



EMORY  
UNIVERSITY

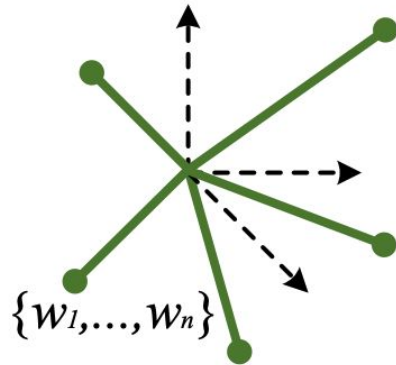


# Regularizing Neural Networks via Minimizing Hyperspherical Energy

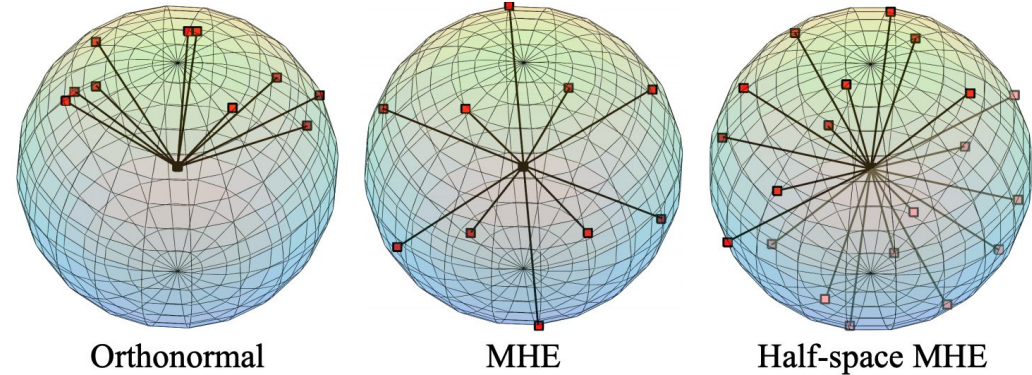
Rongmei Lin, Weiyang Liu, Zhen Liu, Chen Feng, Zhiding Yu,  
James Rehg, Li Xiong, Le Song

# Hyperspherical Energy and Motivation

## Neurons in one layer:



## Minimum hyperspherical energy (MHE):



## Hyperspherical energy: $(\hat{w}_i = \frac{w_i}{\|w_i\|})$

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

Minimizing the hyperspherical energy promotes the **diversity** of neurons on a hypersphere.

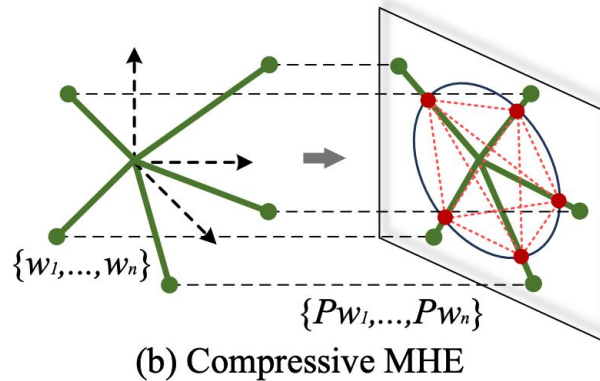
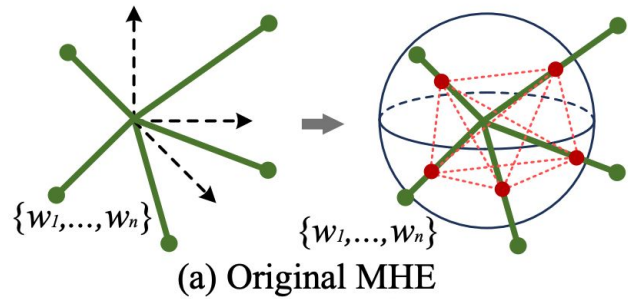
[1] shows that minimum hyperspherical energy leads to **better generalization**.

## Naively minimizing hyperspherical energy in [1]:

- Higher neuron dimension makes the optimization difficult.
- Highly non-linear and non-convex objective leads to many bad local minima.
- Deterministic gradients from naive MHE is sub-optimal to run away from bad local minima.

# Compressive Minimum Hyperspherical Energy (CoMHE)

## Overview of CoMHE:



CoMHE uses projections to reduce the neuron dimension and perform MHE in the projected space.

- Stochastic and dynamic regularization (CoMHE gradients also have stochasticity)
- Low neuron dimension benefits the optimization

## Random Projection CoMHE:

The projection matrices  $P$  are randomly initialized every certain number of iterations:

$$E_s^R(\hat{W}_N) := \frac{1}{C} \sum_{c=1}^C \sum_{i=1}^N \sum_{j=1, j \neq i}^N f_s \left( \left\| \frac{P_c \hat{w}_i}{\|P_c \hat{w}_i\|} - \frac{P_c \hat{w}_j}{\|P_c \hat{w}_j\|} \right\| \right)$$

## Angle-preserving Projection CoMHE:

The projections are learned to preserve angles:

$$E_s^A(\hat{W}_N, P^*) := \sum_{i=1}^N \sum_{j=1, j \neq i}^N f_s \left( \left\| \frac{P^* \hat{w}_i}{\|P^* \hat{w}_i\|} - \frac{P^* \hat{w}_j}{\|P^* \hat{w}_j\|} \right\| \right)$$

$$\text{s.t. } P^* = \arg \min_P \sum_{i \neq j} (\theta(\hat{w}_i, \hat{w}_j) - \theta(P\hat{w}_i, P\hat{w}_j))^2$$

## Adversarial Projection CoMHE:

The projections are learned adversarially:

$$\min_{\hat{W}_N} \max_P E_s^V(\hat{W}_N, P) := \sum_{i=1}^N \sum_{j=1, j \neq i}^N f_s \left( \left\| \frac{P \hat{w}_i}{\|P \hat{w}_i\|} - \frac{P \hat{w}_j}{\|P \hat{w}_j\|} \right\| \right)$$

## Theoretical Guarantees for RP-CoMHE:

$$\frac{\cos(\theta(w_1, w_2)) - \epsilon}{1 + \epsilon} < \cos(\theta(Pw_1, Pw_2)) < \frac{\cos(\theta(w_1, w_2)) + \epsilon}{1 - \epsilon}$$

It holds with probability  $(1 - 2 \exp(-\frac{k\epsilon^2}{8}))^2$

# Experiments and Results

## Convolutional neural networks (CNN):

Method	C-10	C-100
ResNet-110 [1]	6.61	25.16
ResNet-1001 [60]	4.92	<b>22.71</b>
Baseline	5.19	22.87
Orthogonal [29]	5.02	22.36
SRIP [9]	4.75	22.08
MHE [12]	4.72	22.19
HS-MHE [12]	4.66	22.04
RP-CoMHE	4.59	21.82
AP-CoMHE	<b>4.57</b>	<b>21.63</b>

CIFAR-10/100

Method	Res-18	Res-34	Res-50
baseline	32.95	30.04	25.30
Orthogonal [29]	32.65	29.74	25.19
Orthnormal [32]	32.61	29.75	25.21
SRIP [9]	32.53	29.55	24.91
MHE [12]	32.50	29.60	25.02
HS-MHE [12]	32.45	29.50	24.98
RP-CoMHE	31.90	29.38	<b>24.51</b>
AP-CoMHE	<b>31.80</b>	<b>29.32</b>	24.53

ImageNet-2012

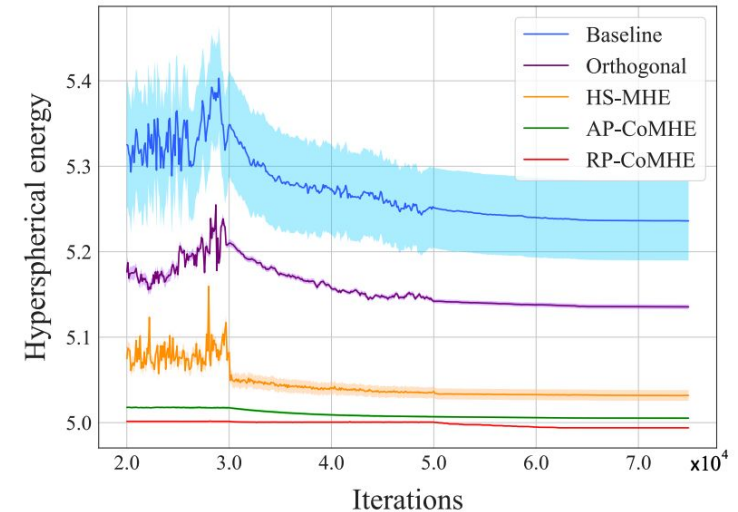
## Point cloud networks (PointNet):

Method	PN	PN (T)	PN++	
Original	87.1	89.20	90.07	ModelNet-40
MHE [12]	87.31	89.33	90.25	
HS-MHE [12]	87.44	89.41	90.31	
RP-CoMHE	87.82	89.69	90.52	
AP-CoMHE	<b>87.85</b>	<b>89.70</b>	<b>90.56</b>	

## Graph convolution networks (GCN):

Method	Citeseer	Cora	Pubmed
GCN Baseline	70.3	81.3	79.0
HS-MHE [12]	71.5	82.0	79.0
RP-CoMHE	<b>72.1</b>	<b>82.7</b>	<b>79.5</b>
AP-CoMHE	72.0	82.6	<b>79.5</b>

## Hyperspherical energy:




## Different network configurations:

Depth	CNN-6	CNN-9	CNN-15
Baseline	32.08	28.13	N/C
MHE [12]	28.16	26.75	26.90
HS-MHE [12]	27.56	25.96	25.84
RP-CoMHE	26.73	24.39	<b>24.21</b>
AP-CoMHE	<b>26.55</b>	<b>24.33</b>	24.55

CoMHE can effectively minimize hyperspherical energy and can improve different types of neural networks. (i.e., CoMHE is **architecture-agnostic.**)

# Thank you!

- ❖ For any question, please feel free to send emails to [rongmei.lin@emory.edu](mailto:rongmei.lin@emory.edu) or [wylu@gatech.edu](mailto:wylu@gatech.edu)
- ❖ Welcome to try CoMHE! The code is made available at <https://github.com/rmlin/CoMHE> 
- ❖ For our full paper and related material, please visit <https://wylu.com/> 