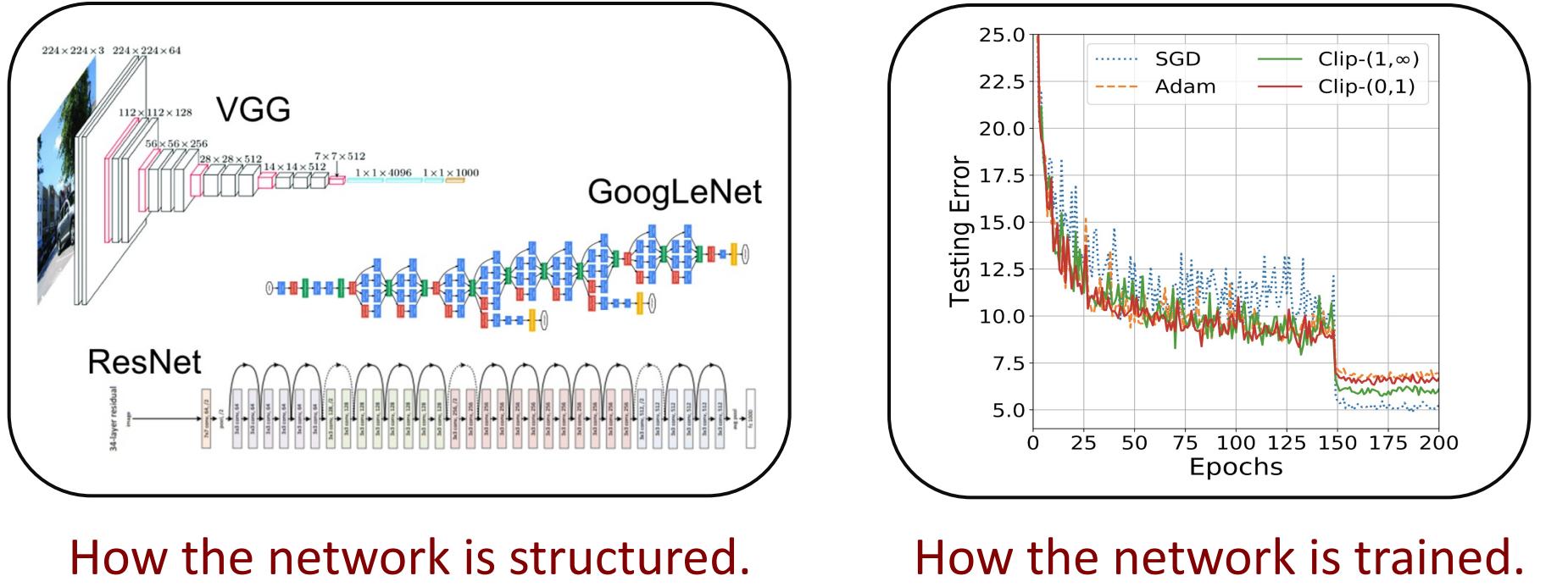


Orthogonal Over-Parameterized Training

Weiyang Liu* Rongmei Lin* Zhen Liu James Rehg Liam Paull Li Xiong Le Song Adrian Weller

Introduction

Empirical Generalization of Neural Networks



Motivation I: Over-parameterization

- Recent theories suggest the importance of over-parameterization in linear neural networks.
- For example,
$$\mathbf{Y} \approx \mathbf{X} = \mathbf{U} \times \mathbf{V}^T$$

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} \|\text{observed}(\mathbf{X}) - \mathbf{y}\|_2^2 \equiv \min_{\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times n}} \|\text{observed}(\mathbf{UV}^T) - \mathbf{y}\|_2^2$$

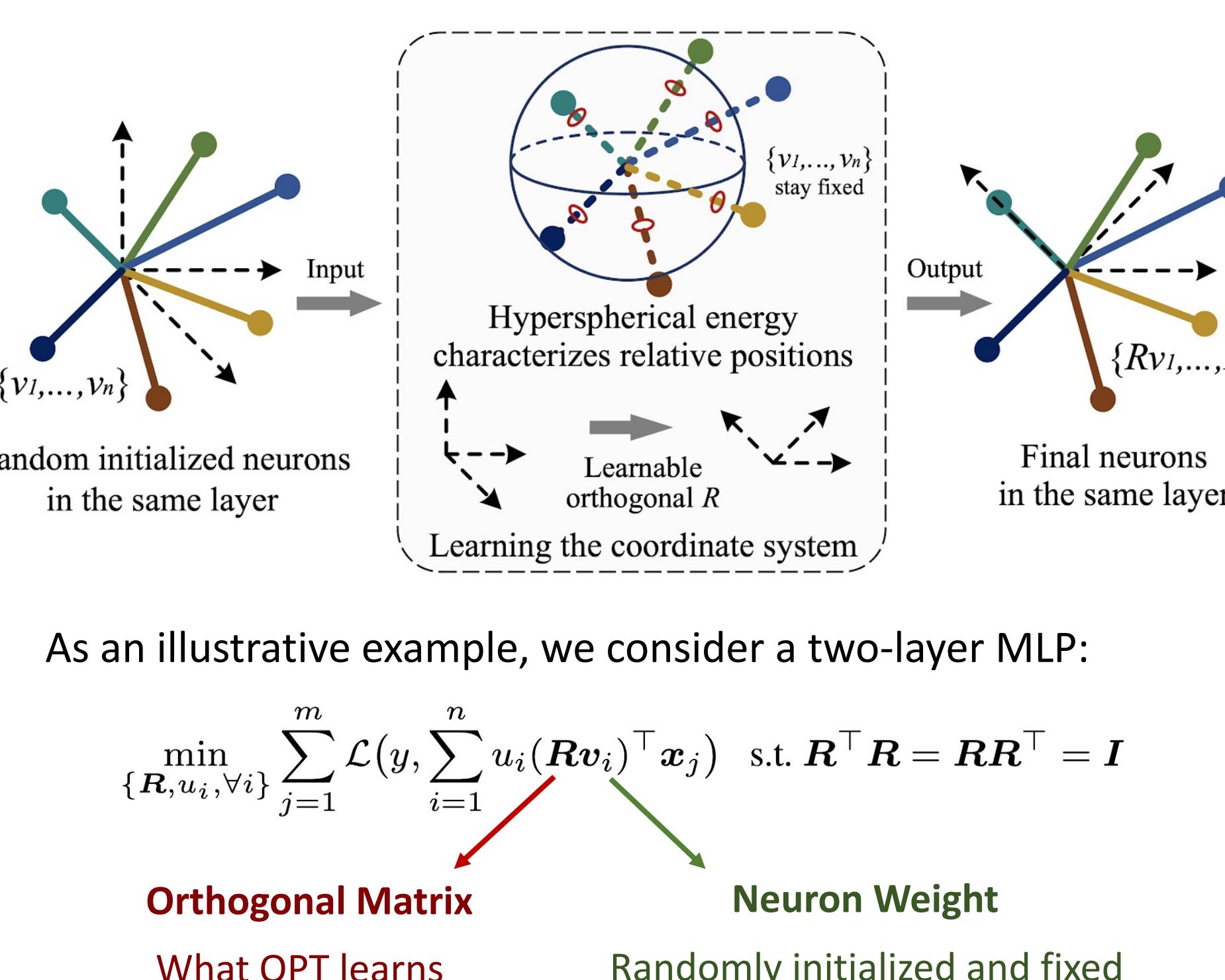
Motivation II: Minimum Hyperspherical Energy (MHE)

- Definition of hyperspherical energy:
$$\min_{\{\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_n \in \mathbb{S}^{d-1}\}} \{E_s(\hat{\mathbf{W}}_n) := \sum_{i=1}^n \sum_{j=1, j \neq i}^n K_s(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j)\}$$
- where $\hat{\mathbf{w}}_i := \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|}$ $K_s(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j) = \begin{cases} \rho(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j)^{-s}, & s > 0 \\ \log(\rho(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j)^{-1}), & s = 0 \\ -\rho(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j)^{-s}, & s < 0 \end{cases}$
- ρ denotes either Euclidean distance or angular distance on the unit hypersphere.
 - Hyperspherical energy characterizes the neuron diversity of the neural network.
 - Previous work shows that lower hyperspherical energy leads to better empirical generalization.

Goal: A Principled Training Framework for Neural Networks

- Make use of the over-parameterization within each neuron
- Naturally guarantee the minimum hyperspherical energy
- Compatible to different network architectures and optimizers

Orthogonal Over-Parameterized Training (OPT)



- As an illustrative example, we consider a two-layer MLP:

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

Orthogonal Matrix

What OPT learns

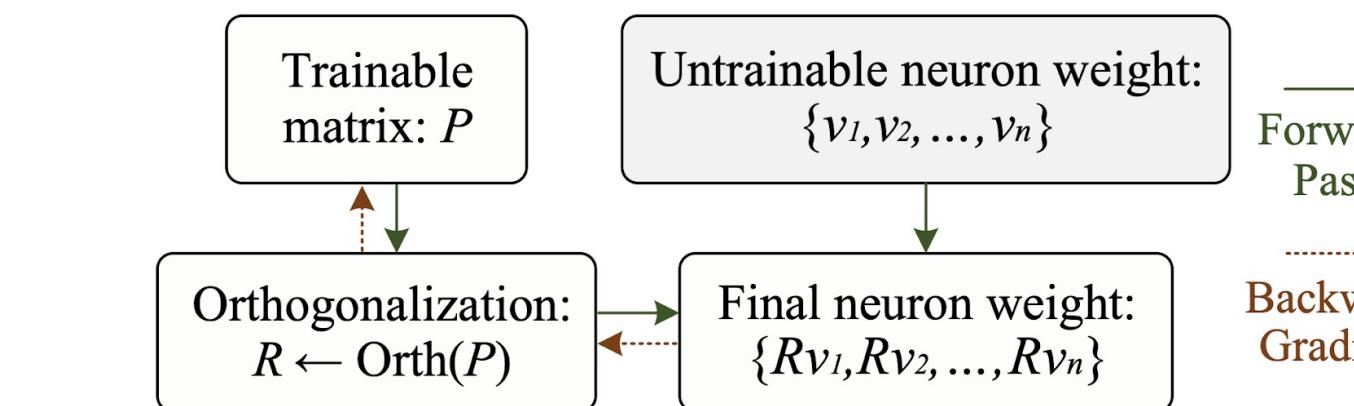
Neuron Weight

Randomly initialized and fixed

Ways to guarantee orthogonality

- Unrolling orthogonalization algorithms:

- Gram-Schmidt Process
- Householder reflection
- Lowdin's Symmetric Orthogonalization



- Orthogonal parameterization:

$$\mathbf{R} = (\mathbf{I} + \mathbf{W})(\mathbf{I} - \mathbf{W})^{-1} \quad \mathbf{W} = -\mathbf{W}^T$$

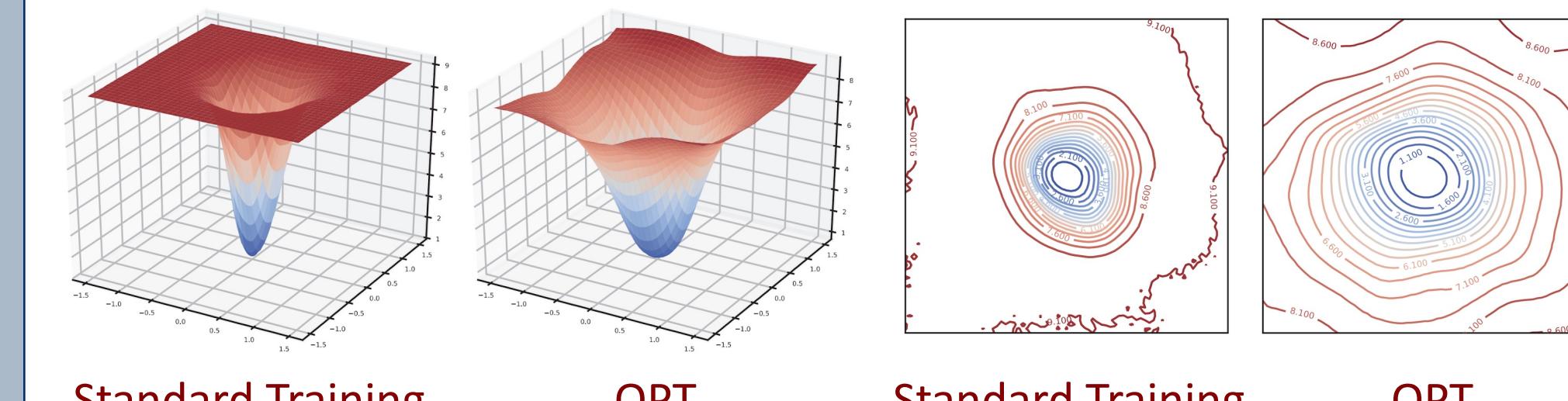
- Orthogonality-preserving gradient descent

- Orthogonality as a regularization:

$$\min_{\mathbf{R}, u_i, \forall i} \sum_{j=1}^m \mathcal{L}(y, \sum_{i=1}^n u_i (\mathbf{R} \mathbf{v}_i)^T \mathbf{x}_j) + \beta \|\mathbf{R}^T \mathbf{R} - \mathbf{I}\|_F^2$$

Intriguing Insights and Extension

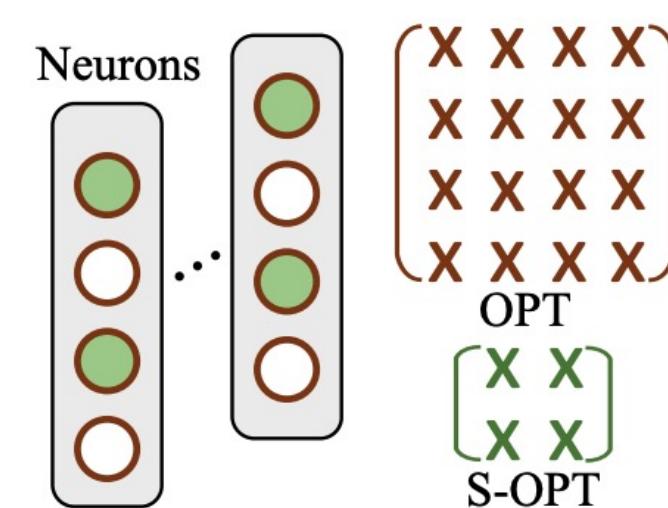
Loss Landscape Comparison



Standard Training OPT Standard Training OPT

Stochastic OPT for Better Scalability

- Approximate a large orthogonal matrix with many small orthogonal ones.
- Similar to the idea of DropOut.
- Randomly selecting a subset of the neuron dimensions in each iteration and perform OPT on this subset.



Experiments and Results

Ablation and exploratory experiments

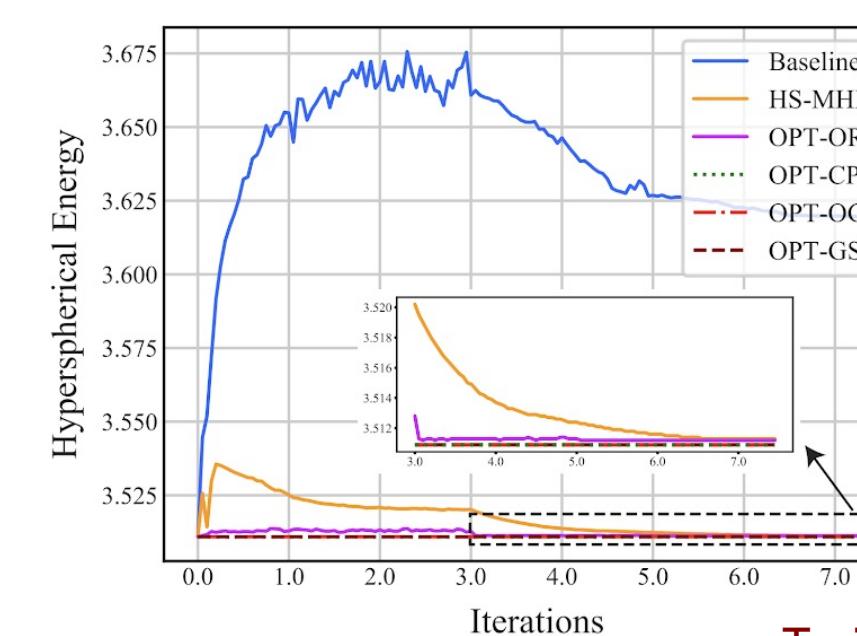
Ablation

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

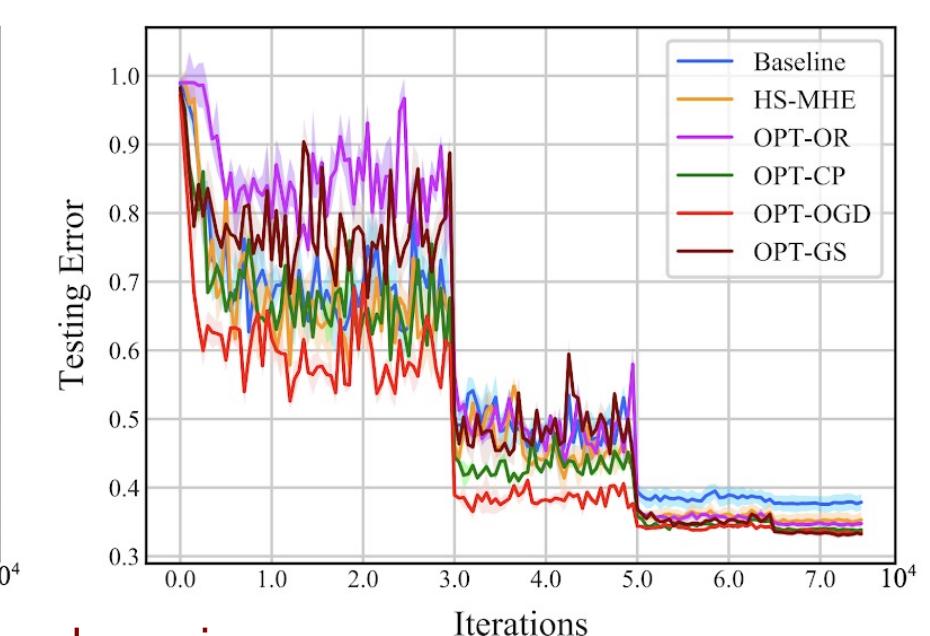
FN: whether neurons are fixed after initialization
LR: whether we enforce orthogonality on R .

	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

CIFAR-100



CIFAR-100



Results on OPT and Stochastic OPT

OPT for PointNet

OPT for MLP, plain CNN and ResNet

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

Method	GCN	Pubmed	MN-40
Baseline	81.3	79.0	87.1
OPT (GS)	81.0	79.4	87.23
OPT (CP)	82.0	79.4	87.81
OPT (OGD)	82.3	79.5	87.86

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

S-OPT for plain CNN and ResNet

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

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

Large Categorical Training

Method	ResNet-18A		ResNet-18B	
Error	Params	Error	Params	

</tbl