled orthogonalization for learning orthogonal transformations for (not necessarily orthogonal) neurons to minimize hyperspherical energy.

Parameterization of Neurons. There are various ways to parameterize a neuron for different applications. [11] over-parameterizes a 2D convolution kernel by combining a 2D kernel of the same size and two additional 1D asymmetric kernels. The resulting convolution kernel has the same effective parameters during testing but more parameters during training. [51] constructs a neuron with a bilinear parameterization and regularizes the bilinear similarity matrix. [79] reparameterizes the neuron matrix with an adaptive fastfood transform to compress model parameters. [30, 48, 73] employ sparse and low-rank structures to construct convolution kernels for a efficient neural network.

Hyperspherical Learning. [54, 52, 72, 10, 71, 47, 50] propose to learn representations on a hypersphere and show that the angular information, in contrast to magnitude information, preserves the most semantic meaning. [49] define the hyperspherical energy that quantifies the diversity of neurons on a hypersphere and shows that the small hyperspherical energy generally improves empirical generalization.

3. Orthogonal Over-Parameterized Training

3.1. General Framework

OPT parameterizes the neuron as the multiplication of an orthogonal matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$ and a neuron weight vector $\mathbf{v} \in \mathbb{R}^d$, and the equivalent neuron weight becomes $\mathbf{w} = \mathbf{R}\mathbf{v}$. The output \hat{y} of this neuron can be represented by $\hat{y} = (\mathbf{R}\mathbf{v})^\top \mathbf{x}$

where $\mathbf{x} \in \mathbb{R}^d$ is the input vector. In OPT, we typically fix the randomly initialized neuron weight \mathbf{v} and only learn the orthogonal matrix \mathbf{R} . In contrast, the standard neuron is directly formulated as $\hat{y} = \mathbf{v}^\top \mathbf{x}$, where the weight vector \mathbf{v} is learned via back-propagation in training.

As an illustrative example, we consider a linear MLP with a loss function \mathcal{L} (e.g., the least squares loss: $\mathcal{L}(e_1, e_2) = (e_1 - e_2)^2$). Specifically, the learning objective of the standard training is $\min_{\{\mathbf{v}_i, u_i, \forall i\}} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i \mathbf{v}_i^\top \mathbf{x}_j)$, while differently, our OPT is formulated as

$$\min_{\{\mathbf{R}, u_i, \forall i\}} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i (\mathbf{R} \mathbf{v}_i)^\top \mathbf{x}_j) \quad \text{s.t. } \mathbf{R}^\top \mathbf{R} = \mathbf{R} \mathbf{R}^\top = \mathbf{I} \quad (1)$$

where $\mathbf{v}_i \in \mathbb{R}^d$ is the i -th neuron in the first layer, and $\mathbf{u} = \{u_1, \dots, u_n\} \in \mathbb{R}^n$ is the output neuron in the second layer. In OPT, each element of \mathbf{v}_i is usually sampled from a zero-mean Gaussian distribution (e.g., both Xavier [14] and Kaiming [20] initializations belong to this class), and is fixed throughout the entire training process. In general, OPT learns an orthogonal matrix that is applied to all the neurons instead of learning the individual neuron weight. Note that, we usually do not apply OPT to neurons in the output layer (e.g., \mathbf{u} in this MLP example, and the final linear classifiers in CNNs), since it makes little sense to fix a set of random linear classifiers. Therefore, the central problem is how to learn these layer-shared orthogonal matrices.

3.2. Hyperspherical Energy Perspective

One of the most important properties of OPT is its invariance to hyperspherical energy. Based on [49], the hyperspherical energy of n neurons is defined as $E(\hat{\mathbf{v}}_i|_{i=1}^n) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n \|\hat{\mathbf{v}}_i - \hat{\mathbf{v}}_j\|^{-1}$ in which $\hat{\mathbf{v}}_i = \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}$ is the i -th neuron weight projected onto the unit hypersphere $\mathbb{S}^{d-1} = \{\mathbf{v} \in \mathbb{R}^d | \|\mathbf{v}\| = 1\}$. Hyperspherical energy is used to characterize the diversity of n neurons on a unit hypersphere. Assume that we have n neurons in one layer, and we have learned an orthogonal matrix \mathbf{R} for these neurons. The hyperspherical energy of these n OPT-trained neurons is

$$E(\hat{\mathbf{R}} \hat{\mathbf{v}}_i|_{i=1}^n) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n \|\mathbf{R} \hat{\mathbf{v}}_i - \mathbf{R} \hat{\mathbf{v}}_j\|^{-1} \quad (2)$$

$$(\text{since } \|\mathbf{R}\|^{-1} = 1) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n \|\hat{\mathbf{v}}_i - \hat{\mathbf{v}}_j\|^{-1} = E(\hat{\mathbf{v}}_i|_{i=1}^n)$$

which verifies that the hyperspherical energy does not change in OPT. Moreover, [49] proves that minimum hyperspherical energy corresponds to the uniform distribution over the hypersphere. As a result, if the initialization of the neurons in the same layer follows the uniform distribution over the hypersphere, then we can guarantee that the hyperspherical energy is minimal in a probabilistic sense.

Theorem 1. For the neuron $\mathbf{h} = \{h_1, \dots, h_d\}$ where $h_i, \forall i$ are initialized i.i.d. following a zero-mean Gaussian distribution (i.e., $h_i \sim N(0, \sigma^2)$), the projections onto a unit hypersphere $\hat{\mathbf{h}} = \mathbf{h} / \|\mathbf{h}\|$ where $\|\mathbf{h}\| = (\sum_{i=1}^d h_i^2)^{1/2}$ are uniformly

distributed on the unit hypersphere \mathbb{S}^{d-1} . The neurons with minimum hyperspherical energy attained asymptotically approach the uniform distribution on \mathbb{S}^{d-1} .

Theorem 1 proves that, as long as we initialize the neurons in the same layer with zero-mean Gaussian distribution, the resulting hyperspherical energy is guaranteed to be small (i.e., the expected energy is minimal). It is because the neurons are uniformly distributed on the unit hypersphere and hyperspherical energy quantifies the uniformity on the hypersphere in some sense. More importantly, prevailing neuron initializations such as [14] and [20] are zero-mean Gaussian distribution. Therefore, our neurons naturally have low hyperspherical energy from the beginning. Appendix L gives geometric properties of the random initialized neurons.

3.3. Unrolling Orthogonalization Algorithms

In order to learn the orthogonal transformation, we unroll classic orthogonalization algorithms and embed them into the

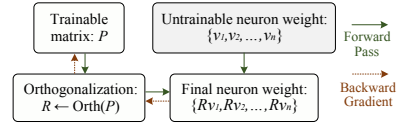


Figure 2: Unrolled orthogonalization.

neural network such that the training can be performed in an end-to-end fashion. We need to make every step of the orthogonalization algorithm differentiable, as shown in Fig. 2.

Gram-Schmidt Process. This method takes a linearly independent set and eventually produces an orthogonal set based on it. The Gram-Schmidt Process (GS) usually takes the following steps to orthogonalize a set of vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_n\} \in \mathbb{R}^{n \times n}$ and obtain an orthonormal set $\{\mathbf{e}_1, \dots, \mathbf{e}_i, \dots, \mathbf{e}_n\} \in \mathbb{R}^{n \times n}$. First, when $i=1$, we have $\mathbf{e}_1 = \frac{\tilde{\mathbf{e}}_1}{\|\tilde{\mathbf{e}}_1\|}$ where $\tilde{\mathbf{e}}_1 = \mathbf{u}_1$. Then, when $n \geq i \geq 2$, we have $\mathbf{e}_i = \frac{\tilde{\mathbf{e}}_i}{\|\tilde{\mathbf{e}}_i\|}$ where $\tilde{\mathbf{e}}_i = \mathbf{u}_i - \sum_{j=1}^{i-1} \text{Proj}_{\mathbf{e}_j}(\mathbf{u}_i)$. Note that, $\text{Proj}_{\mathbf{b}}(\mathbf{a}) = \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\langle \mathbf{b}, \mathbf{b} \rangle} \mathbf{b}$ is defined as the projection operator.

Householder Reflection. A Householder reflector is defined as $\mathbf{H} = \mathbf{I} - 2 \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2}$ where \mathbf{u} is perpendicular to the reflection hyperplane. In QR factorization, Householder reflection (HR) is used to transform a (non-singular) square matrix into an orthogonal matrix and an upper triangular matrix. Given a matrix $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\} \in \mathbb{R}^{n \times n}$, we consider the first column vector \mathbf{u}_1 . We use Householder reflector to transform \mathbf{u}_1 to $\mathbf{e}_1 = \{1, 0, \dots, 0\}$. Specifically, we construct an orthogonal matrix \mathbf{H}_1 with $\mathbf{H}_1 = \mathbf{I} - 2 \frac{(\mathbf{u}_1 - \|\mathbf{u}_1\| \mathbf{e}_1)(\mathbf{u}_1 - \|\mathbf{u}_1\| \mathbf{e}_1)^\top}{\|\mathbf{u}_1 - \|\mathbf{u}_1\| \mathbf{e}_1\|^2}$. The first column of $\mathbf{H}_1 \mathbf{U}$ becomes $\{\|\mathbf{u}_1\|, 0, \dots, 0\}$. At the k -th step, we can view the sub-matrix $\mathbf{U}_{(k:n, k:n)}$ as a new \mathbf{U} , and use the same procedure to construct the Householder transformation $\tilde{\mathbf{H}}_k \in \mathbb{R}^{(n-k) \times (n-k)}$. We construct the final Householder transformation as $\mathbf{H}_k = \text{Diag}(\mathbf{I}_k, \tilde{\mathbf{H}}_k)$. Now we can gradually transform \mathbf{U} to an upper triangular matrix with n Householder reflections. Therefore, we have that $\mathbf{H}_n \cdots \mathbf{H}_2 \mathbf{H}_1 \mathbf{U} = \mathbf{R}^{\text{up}}$ where \mathbf{R}^{up} is an upper triangular matrix and the obtained

orthogonal set is $Q^\top = H_n \cdots H_2 H_1$.

Löwdin’s Symmetric Orthogonalization. Let the matrix $U = \{u_1, \dots, u_n\} \in \mathbb{R}^{n \times n}$ be a given set of linearly independent vectors in an n -dimensional space. A non-singular linear transformation A can transform the basis U to an orthogonal basis R : $R = UA$. The matrix R will be orthogonal if $R^\top R = (UA)^\top UA = A^\top MA = I$ where $M = U^\top U$ is the Gram matrix of the given set U . We obtain a general solution to the orthogonalization problem via the substitution: $A = M^{-\frac{1}{2}} B$ where B is an arbitrary unitary matrix. The specific choice $B = I$ gives the Löwdin’s symmetric orthogonalization (LS): $R = UM^{-\frac{1}{2}}$. We can analytically obtain the symmetric orthogonalization from the singular value decomposition: $U = W\Sigma V^\top$. Then LS gives $R = WV^\top$ as the orthogonal set for U . LS has a unique property which the other orthogonalizations do not have. The orthogonal set resembles the original set in a nearest-neighbour sense. More specifically, LS guarantees that $\sum_i \|R_i - U_i\|^2$ (where R_i and U_i are the i -th column of R and U , respectively) is minimized. Intuitively, LS indicates the gentlest pushing of the directions of the vectors in order to get them orthogonal to each other.

Discussion. These orthogonalization algorithms are fully differentiable and end-to-end trainable. For accurate orthogonality, these algorithms can be used repeatedly and unrolled with multiple steps. Empirically, one-step unrolling already works well. Givens rotations can also construct the orthogonal matrix, but it requires traversing all lower triangular elements in the original set U , which takes $\mathcal{O}(n^2)$ complexity and is too costly. Interestingly, each orthogonalization encodes a unique inductive bias to the neurons by imposing implicit regularizations (e.g., least distance in Hilbert space for LS). Details about these orthogonalizations are in Appendix A. Unrolling orthogonalization has been considered in different scenarios [27, 69, 56]. More orthogonalization methods [41] can be applied in OPT, but exhaustively applying them to OPT is out of the scope of this paper.

3.4. Orthogonal Parameterization

A convenient way to ensure orthogonality while learning the matrix R is to use a special parameterization that inherently guarantees orthogonality. The exponential parameterization use $R = \exp(W)$ (where $\exp(\cdot)$ denotes the matrix exponential) to represent an orthogonal matrix from a skew-symmetric matrix W . The Cayley parameterization (CP) is a Padé approximation of the exponential parameterization, and is a more natural choice due to its simplicity. CP uses the following transform to construct an orthogonal matrix R from a skew-symmetric matrix W : $R = (I + W)(I - W)^{-1}$ where $W = -W^\top$. We note that CP only produces the orthogonal matrices with determinant 1, which belong to the special orthogonal group and thus $R \in SO(n)$. Specifically, it suffices to learn the upper or lower triangular of the matrix W with unconstrained optimization to obtain a

desired orthogonal matrix R . Cayley parameterization does not cover the entire orthogonal group and is less flexible in terms of representation power, which serves as an explicit regularization for the neurons.

3.5. Orthogonality-Preserving Gradient Descent

An alternative way to guarantee orthogonality is to modify the gradient update for the matrix R . The idea is to initialize R with an arbitrary orthogonal matrix and then ensure each gradient update is to apply an orthogonal transformation to R . It is essentially conducting gradient descent on the Stiefel manifold [44, 74, 75, 42, 3, 22, 33]. Given a matrix $U_{(0)} \in \mathbb{R}^{n \times n}$ that is initialized as an orthogonal matrix, we aim to construct an orthogonal transformation as the gradient update. We use the Cayley transform to compute a parametric curve on the Stiefel manifold $\mathcal{M}_s = \{U \in \mathbb{R}^{n \times n} : U^\top U = I\}$ with a specific metric via a skew-symmetric matrix W and use it as the update rule:

$$Y(\lambda) = (I - \frac{\lambda}{2}W)^{-1}(I + \frac{\lambda}{2}W)U_{(i)}, \quad U_{(i+1)} = Y(\lambda) \quad (3)$$

where $\hat{W} = \nabla f(U_{(i)})U_{(i)}^\top - \frac{1}{2}U_{(i)}(U_{(i)}^\top \nabla f(U_{(i)}U_{(i)}^\top))$ and $W = \hat{W} - \hat{W}^\top$. $U_{(i)}$ denotes the orthogonal matrix in the i -th iteration. $\nabla f(U_{(i)})$ denotes the original gradient of the loss function w.r.t. $U_{(i)}$. We term this gradient update as orthogonal-preserving gradient descent (OGD). To reduce the computational cost of the matrix inverse in Eq. 3, we use an iterative method [44] to approximate the Cayley transform without matrix inverse. We arrive at the fixed-point iteration:

$$Y(\lambda) = U_{(i)} + \frac{\lambda}{2}W(U_{(i)} + Y(\lambda)) \quad (4)$$

which converges to the closed-form Cayley transform with a rate of $\mathcal{O}(\lambda^{2+n})$ (n is the iteration number). In practice, two iterations suffice for a reasonable approximation accuracy.

3.6. Relaxation to Orthogonal Regularization

Alternatively, we also consider relaxing the original optimization with an orthogonality constraint to an unconstrained optimization with orthogonality regularization (OR). Specifically, we remove the orthogonality constraint, and adopt an orthogonality regularization for R , i.e., $\|R^\top R - I\|_F^2$. However, OR cannot guarantee the energy stays unchanged. Taking Eq. 1 as an example, the objective becomes

$$\min_{R, u_i, v_i} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i (Rv_i)^\top x_j) + \beta \|R^\top R - I\|_F^2 \quad (5)$$

where β is a hyperparameter. This serves as an relaxation of the original OPT objective. Note that, OR is imposed to R instead of neurons and is quite different from the existing orthogonality regularization on neurons [54, 4, 27, 78, 7].

3.7. Refining the Initialization as Preprocessing

Minimizing the energy beforehand. Because we randomly initialize the neurons $\{v_1, \dots, v_n\}$, there exists a variance that makes the hyperspherical energy deviate from

the minima even if the hyperspherical energy is minimal in a probabilistic sense. To further reduce the hyperspherical energy, we propose to refine the random initialization by minimizing its hyperspherical energy as a preprocessing step before the OPT training. Specifically, before feeding these neurons to OPT, we first minimize the hyperspherical energy of the initialized neurons with gradient descent (without fitting the training data). Moreover, since the randomly initialized neurons cannot guarantee to get rid of the collinearity redundancy as shown in [49] (*i.e.*, two neurons are on the same line but have opposite directions), we can perform the half-space hyperspherical energy minimization [49].

Normalizing the neurons. The norm of the randomly initialized neurons may have some influence on OPT, serving a role similar to weighting the importance of different neurons. Moreover, the norm makes the hyperspherical energy less expressive to characterize the diversity of neurons, as discussed in Section 5.3. To make the coordinate frame (*i.e.* the rotation matrix R) truly independent of the relative positions of the neurons, we propose to normalize the neuron weights such that each neuron has unit norm. Because the weights of the neurons $\{v_1, \dots, v_n\}$ are fixed during training and orthogonal matrices will not change the norm of the neurons, we only need to normalize the randomly initialized neuron weights as a preprocessing before the OPT training.

We have comprehensively evaluated both refinement strategies in Section 6.2 and verified their effectiveness. Note that the effectiveness of OPT is not dependent on these refinements. Our experiments do not use these refinements by default and the results show that OPT still performs well.

4. Towards Better Scalability for OPT

If the dimension of neurons becomes extremely large, then the orthogonal matrix to transform the neurons will also be large. Therefore, it may take large GPU memory and time to train the neural networks with the original OPT. To address this, we propose a scalable variant – stochastic OPT (S-OPT). The key idea of S-OPT is to randomly select some dimensions from the neurons in the same layer and construct a small orthogonal matrix to transform these dimensions together. The selection of dimensions is stochastic in each outer iteration, so a small orthogonal matrix is sufficient to cover all the neuron dimensions. S-OPT aims to approximate a large orthogonal transformation for all the neuron dimensions with many small orthogonal transformations for random subsets of these dimensions, which shares similar spirits with Givens rotation. The approximation will be more accurate when the procedure is randomized over many times. Fig. 3 compares the size of the orthogonal matrix in OPT and S-OPT. The orthogonal matrix in OPT is of size $d \times d$, while the orthogonal matrix in S-OPT is of size $p \times p$ where p is usually much smaller than d . Most

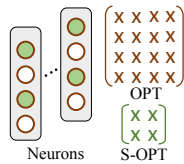


Figure 3: S-OPT.

importantly, S-OPT can still preserve the low hyperspherical energy of neurons because of the following result.

Theorem 2. *For n d -dimensional neurons, selecting any p ($p \leq d$) dimensions and applying an shared orthogonal transformation ($p \times p$ orthogonal matrix) to these p dimensions of all neurons will not change the hyperspherical energy.*

A description of S-OPT is given in Algorithm 1. S-OPT has outer and inner iterations. In each inner iteration, the training is almost the same as OPT, except that the orthogonal matrix transforms a subset of the dimensions and the learnable orthogonal matrix has to be re-initialized to an identity matrix. The selection of neuron dimension is randomized in every outer iteration such that all neuron dimensions can be sufficiently covered as the number of outer iterations increases. Therefore, given sufficient number of iterations, S-OPT will perform comparably to OPT, as empirically verified in Section 6.3. As a parallel direction to improve the scalability, we further propose a parameter-efficient OPT in Appendix I. This OPT variant explores structure priors in R to improve parameter efficiency.

Algorithm 1 Stochastic OPT

```

for  $i = 1, 2, \dots, N_{\text{out}}$  do
  for  $j = 1, 2, \dots, N_{\text{in}}$  do
    1. Randomly select  $p$  dimensions from  $d$ -dimensional neurons in the same layer.
    2. Construct an orthogonal matrix  $R_p \in \mathbb{R}^{p \times p}$  and initialize it as identity matrix.
    3. Update  $R_p$  by applying OPT with one iteration.
  end
  4. Multiply  $R_p$  back to the  $p$ -dim sub-vectors from the  $d$ -dim neurons to transform these neurons.
end

```

5. Intriguing Insights and Discussions

5.1. Local Landscape

We follow [43] to visualize the loss landscapes of both standard training and OPT in Fig. 4. For standard training, we perturb the parameter space of all the neurons (*i.e.*, filters). For OPT, we perturb the parameter space of all the trainable matrices (*i.e.*, P in Fig. 2), because OPT

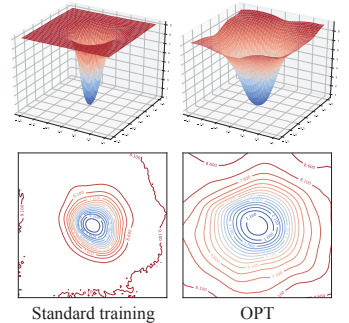


Figure 4: Training loss landscapes.

does not directly learn neuron weights. The general idea is to use two random vectors (*e.g.*, normal distribution) to perturb the parameter space and obtain the loss value with the perturbed network parameters. Details and full results about the visualization are given in Appendix E. The loss landscape of standard training has extremely sharp minima. The red region is very flat, leading to small gradients. In contrast, the loss landscape of OPT is much more smooth and convex with flatter minima, well matching the finding that flat minimizers generalize well [23, 8, 29]. Additional loss landscape visualization results in Appendix F (with uniform perturbation distributions) also support the same argument.

We also show the landscape of testing error on CIFAR-100 in Fig. 5. Full results and details are in Appendix E. Compared to standard training, the testing error of OPT increases more slowly and smoothly while the network parameters move away from the minima, which indicates that the parameter space of OPT yields better robustness to perturbations.

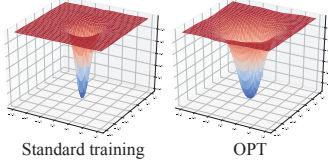


Figure 5: Testing error landscapes.

5.2. Optimization and Generalization

We discuss why OPT may improve optimization and generalization. On one hand, [77] proves that once the neurons are hyperspherically diverse enough in a one-hidden-layer network, the training loss is on the order of the square norm of the gradient and the generalization error will have an additional term $\tilde{O}(1/\sqrt{m})$ where m is the number of samples. This suggests that SGD-optimized networks with minimum hyperspherical energy (MHE) attained have no spurious local minima. Since OPT is guaranteed to achieve MHE in expectation, OPT-trained networks enjoy the inductive bias induced by MHE. On the other hand, [34, 1, 12, 45, 16] shows that over-parameterization in neural networks improves the first-order optimization, leads to better generalization, and imposes implicit regularizations. In the light of this, OPT also introduces over-parameterization to each neuron, which shares similar spirits with [46]. Specifically, one d -dimensional neuron has $d^2 + d$ parameters in OPT (with d^2 being layer-shared), compared to d parameters in a standard neuron. Although OPT uses more parameters for a neuron in training, the equivalent number of parameters for a neuron stays unchanged and it will not affect testing speed.

5.3. Discussions

Over-parameterization. We delve deeper into the over-parameterization in the context of OPT. Its definition varies in different cases. OPT is over-parameterized in terms of training in the following sense. Although OPT-trained networks have the same effective number of parameters as the standard networks in testing, the OPT neuron is decomposed into two sets of parameters in training: orthogonal matrix and neuron weights. It means that the same set of parameters in a neural network can be represented by different sets of training parameters in OPT (*i.e.*, different combinations of orthogonal matrices and neuron weights can lead to the same neural network). OPT is still over-parameterized even if we only count the number of learnable parameters. For a layer of n d -dimensional neurons, the number of learnable parameters in vanilla OPT is $\frac{d(d-1)}{2}$ in contrast to nd in standard training. In prevailing architectures (*e.g.*, ResNet [21]), the neuron dimension is far larger than the number of neurons.

Coordinate system and relative position. OPT shows that learning the coordinate system yields better general-

ization than learning neuron weights directly. This implies that the coordinate system is crucial to generalization. However, the relative position does not matter only when the hyperspherical energy is sufficiently low, indicating that the neurons need to be diverse enough on the unit hypersphere.

The effects of neuron norm. Because we will normalize the neuron norm when computing the hyperspherical energy, the effects of neuron norm will not be taken into consideration. Moreover, simply learning the orthogonal matrices will not change the neuron norm. Therefore, the neuron norm may affect the training. We use an extreme example to demonstrate the effects. Assume that one of the neurons has norm 1000 and the other neurons have norm 0.01. Then no matter what orthogonal matrices we have learned, the final performance will be bad. In this case, the hyperspherical energy can still be minimized to a very low value, but it can not capture the norm distribution. Fortunately, such an extreme case is unlikely to happen, because we are using zero-mean Gaussian distribution to initialize the neuron and every neuron also has the same expected value for the norm. To eliminate the effects of norms, we can normalize the neuron weights in training, as proposed in Section 3.7.

6. Applications and Experimental Results

We put all the experimental settings and many additional results in Appendix D and Appendix I, K, respectively.

6.1. Ablation Study and Exploratory Experiments

Orthogonality. We evaluate whether orthogonality in OPT is necessary. We use 6-layer and 9-layer CNN (Appendix D) on CIFAR-100. Then we compare OPT with unconstrained over-

| Method | FN | LR | CNN-6 | CNN-9 |
|----------|----|----|--------------|--------------|
| Baseline | - | - | 37.59 | 33.55 |
| UPT | ✗ | U | 48.47 | 46.72 |
| UPT | ✓ | U | 42.61 | 39.38 |
| OPT | ✗ | GS | 37.24 | 32.95 |
| OPT | ✓ | GS | 33.02 | 31.03 |

Table 1: Error (%) on C-100.

parameterized training (UPT) which learns an unconstrained matrix \mathbf{R} (with weight decay) using the same network. In Table 1, “FN” denotes whether the randomly initialized neuron weights are fixed in training. “LR” denotes whether the learnable matrix \mathbf{R} is unconstrained (“U”) or orthogonal (“GS” for Gram-Schmidt process). Table 1 shows that without orthogonality, UPT performs much worse than OPT.

Fixed or learnable weights. From Table 1, we can see that using fixed neuron weights is consistently better than learnable neuron weights in both UPT and OPT. It indicates that fixing the neuron weights can well maintain low hyperspherical energy and is beneficial to empirical generalization.

Refining initialization. We evaluate two refinement methods in Section 3.7 for neuron initialization. First, we

| Method | Original | MHE | HS-MHE | CoMHE |
|----------|----------|--------|--------|---------------|
| OPT (GS) | 33.02 | 32.99 | 32.78 | 32.69 |
| OPT (LS) | 34.48 | 34.43 | 34.37 | 34.15 |
| OPT (CP) | 33.53 | 33.50 | 33.42 | 33.27 |
| Energy | 3.5109 | 3.5003 | 3.4976 | 3.4954 |

Table 2: Refining initialized energy.

consider the hyperspherical energy minimization as a preprocessing for the neuron weights. Our experiment uses CNN-6 on CIFAR-100. Specifically, we run gradient descent for 5k

iterations to minimize the objective of MHE/HS-MHE [49] or CoMHE [47] before the training starts. Table 2 shows the hyperspherical energy before and after the preprocessing. All methods start with the same random initialization, so all hyperspherical energies start at 3.5109. Testing errors (%) in Table 2 show that the refinement well improves OPT. Although using advanced regularizations such as CoMHE as pre-processing can improve the performance significantly, we do not use them in the other experiments in order to keep our comparison fair and clean. More different ways to minimize the hyperspherical energy can also be considered [50].

We evaluate the second refinement strategy, *i.e.*, neuron weight normalization. Section 5.3 has explained why normalizing the neuron weights may be useful. After initialization, we normalize all the neuron weights to 1. Since OPT does not change the neuron norm, the neuron will keep the norm as 1. More importantly, the hyperspherical energy will not be affected by the neuron normalization. We conduct classification with CNN-6 on CIFAR-100. Testing errors in Table 3 show that normalizing the neurons greatly improves OPT, validating our previous analysis. Note that, these two refinements are not used by default in other experiments.

High vs. low energy. We validate that high hyperspherical energy corresponds to inferior empirical generalization. To initialize high energy neurons, we use [20] and set the mean as 1e-3, 1e-2, 2e-2, and 3e-2. We experiment on CIFAR-100 with CNN-6. Table 4 (“N/C” denotes not converged) show that higher energy generalizes worse and also leads to difficulty in convergence. We see that a small change in energy can lead to a dramatic generalization gap.

No BatchNorm. We evaluate how OPT performs without BatchNorm (BN) [28]. We perform classification on CIFAR-100 with CNN-6. In Table 5, we see that all OPT variants consistently outperform both the baseline and HS-MHE [49] by a significant margin, validating that OPT can work well without BN. CP achieves the best error with more than 4% lower than standard training.

6.2. Empirical Evaluation on OPT

Multi-layer perceptrons. We evaluate OPT on MNIST with a 3-layer MLP. Appendix D gives specific settings. Table 6 shows the testing error with normal initialization (MLP-N) or Xavier initialization [14] (MLP-X). GS/HR/LS denote different orthogonalization unrolling. CP denotes Cayley parameterization. OGD denotes orthogonal-preserving gradient descent. OR denotes relaxed orthogonal regularization. All OPT variants outperform the others by a large margin.

| Method | w/o Norm | w/ Norm |
|-----------|----------|--------------|
| Baseline | 37.59 | 36.05 |
| OPT (GS) | 33.02 | 32.54 |
| OPT (HR) | 35.67 | 35.30 |
| OPT (LS) | 34.48 | 32.11 |
| OPT (CP) | 33.53 | 32.49 |
| OPT (OGD) | 33.37 | 32.70 |
| OPT (OR) | 34.70 | 33.27 |

Table 3: Normalization (%).

| | Mean Energy | Error (%) |
|------|---------------|--------------|
| 0 | 3.5109 | 32.49 |
| 1e-3 | 3.5117 | 33.11 |
| 1e-2 | 3.5160 | 39.51 |
| 2e-2 | 3.5531 | 53.89 |
| 3e-2 | 3.6761 | N/C |

Table 4: Initial energy.

| Method | Error (%) |
|-----------|--------------|
| Baseline | 38.95 |
| HS-MHE | 36.90 |
| OPT (GS) | 35.61 |
| OPT (HR) | 37.51 |
| OPT (LS) | 35.83 |
| OPT (CP) | 34.88 |
| OPT (OGD) | 35.38 |

Table 5: No BN.

| Method | MNIST | | CIFAR-100 | | | |
|----------------|-------------|-------------|--------------|--------------|--------------|--------------|
| | MLP-N | MLP-X | CNN-6 | CNN-9 | ResNet-20 | ResNet-32 |
| Baseline | 6.05 | 2.14 | 37.59 | 33.55 | 31.11 | 30.16 |
| Orthogonal [7] | 5.78 | 1.93 | 36.32 | 33.24 | 31.06 | 30.05 |
| SRIP [4] | - | - | 34.82 | 32.72 | 30.89 | 29.70 |
| HS-MHE [49] | 5.57 | 1.88 | 34.97 | 32.87 | 30.98 | 29.76 |
| OPT (GS) | 5.11 | 1.45 | 33.02 | 31.03 | 30.49 | 29.34 |
| OPT (HR) | 5.31 | 1.60 | 35.67 | 32.75 | 30.73 | 29.56 |
| OPT (LS) | 5.32 | 1.54 | 34.48 | 31.22 | 30.51 | 29.42 |
| OPT (CP) | 5.14 | 1.49 | 33.53 | 31.28 | 30.47 | 29.31 |
| OPT (OGD) | 5.38 | 1.56 | 33.33 | 31.47 | 30.50 | 29.39 |
| OPT (OR) | 5.41 | 1.78 | 34.70 | 32.63 | 30.66 | 29.47 |

Table 6: Testing error (%) of OPT for MLPs and CNNs.

Convolutional networks. We evaluate OPT with 6/9-layer plain CNNs and ResNet-20/32 [21] on CIFAR-100. Detailed settings are in Appendix D. All neurons (*i.e.*, convolution kernels) are initialized by [20]. BatchNorm is used by default. Table 6 shows that all OPT variants outperform both baseline and HS-MHE by a large margin. HS-MHE puts the hyperspherical energy into the loss function and naively minimizes it along with the CNN. We observe that OPT (HR) performs the worse among all OPT variants partially because of its intensive unrolling computation. OPT (GS) achieves the best testing error on CNN-6/9, while OPT (CP) achieves the best testing error on ResNet-20/34, implying that different OPT encodes different inductive bias.

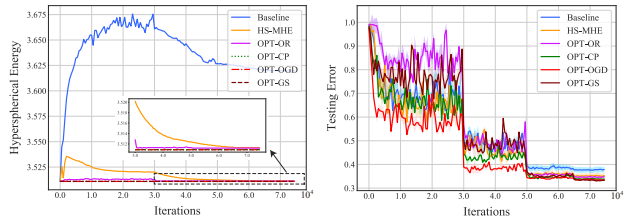


Figure 6: Training dynamics on CIFAR-100. Left: Hyperspherical energy vs. iteration. Right: Testing error vs. iteration.

Training dynamics. We look into how hyperspherical energy and testing error changes in OPT. Fig. 6 shows that the energy of the baseline will increase dramatically at the beginning and then gradually go down, but it still stays in a high value in the end. HS-MHE well reduces the energy at the end of the training. In contrast, OPT variants always maintain very small energy in training. OPT with GS, CP and OGD keep exactly the same energy as the random initialization, while OPT (OR) slightly increases the energy due to relaxation. All OPT variants converge efficiently and stably.

Large-scale learning. To see how OPT performs in large-scale settings, we evaluate OPT on the large-scale ImageNet-2012 [62]. Specifically, we use OPT with OGD and CP to train a plain 10-layer CNN (Appendix D) on ImageNet. Note that, our purpose is to validate the superiority of OPT over the corresponding baseline rather than achieving state-of-the-art results. Table 7 shows that OPT (CP) reduces top-1 and top-5 error for the baseline by $\sim 0.7\%$ and $\sim 0.9\%$, respectively.

| Method | Top-1 | Top-5 |
|----------------|--------------|--------------|
| Baseline | 44.32 | 21.13 |
| Orthogonal [7] | 44.13 | 20.97 |
| HS-MHE [49] | 43.92 | 20.85 |
| OPT (OGD) | 43.81 | 20.49 |
| OPT (CP) | 43.67 | 20.26 |

Table 7: ImageNet (%).

Few-shot recognition. For evaluating OPT on cross-task generalization, we perform the few-shot recognition on Mini-ImageNet, following the same setup as [9]. Appendix D gives more detailed settings. We apply OPT with CP to train the baseline and baseline++ in [9], and immediately obtain improvements. Therefore, OPT-trained networks generalize well in this challenging scenario.

Geometric learning. We apply OPT to graph convolution network (GCN) [37] and point cloud network (PointNet) [57] for graph node and point cloud classification, respectively. The training of GCN and PointNet is conceptually similar to MLP, and the detailed training procedures are given in Appendix D. For GCN, we evaluate OPT on Cora and Pubmed datasets [63]. For PointNet, we conduct experiments on ModelNet-40 dataset [76]. Table 9 shows that OPT effectively improves both GCN and PointNet.

6.3. Empirical Evaluation on S-OPT

| Method | CIFAR-100 | | | | ImageNet | |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | CNN-6 | Params | Wide CNN-9 | Params | ResNet-18 | Params |
| Baseline | 37.59 | 258K | 28.03 | 2.99M | 32.95 | 11.7M |
| HS-MHE [49] | 34.97 | 258K | 25.96 | 2.99M | 32.50 | 11.7M |
| OPT (GS) | 33.02 | 1.36M | OOM | 16.2M | OOM | 46.5M |
| S-OPT (GS) | 33.70 | 90.9K | 25.59 | 1.04M | 32.26 | 3.39M |

Table 10: OPT vs. S-OPT on CIFAR-100 & ImageNet.

Convolutional networks. S-OPT is a scalable OPT variant, and we evaluate its performance in terms of number of *trainable parameters* and testing error. Training parameters are learnable variables in training, and are different from model parameters in testing. In testing, all methods have the same number of model parameters. We perform classification on CIFAR-100 with CNN-6 and wide CNN-9. We also evaluate S-OPT with standard ResNet-18 on ImageNet. Detailed settings are in Appendix D. For S-OPT, we set the sampling dimension as 25% of the original neuron dimension in each layer. Table 10 shows that S-OPT achieves a good trade-off between accuracy and scalability. More importantly, S-OPT can be applied to large neural networks, making OPT more useful in practice. Additionally, Appendix I discusses an efficient parameter sharing for OPT.

Sampling dimensions. We study how the sampling dimension p affect the performance by performing classification with wide CNN-9 on CIFAR-100. In Table 11, $p = d/4$ means that we randomly sample 1/4 of the original neuron dimension in each layer, so p may vary in different layer. $p = 16$ means that we sample 16 dimensions in each layer. Note that there are 25.6K pa-

| Method | 5-shot Acc. (%) |
|-------------------|------------------------------------|
| MAML [13] | 62.71 \pm 0.71 |
| MatchingNet [70] | 63.48 \pm 0.66 |
| ProtoNet [65] | 64.24 \pm 0.72 |
| Baseline [9] | 62.53 \pm 0.69 |
| Baseline w/ OPT | 63.27 \pm 0.68 |
| Baseline++ [9] | 66.43 \pm 0.63 |
| Baseline++ w/ OPT | 66.82 \pm 0.62 |

Table 8: Few-shot learning.

| Method | GCN | | PointNet |
|-----------|-------------|-------------|--------------|
| | Cora | Pubmed | MN-40 |
| Baseline | 81.3 | 79.0 | 87.1 |
| OPT (GS) | 81.9 | 79.4 | 87.23 |
| OPT (CP) | 82.0 | 79.4 | 87.81 |
| OPT (OGD) | 82.3 | 79.5 | 87.86 |

Table 9: Geometric networks.

rameters used for the final classification layer, which can not be saved in S-OPT. Table 11 shows that S-OPT can achieve highly competitive accuracy with a reasonably large p .

6.4. Large Categorical Training

Previously, OPT is not applied to the final classification layer, since it makes little sense to fix random classifiers and learn an orthogonal matrix to transform them. However, learning the classification layer can be costly with large number of classes. The number of trainable parameters of the classification layer grows linearly with the number of classes. To address this, OPT can be used to learn the classification layer, because its number of trainable parameters only depends on the classifier dimension. To be fair, we *only* learn the last classification layer with OPT and the other layers are normally learned (CLS-OPT). The oracle learns the entire network normally. Experimental details are in Appendix D.

We intuitively compare the oracle and CLS-OPT by visualizing the deep MNIST features following [53]. The features are the direct outputs of CNN by setting the output dimension as 3. Fig. 7 show

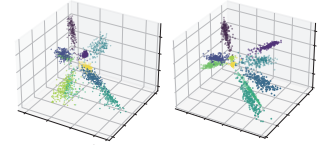


Figure 7: Feature visualization.

that even if CLS-OPT fixes randomly initialized classifiers, it can still learn discriminative and separable deep features.

We evaluate its performance on ImageNet with 1K classes. We use ResNet-18 with different output dimensions (A:128, B:512).

| Method | ResNet-18A | | ResNet-18B | |
|---------|--------------|--------------|--------------|-------------|
| | Error | Params | Error | Params |
| Oracle | 18.08 | 64.0K | 12.12 | 512K |
| CLS-OPT | 21.12 | 8.13K | 12.05 | 131K |

Table 12: CLS-OPT on ImageNet.

Table 12 gives the top-5 test error (%) and “Params” denotes the number of trainable parameters in the classification layer. CLS-OPT performs well with far less trainable parameters.

Since face datasets usually contain large number of identities [17], it is natural to apply CLS-OPT to learn face embeddings. We train

| Method | 512 Dim. | | 1024 Dim. | |
|---------|-------------|-------------|-------------|-------------|
| | Acc. | Params | Acc. | Params |
| Oracle | 95.7 | 5.41M | 96.4 | 10.83M |
| CLS-OPT | 94.9 | 131K | 95.8 | 524K |

Table 13: Verification (%) on LFW.

on CASIA [80] which has 0.5M face images of 10,572 identities, and test on LFW [25]. Since the training and testing sets do not overlap, the task well evaluates the generalizability of learned features. All methods use CNN-20 [52] and standard softmax loss. We set the output feature dimension as 512 or 1024. Table 13 validates CLS-OPT’s effectiveness.

7. Concluding Remarks

We propose a novel training framework for neural networks. By parameterizing neurons with weights and a shared orthogonal matrix, OPT can provably achieve small hyperspherical energy and yield superior generalizability.

Acknowledgements. Weiyang Liu and Adrian Weller are supported by DeepMind and the Leverhulme Trust via CFI. Adrian Weller acknowledges support from the Alan Turing Institute under EPSRC grant EP/N510129/1 and U/B/000074. Rongmei Lin and Li Xiong are supported by NSF under CNS-1952192, IIS-1838200.

| $p =$ | Error (%) | Params |
|--------|--------------|--------------|
| d | OOM | 16.2M |
| $d/4$ | 25.59 | 1.04M |
| $d/8$ | 28.61 | 278K |
| $d/16$ | 32.52 | 88.7K |
| 16 | 33.03 | 27.0K |
| 3 | 45.22 | 26.0K |
| 0 | 60.64 | 25.6K |

Table 11: Sampling dim.

References

- [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018. 6, 36
- [2] Ramesh Naidu Annavarapu. Singular value decomposition and the centrality of löwdin orthogonalizations. *American Journal of Computational and Applied Mathematics*, 3(1):33–35, 2013. 16
- [3] Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In *ICML*, 2016. 2, 4
- [4] Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regularizations in training deep cnns? In *NeurIPS*, 2018. 2, 4, 7
- [5] Johann S Brauchart, Alexander B Reznikov, Edward B Saff, Ian H Sloan, Yu Guang Wang, and Robert S Womersley. Random point sets on the sphere—hole radii, covering, and separation. *Experimental Mathematics*, 27(1):62–81, 2018. 42
- [6] Anna Breger, Martin Ehler, and Manuel Gräf. Points on manifolds with asymptotically optimal covering radius. *Journal of Complexity*, 48:1–14, 2018. 42
- [7] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. In *ICLR*, 2017. 2, 4, 7
- [8] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient descent into wide valleys. In *ICLR*, 2017. 5
- [9] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019. 8, 20
- [10] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 2
- [11] Xiaohan Ding, Yuchen Guo, Guiguang Ding, and Jungong Han. Acnet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks. In *ICCV*, 2019. 1, 2
- [12] Simon S Du, Jason D Lee, Yuandong Tian, Barnabas Poczos, and Aarti Singh. Gradient descent learns one-hidden-layer cnn: Don’t be afraid of spurious local minima. *arXiv preprint arXiv:1712.00779*, 2017. 6, 36
- [13] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 8
- [14] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. 1, 3, 7
- [15] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In *NeurIPS*, 2018. 1
- [16] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. In *NeurIPS*, 2017. 1, 6
- [17] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *ECCV*, 2016. 8
- [18] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. 40
- [19] DP Hardin and EB Saff. Minimal riesz energy point configurations for rectifiable d-dimensional manifolds. *Advances in Mathematics*, 193(1):174–204, 2005. 18
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015. 1, 3, 7
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6, 7, 20, 21
- [22] Mikael Henaff, Arthur Szlam, and Yann LeCun. Recurrent orthogonal networks and long-memory tasks. *arXiv preprint arXiv:1602.06662*, 2016. 4
- [23] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Computation*, 1997. 5
- [24] Walter Hoffmann. Iterative algorithms for gram-schmidt orthogonalization. *Computing*, 41(4):335–348, 1989. 14
- [25] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. *Technical Report*, 2008. 8
- [26] Lei Huang, Li Liu, Fan Zhu, Diwen Wan, Zehuan Yuan, Bo Li, and Ling Shao. Controllable orthogonalization in training dnns. In *CVPR*, 2020. 2
- [27] Lei Huang, Xianglong Liu, Bo Lang, Adams Wei Yu, Yongliang Wang, and Bo Li. Orthogonal weight normalization: Solution to optimization over multiple dependent stiefel manifolds in deep neural networks. In *AAAI*, 2018. 2, 4
- [28] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 7
- [29] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. In *UAI*, 2018. 5
- [30] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. In *BMVC*, 2014. 2
- [31] Kui Jia, Shuai Li, Yuxin Wen, Tongliang Liu, and Dacheng Tao. Orthogonal deep neural networks. *TPAMI*, 2019. 2
- [32] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *NeurIPS*, 2016. 40
- [33] Li Jing, Yichen Shen, Tena Dubcek, John Peurifoy, Scott Skirlo, Yann LeCun, Max Tegmark, and Marin Soljačić. Tunable efficient unitary neural networks (eunn) and their application to rnns. In *ICML*, 2017. 4
- [34] Kenji Kawaguchi. Deep learning without poor local minima. In *NeurIPS*, 2016. 1, 6, 36
- [35] Kenji Kawaguchi, Bo Xie, and Le Song. Deep semi-random features for nonlinear function approximation. In *AAAI*, 2018. 37
- [36] Nitish Shirish Keskar and Richard Socher. Improving generalization performance by switching from adam to sgd. *arXiv preprint arXiv:1712.07628*, 2017. 1
- [37] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 8
- [38] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 21
- [39] N.S. Landkof. *Foundations of modern potential theory*. Springer-Verlag, 1972. 18
- [40] Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht. Gradient descent only converges to minimizers. In *COLT*, 2016. 36
- [41] Mario Lezcano-Casado. Trivializations for gradient-based optimization on manifolds. In *NeurIPS*, 2019. 4
- [42] Mario Lezcano-Casado and David Martínez-Rubio. Cheap orthogonal constraints in neural networks: A simple parametrization of the orthogonal and unitary group. In *ICML*, 2019. 2, 4
- [43] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *NeurIPS*, 2018. 5, 23, 30
- [44] Jun Li, Fuxin Li, and Sinisa Todorovic. Efficient riemannian optimization on the stiefel manifold via the cayley transform. In *ICLR*, 2020. 4
- [45] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *COLT*, 2018. 1, 6
- [46] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013. 6, 40

- [47] Rongmei Lin, Weiyang Liu, Zhen Liu, Chen Feng, Zhiding Yu, James M. Rehg, Li Xiong, and Le Song. Regularizing neural networks via minimizing hyperspherical energy. In *CVPR*, 2020. 2, 7
- [48] Baoyuan Liu, Min Wang, Hassan Foroosh, Marshall Tappen, and Marianna Pinsky. Sparse convolutional neural networks. In *CVPR*, 2015. 2
- [49] Weiyang Liu, Rongmei Lin, Zhen Liu, Lixin Liu, Zhiding Yu, Bo Dai, and Le Song. Learning towards minimum hyperspherical energy. In *NeurIPS*, 2018. 2, 3, 5, 7, 8, 23, 30, 36, 43
- [50] Weiyang Liu, Rongmei Lin, Zhen Liu, Li Xiong, Bernhard Schölkopf, and Adrian Weller. Learning with hyperspherical uniformity. In *AISTATS*, 2021. 2, 7
- [51] Weiyang Liu, Zhen Liu, James M Rehg, and Le Song. Neural similarity learning. In *NeurIPS*, 2019. 1, 2, 40
- [52] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *CVPR*, 2017. 2, 8, 22
- [53] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks. In *ICML*, 2016. 8
- [54] Weiyang Liu, Yan-Ming Zhang, Xingguo Li, Zhiding Yu, Bo Dai, Tuo Zhao, and Le Song. Deep hyperspherical learning. In *NeurIPS*, 2017. 2, 4, 21
- [55] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *ECCV*, 2018. 37
- [56] Zakaria Mhammedi, Andrew Hellicar, Ashfaqur Rahman, and James Bailey. Efficient orthogonal parametrisation of recurrent neural networks using householder reflections. In *ICML*, 2017. 4
- [57] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017. 8, 21
- [58] Haozhi Qi, Chong You, Xiaolong Wang, Yi Ma, and Jitendra Malik. Deep isometric learning for visual recognition. *arXiv preprint arXiv:2006.16992*, 2020. 2
- [59] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NeurIPS*, 2008. 37
- [60] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019. 1
- [61] Matthias Reitzner. Stochastic approximation of smooth convex bodies. *Mathematika*, 51(1-2):11–29, 2004. 42
- [62] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 7
- [63] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 2008. 8, 21
- [64] Ian H Sloan and Robert S Womersley. Extremal systems of points and numerical integration on the sphere. *Advances in Computational Mathematics*, 21(1-2):107–125, 2004. 43
- [65] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, 2017. 8
- [66] Daniel Soudry and Yair Carmon. No bad local minima: Data independent training error guarantees for multilayer neural networks. *arXiv preprint arXiv:1605.08361*, 2016. 36
- [67] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 2014. 37
- [68] Vipin Srivastava. A unified view of the orthogonalization methods. *Journal of Physics A: Mathematical and General*, 33(35):6219, 2000. 16
- [69] Yun Tang, Jing Huang, Guangtao Wang, Xiaodong He, and Bowen Zhou. Orthogonal relation transforms with graph context modeling for knowledge graph embedding. In *ACL*, 2020. 4
- [70] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *NeurIPS*, 2016. 8
- [71] Feng Wang, Weiyang Liu, Haijun Liu, and Jian Cheng. Additive margin softmax for face verification. *arXiv preprint arXiv:1801.05599*, 2018. 2
- [72] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *CVPR*, 2018. 2
- [73] Min Wang, Baoyuan Liu, and Hassan Foroosh. Factorized convolutional neural networks. In *ICCV Workshops*, 2017. 2
- [74] Zaiwen Wen and Wotao Yin. A feasible method for optimization with orthogonality constraints. *Mathematical Programming*, 2013. 4
- [75] Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-capacity unitary recurrent neural networks. In *NeurIPS*, 2016. 2, 4
- [76] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, 2015. 8, 21
- [77] Bo Xie, Yingyu Liang, and Le Song. Diverse neural network learns true target functions. In *AISTATS*, 2017. 6, 36, 37
- [78] Di Xie, Jiang Xiong, and Shiliang Pu. All you need is beyond a good init: Exploring better solution for training extremely deep convolutional neural networks with orthonormality and modulation. In *CVPR*, 2017. 2, 4
- [79] Zichao Yang, Marcin Moczulski, Misha Denil, Nando de Freitas, Alex Smola, Le Song, and Ziyu Wang. Deep fried convnets. In *ICCV*, 2015. 2
- [80] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014. 8

Appendix

Table of Contents

| | |
|--|-----------|
| 1. Introduction | 1 |
| 2. Related Work | 2 |
| 3. Orthogonal Over-Parameterized Training | 2 |
| 3.1. General Framework | 2 |
| 3.2. Hyperspherical Energy Perspective | 3 |
| 3.3. Unrolling Orthogonalization Algorithms | 3 |
| 3.4. Orthogonal Parameterization | 4 |
| 3.5. Orthogonality-Preserving Gradient Descent | 4 |
| 3.6. Relaxation to Orthogonal Regularization | 4 |
| 3.7. Refining the Initialization as Preprocessing | 4 |
| 4. Towards Better Scalability for OPT | 5 |
| 5. Intriguing Insights and Discussions | 5 |
| 5.1. Local Landscape | 5 |
| 5.2. Optimization and Generalization | 6 |
| 5.3. Discussions | 6 |
| 6. Applications and Experimental Results | 6 |
| 6.1. Ablation Study and Exploratory Experiments | 6 |
| 6.2. Empirical Evaluation on OPT | 7 |
| 6.3. Empirical Evaluation on S-OPT | 8 |
| 6.4. Large Categorical Training | 8 |
| 7. Concluding Remarks | 8 |
| A Details of Unrolled Orthogonalization Algorithms | 13 |
| A.1. Gram-Schmidt Process | 13 |
| A.2. Householder Reflection | 14 |
| A.3. Löwdin’s Symmetric Orthogonalization | 16 |
| B Proof of Theorem 1 | 17 |
| C Proof of Theorem 2 | 19 |
| D Experimental Settings | 20 |
| E Loss Landscape Visualization (Normal Distribution Perturbation) | 23 |
| E.1. Visualization Procedure | 23 |
| E.2. Experimental Details | 23 |
| E.3. Full Visualization Results for the Main Paper | 23 |
| E.4. Loss Landscape Visualization of Different Neural Network | 26 |
| E.5. Loss Landscape Visualization of Different Dataset | 28 |
| F. Loss Landscape Visualization (Uniform Distribution Perturbation) | 30 |
| F.1. Visualization Procedure | 30 |

| | |
|--|-----------|
| F.2. Experimental Details | 30 |
| F.3. Full Visualization Results | 30 |
| F.4. Loss Landscape Visualization of Different Neural Network | 33 |
| F.5. Loss Landscape Visualization on Different Dataset | 34 |
| G Theoretical Discussion on Optimization and Generalization | 36 |
| H More Discussions | 37 |
| I. On Parameter-Efficient OPT | 38 |
| I.1 . Formulation | 38 |
| I.2 . Experiments and Results | 39 |
| J. On Generalizing OPT: Over-Parameterized Training with Constraint | 40 |
| K Hyperspherical Energy Training Dynamics of Individual Layers | 41 |
| L Geometric Properties of Randomly Initialized Neurons | 42 |

A. Details of Unrolled Orthogonalization Algorithms

A.1. Gram-Schmidt Process

Gram-Schmidt Process. GS process is a method for orthonormalizing a set of vectors in an inner product space, *i.e.*, the Euclidean space \mathbb{R}^n equipped with the standard inner product. Specifically, GS process performs the following operations to orthogonalize a set of vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_n\} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned}
 \text{Step 1: } \tilde{\mathbf{e}}_1 &= \mathbf{u}_1, \quad \mathbf{e}_1 = \frac{\tilde{\mathbf{e}}_1}{\|\tilde{\mathbf{e}}_1\|} \\
 \text{Step 2: } \tilde{\mathbf{e}}_2 &= \mathbf{u}_2 - \text{Proj}_{\tilde{\mathbf{e}}_1}(\mathbf{u}_2), \quad \mathbf{e}_2 = \frac{\tilde{\mathbf{e}}_2}{\|\tilde{\mathbf{e}}_2\|} \\
 \text{Step 3: } \tilde{\mathbf{e}}_3 &= \mathbf{u}_3 - \text{Proj}_{\tilde{\mathbf{e}}_1}(\mathbf{u}_3) - \text{Proj}_{\tilde{\mathbf{e}}_2}(\mathbf{u}_3), \quad \mathbf{e}_3 = \frac{\tilde{\mathbf{e}}_3}{\|\tilde{\mathbf{e}}_3\|} \\
 \text{Step 4: } \tilde{\mathbf{e}}_4 &= \mathbf{u}_4 - \text{Proj}_{\tilde{\mathbf{e}}_1}(\mathbf{u}_4) - \text{Proj}_{\tilde{\mathbf{e}}_2}(\mathbf{u}_4) - \text{Proj}_{\tilde{\mathbf{e}}_3}(\mathbf{u}_4), \quad \mathbf{e}_4 = \frac{\tilde{\mathbf{e}}_4}{\|\tilde{\mathbf{e}}_4\|} \\
 &\vdots \\
 \text{Step n: } \tilde{\mathbf{e}}_n &= \mathbf{u}_n - \text{Proj}_{\tilde{\mathbf{e}}_1}(\mathbf{u}_n) - \text{Proj}_{\tilde{\mathbf{e}}_2}(\mathbf{u}_n) - \text{Proj}_{\tilde{\mathbf{e}}_3}(\mathbf{u}_n) - \dots - \text{Proj}_{\tilde{\mathbf{e}}_{n-1}}(\mathbf{u}_n), \quad \mathbf{e}_n = \frac{\tilde{\mathbf{e}}_n}{\|\tilde{\mathbf{e}}_n\|}
 \end{aligned} \tag{6}$$

where $\text{Proj}_{\mathbf{a}}(\mathbf{b}) = \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\langle \mathbf{a}, \mathbf{a} \rangle} \mathbf{a}$ denotes the projection of the vector \mathbf{b} onto the vector \mathbf{a} . The set $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$ denotes the output orthonormal set. The algorithm flowchart can be described as follows:

Algorithm 2 Gram-Schmidt Process

Input: $U = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\} \in \mathbb{R}^{n \times n}$
Output: $R = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \in \mathbb{R}^{n \times n}$
 $R = \mathbf{0}$
for $j = 1, 2, \dots, n$ **do**
 $\mathbf{q}_j = R^\top \mathbf{u}_j$
 $\mathbf{t} = \mathbf{u}_j - R \mathbf{q}_j$
 $q_{jj} = \|\mathbf{t}\|_2$
 $\mathbf{e}_j = \frac{\mathbf{t}}{q_{jj}}$
end

The vectors $\mathbf{q}_j, \forall j$ in the algorithm above are used to compute the QR factorization, which is not useful in orthogonalization and therefore does not need to be stored. When the GS process is implemented on a finite-precision computer, the vectors $\mathbf{e}_j, \forall j$ are often not quite orthogonal, because of rounding errors. Besides the standard GS process, there is a modified Gram-Schmidt (MGS) algorithm which enjoys better numerical stability. This approach gives the same result as the original formula in exact arithmetic and introduces smaller errors in finite-precision arithmetic. Specifically, GS computes the following formula:

$$\begin{aligned}
 \tilde{\mathbf{e}}_j &= \mathbf{u}_j - \sum_{k=1}^{j-1} \text{Proj}_{\tilde{\mathbf{e}}_k}(\mathbf{u}_j) \\
 \mathbf{e}_j &= \frac{\tilde{\mathbf{e}}_j}{\|\tilde{\mathbf{e}}_j\|}
 \end{aligned} \tag{7}$$

Instead of computing the vector \mathbf{e}_j as in Eq. 7, MGS computes the orthogonal basis differently. MGS does not subtract the projections of the original vector set, and instead remove the projection of the previously constructed orthogonal basis.

Specifically, MGS computes the following series of formulas:

$$\begin{aligned}
\tilde{e}_j^{(1)} &= \mathbf{u}_j - \text{Proj}_{\tilde{e}_1}(\mathbf{u}_j) \\
\tilde{e}_j^{(2)} &= \tilde{e}_j^{(1)} - \text{Proj}_{\tilde{e}_2}(\tilde{e}_j^{(1)}) \\
&\vdots \\
\tilde{e}_j^{(j-2)} &= \tilde{e}_j^{(j-3)} - \text{Proj}_{\tilde{e}_{j-2}}(\tilde{e}_j^{(j-3)}) \\
\tilde{e}_j^{(j-1)} &= \tilde{e}_j^{(j-2)} - \text{Proj}_{\tilde{e}_{j-1}}(\tilde{e}_j^{(j-2)}) \\
\mathbf{e}_j &= \frac{\tilde{e}_j^{(j-1)}}{\|\tilde{e}_j^{(j-1)}\|}
\end{aligned} \tag{8}$$

where each step finds a vector $\tilde{e}_j^{(i)}$ that is orthogonal to $\tilde{e}_j^{(i-1)}$. Therefore, $\tilde{e}_j^{(i)}$ is also orthogonalized against any errors brought by the computation of $\tilde{e}_j^{(i-1)}$. In practice, although MGS enjoys better numerical stability, we find the empirical performance of GS and MGS is almost the same in OPT. However, MGS takes longer time to complete since the computation of each orthogonal basis is an iterative process. Therefore, we usually stick to classic GS for OPT.

Iterative Gram-Schmidt Process. Iterative Gram-Schmidt (IGS) process is an iterative version of the GS process. It is shown in [24] that GS process can be carried out iteratively to obtain a basis matrix that is orthogonal in almost full working precision. The IGS algorithm is given as follows:

Algorithm 3 Iterative Gram-Schmidt Process

Input: $U = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\} \in \mathbb{R}^{n \times n}$
Output: $R = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \in \mathbb{R}^{n \times n}$
 $R = 0$
for $j = 1, 2, \dots, n$ **do**
 $\mathbf{q}_j = 0$
 $\mathbf{t} = \mathbf{u}_j$
 while $\mathbf{t} \perp \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_{j-1})$ **is False** **do**
 $\mathbf{p} = \mathbf{t}$
 $\mathbf{s} = R^\top \mathbf{p}$
 $\mathbf{v} = R\mathbf{s}$
 $\mathbf{t} = \mathbf{p} - \mathbf{v}$
 $\mathbf{q}_j \leftarrow \mathbf{q}_j + \mathbf{s}$
 end
 $q_{jj} = \|\mathbf{t}\|_2$
 $\mathbf{e}_j = \frac{\mathbf{t}}{q_{jj}}$
end

The vectors $\mathbf{q}_j, \forall j$ in the algorithm above are used to compute the QR factorization, which is not useful in orthogonalization and therefore does not need to be explicitly computed. The while loop in IGS is an iterative procedure. In practice, we can unroll a fixed number of steps for the while loop in order to improve the orthogonality. The resulting \mathbf{q}_j in the j -th step corresponds to the solution of the equation $\tilde{R}^\top \tilde{R}\mathbf{q}_j = \tilde{R}^\top \mathbf{u}_j$ where $\tilde{R} = \{\mathbf{e}_1, \dots, \mathbf{e}_{j-1}\}$. The IGS process corresponds to the Gauss-Jacobi iteration for solving this equation.

Both GS and IGS are easy to be embedded in the neural networks, since they are both differentiable. In our experiments, we find that the performance gain of unrolling multiple steps in IGS over GS is not very obvious (partially because GS has already achieved nearly perfect orthogonality), but IGS costs longer training time. Therefore, we unroll the classic GS process by default.

A.2. Householder Reflection

Let $\mathbf{v} \in \mathbb{R}^n$ be a non-zero vector. A matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ of the form

$$\mathbf{H} = \mathbf{I} - \frac{2\mathbf{v}\mathbf{v}^\top}{\mathbf{v}^\top \mathbf{v}} \tag{9}$$

is a Householder reflection. The vector \mathbf{v} is the Householder vector. If a vector \mathbf{x} is multiplied by the matrix \mathbf{H} , then it will be reflected in the hyperplane $\text{span}(\mathbf{v})^\perp$. Householder matrices are symmetric and orthogonal.

For a vector $\mathbf{x} \in \mathbb{R}^n$, we let $\mathbf{v} = \mathbf{x} \pm \|\mathbf{x}\|_2 \mathbf{e}_1$ where \mathbf{e}_1 is a vector of $\{1, 0, \dots, 0\}$ (the first element is 1 and the remaining elements are 0). Then we construct the Householder reflection matrix with \mathbf{v} and multiply it to \mathbf{x} :

$$\mathbf{H}\mathbf{x} = \left(\mathbf{I} - 2 \frac{\mathbf{v}\mathbf{v}^\top}{\mathbf{v}^\top \mathbf{v}} \right) \mathbf{x} = \mp \|\mathbf{x}\|_2 \mathbf{e}_1 \quad (10)$$

which indicates that we can make any non-zero vector become $\alpha \mathbf{e}_1$ where α is some constant by using Householder reflection. By left-multiplying a reflection we can turn a dense vector \mathbf{x} into a vector with the same length and with only a single nonzero entry. Repeating this n times gives us the Householder QR factorization, which also orthogonalizes the original input matrix. Householder reflection orthogonalizes a matrix $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ by triangularizing it:

$$\mathbf{U} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_n \mathbf{R} \quad (11)$$

where \mathbf{R} is a upper-triangular matrix in the QR factorization. $\mathbf{H}_j, j \geq 2$ is constructed by $\text{Diag}(\mathbf{I}_{j-1}, \tilde{\mathbf{H}}_{n-j+1})$ where $\tilde{\mathbf{H}}_{n-j+1} \in \mathbb{R}^{(n-j+1) \times (n-j+1)}$ is the Householder reflection that is performed on the vector $\mathbf{U}_{(j:n,j)}$. The algorithm flowchart is given as follows:

Algorithm 4 Householder Reflection Orthogonalization

Input: $\mathbf{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\} \in \mathbb{R}^{n \times n}$

Output: $\mathbf{U} = \mathbf{Q}\mathbf{R}$, where $\mathbf{Q} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \in \mathbb{R}^{n \times n}$ is the orthogonal matrix and $\mathbf{R} \in \mathbb{R}^{n \times n}$ is a upper triangular matrix

```

for  $j = 1, 2, \dots, n-1$  do
     $\{\mathbf{v}, \beta\} = \text{Householder}(\mathbf{U}_{j:n,j})$ 
     $\mathbf{U}_{j:n,j:n} \leftarrow \mathbf{U}_{j:n,j:n} - \beta \mathbf{v}(\mathbf{v}^\top \mathbf{U}_{j:n,j:n})$ 
     $\mathbf{U}_{j+1:n,j} \leftarrow \mathbf{v}_{(2:\text{end})}$ 
end
function  $\{\mathbf{v}, \beta\} = \text{Householder}(\mathbf{x})$ 
     $\sigma^2 = \|\mathbf{x}_{2:\text{end}}\|_2^2$ 
     $\mathbf{v} \leftarrow \begin{bmatrix} 1 \\ \mathbf{x}_{2:\text{end}} \end{bmatrix}$ 
    if  $\sigma^2 = 0$  then
         $\beta = 0$ 
    else
        if  $x_1 \leq 0$  then
             $\mathbf{v}_1 = x_1 - \sqrt{x_1^2 + \sigma^2}$ 
        else
             $\mathbf{v}_1 = -\frac{\sigma^2}{x_1 + \sqrt{x_1^2 + \sigma^2}}$ 
        end
         $\beta = \frac{2\mathbf{v}_1^2}{\sigma^2 + \mathbf{v}_1^2}$ 
         $\mathbf{v} \leftarrow \frac{\mathbf{v}}{\mathbf{v}_1}$ 
    end
end function

```

The algorithm follows the Matlab notation where $\mathbf{U}_{j:n,j:n}$ denotes the submatrix of \mathbf{U} from the j -th column to the n -th column and from the j -th row to the n -th row. Note that, there are a number of variants for the Householder reflection orthogonalization, such as the implicit variant where we do not store each reflection \mathbf{H}_j explicitly. Here \mathbf{Q} is the final orthogonal matrix we need.

A.3. Löwdin's Symmetric Orthogonalization

Let $U = \{u_1, u_2, \dots, u_n\}$ be a set of linearly independent vectors in a n -dimensional space. We define a general non-singular linear transformation A that can transform the basis U to a new basis R :

$$R = UA \quad (12)$$

where the basis R will be orthonormal if (the transpose will become conjugate transpose in complex space)

$$R^\top R = (UA)^\top (UA) = A^\top U^\top U A = A^\top M A = I \quad (13)$$

where $M = U^\top U$ is the gram matrix of the given basis U .

A general solution to this orthogonalization problem can be obtained via the substitution:

$$A = M^{-1} B \quad (14)$$

in which B is an arbitrary orthogonal (or unitary) matrix. When $B = I$, we will have the symmetric orthogonalization, namely

$$R := \Phi = U M^{-\frac{1}{2}} \quad (15)$$

When $B = V$ in which V diagonalizes M , then we have the canonical orthogonalization, namely

$$\Lambda = UV d^{-\frac{1}{2}}. \quad (16)$$

Because V diagonalizes M , we have that $M = V d V^\top$. Therefore, we have the $M^{-\frac{1}{2}}$ transformation as $M^{-\frac{1}{2}} = V d^{-\frac{1}{2}} V^\top$. This is essentially an eigenvalue decomposition of the symmetric matrix $M = U^\top U$.

In order to compute the Löwdin's symmetric orthogonalized basis sets, we can use singular value decomposition. Specifically, SVD of the original basis set U is given by

$$U = W \Sigma V^\top \quad (17)$$

where both $W \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices. Σ is the diagonal matrix of singular values. Therefore, we have that

$$\begin{aligned} R &= U M^{-\frac{1}{2}} \\ &= W \Sigma V^\top V d^{-\frac{1}{2}} V^\top \\ &= W \Sigma d^{-\frac{1}{2}} V^\top \end{aligned} \quad (18)$$

where we have $\Sigma = d^{\frac{1}{2}}$ due to the connections between eigenvalue decomposition and SVD. Therefore, we end up with

$$R = W V^\top \quad (19)$$

which is the output orthogonal matrix for Löwdin's symmetric orthogonalization.

An interesting feature of the symmetric orthogonalization is to ensure that

$$R = \arg \min_{P \in \text{orth}(U)} \sum_i \|P_i - U_i\| \quad (20)$$

where P_i and U_i are the i -th column vectors of $P \in \mathbb{R}^{n \times n}$ and U , respectively. $\text{orth}(U)$ denotes the set of all possible orthonormal sets in the range of U . This means that the symmetric orthogonalization functions R_i (or Φ_i) are the least distant in the Hilbert space from the original functions U_i . Therefore, symmetric orthogonalization indicates the gentlest pushing of the directions of the vectors in order to make them orthogonal.

More interestingly, the symmetric orthogonalized basis sets has unique geometric properties [68, 2] if we consider the Schweinler-Wigner matrix in terms of the sum of squared projections.

B. Proof of Theorem 1

To be more specific, neurons with each element initialized by a zero-mean Gaussian distribution are uniformly distributed on a hypersphere. We show this argument with the following theorem.

Theorem 3. *The normalized vector of Gaussian variables is uniformly distributed on the sphere. Formally, let $x_1, x_2, \dots, x_n \sim \mathcal{N}(0, 1)$ and be independent. Then the vector*

$$\mathbf{x} = \left[\frac{x_1}{z}, \frac{x_2}{z}, \dots, \frac{x_n}{z} \right] \quad (21)$$

follows the uniform distribution on \mathbb{S}^{n-1} , where $z = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$ is a normalization factor.

Proof. A random variable has distribution $\mathcal{N}(0, 1)$ if it has the density function

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}. \quad (22)$$

A n -dimensional random vector \mathbf{x} has distribution $\mathcal{N}(0, 1)$ if the components are independent and have distribution $\mathcal{N}(0, 1)$ each. Then the density of \mathbf{x} is given by

$$f(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}\langle \mathbf{x}, \mathbf{x} \rangle}. \quad (23)$$

Then we introduce the following lemma (Lemma 1) about the orthogonal-invariance of the normal distribution.

Lemma 1. *Let \mathbf{x} be a n -dimensional random vector with distribution $\mathcal{N}(0, 1)$ and $\mathbf{U} \in \mathbb{R}^{n \times n}$ be an orthogonal matrix ($\mathbf{U}\mathbf{U}^\top = \mathbf{U}^\top\mathbf{U} = \mathbf{I}$). Then $\mathbf{Y} = \mathbf{U}\mathbf{x}$ also has the distribution of $\mathcal{N}(0, 1)$.*

Proof. For any measurable set $A \subset \mathbb{R}^n$, we have that

$$\begin{aligned} P(\mathbf{Y} \in A) &= P(\mathbf{x} \in \mathbf{U}^\top A) \\ &= \int_{\mathbf{U}^\top A} \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}\langle \mathbf{x}, \mathbf{x} \rangle} \\ &= \int_A \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}\langle \mathbf{U}\mathbf{x}, \mathbf{U}\mathbf{x} \rangle} \\ &= \int_A \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}\langle \mathbf{x}, \mathbf{x} \rangle} \end{aligned} \quad (24)$$

because of orthogonality of \mathbf{U} . Therefore the lemma holds. \square

Because any rotation is just a multiplication with some orthogonal matrix, we know that normally distributed random vectors are invariant to rotation. As a result, generating $\mathbf{x} \in \mathbb{R}^n$ with distribution $\mathcal{N}(0, 1)$ and then projecting it onto the hypersphere \mathbb{S}^{n-1} produces random vectors $\mathbf{U} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$ that are uniformly distributed on the hypersphere. Therefore the theorem holds. \square

Then we show the normalized vector \mathbf{y} where each element follows a zero-mean Gaussian distribution with some constant variance σ^2 :

$$\mathbf{y} = \left[\frac{y_1}{r}, \frac{y_2}{r}, \dots, \frac{y_n}{r} \right] \quad (25)$$

where $r = \sqrt{y_1^2 + y_2^2 + \dots + y_n^2}$. Because we have that $\frac{y_i}{\sigma} \sim \mathcal{N}(0, 1)$, we can rewrite \mathbf{y} as the following random vector:

$$\mathbf{y} = \left[\frac{y_1/\sigma}{r/\sigma}, \frac{y_2/\sigma}{r/\sigma}, \dots, \frac{y_n/\sigma}{r/\sigma} \right] \quad (26)$$

where $r/\sigma = \sqrt{(y_1/\sigma)^2 + (y_2/\sigma)^2 + \dots + (y_n/\sigma)^2}$. Therefore, we directly can apply Theorem 3 and conclude that \mathbf{y} also follows the uniform distribution on \mathbb{S}^{n-1} . Now we obtain that any random vector with each element following a zero-mean Gaussian distribution with some constant variance follows the uniform distribution on \mathbb{S}^{n-1} .

Then we show that the minimum hyperspherical energy asymptotically corresponds to the uniform distribution over the unit hypersphere. We first write down the hyperspherical energy of N neurons $\{\mathbf{w}_1, \dots, \mathbf{w}_N \in \mathbb{R}^{d+1}\}$ (we also define that $\hat{\mathbf{w}}_i = \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|} \in \mathbb{S}^d$):

$$E_{s,d}(\hat{\mathbf{w}}_i|_{i=1}^N) = \sum_{i=1}^N \sum_{j=1, j \neq i}^N f_s(\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\|) = \begin{cases} \sum_{i \neq j} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\|^{-s}, & s > 0 \\ \sum_{i \neq j} \log(\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\|^{-1}), & s = 0 \end{cases} \quad (27)$$

where s is a hyperparameter that controls the behavior of hyperspherical energy. We then define a N -point minimal hyperspherical s -energy over \mathbf{A} with

$$\varepsilon_{s,d}(\mathbf{A}, \hat{\mathbf{W}}_N) := \inf_{\hat{\mathbf{W}}_N \subset \mathbf{A}} E_{s,d}(\hat{\mathbf{w}}_i|_{i=1}^N) \quad (28)$$

where we denote that $\hat{\mathbf{W}}_N = \{\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_N\}$. Typically, we will assume that \mathbf{A} is compact. Based on [19], we discuss the asymptotic behavior (as $N \rightarrow \infty$) of $\varepsilon_{s,d}(\mathbf{A}, \hat{\mathbf{W}}_N)$ in three different scenarios: (1) $0 < s < d$; (2) $s = d$; and $s > d$. The reason behind is the behavior of the following energy integral:

$$I_s(\mu) = \iint_{\mathbb{S}^d \times \mathbb{S}^d} \|\mathbf{u} - \mathbf{v}\|^{-s} d\mu(\mathbf{u}) d\mu(\mathbf{v}), \quad (29)$$

is quite different under these three scenarios. In scenario (1), Eq. 29 that is taken over all probability measures μ supported on \mathbb{S}^d will be minimal for normalized Lebesgue measure $\frac{\mathcal{H}_d(\cdot)|_{\mathbb{S}^d}}{\mathcal{H}_d(\mathbb{S}^d)}$ on \mathbb{S}^d . In the case of $s \geq d$, we will have that $I_s(\mu)$ is positive infinity for all such measures μ . Therefore, the behaviour of the minimum hyperspherical energy is different in these three cases. In general, as the parameter s increases, there is a transition from the global effects to the more local influences (from nearest neighbors). The transition happens when $s = d$. However, we typically have $0 < s < d$ in the neural networks. Therefore, we will mostly study the case of $0 < s < d$ and the theoretical asymptotic behavior is quite standard results from the potential theory [39]. From the classic potential theory, we have the following known lemma:

Lemma 2. *If $0 < s < d$, we have that*

$$\lim_{N \rightarrow \infty} \frac{\varepsilon_{s,d}(\mathbf{S}^d, \hat{\mathbf{W}}_N)}{N^2} = I_s\left(\frac{\mathcal{H}_d(\cdot)|_{\mathbb{S}^d}}{\mathcal{H}_d(\mathbb{S}^d)}\right) \quad (30)$$

Moreover, any sequence of s -energy configuration of minimal hyperspherical energy $((\hat{\mathbf{W}}_N^)_{N=1}^\infty \subset \mathbb{S}^d)$ is asymptotically uniformly distributed in the sense that for the weak-star topology of measures,*

$$\frac{1}{N} \sum_{\mathbf{v} \in \hat{\mathbf{W}}_N^*} \delta_{\mathbf{v}} \rightarrow \frac{\mathcal{H}_d(\cdot)|_{\mathbb{S}^d}}{\mathcal{H}_d(\mathbb{S}^d)} \quad \text{as } N \rightarrow \infty \quad (31)$$

where $\delta_{\mathbf{v}}$ denotes the unit point mass at \mathbf{v} .

The lemma above concludes that the neuron configuration with minimal hyperspherical energy asymptotically corresponds to the uniform distribution on \mathbb{S}^d when $0 < s < d$. From [19], we also have the following lemma that shows the same conclusion holds for the the case of $s = d$ and $s > d$:

Lemma 3. *Let $\mathcal{B}^d := \bar{B}(0, 1)$ denote the closed unit ball in \mathbb{R}^d . For the case of $s = d$, we have that*

$$\lim_{N \rightarrow \infty} \frac{\varepsilon_{s,d}(\mathbf{S}^d, \hat{\mathbf{W}}_N)}{N^2 \log N} = \frac{\mathcal{H}_d(\mathcal{B}^d)}{\mathcal{H}_d(\mathbb{S}^d)} = \frac{1}{d} \frac{\Gamma(\frac{d+1}{2})}{\sqrt{\pi} \Gamma(\frac{d}{2})} \quad (32)$$

and any sequence $(\hat{\mathbf{W}}_N^) \subset \mathbb{S}^d$ of minimal s -energy configurations satisfies Eq. 31.*

The lemma above shows that the same conclusion holds for $s = d$. For the case of $s > d$, the theoretical analysis is more involved, but the conclusion that the neuron configuration with minimal hyperspherical energy asymptotically corresponds to the uniform distribution on \mathbb{S}^d still holds. Note that, we usually will not have the case of $s > d$ in our applications.

C. Proof of Theorem 2

We consider a set of n d -dimensional neurons $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_n\} \in \mathbb{R}^{d \times n}$. The hyperspherical energy of the original set of neurons can be written as:

$$E(\mathbf{w}_i|_{i=1}^n) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n \left\| \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|_2} - \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\|^{-1} \quad (33)$$

which means that if the pairwise angle between any two neurons stays unchanged, then the hyperspherical energy will also stay unchanged. Now we consider the cosine value of the angle $\theta_{(\mathbf{w}_i, \mathbf{w}_j)}$ between any two neuron \mathbf{w}_i and \mathbf{w}_j :

$$\cos(\theta_{(\mathbf{w}_i, \mathbf{w}_j)}) = \frac{\mathbf{w}_i^\top \mathbf{w}_j}{\|\mathbf{w}_i\| \cdot \|\mathbf{w}_j\|} = \frac{\sum_{k=1}^d w_{ik} \cdot w_{jk}}{\|\mathbf{w}_i\| \cdot \|\mathbf{w}_j\|} \quad (34)$$

where w_{ik} is the k -th element of the neuron \mathbf{w}_i . From the equation above, we can observe that permuting the order of the elements in the neurons together will not change the angle. For example, switching the i -th and j -th element in all the neurons will not change the hyperspherical energy. Assume that we randomly select p dimensions from the d dimensions and denote the set of p dimension as $\mathbf{s} = \{s_1, \dots, s_p\} \in \mathbb{R}^p$. Therefore we can construct a new set of neurons $\tilde{\mathbf{W}} = \{\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_n\}$ by permuting the p dimensions in \mathbf{s} to become the first p elements for all the neurons. Essentially, we use permutation to make $\tilde{\mathbf{w}}_i = \mathbf{w}_{s_i}$ for $i \in [1, p]$. Therefore, we can have the following equation:

$$E(\mathbf{w}_i|_{i=1}^n) = E(\tilde{\mathbf{w}}_i|_{i=1}^n) \quad (35)$$

Then we consider an orthogonal matrix $\mathbf{R}_p \in \mathbb{R}^{d \times d}$ that is used to transform the p dimension in the neurons. The equivalent orthogonal transformation for the d -dimensional neurons $\tilde{\mathbf{W}}$ is

$$\tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n-p} \end{bmatrix} = \begin{bmatrix} \mathbf{R}_p & 0 & \dots & 0 \\ 0 & 1 & \ddots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & 1 \end{bmatrix} \quad (36)$$

where \mathbf{I}_{n-p} is an identity matrix of size $(n-p) \times (n-p)$. It is easy to verify that $\tilde{\mathbf{R}}$ is also an orthogonal matrix: $\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}} = \mathbf{I}_n$. Then we permute the order of $\tilde{\mathbf{W}}$ back to the original neuron set \mathbf{W} and obtain a new set of neurons $\mathbf{W}^t = \{\mathbf{w}_1^t, \dots, \mathbf{w}_n^t\}$. \mathbf{W}^t is in fact the result of directly performing orthogonal transformation to the p dimensions in \mathbf{W} . Because any order permutation of elements in neurons does not change the hyperspherical energy, we have the following equation

$$E(\mathbf{w}_i|_{i=1}^n) = E(\tilde{\mathbf{w}}_i|_{i=1}^n) = E(\tilde{\mathbf{R}}\tilde{\mathbf{w}}_i|_{i=1}^n) = E(\tilde{\mathbf{R}}\mathbf{w}_i^t|_{i=1}^n) \quad (37)$$

which concludes our proof.

D. Experimental Settings

| Layer | CNN-6 (CIFAR-100) | CNN-9 (CIFAR-100) | CNN-10 (ImageNet-2012) |
|-----------------|------------------------------------|-----------------------------|--|
| Conv1.x | $[3 \times 3, 64] \times 2$ | $[3 \times 3, 64] \times 3$ | $[7 \times 7, 64]$, Stride 2 3×3 , Max Pooling, Stride 2 $[3 \times 3, 64] \times 3$ |
| Pool1 | 2×2 Max Pooling, Stride 2 | | |
| Conv2.x | $[3 \times 3, 64] \times 2$ | $[3 \times 3, 64] \times 3$ | $[3 \times 3, 128] \times 3$ |
| Pool2 | 2×2 Max Pooling, Stride 2 | | |
| Conv3.x | $[3 \times 3, 64] \times 2$ | $[3 \times 3, 64] \times 3$ | $[3 \times 3, 256] \times 3$ |
| Pool3 | 2×2 Max Pooling, Stride 2 | | |
| Fully Connected | 64 | 64 | 256 |

Table 14: Our plain CNN architectures with different convolutional layers. Conv1.x, Conv2.x and Conv3.x denote convolution units that may contain multiple convolution layers. E.g., $[3 \times 3, 64] \times 3$ denotes 3 cascaded convolution layers with 64 filters of size 3×3 .

| Layer | ResNet-20 (CIFAR-100) | ResNet-32 (CIFAR-100) |
|-----------------|--|--|
| Conv1.x | $[3 \times 3, 16] \times 1$ $[3 \times 3, 16] \times 3$ | $[3 \times 3, 16] \times 1$ $[3 \times 3, 16] \times 5$ |
| Conv2.x | $[3 \times 3, 32] \times 3$ | $[3 \times 3, 32] \times 5$ |
| Conv3.x | $[3 \times 3, 64] \times 3$ | $[3 \times 3, 64] \times 5$ |
| Average Pooling | | |

Table 15: Our ResNet architectures with different convolutional layers. Conv0.x, Conv1.x, Conv2.x, Conv3.x and Conv4.x denote convolution units that may contain multiple convolutional layers, and residual units are shown in double-column brackets. Conv1.x, Conv2.x and Conv3.x usually operate on different size feature maps. These networks are essentially the same as [21], but some may have a different number of filters in each layer. The downsampling is performed by convolutions with a stride of 2. E.g., $[3 \times 3, 64] \times 4$ denotes 4 cascaded convolution layers with 64 filters of size 3×3 , S2 denotes stride 2.

Reported Results. For all the experiments on MLPs and CNNs (except CNNs in the few-shot learning), we report testing error rates. For the few-shot learning experiment, we report testing accuracy. For all the experiments on both GCNs and PointNets, we report testing accuracy. All results are averaged over 10 runs of the model.

Multilayer perceptron. We conduct digit classification task on MNIST with a three-layer multilayer perceptron following this repository¹. The input dimension of each MNIST digit is 28×28 , which is 784 dimensions after flattened. Our two hidden layers have 256 output dimensions, *i.e.*, 256 neurons. The output layer will output 10 logits for classification. We use a cross-entropy loss with softmax function. For the optimization, we use a momentum SGD with learning rate 0.01, momentum 0.9 and batch size 100. The training stops at 100 epochs.

Convolutional neural networks. The network architectures used in the paper are elaborated in Table 14 and Table 15. For CIFAR-100, we use 128 as the mini-batch size. We use momentum SGD with momentum 0.9 and the learning rate starts with 0.1, divided by 10 when the performance is saturated. For ImageNet-2012, we use batch size 128 and start with learning rate 0.1. The learning rate is divided by 10 when the performance is saturated, and the training is terminated at 700k iterations. For ResNet-20 and ResNet-32 on CIFAR-100, we use exactly the same architecture used on CIFAR-10 as [21]. The rotation matrix is initialized with random normal distribution (mean is 0 and variance is 1). Note that, for all the compared methods, we always use the best possible hyperparameters to make sure that the comparison is fair. The baseline has exactly the same architecture and training settings as the one that OPT uses. If not otherwise specified, standard ℓ_2 weight decay ($5e-4$) is applied to all the neural network including baselines and the networks that use OPT training.

Few-shot learning. The network architecture (Table 16) we used for few-shot learning experiments is the same as that used in [9]. In our experiments, we show comparison of our OPT training with standard training on ‘baseline’ and ‘baseline++’ settings in [9]. In ‘baseline’ setting, a standard CNN model is pretrained on the whole meta-train dataset (standard non-MAML supervised training) and later only the classifier layer is finetuned on few-shot dataset. ‘baseline++’ differs from ‘baseline’ on the classifier: in ‘baseline’, each output dimension of the classifier is computed as the inner product between weight w and input x , *i.e.* $w \cdot x$; while in ‘baseline++’ it becomes the scaled cosine distance $c \frac{w \cdot x}{\|w\| \|x\|}$ where c is a positive scalar. Following [9], we set $c = 2$.

During pretraining, the model is trained for 200 epochs on the meta-train set of mini-ImageNet with an Adam optimizer (learning rate $1e-3$, weight decay $5e-4$) and the classifier is discarded after pretraining. The model is later finetuned, with

¹https://github.com/hwalsuklee/tensorflow-mnist-MLP-batch_normalization-weight_initializers

| Layer | CNN-4 |
|-------------------|------------------------------------|
| Conv1 | $3 \times 3, 64$ |
| Pool1 | 2×2 Max Pooling, Stride 2 |
| Conv2 | $3 \times 3, 64$ |
| Pool2 | 2×2 Max Pooling, Stride 2 |
| Conv3 | $3 \times 3, 64$ |
| Pool3 | 2×2 Max Pooling, Stride 2 |
| Conv4 | $3 \times 3, 64$ |
| Pool4 | 2×2 Max Pooling, Stride 2 |
| Linear Classifier | number of classes |

Table 16: Architecture for few-shot learning. The number of classes is different for pretraining and finetuning.

| Layer | Wide CNN-9 (CIFAR-100) | ResNet-18 (ImageNet-2012) |
|---------|--|---|
| Conv0.x | N/A | $[7 \times 7, 64]$, Stride 2 3×3 , Max Pooling, Stride 2 |
| Conv1.x | $[3 \times 3, 64] \times 3$ 2×2 Max Pooling, Stride 2 | $3 \times 3, 64$ $3 \times 3, 64$ $\times 2$ |
| Conv2.x | $[3 \times 3, 128] \times 3$ 2×2 Max Pooling, Stride 2 | $3 \times 3, 128$ $3 \times 3, 128$ $\times 2$ |
| Conv3.x | $[3 \times 3, 256] \times 3$ 2×2 Max Pooling, Stride 2 | $3 \times 3, 256$ $3 \times 3, 256$ $\times 2$ |
| Conv4.x | N/A | $3 \times 3, 512$ $3 \times 3, 512$ $\times 2$ |
| Final | 256-dim Fully Connected | Average Pooling |

Table 17: Our wide CNN-9 and wide ResNet-18 architectures with different convolutional layers.

a new classifier, on the few-shot samples (5 way, support size 5) with a momentum SGD optimizer (learning rate $1e - 2$, momentum 0.9, dampening 0.9, weight decay $1e - 3$, batch size 4) for 100 epochs. We re-initialize the classifier for each few-shot sample.

Graph neural networks. We implement the OPT training for GCN in the official repository². The experimental settings also follow the official repository to ensure a fair comparison. For OPT (CP) method, we use the original hyperparameters and experimental setup except the added rotation matrix. For OPT (OGD) method, we use our own OGD optimizer in Tensorflow to train the rotation matrix in order to maintain orthogonality and use the original optimizer to train the other variables.

Training a GCN with OPT is not that straightforward. Specifically, the forward model of GCN is $\mathbf{Z} = \text{Softmax}(\hat{\mathbf{A}} \cdot \text{ReLU}(\hat{\mathbf{A}} \cdot \mathbf{X} \cdot \mathbf{W}_0) \cdot \mathbf{W}_1)$ where $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$. We note that \mathbf{A} is the adjacency matrix of the graph, $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ (\mathbf{I} is an identity matrix), and $\tilde{\mathbf{D}} = \sum_j \tilde{\mathbf{A}}_{ij}$. $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the feature matrix of n nodes in the graph (feature dimension is d). \mathbf{W}_1 is the weights of the classifiers. \mathbf{W}_0 is the weight matrix of size $d \times h$ where h is the dimension of the hidden space. We treat each column vector of \mathbf{W}_0 as a neuron, so there are h neurons in total. Then we apply OPT to train these h neurons of dimension d in GCN. We conduct experiments on Cora and Pubmed datasets [63]. We aim to verify the effectiveness of OPT on GCN instead of achieving state-of-the-art performance on this task.

Point cloud recognition. To simplify the comparison and remove all the bells and whistles, we use a vanilla PointNet (without T-Net) as our backbone network. We apply OPT to train the MLPs in PointNet. We follow the same experimental settings as [57] and evaluate on the ModelNet-40 dataset [76]. We exactly follow the same setting in the original paper [57] and the official repositories³. Specifically, we multiply the rotation matrix to the original fixed neurons in all the 1×1 convolution layers and the fully connected layer except the final classifier. All the rotation matrix is initialized with random normal distribution. For the experiments, we use point number 1024, batch size 32 and Adam optimizer with initial learning rate 0.001. The learning rate will decay by 0.7 every 200k iterations, and the training is terminated at 250 epochs.

Experimental settings for S-OPT. For the experiment of S-OPT, the architecture of wide CNN-9 and wide ResNet-18 is given in Table 17. CNN-6 is the same as the one in Table 14. We use standard data augmentation for CIFAR-100, following [54]. For ImageNet-2012, we use the same data augmentation in [38, 54]. This data augmentation does not contain as many transformation as the one in [21], so the final performance may be worse than [21]. However, all the compared methods use the same data augmentation in our experiments, so the experiment is still a fair comparison. For CIFAR-100, we use $N_{\text{out}} = 300$ and $N_{\text{in}} = 750$. For ImageNet, we use $N_{\text{out}} = 700$ and $N_{\text{in}} = 1000$. For S-OPT, we directly use the original OPT for the first layer, as its neuron dimension is typically very small. We decrease the learning rate by a factor of 10 when the performance is saturated in the outer iteration.

²<https://github.com/tkipf/gcn>

³<https://github.com/charlesq34/pointnet>

| Layer | MNIST Visualization | ResNet-18A | ResNet-18B | CNN-20 for Face Embeddings |
|---------|--|---|---|--|
| Conv0.x | N/A | $[7 \times 7, 64]$, Stride 2 3×3 , Max Pooling, Stride 2 | $[7 \times 7, 64]$, Stride 2 3×3 , Max Pooling, Stride 2 | N/A |
| Conv1.x | $[3 \times 3, 32] \times 2$ 3×3 , Max Pooling, Stride 2 | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$ | $[3 \times 3, 64] \times 1$, Stride 2 $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 1$ |
| Conv2.x | $[3 \times 3, 64] \times 2$ 3×3 , Max Pooling, Stride 2 | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $[3 \times 3, 128] \times 1$, Stride 2 $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ |
| Conv3.x | $[3 \times 3, 128] \times 2$ 3×3 , Max Pooling, Stride 2 | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$ | $[3 \times 3, 256] \times 1$, Stride 2 $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 4$ |
| Conv4.x | N/A | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$ | $[3 \times 3, 512] \times 1$, Stride 2 $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 1$ |
| Final | Fully Connected (3-dim) | Average Pooling (128-dim) | Average Pooling (512-dim) | Fully Connected (512-dim) |

Table 18: Our network architectures for large categorical training.

Experimental settings for large categorical training. All the network architectures used in the large categorical training are specified in Table 18. For the visualization on MNIST, we simply set the output dimension as 3 and directly plot the 3-dimensional features. In Fig. 7, each color denotes a class of digits, and each dot point denotes 3-dimensional features for a digit image. The experiments on ImageNet follows the same setting as the previous section. For the open-set face recognition experiments, we generally follow the same training configuration as SphereFace [52]. For all the methods used in face recognition, we use the 20-layer residual network as described in Table 18. Since OPT is originally implemented in TensorFlow, we re-implement the CNN-20 for deep face recognition in TensorFlow, which yields an accuracy gap compared to [52]. This is due to some mis-match in data augmentation and optimizations. However, since both the baseline and our CLS-OPT use the same network implementation in TensorFlow and achieving state-of-the-art results is not our major focus, it is still a fair and valid comparison. We expect CLS-OPT can also be generally useful for large categorical training of deep face recognition.

E. Loss Landscape Visualization (Normal Distribution Perturbation)

E.1. Visualization Procedure

We generally follow the visualization procedure in [43]. However, since OPT has a different training process, we use a modified visualization method but still make it comparable to the baseline.

Specifically, if we want to plot the loss landscape of a neural network with loss $\mathcal{L}(\theta)$ where θ is the learnable model parameters, we need to first choose pretrained model parameters θ^* as a center point. Then we choose two random direction vectors δ and η . The 2D plot $f(\alpha, \beta)$ is defined as

$$f(\alpha, \beta) = \mathcal{L}(\theta^* + \alpha\delta + \beta\eta) \quad (38)$$

which can be used as a 2D surface visualization. Note that, after we randomly initialize the direction vectors δ and η (with normal distribution), we need to perform the filter normalization [43]. Specifically, we normalize each filter in δ and η to have the same norm as the corresponding filter in θ^* . The loss landscape of our baseline is plotted using this visualization approach.

In contrast, the learnable parameters in OPT are no longer the weights of neurons. Instead, the learnable parameters are the orthogonal matrices. More precisely, the trainable matrices are used to perform orthogonalization in the neural networks (*i.e.*, P in Fig. 2). We denote the combination of all the trainable matrices as \tilde{R} , and the corresponding pretrained matrices as \tilde{R}^* . Then the 2D visualization of OPT is

$$f(\alpha, \beta) = \mathcal{L}(\tilde{R}^* + \alpha\gamma + \beta\kappa) \quad (39)$$

where γ and κ are two random direction vectors (which follow the normal distribution) to perturb \tilde{R}^* . The visualization procedures of baseline and OPT are essentially the same except that the trainable variables are different. Therefore, their loss landscapes are comparable.

E.2. Experimental Details

In Fig. 4, we vary α and β from -1.5 to 1.5 for both baseline and OPT, and then plot the surface of 2D function f . We use the CNN-6 (as specified in Appendix D) on CIFAR-100. We use the same data augmentation as [49]. We train the network with SGD with momentum 0.9 and batch size 128. We start with learning rate 0.1, divide it by 10 at 30k, 50k and 64k iterations, and terminate training at 75k iterations. The training details basically follows [49]. We mostly use CP for OPT due to efficiency. Note that, the other orthogonalization methods in OPT yields similar loss landscapes in general. The pretrained model for standard training yields 37.59% testing error on CIFAR-100, while the pretrained model for OPT yields 33.53% error. This is also reported in Section 6.

E.3. Full Visualization Results for the Main Paper

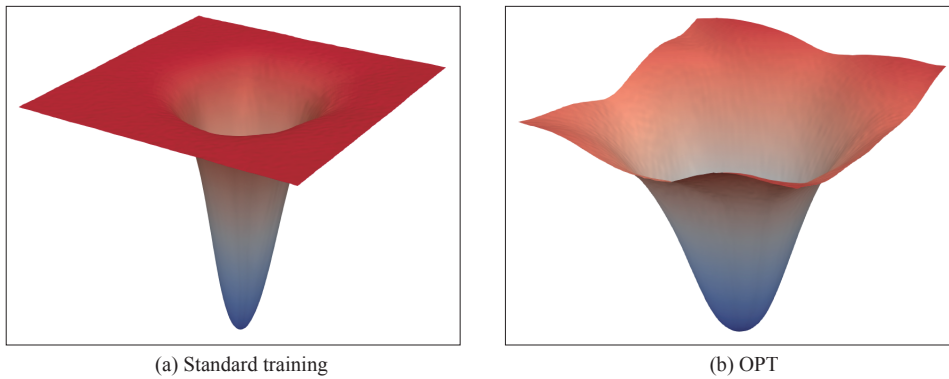


Figure 8: High-quality rendered loss landscapes of standard training and OPT.

Following the same experimental settings in Appendix E.2, we render the 3D loss landscapes with some color and lighting effects for Fig. 8. The visualization data is exactly the same as Fig. 4, and we simply use ParaView to plot the figure. The rendered loss landscape better reflects that OPT yields a much more smooth loss geometry.

We also give the large and full version of Fig. 4(b) (in the main paper) in the following figure. Fig. 9 is identical to Fig. 4(b) in the main paper except that Fig. 9 has larger size. From Fig. 9, we can better observe the dramatically different loss landscape

between standard training and OPT. From the contour plots, we can better see that the red region of standard training is extremely flat and is highly non-convex around the edge. In comparison to standard training, OPT has very smooth and relatively convex contour shape.

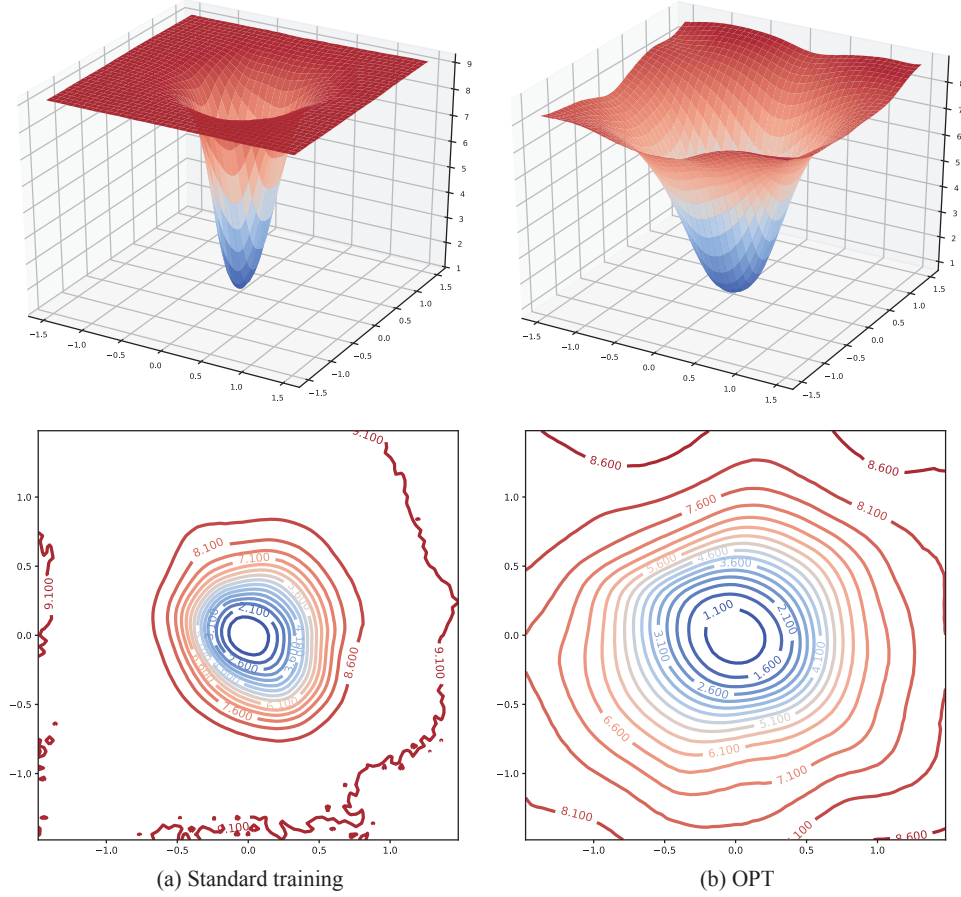
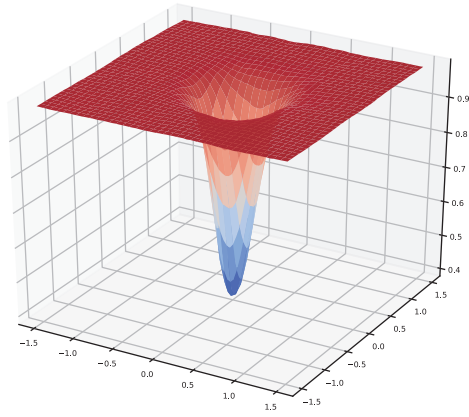
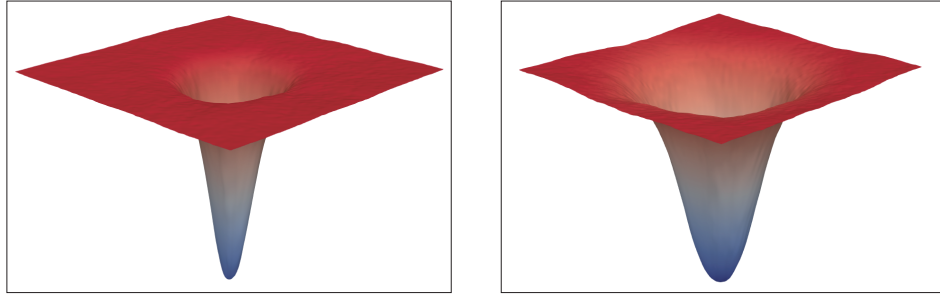
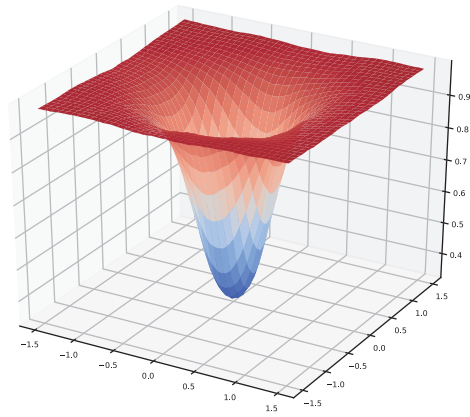


Figure 9: Comparison of loss landscapes between standard training and OPT (full results of Fig. 4(b) in the main paper). Top row: loss landscape visualization with Cartesian coordinate system; Bottom row: loss contour visualization.

Then we visualize the testing error landscape. Additionally, we include a high-quality visualization using ParaView. We render the plot with lighting and color effects in order to better demonstrate the difference between OPT and standard training. The visualization results are given in Fig. 10. The comparison of testing error landscape shows that the parameter space of OPT is more robust than standard training, because the testing error of OPT is increased in a slower speed while the model is perturbed away from the pretrained parameters. In other words, the parameter space of OPT is more robust to the random perturbation than standard training. Combining the visualization of loss landscape, it is well justified that OPT can significantly alleviate training difficulty and improve generalization.



(a) Standard training



(b) OPT

Figure 10: Comparison of testing error landscapes between standard training and OPT. Top row: high-quality rendered testing error landscape visualization with lighting effects; Bottom row: testing error landscape visualization with Cartesian coordinate system.

E.4. Loss Landscape Visualization of Different Neural Network

In order to show that the loss landscape of OPT is quite general and consistent across different neural network architectures, we also visualize the loss landscape using a deep network (CNN-9 as specified in Appendix D). The experiments are conducted on CIFAR-100. The results are given in Fig. 11. We can see that the loss landscape of OPT is much more smooth than standard training, similar to the observation for CNN-6. Therefore, the loss landscape difference between OPT and standard training is consistent across different network architectures, and OPT consistently shows better optimization landscape.

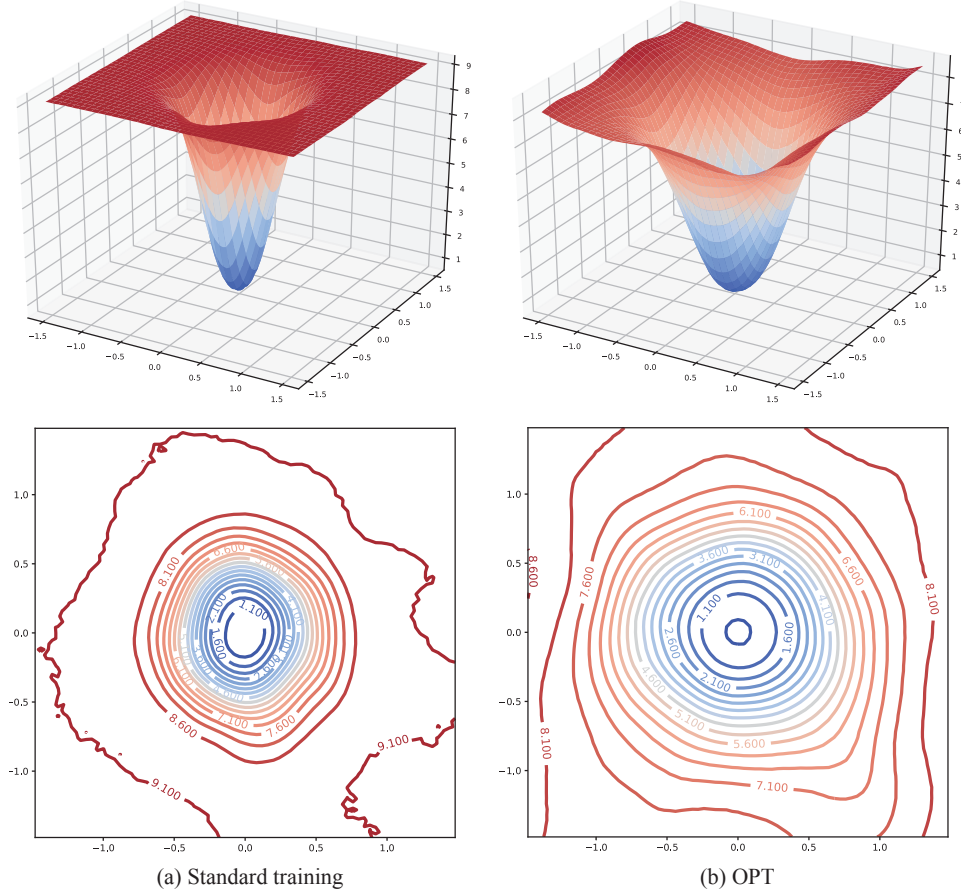
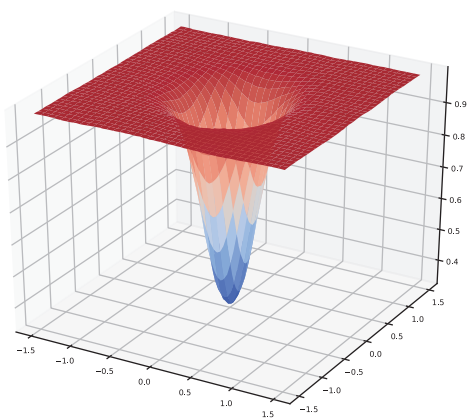
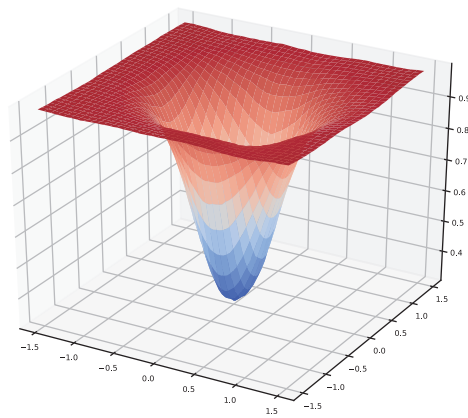


Figure 11: Comparison of loss landscapes between standard training and OPT on CIFAR-100 (CNN-9). Top row: loss landscape visualization with Cartesian coordinate system; Bottom row: loss contour visualization.

Then we show the landscape of testing error for CNN-9 on CIFAR-100. The results are given in Fig. 12. Similar to CNN-6, the testing error landscape of OPT is more smooth and convex than standard training. Moreover, OPT has a more flat local minima of testing error, while standard training has a sharp local minima. The testing error landscape in Fig. 12 generally follows the same pattern as the loss landscape in Fig. 11. The visualization further verifies the superiority of OPT is very consistent across different network architectures.



(a) Standard training



(b) OPT

Figure 12: Comparison of testing error landscapes between standard training and OPT on CIFAR-100 (CNN-9).

E.5. Loss Landscape Visualization of Different Dataset

Similar to Appendix F.5, we also visualize the loss landscape of OPT and standard training on a different dataset (CIFAR-10) with CNN-6. The loss landscape visualization is given in Fig. 13. Although the loss landscape on CIFAR-10 is quite different from the one on CIFAR-100, we can still observe that the loss landscape of OPT has a very flat local minima and the loss values are increasing smoothly and slowly. In contrast, the loss landscape of standard training has a sharp local minima and the loss values quickly increase to a large value. The red region of standard training will lead to very small gradient, potentially affecting the training. From the contour plots, the comparison apparently shows that the loss landscape of OPT is much more smooth than standard training.

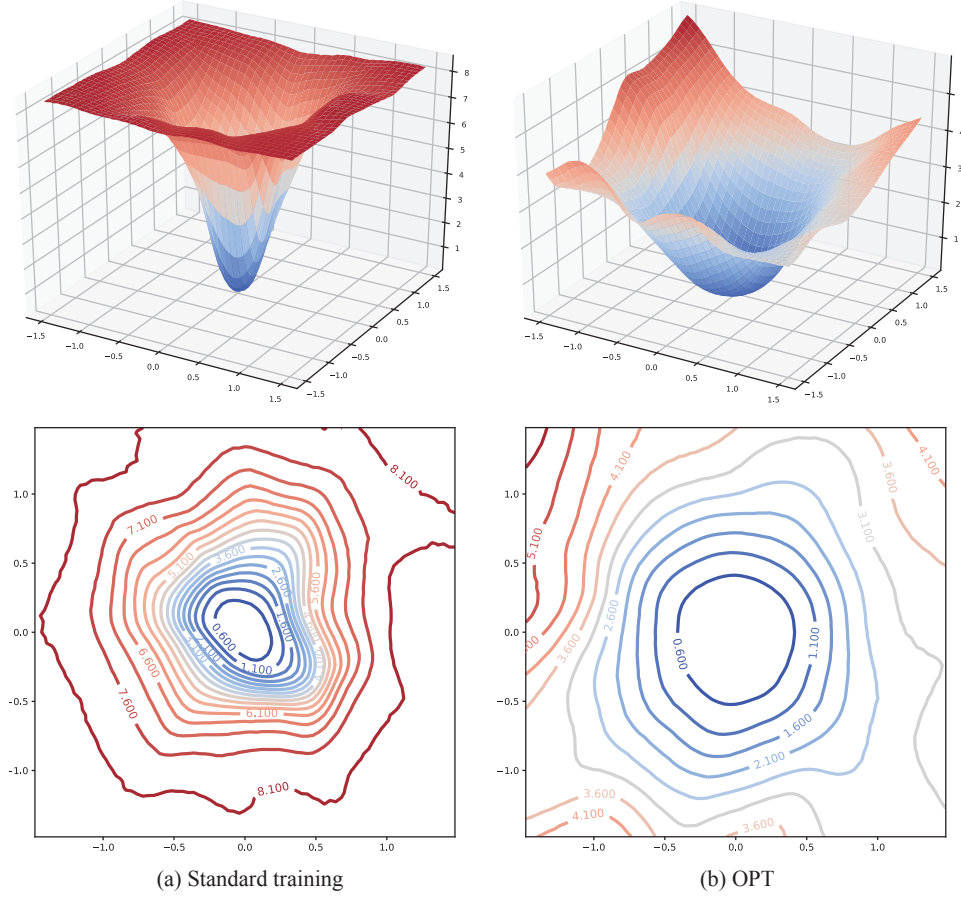
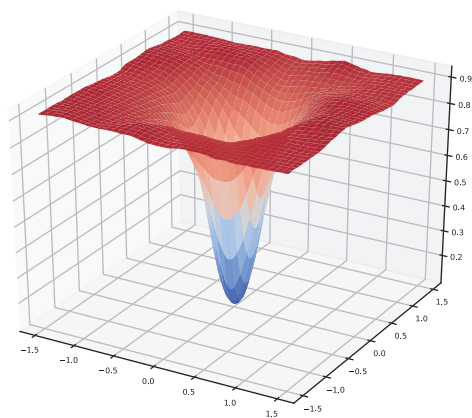
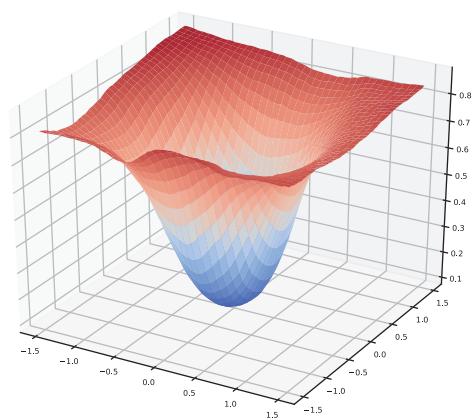


Figure 13: Comparison of loss landscapes between standard training and OPT on CIFAR-10 (CNN-6). Top row: loss landscape visualization with Cartesian coordinate system; Bottom row: loss contour visualization.

We also visualize the landscape of testing error in Fig. 14. The testing error landscape generally follows the pattern in the loss landscape. One can easily observe that the parameter space of standard training is very sensitive to perturbations. A small perturbation can make the model parameters completely fail (*i.e.*, the testing error dramatically increase). Differently, the parameter space of OPT is much more robust to perturbations. The model parameter can still work well with a small perturbation. Both Fig. 13 and Fig. 14 validate that the superiority of OPT is consistent across different training datasets.



(a) Standard training



(b) OPT

Figure 14: Comparison of testing error landscapes between standard training and OPT on CIFAR-10 (CNN-6).

F. Loss Landscape Visualization (Uniform Distribution Perturbation)

F.1. Visualization Procedure

Different from Appendix E, we choose two random direction vectors δ and η based on $[0, 1]$ uniform distribution to further justify the effectiveness of OPT. The 2D plot $f(\alpha, \beta)$ is defined as

$$f(\alpha, \beta) = \mathcal{L}(\theta^* + \alpha\delta + \beta\eta) \quad (40)$$

which can be used as a 2D surface visualization. Note that, after we randomly initialize the direction vectors δ and η with $[0, 1]$ normal distribution, we need to perform the filter normalization [43]. Specifically, we normalize each filter in δ and η to have the same norm as the corresponding filter in θ^* . The loss landscape of our baseline is plotted using this visualization approach.

In contrast, the learnable parameters in OPT are no longer the weights of neurons. Instead, the learnable parameters are the orthogonal matrices. More precisely, the trainable matrices are used to perform orthogonalization in the neural networks (*i.e.*, P in Fig. 2). We denote the combination of all the trainable matrices as \tilde{R} , and the corresponding pretrained matrices as \tilde{R}^* . Then the 2D visualization of OPT is

$$f(\alpha, \beta) = \mathcal{L}(\tilde{R}^* + \alpha\gamma + \beta\kappa) \quad (41)$$

where γ and κ are two random direction vectors (which follow the $[0, 1]$ uniform distribution) to perturb \tilde{R}^* . The visualization procedures of baseline and OPT are essentially the same except that the trainable variables are different. Therefore, their loss landscapes are comparable.

F.2. Experimental Details

In Fig. 4, we vary α and β from -1 to 1 for both baseline and OPT, and then plot the surface of 2D function f . We use the CNN-6 (as specified in Appendix D) on CIFAR-100. We use the same data augmentation as [49]. We train the network with SGD with momentum 0.9 and batch size 128. We start with learning rate 0.1, divide it by 10 at 30k, 50k and 64k iterations, and terminate training at 75k iterations. The training details basically follows [49]. We mostly use CP for OPT due to efficiency. Note that, the other orthogonalization methods in OPT yields similar loss landscapes in general. The pretrained model for standard training yields 37.59% testing error on CIFAR-100, while the pretrained model for OPT yields 33.53% error. This is also reported in Section 6.

F.3. Full Visualization Results

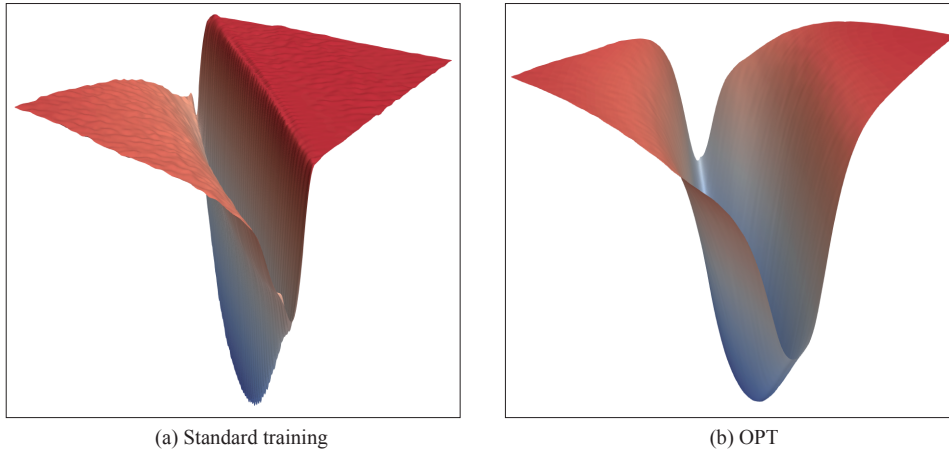


Figure 15: High-quality rendered loss landscapes of standard training and OPT.

Following the same experimental settings in Appendix F.2, we render the 3D loss landscapes with some color and lighting effects for Fig. 15. We first use ParaView to plot a high-quality loss landscape comparison between standard training and OPT. As expected, the loss landscape of OPT is much more smooth than standard training. Note that, for the flat red region in standard training, we can still observe numerous small local minima, while the red region of OPT is very smooth. Fig. 15 better validates our analysis and discussion in Section 5.1, and also shows the superiority of OPT.

We provide the visualization results in the rest of the subsection. From Fig. 16, we can better observe the dramatically different loss landscape between standard training and OPT.

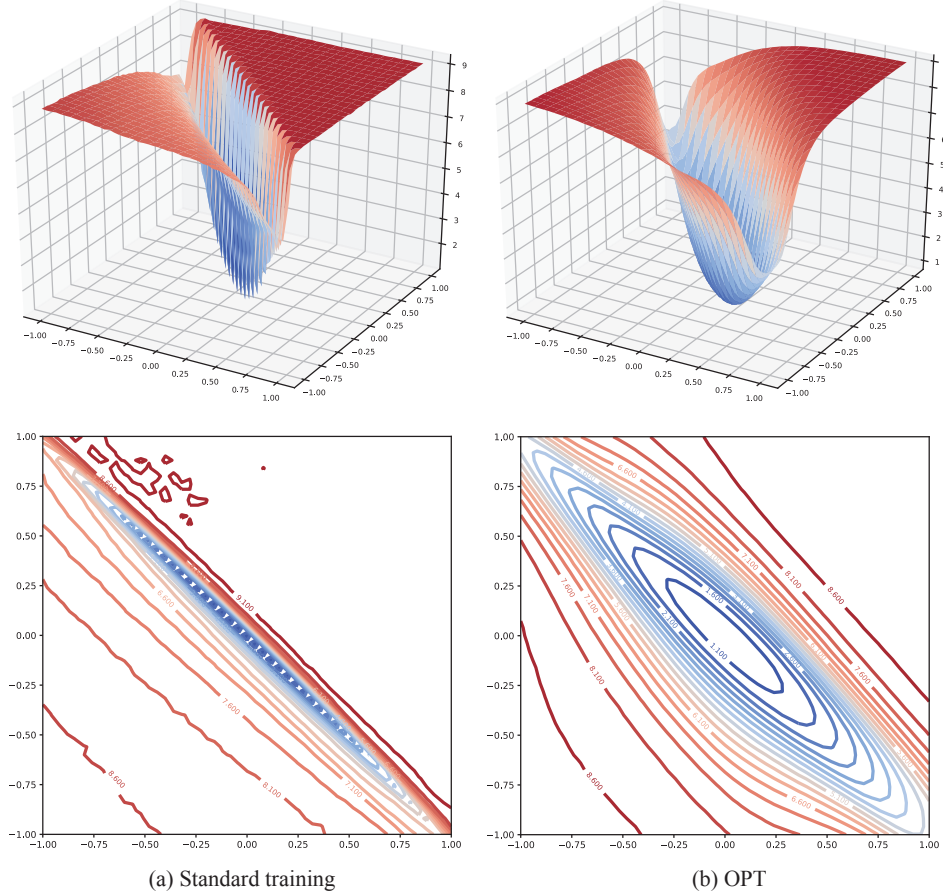


Figure 16: Comparison of loss landscapes between standard training and OPT (full results of Fig. 4(a) in the main paper). Top row: loss landscape visualization with Cartesian coordinate system; Bottom row: loss contour visualization.

To better understand the difference of the training dynamics between standard training and OPT, we also plot the testing error landscapes in Fig. 17 for both methods. The testing error is computed on the testing set of CIFAR-100 with the perturbed pretrained model (α and β are the perturbation parameters). From the testing error landscape comparison in Fig. 17, we can see that once the baseline pretrained model is slightly perturbed, the testing error will immediately increase to 99.99% which is random selection-level testing error (because we have 100 balanced classes in total, randomly picking a class leads to 0.01% accuracy). In contrast, the testing error landscape of OPT is much more smooth. Even if we perturb the OPT pretrained model, we still end up with a reasonably low testing error, show that the parameter space of OPT is more smooth and continuous. All these evidences suggest that OPT is a better training framework for neural networks and can significantly alleviate the optimization difficulty. In this following subsections, we aim to show that the loss and testing error landscape difference between standard training and OPT is not a coincidence. We will show that the improvement of OPT on the loss and testing error landscape is both dataset-agnostic and architecture-agnostic.

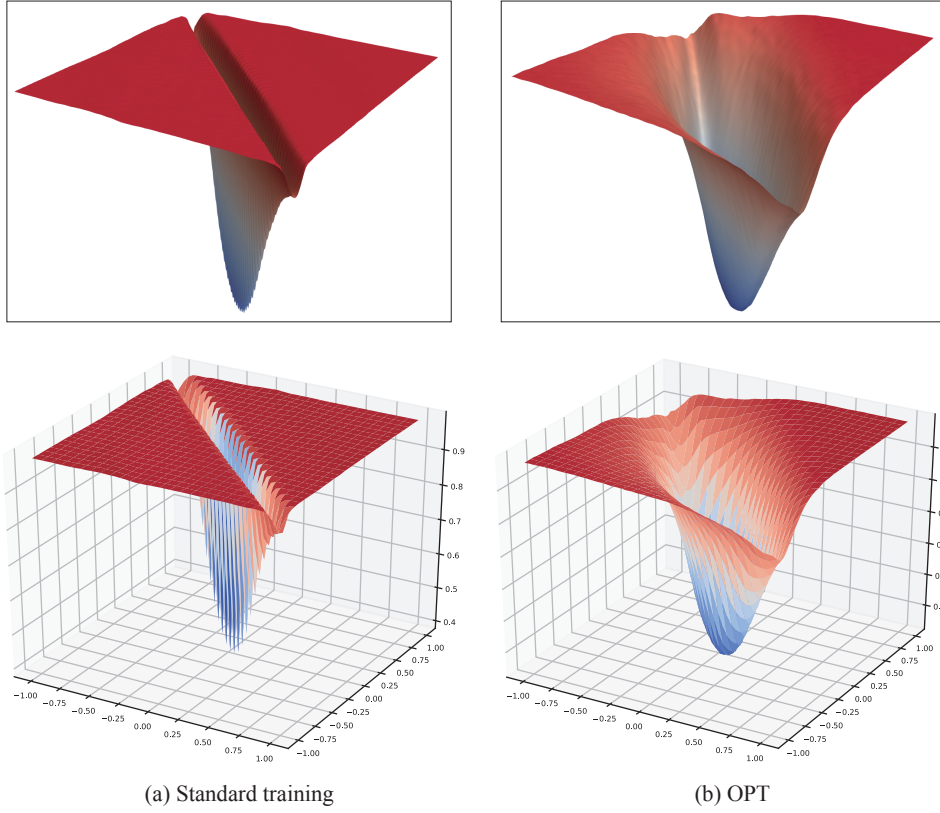


Figure 17: Comparison of testing error landscapes between standard training and OPT. Top row: high-quality rendered testing error landscape visualization with lighting effects; Bottom row: testing error landscape visualization with Cartesian coordinate system.

F.4. Loss Landscape Visualization of Different Neural Network

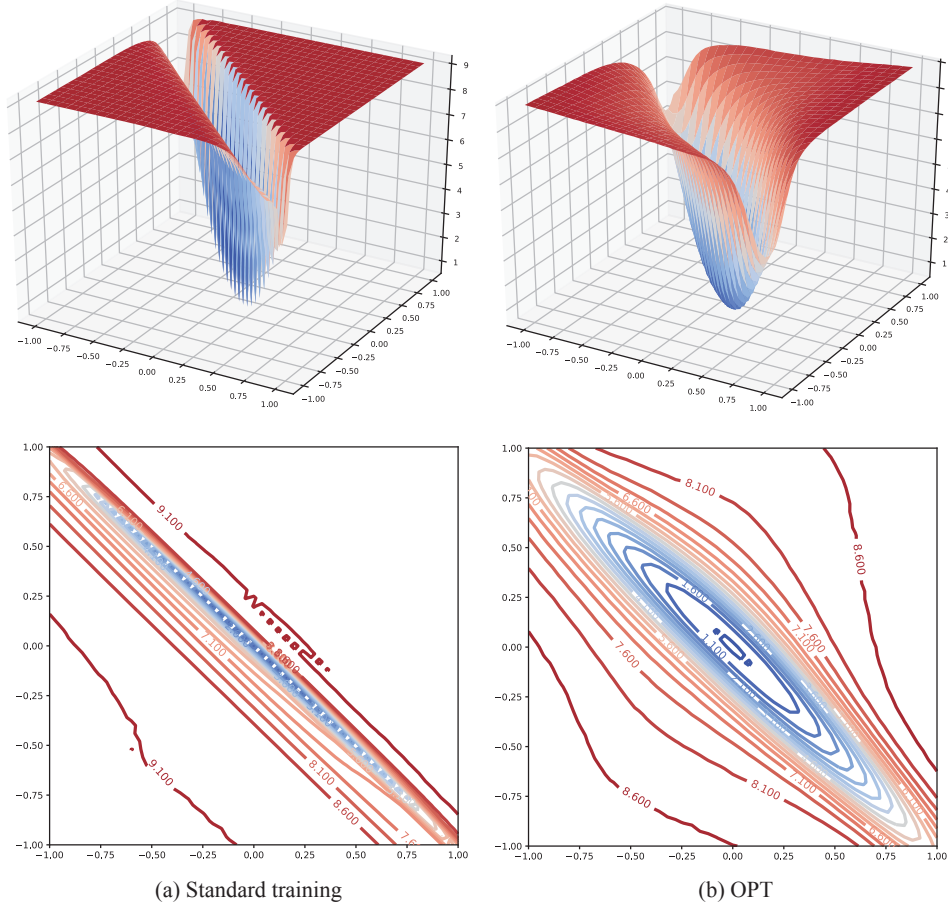


Figure 18: Comparison of loss landscapes between standard training and OPT on CIFAR-100 (CNN-9). Top row: loss landscape visualization with Cartesian coordinate system; Bottom row: loss contour visualization.

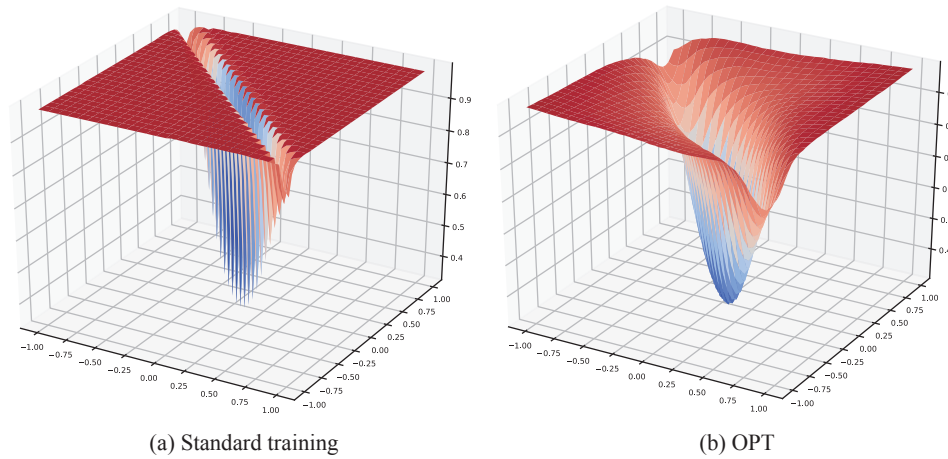


Figure 19: Comparison of testing error landscapes between standard training and OPT on CIFAR-100 (CNN-9).

To show that the difference of the loss landscape between standard training and OPT is consistent across different neural networks. We use a deeper CNN-9 (as specified in Appendix D) to visualize the loss landscape. The experimental settings generally follow Appendix F.2. We use CNN-9 as our backbone architecture and train it on CIFAR-100. The visualization of

the loss landscapes is given in Fig. 18. We observe that the loss landscapes of OPT with CNN-9 and CNN-6 are very similar. In general, the conclusion that OPT yields better loss landscape still holds for a deeper neural network, showing that the effectiveness of OPT is not architecture-dependent.

Moreover, we also visualize the landscape of testing error in Fig. 19. We can observe that the testing error landscapes are somewhat similar to the loss landscape in Fig. 18. The results further validate the superiority of OPT. We can observe in Fig. 19 that OPT has more smooth testing error landscape and can make the training parameter space of the neural network less sensitive to perturbations.

F.5. Loss Landscape Visualization on Different Dataset

We also perform the same loss and testing error landscape visualization on CIFAR-10. The training details basically follows Appendix F.2. For CIFAR-10, we use the same data augmentation as in Appendix D. The results are given in Fig. 20. From Fig. 20, we can observe even more dramatic difference of the loss landscape between standard training and OPT. In standard training, the loss landscape exhibits highly non-convex and non-smooth behavior. There are countless local minima in the loss landscape. Different from the results in Fig. 18, the loss landscape of standard training on CIFAR-10 has some huge local minima that are hard to escape from. In contrast, the loss landscape of OPT on CIFAR-10 does not show obvious and huge local minima and is far more convex and smooth than standard training. The contour maps show more significant difference between standard training and OPT. The contour map of OPT shows the shape of a single symmetric and convex valley, while the contour map of standard training presents the shape of multiple highly irregular valleys. The visualization further validate that the improvement of OPT on optimization landscape is very consistent across different training datasets.

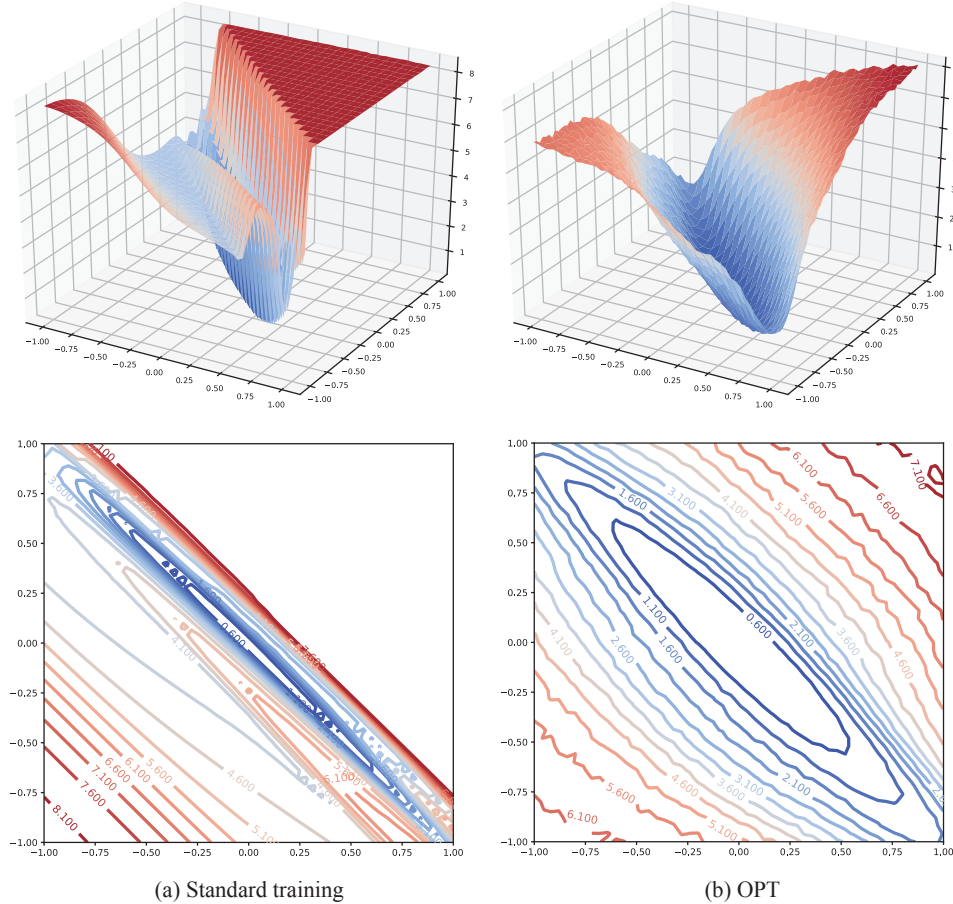


Figure 20: Comparison of loss landscapes between standard training and OPT on CIFAR-10 (CNN-6). Top row: loss landscape visualization with Cartesian coordinate system; Bottom row: loss contour visualization.

Then we also visualize the testing error landscape of standard training and OPT in Fig. 21. As expected, the testing error landscape of standard training shows that changing the pretrained model parameters with a very small perturbation could lead to a dramatic increase in testing error. It indicates that the parameter space of standard training is very sensitive to even a small perturbation. In comparison, the testing error landscape of OPT shows similar shape to the training loss landscape which is the shape of a single regular, smooth and convex valley. We can conclude that OPT has huge advantages over standard training in terms of the optimization landscape. Although the conclusion is drawn from a simple visualization method, it can still shed some light on why OPT yields better training dynamics and generalization ability.

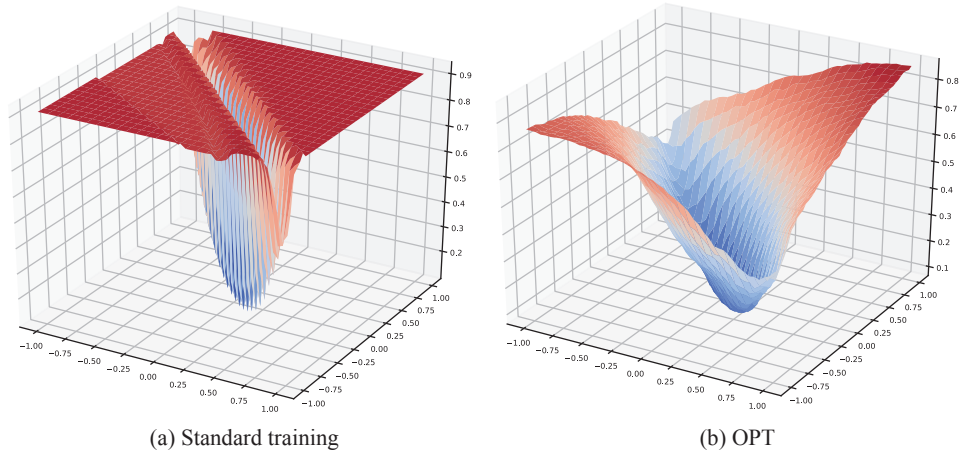


Figure 21: Comparison of testing error landscapes between standard training and OPT on CIFAR-10 (CNN-6).

G. Theoretical Discussion on Optimization and Generalization

The key problem we discuss in this section is why OPT may lead to easier optimization and better generalization. We have already shown that OPT can guarantee the minimum hyperspherical energy (MHE) in a probabilistic sense. Although empirical evidences [49] have shown significant and consistent performance gain by minimizing hyperspherical energy, why lower hyperspherical energy will lead to better generalization is still unclear. We argue that OPT leads to better generalization from two aspects: how OPT may affect the training and generalization, and why minimum hyperspherical energy serves as a good inductive bias. We note that rigorously proving that OPT generalizes better is out of the scope of this paper and remains our future work. The section serves as a very preliminary discussion for this topic, and hopefully the discussion can inspire more theoretical studies about OPT.

Our goal here is to leverage and apply existing theoretical results [34, 77, 66, 40, 12, 1] to explain the role that MHE plays rather than proving sharp and novel generalization bounds. We emphasize that our paper is *NOT* targeted as a theoretical one that proves novel generalization bounds.

We simply consider one-hidden-layer networks as the hypothesis class:

$$\mathcal{F} = \{f(x) = \sum_{j=1}^n v_j \sigma(\mathbf{w}_j^\top \mathbf{x}) : v_j \in \{\pm 1\}, \sum_{j=1}^n \|\mathbf{w}_j\| \leq C_w\} \quad (42)$$

where $\sigma(\cdot) = \max(0, \cdot)$ is ReLU. Since the magnitude of v_j can be scaled into \mathbf{w}_j , we can restrict v_j to be ± 1 . Given a set of *i.i.d.* training sample $\{\mathbf{x}_i, y_i\}_{i=1}^m$ where $\mathbf{x} \in \mathbb{R}^d$ is drawn uniformly from the unit hypersphere, we minimize the least square loss $\mathcal{L} = \frac{1}{2m} \sum_{i=1}^m (y_i - f(\mathbf{x}_i))^2$. The gradient w.r.t. \mathbf{w}_i is

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_j} = \frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i) v_j \sigma'(\mathbf{w}_j^\top \mathbf{x}_i) \mathbf{x}_i. \quad (43)$$

Let $\mathbf{W} := \{\mathbf{w}_1^\top, \dots, \mathbf{w}_n^\top\}^\top$ be the column concatenation of neuron weights. We aim to identify the conditions under which there are no spurious local minima. We rewrite that

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}_1}, \dots, \frac{\partial \mathcal{L}}{\partial \mathbf{w}_n} \right)^\top = \mathbf{D} \mathbf{r} \quad (44)$$

where $\mathbf{r} \in \mathbb{R}^m$ $\mathbf{r}_i = \frac{1}{m} f(\mathbf{x}_i) - y_i$, $\mathbf{D} \in \mathbb{R}^{n \times m}$, and $\mathbf{D}_{ij} = v_i \sigma'(\mathbf{w}_i^\top \mathbf{x}_j) \mathbf{x}_j$. Therefore, we can obtain the following inequality:

$$\|\mathbf{r}\| \leq \frac{1}{s_m(\mathbf{D})} \left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}} \right\| \quad (45)$$

where $\|\mathbf{r}\|$ is the training error and $s_m(\mathbf{D})$ is the minimum singular value of \mathbf{D} . If we need the training error to be small, then we have to lower bound $s_m(\mathbf{D})$ away from zero. Therefore, the essential problem now becomes the relationship between MHE and the lower bound of $s_m(\mathbf{D})$. We have the following result from [77]:

Lemma 4. *With probability larger than $1 - m \exp(-m\gamma_m/8) - 2m^2 \exp(-4 \log^2 d) - \delta$, we will have that $s_m(\mathbf{D})^2 \geq \frac{1}{2} nm\gamma_m - cn\rho(\mathbf{W})$ where*

$$\begin{aligned} \rho(\mathbf{W}) \leq & \frac{\log d}{\sqrt{d}} \sqrt{2L_2(\mathbf{W})} m \left(\frac{4}{m} \log \frac{1}{\delta} \right)^{1/4} \\ & + \frac{2 \log d}{\sqrt{d}} m \sqrt{\frac{4}{3m} \log \frac{1}{\delta}} + \frac{\log d}{\sqrt{d}} m L_2(\mathbf{W}) + 2, \end{aligned} \quad (46)$$

and $L_2(\mathbf{W}) = \frac{1}{n^2} \sum_{i,j=1}^n k(\mathbf{w}_i, \mathbf{w}_j)^2 - \mathbb{E}_{\mathbf{u}, \mathbf{v}}[k(\mathbf{u}, \mathbf{v})^2]$. The kernel function $k(\mathbf{u}, \mathbf{v})$ is $\frac{1}{2} - \frac{1}{2\pi} \arccos(\frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\| \|\mathbf{v}\|})$.

Once MHE is achieved, the neurons will be uniformly distributed on the unit hypersphere. From Lemma 4, we can see that if the neurons are uniformly distributed on the unit hypersphere, $L_2(\mathbf{W})$ will be very small and close to zero. Then $\rho(\mathbf{W})$ will also be small, leading to large lower bound for $s_m(\mathbf{D})$. Therefore, MHE can result in small training error once the gradient norm $\left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}} \right\|$ is small. The result implies no spurious local minima if we use OPT for training.

Furthermore, suppose that $\|\frac{\partial \mathcal{L}}{\partial \mathbf{W}}\|^2 \leq \epsilon$, [77] also proves a training error bound $\tilde{\mathcal{O}}(\epsilon)$ and a generalization bound $\tilde{\mathcal{O}}(\epsilon + \frac{1}{\sqrt{m}})$ based on the assumption that \mathbf{W} belongs to a specific set $\mathcal{G}_{\mathbf{W}}$ (for the definition of $\mathcal{G}_{\mathbf{W}}$, please refer to [77]). Therefore, MHE is also connected to the training and generalization error. Note that, the analysis here is highly simplified and the purpose here is to give some justifications rather than rigorously proving any bound.

We further argue that MHE induced by OPT serves as an important inductive bias for neural networks. As the standard regularizer for neural networks, weight decay controls the norm of the neuron weights, regularizing essentially one dimension of the weight. In contrast, MHE completes an important missing piece by regularizing the remaining dimensions of the weight. MHE encourages minimum hyperspherical redundancy between neurons. In the linear classifier case, MHE impose a prior of maximal inter-class separability.

H. More Discussions

Semi-randomness. OPT fixes the randomly initialized neuron weight vectors and simply learns layer-shared orthogonal matrices, so OPT naturally imposes strong randomness to the neurons. OPT well combines the good generalizability from randomness and the strong approximation power from neural networks. Such randomness suggests that the specific configuration of relative position among neurons does not matter that much, and the coordinate system is more crucial for generalization. [35, 59, 67] also show that randomness can be beneficial to generalization.

Flexible training. First, OPT can be used in multi-task training [55] where each set of orthogonal matrices represent one task. OPT can learn different set of orthogonal matrices for different tasks with the neuron weights remain the same. Second, we can perform progressive training with OPT. For example, after learning a set of orthogonal matrices on a large coarse-grained dataset (*i.e.*, pretraining), we can multiple the orthogonal matrices back to the neuron weights and construct a new set of neuron weights. Then we can use the new neuron weights as a starting point and apply OPT to train on a small fine-grained dataset (*i.e.*, finetuning).

Limitations and open problems The limitations of OPT include more GPU memory consumption and heavy computation during training, more numerical issues when ensuring orthogonality and weak scalability for ultra wide neural networks. Therefore, there will be plenty of open problems in OPT, such as scalable and efficient training. Most significantly, OPT opens up a new possibility for studying theoretical generalization of deep networks. With the decomposition to hyperspherical energy and coordinate system, OPT provides a new perspective for future research.

I. On Parameter-Efficient OPT

I.1. Formulation

Since OPT over-parameterizes the neurons, it will consume more GPU memory in training (note that, the number of parameters will not increase in testing). For a d -dimensional neuron, OPT will learn an orthogonal matrix of size $d \times d$ that applies to the neuron. Therefore, we will need d^2 extra parameters for one layer of neurons, making the training more expensive in terms of the GPU memory. Although the extra training overhead in OPT will not affect the inference speed of the trained neural networks, we still desire to achieve better parameter efficiency in OPT. To this end, we discuss some design possibilities for the *parameter-efficient OPT* (PE-OPT) in this section.

Original OPT over-parameterize a neuron $\mathbf{v} \in \mathbb{R}^{n \times n}$ with $\mathbf{R}\mathbf{v}$ where \mathbf{R} is a layer-shared orthogonal matrix of size $d \times d$. We aim to reduce the effective parameters of this $d \times d$ orthogonal matrix. We incorporate a block-diagonal structure to the orthogonal matrix \mathbf{R} . Specifically, we formulate \mathbf{R} as $\text{Diag}(\mathbf{R}^{(1)}, \mathbf{R}^{(2)}, \dots, \mathbf{R}^{(k)})$ where $\mathbf{R}^{(i)}$ is an orthogonal matrix with size $d_i \times d_i$ (it is easy to see that we need $d = \sum_i d_i$). As an example, we only consider the case where all $\mathbf{R}^{(i)}$ are of the same size (i.e., $d_1 = d_2 = \dots = d_k = \frac{d}{k}$). It is also obvious that as long as each block is an orthogonal matrix, then the overall matrix \mathbf{R} remains an orthogonal matrix.

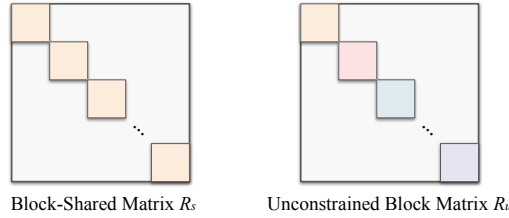


Figure 22: Comparison between the block-shared matrix \mathbf{R}_s and the unconstrained block matrix \mathbf{R}_u .

First, we consider that all the block matrices on the diagonal of the orthogonal matrix \mathbf{R} are shared, meaning that $\mathbf{R} = \text{Diag}(\mathbf{R}^{(1)}, \mathbf{R}^{(1)}, \dots, \mathbf{R}^{(1)})$ (i.e., $\mathbf{R}^{(1)} = \mathbf{R}^{(2)} = \dots = \mathbf{R}^{(k)}$). Therefore, we have a block-diagonal matrix \mathbf{R}_s with shared block $\mathbf{R}^{(1)}$ as the final orthogonal matrix for the neuron \mathbf{v} :

$$\mathbf{R}_s = \begin{bmatrix} \mathbf{R}^{(1)} & 0 & \dots & 0 \\ 0 & \mathbf{R}^{(1)} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mathbf{R}^{(1)} \end{bmatrix} \quad (47)$$

where $\mathbf{R}^{(1)} \in \mathbb{R}^{\frac{d}{k} \times \frac{d}{k}}$. The effective number of parameters for the orthogonal matrix \mathbf{R}_s immediately reduces to $\frac{d^2}{k}$. The left figure in Fig. 22 gives an intuitive illustration for the block-shared matrix \mathbf{R}_s . Therefore, PE-OPT only needs to learn $\mathbf{R}^{(1)}$ in order to construct the orthogonal matrix of size $d \times d$.

Second, we consider that all the diagonal block matrices are independent, indicating that $\mathbf{R} = \text{Diag}(\mathbf{R}^{(1)}, \mathbf{R}^{(2)}, \dots, \mathbf{R}^{(k)})$ where $\mathbf{R}^{(i)}, \forall i$ are different orthogonal matrices in general. We term such matrix \mathbf{R} as unconstrained block matrix. Therefore, we have the unconstrained block diagonal matrix \mathbf{R}_u as

$$\mathbf{R}_u = \begin{bmatrix} \mathbf{R}^{(1)} & 0 & \dots & 0 \\ 0 & \mathbf{R}^{(2)} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mathbf{R}^{(k)} \end{bmatrix} \quad (48)$$

where the orthogonal matrices $\mathbf{R}^{(i)}, \forall i$ will be learned independently. The effective number of parameters for the orthogonal matrix \mathbf{R}_u is $\frac{d^2}{k}$, making it more flexible than the block-shared matrix \mathbf{R}_s .

Let's consider a convolution neuron (i.e., convolution filter) $\mathbf{v} \in \mathbb{R}^{c_1 \times c_2 \times c_3}$ (e.g., a typical convolution neuron is of size $3 \times 3 \times 64$) as an example. The orthogonal matrix \mathbf{R} for the convolution neuron is of size $(c_1 c_2 c_3) \times (c_1 c_2 c_3)$. Typically, we will divide the neuron into k sub-neuron along the c_3 -axis, each with size $c_1 \times c_2 \times \frac{c_3}{k}$. Then in order to learn a block-shared

orthogonal matrix \mathbf{R}_s , we will essentially learn a shared orthogonal matrix of size $(\frac{1}{k}c_1c_2c_3) \times (\frac{1}{k}c_1c_2c_3)$ that applies to each sub-neuron (there are k sub-neurons of size $c_1 \times c_2 \times \frac{c_3}{k}$ in total). For the case of learning a unconstrained block-diagonal orthogonal matrix \mathbf{R}_u , we simply learn different orthogonal matrices for different sub-neurons.

I.2. Experiments and Results

We conduct the image recognition experiments on CIFAR-100 with CNN-6 described in Table 14. The setting is exactly the same as Section 6.2. For the convolution filter, we use the size of $3 \times 3 \times 64$, *i.e.*, $c_1 = 3, c_2 = 3, c_3 = 64$. The results are given in Table 19 and Table 20. “# Parameters” in both tables denote the number of effective parameters for the orthogonal matrix \mathbf{R} in a single layer. The baseline with fixed neurons is only to train the final classifiers with the randomly initialized neuron weights staying fixed. It means that this baseline basically removes the learnable orthogonal matrices but still fixes the neuron weights, so it only achieves 73.81% testing error. As expected, as the number of effective parameters goes down, the performance of PE-OPT generally decreases. One can also observe that using separate orthogonal matrices generally yields better performance than shared orthogonal matrices. $k = 2$ and $k = 4$ seems to be a reasonable trade-off between better accuracy and less parameters.

When k becomes larger (*i.e.*, the number of parameters become less) in the case of block-shared orthogonal matrices, we find that PE-OPT (LS) performs the best among all the variants. When k becomes larger (*i.e.*, the number of parameters become less) in the case of unconstrained block orthogonal matrices, we can see that both PE-OPT (GS) and PE-OPT (LS) performs better than the other variants.

| Method | # Parameters | PE-OPT (CP) | PE-OPT (GS) | PE-OPT (HR) | PE-OPT (LS) | PE-OPT (OGD) |
|--|--------------|-------------|--------------|-------------|--------------|--------------|
| $c_3/k = 64$ ($k = 1$) (<i>i.e.</i> , Original OPT) | 331.7K | 33.53 | 33.02 | 35.67 | 34.48 | 33.33 |
| $c_3/k = 32$ ($k = 2$) | 82.9K | 34.93 | 34.39 | 35.83 | 34.50 | 35.06 |
| $c_3/k = 16$ ($k = 4$) | 20.7K | 39.40 | 39.13 | 39.67 | 37.58 | 39.80 |
| $c_3/k = 8$ ($k = 8$) | 5.2K | 47.77 | 46.65 | 46.69 | 45.62 | 47.43 |
| $c_3/k = 4$ ($k = 16$) | 1.3K | 56.65 | 55.91 | 55.69 | 54.75 | 57.15 |
| $c_3/k = 2$ ($k = 32$) | 0.3K | 63.46 | 62.65 | 62.38 | 61.60 | 62.46 |
| $c_3/k = 1$ ($k = 64$) | 0.1K | 67.36 | 67.11 | 67.05 | 66.61 | 67.23 |
| Baseline | - | | | 37.59 | | |
| Baseline with fixed random neurons | - | | | 73.81 | | |

Table 19: Testing error (%) on CIFAR-100 with different settings of PE-OPT (with block-shared orthogonal matrix \mathbf{R}_s).

| Method | # Parameters | PE-OPT (CP) | PE-OPT (GS) | PE-OPT (HR) | PE-OPT (LS) | PE-OPT (OGD) |
|--|--------------|-------------|--------------|-------------|--------------|--------------|
| $c_3/k = 64$ ($k = 1$) (<i>i.e.</i> , Original OPT) | 331.7K | 33.53 | 33.02 | 35.67 | 34.48 | 33.33 |
| $c_3/k = 32$ ($k = 2$) | 165.9K | 33.54 | 33.15 | 35.65 | 34.09 | 34.27 |
| $c_3/k = 16$ ($k = 4$) | 82.9K | 34.77 | 34.50 | 35.71 | 34.96 | 35.97 |
| $c_3/k = 8$ ($k = 8$) | 41.5K | 37.25 | 36.43 | 36.40 | 36.17 | 39.75 |
| $c_3/k = 4$ ($k = 16$) | 20.7K | 40.74 | 39.89 | 39.98 | 39.93 | 43.43 |
| $c_3/k = 2$ ($k = 32$) | 10.4K | 45.36 | 44.77 | 44.83 | 44.61 | 48.98 |
| $c_3/k = 1$ ($k = 64$) | 5.2K | 50.94 | 49.16 | 49.57 | 49.23 | 54.93 |
| Baseline | - | | | 37.59 | | |
| Baseline with fixed random neurons | - | | | 73.81 | | |

Table 20: Testing error (%) on CIFAR-100 with different settings of PE-OPT (with unconstrained block orthogonal matrix \mathbf{R}_u).

J. On Generalizing OPT: Over-Parameterized Training with Constraint

OPT opens many new possibilities in training neural networks. We consider a simple generalization to OPT in this section to showcase the great potential of OPT. Instead of constraining the over-parameterization matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$ in Eq. 1 to be orthogonal, we can use any meaningful structural constraints for this matrix, and even regularize it in a task-driven way. Furthermore, instead of a linear over-parameterization (*i.e.*, multiplying a matrix \mathbf{R}) to the neuron, we can also consider nonlinear mapping. We come up with the following straightforward generalization to OPT (the settings and notations exactly follow Eq. 1):

$$\begin{aligned}
&\text{Standard: } \min_{\mathbf{v}_i, u_i, \forall i} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i \mathbf{v}_i^\top \mathbf{x}_j) \\
&\text{Original OPT: } \min_{\mathbf{R}, u_i, \forall i} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i (\mathbf{R} \mathbf{v}_i)^\top \mathbf{x}_j) \\
&\quad \text{s.t. } \mathbf{R}^\top \mathbf{R} = \mathbf{R} \mathbf{R}^\top = \mathbf{I} \\
&\text{Generalized OPT: } \min_{\mathbf{R}, u_i, \forall i} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i (\mathcal{T}(\mathbf{v}_i))^\top \mathbf{x}_j) \\
&\quad \text{s.t. Some constraints on } \mathcal{T}(\cdot)
\end{aligned} \tag{49}$$

where $\mathcal{T}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ denotes some transformation (including both linear and nonlinear). Notice that the generalized OPT (G-OPT) no longer requires orthogonality. Such formulation of G-OPT can immediately inspire a number of instances. We will discuss some obvious ones here.

If we consider $\mathcal{T}(\cdot)$ to be a linear mapping, we may constrain \mathbf{R} to be symmetric positive definite other than orthogonal. A simple way to achieve that is to use Cholesky factorization $\mathbf{L} \mathbf{L}^\top$ where \mathbf{L} is a lower triangular matrix to parameterize the matrix \mathbf{R} . Essentially, we learn a lower triangular matrix \mathbf{L} and use $\mathbf{L} \mathbf{L}^\top$ to replace \mathbf{R} in OPT. The positive definiteness provides the transformation \mathbf{R} with some geometric constraint. Specifically, a positive definite \mathbf{R} only transforms the neuron weight \mathbf{v} to the direction that has the angle less than $\frac{\pi}{2}$ to \mathbf{v} , because $\mathbf{v}^\top \mathbf{R} \mathbf{v} > 0$. Moreover, we can also require the transformation to have structural constraints on \mathbf{R} . For example, \mathbf{R} can be upper (lower) triangular, banded, symmetric, skew-symmetric, upper (lower) Hessenberg, etc.

We can also consider $\mathcal{T}(\cdot)$ to be a nonlinear mapping. A obvious example is to use a neural network (*e.g.*, MLP, CNN) as $\mathcal{T}(\cdot)$. Then the nonlinear G-OPT will share some similarities with HyperNetworks [18] and Network-in-Network [46]. If we further consider $\mathcal{T}(\cdot)$ to be dependent on the input, then the nonlinear G-OPT will have close connections to dynamic neural networks [32, 51].

To summarize, OPT provides a novel and effective framework to train neural networks and may inspire many different threads of future research.

K. Hyperspherical Energy Training Dynamics of Individual Layers

We also plot the hyperspherical energy ($E(\hat{\mathbf{v}}_i|_{i=1}^n) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n \|\hat{\mathbf{v}}_i - \hat{\mathbf{v}}_j\|^{-1}$ in which $\hat{\mathbf{v}}_i = \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}$ is the i -th neuron weight projected onto the unit hypersphere.) in every layer of CNN-6 during training to show how these hyperspherical energies are being minimized. From Fig. 23, we can observe that OPT can always maintain the minimum hyperspherical energy during the entire training process, while the MHE regularization cannot. Moreover, the hyperspherical energy of the baseline will also decrease as the training proceeds, but it is still much higher than the OPT training.

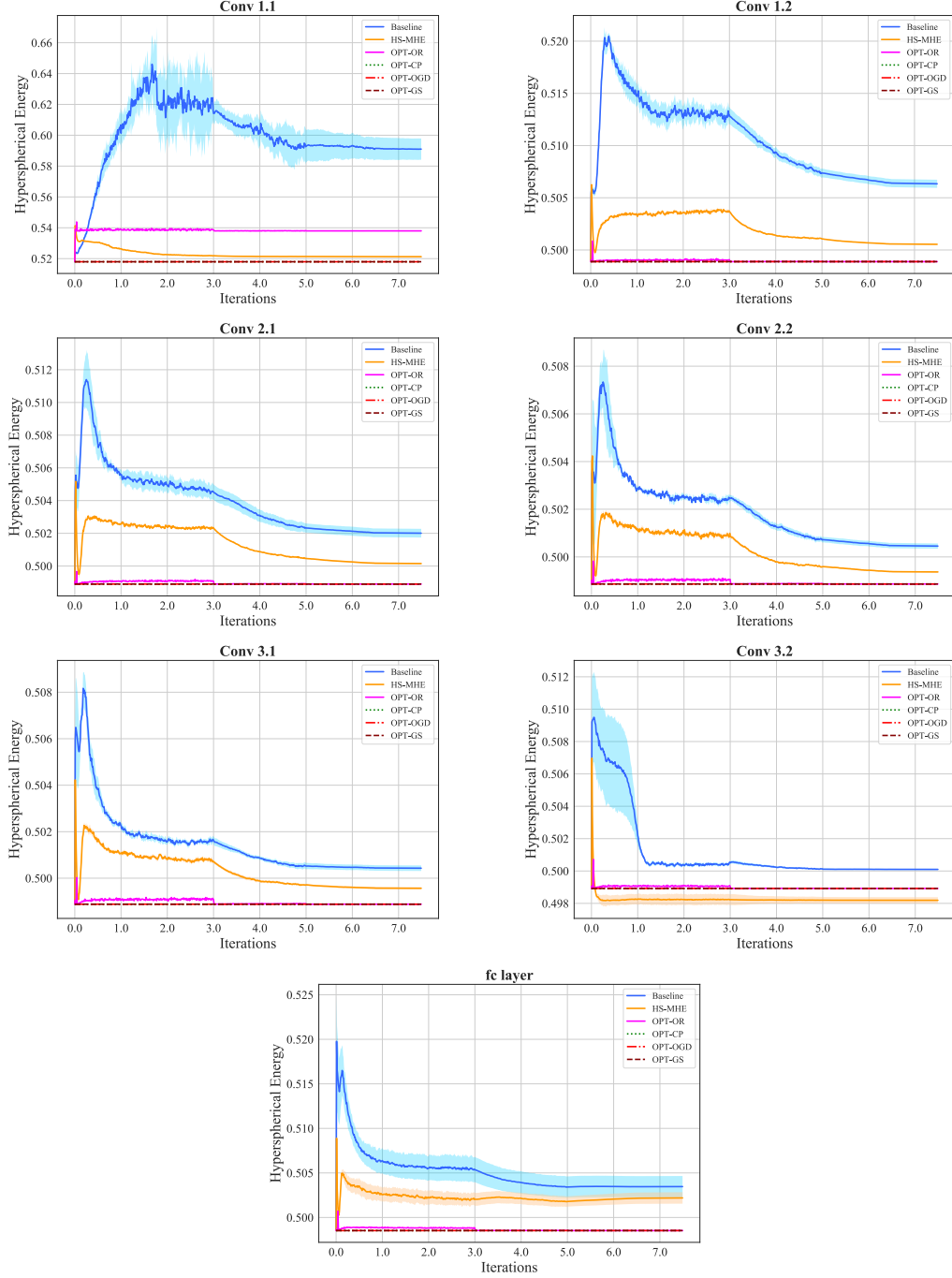


Figure 23: Training dynamics of hyperspherical energy in each layer of CNN-6. We average results with 10 runs.

L. Geometric Properties of Randomly Initialized Neurons

There are many interesting geometric properties [5, 6] of random points distributed independently and uniformly on the unit hypersphere. We summarize a few of them that make randomly initialized neurons distinct from any deterministic neuron configuration. Note that, there exist many deterministic neuron configurations that can also achieve very low hyperspherical energy, and this section aims to describe a few unique geometric properties of randomly initialized neurons.

There are two widely used geometric properties corresponding to a neuron configuration (*i.e.*, a set of neurons) $\hat{\mathbf{W}}_N = \{\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_N \in \mathbb{S}^d\}$. In the main paper, we define neurons on \mathbb{S}^{d-1} , but without loss of generality we define neurons on \mathbb{S}^d here for convenience. The first one is the *covering radius*:

$$\alpha(\hat{\mathbf{W}}_N) := \alpha(\hat{\mathbf{W}}_N; \mathbb{S}^{d-1}) := \max_{\mathbf{u} \in \mathbb{S}^d} \min_{1 \leq i \leq N} \arccos(\mathbf{u}, \hat{\mathbf{w}}_i) \quad (50)$$

which is the biggest geodesic distance from a neuron in \mathbb{S}^d to the nearest point in $\hat{\mathbf{W}}_N$. The second one is the *separation distance*:

$$\psi(\hat{\mathbf{W}}_N) := \min_{1 \leq i, j \leq N, i \neq j} \arccos(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j) \quad (51)$$

which gives the least geodesic distance between arbitrary two points in $\hat{\mathbf{W}}_N$. Random points (*i.e.*, randomly initialized neurons) typically have poor separation properties, since the separation is sensitive to the specific placement of points. [5] shows an example on \mathbb{S}^1 to illustrate this observation.

[5] considers a different but related quantity, *i.e.*, the sum of powers of the “hole radii”. A set of neurons $\hat{\mathbf{W}}_N$ on \mathbb{S}^d uniquely defines a convex polytope, which can be viewed as the convex hull of the neuron configuration. Each facet of the polytope defines a “hole”. Such a hole denotes the maximal spherical cap for a facet that contains neurons of $\hat{\mathbf{W}}_N$ only on the boundary. It is easy to see that the geodesic radius of the largest hole is the covering radius $\alpha(\hat{\mathbf{W}}_N)$. We assume that for the set of neurons $\hat{\mathbf{W}}_N$, there are f_d holes (*i.e.*, facets) in total. Therefore, the i -th hole radius is defined as $\rho_i = \rho_i(\hat{\mathbf{W}}_N)$ which is the Euclidean distance in \mathbb{R}^{d+1} between the center of the i -th spherical cap and the boundary. The i -th spherical cap is located on the sphere corresponding to the i -th facet. We have that $\rho_i = 2 \sin(\frac{\alpha_i}{2})$ where α_i is the geodesic radius of the i -th spherical cap. We are interested in the sums of the p -th powers of the hole radii, *i.e.*,

$$\mathcal{P} = \sum_{i=1}^{f_d} (\rho_i)^p \quad (52)$$

where p is larger than zero. For large p , the largest hole dominates:

$$\lim_{p \rightarrow \infty} (\mathcal{P})^{\frac{1}{p}} = \lim_{p \rightarrow \infty} \left(\sum_{i=1}^{f_d} (\rho_i)^p \right)^{\frac{1}{p}} = \max_{1 \leq i \leq f_d} \rho_i = 2 \sin\left(\frac{\alpha(\hat{\mathbf{W}}_N)}{2}\right) \quad (53)$$

where $\rho(\hat{\mathbf{W}}_N) := \max_{1 \leq i \leq f_d} \rho_i$. Then we introduce some useful notations to state the geometric properties. Let ψ_d be the surface area of \mathbb{S}^d , and we have that

$$\psi_d = \frac{2\pi^{\frac{d+1}{2}}}{\Gamma(\frac{d+1}{2})}, \quad (54)$$

and we also define the following quantities (with $\psi_0 = 2$):

$$\begin{aligned} \kappa_d &:= \frac{1}{d} \frac{\psi_{d-1}}{\psi_d} = \frac{1}{d} \frac{\Gamma(\frac{d+1}{2})}{\sqrt{\pi} \Gamma(\frac{d}{2})} \\ B_d &:= \frac{2}{d+1} \frac{\kappa_{d^2}}{(\kappa_d)^d} \end{aligned} \quad (55)$$

where κ_d can be alternatively defined with the recursion: $\kappa_1 = \frac{1}{\pi}$ and $\kappa_d = \frac{1}{2\pi d \kappa_{d-1}}$. [61] gives the expected number of facets constructed from N random neurons that are independently and uniformly distributed on the unit hypersphere \mathbb{S}^d :

$$\mathbb{E}[f_d] = B_d N (1 + o(1)) \quad (56)$$

where $N \rightarrow \infty$. Then we introduce the main results of [5] (asymptotics for the expected moments of the hole radii) in the following lemma:

Lemma 5. *If $p \geq 0$ and $\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_N$ are N neurons on \mathbb{S}^d that are independently and randomly distributed with respect to the normalized surface area measure σ_d on \mathbb{S}^d , then we have that*

$$\begin{aligned}\mathbb{E}[\mathcal{P}] &= B_d(\kappa_d)^{-\frac{p}{d}} \frac{\Gamma(d + \frac{p}{d})\Gamma(N+1)}{\Gamma(d)\Gamma(N + \frac{p}{d})} (1 + \mathcal{O}(N^{-\frac{2}{d}})) \\ &= c_{d,p} N^{1-\frac{p}{d}} (1 + \mathcal{O}(N^{-\frac{2}{d}}))\end{aligned}\tag{57}$$

as $N \rightarrow \infty$, where $\rho_i = \rho_{i,N}$ is the Euclidean hole radius associated with the i -th facet of the convex hull of $\hat{\mathbf{W}}_N$, $c_{d,p} := B_d B_{d,p}$, and $B_{d,p} := \frac{\Gamma(d+\frac{p}{d})}{\Gamma(d)} (\kappa_d)^{-\frac{p}{d}}$. The \mathcal{O} -terms above depend on d and p .

As we mentioned, there are many deterministic point (*i.e.*, neuron) configurations such as minimizing hyperspherical energy (*i.e.*, Riesz s -energy) [49] (as $s \rightarrow \infty$, the minimal s -energy points approach the best separation), maximizing the determinant for polynomial interpolation [64], Fibonacci points, spherical t -designs, minimizing covering radius (*i.e.*, best covering problem), maximizing the separation (*i.e.*, best packing problem) and maximizing the s -polarization, etc. We note that randomly initialized neurons are quite different from these deterministic neuron configurations and have unique geometric properties.