

# Decoupled Networks

Weiyang Liu, Zhen Liu, Zhiding Yu, Bo Dai, Rongmei Lin, Yisen Wang, James M.Rehg, Le Song

Georgia Institute of Technology, NVIDIA, Emory University, Tsinghua University, Ant Financial

**CVPR 2018**

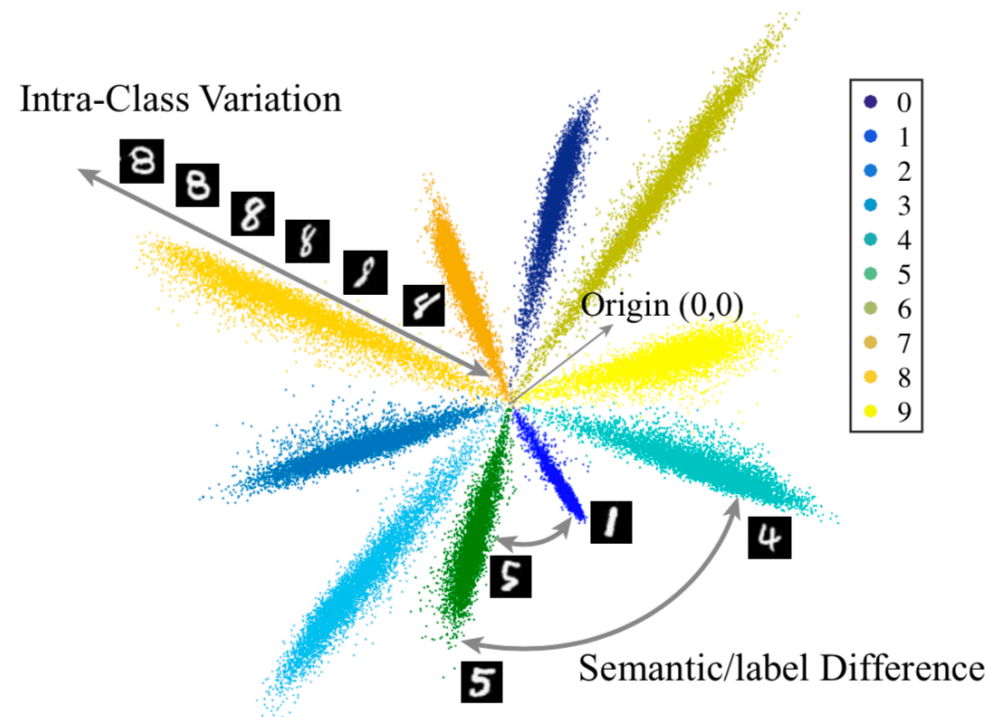
**Sungman, Cho.**

# Introduction

Dive into FR(Face Recognition)

# Abstract

- Propose a generic **decoupled learning framework** which **models the intra-class variation and semantic difference independently**.
- Reparametrize the inner product to a decoupled form
- Decoupled reparameterization renders significant performance gain with easier convergence and stronger robustness.



# Introduction

- $\langle w, x \rangle = w^T x$  couples the semantic difference and the intra-class variation in one unified measure.
- When the inner product between two samples is large, **one can not tell whether the two samples have large semantic difference or have large intra-class variation.**
- Propose a novel decoupled network (DCNet) by generalizing traditional Conv. operation.

$$f(\boldsymbol{w}, \boldsymbol{x}) = \langle \boldsymbol{w}, \boldsymbol{x} \rangle = \boldsymbol{w}^\top \boldsymbol{x},$$

$$f(\boldsymbol{w}, \boldsymbol{x}) = \|\boldsymbol{w}\| \|\boldsymbol{x}\| \cos(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

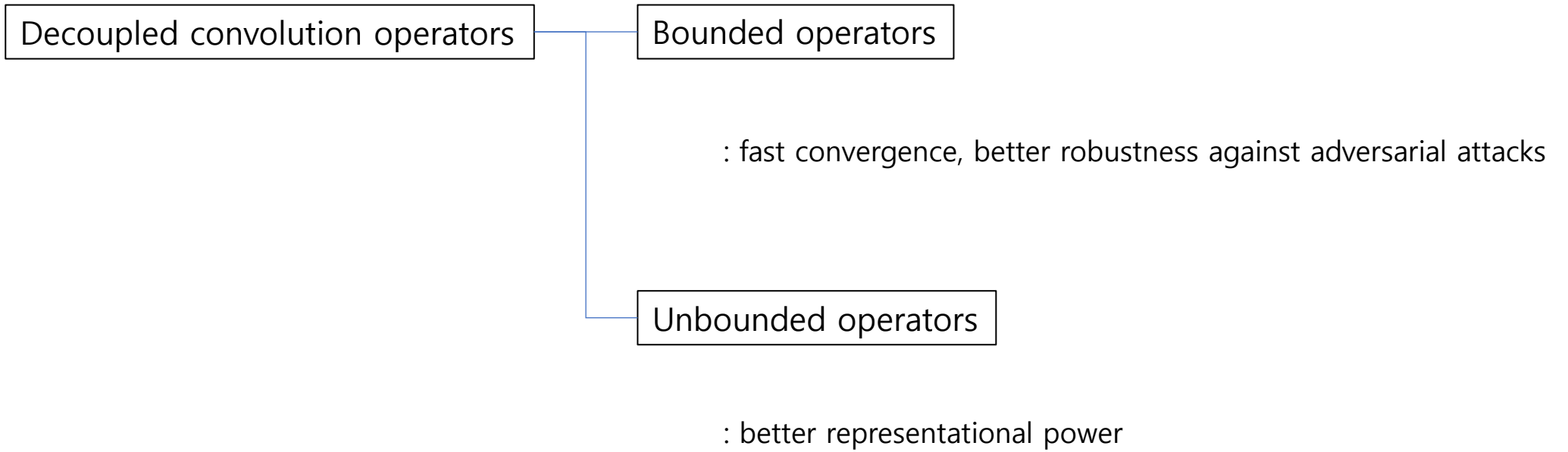
# Introduction

magnitude function      angular function

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

- magnitude function : models the intra-class variation
- angular function : models the semantic difference

# Introduction



# Contributions

- Propose an **explicitly decoupled framework** to study neural networks.
- **By decoupling** the inner product, we are able to **design more effective magnitude and angular functions** rather than the original convolution.
- DCNets have **easier convergence, better accuracy and stronger robustness**.

# Related Works

- Large-margin softmax loss for convolutional neural networks, ICML, 2016
- Additive margin softmax for face verification, arXiv:1801.05599, 2018
- Improving face verification and person re-identification accuracy using hyperplane similarity, ICCV, 2017
- Deep hypersphere embedding for face recognition, CVPR, 2017

- 1. normalizes the weights of the last fully connected layer**
- 2. specific area. (face recognition)**



# Decoupled Networks

# Decoupled Networks

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

The design of angular function  $g(\cdot)$  is relatively easy but restricted, because it only takes the angle as input.

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

$h(\cdot)$  takes the norm of  $w$  and the norm of  $x$  as two inputs, therefore it is more complicated to design.

# Decoupled Networks

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

$\|\boldsymbol{w}\|$  is the **intrinsic property of a kernel itself**, corresponding to the importance of the kernel rather than the intra-class variation of the inputs.

Therefore, we tend not to include  $\|\boldsymbol{w}\|$  into the magnitude function.

# Decoupled Networks

All kernels are assigned with equal importance.

$$f_d(\mathbf{w}, \mathbf{x}) = h(\|\mathbf{w}\|, \|\mathbf{x}\|) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}),$$

Encourages the network to make decision based on as many kernels as possible and therefore may make the **network generalize better**.

However, *weighted decoupled operators* improve the **representational power** and may be useful when dealing with a large-scale dataset with numerous categories.

*weighted decoupled operators :*  
combining  $\|\mathbf{w}\|$  back to  $h(\cdot)$

# Bounded Decoupled Operators

Bounded decoupled operators :  $|f_d(\boldsymbol{w}, \boldsymbol{x})| \leq c$

## Hyperspherical Convolution

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$



$$h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) = \alpha$$

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = \alpha \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

# Bounded Decoupled Operators

Bounded decoupled operators :  $|f_d(\boldsymbol{w}, \boldsymbol{x})| \leq c$

## Hyperspherical Convolution

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$



$$h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) = \alpha$$

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = \boxed{\alpha} \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

controls the output scale

# Bounded Decoupled Operators

Bounded decoupled operators :  $|f_d(\mathbf{w}, \mathbf{x})| \leq c$

## Hyperspherical Convolution

$$f_d(\mathbf{w}, \mathbf{x}) = h(\|\mathbf{w}\|, \|\mathbf{x}\|) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}),$$



$$h(\|\mathbf{w}\|, \|\mathbf{x}\|) = \alpha$$

$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}),$$

geodesic distance on the unit hypersphere and typically outputs value from -1 to 1

# Bounded Decoupled Operators

Bounded decoupled operators :  $|f_d(\boldsymbol{w}, \boldsymbol{x})| \leq c$

## Hyperspherical Convolution

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$



$$h(\|\boldsymbol{w}\|, \|\boldsymbol{x}\|) = \alpha$$

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = \alpha \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

output :  $[-\alpha, \alpha]$

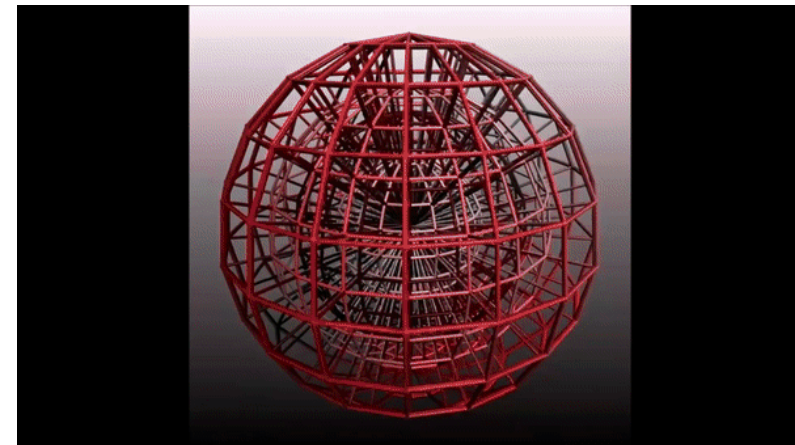


# Bounded Decoupled Operators

Bounded decoupled operators :  $|f_d(\mathbf{w}, \mathbf{x})| \leq c$

## Hyperspherical Convolution

$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}),$$



HyperSphere

Usually, we can use  $\alpha = 1$ , which reduces to SphereConv in this case.

Projecting  $w$  and  $x$  to a hypersphere and then performing inner product. (if  $g(\theta) = \cos(\theta)$ )

**Improves the problem conditioning in neural networks, making the network converge better.**

# Bounded Decoupled Operators

Bounded decoupled operators :  $|f_d(\boldsymbol{w}, \boldsymbol{x})| \leq c$

## Hyperball Convolution

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = \alpha \cdot \frac{\min(\|\boldsymbol{x}\|, \rho)}{\rho} \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

$\rho$  controls the saturation threshold for the input norm  $\|\boldsymbol{x}\|$

Projecting  $w$  to a hypersphere,  $x$  to a hyperball, then performing the inner product. (if  $g(\theta) = \cos(\theta)$ )

**More robust and flexible than SphereConv**

# Bounded Decoupled Operators

Is BallConv more robust and flexible than SphereConv ?

$$\text{SphereConv} : f_d(\boldsymbol{w}, \boldsymbol{x}) = \alpha \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

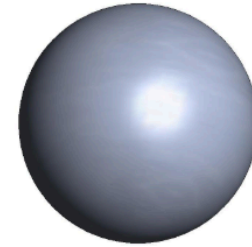
$x$  with small  $\|x\|$  and the same direction as  $w$  could still produce maximum output

# Appendix. Ball vs Sphere

## Ball:

A ball is the shape of a circle. It is used in all the common ball games these days such as the sports of football, basketball, cricket and squash. Hockey is also a game in which ball is used and it is totally round in appearance.

Image Courtesy: istockphoto.com



## Sphere:

A sphere is a shape almost the same as a ball, but it is a bit flatter than the ball and is most commonly used in the game of rugby these days.

Image Courtesy: istockphoto.com



# Bounded Decoupled Operators

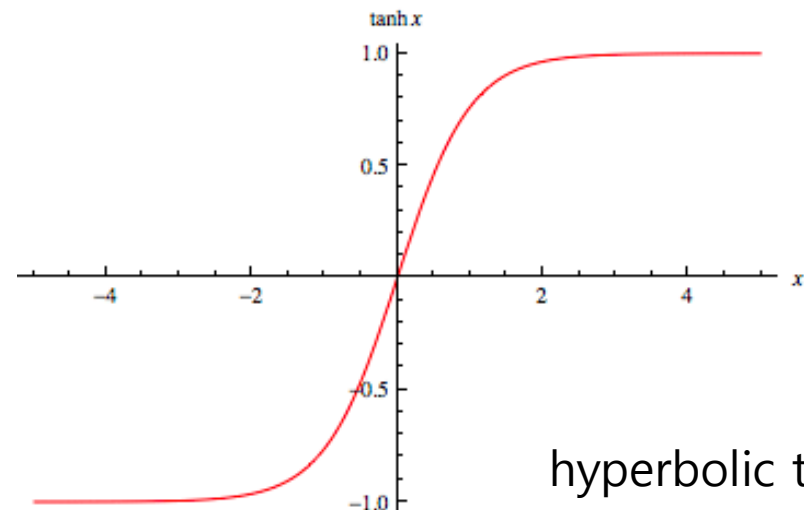
Bounded decoupled operators :  $|f_d(\mathbf{w}, \mathbf{x})| \leq c$

## Hyperbolic Tangent Convolution

$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \tanh\left(\frac{\|\mathbf{x}\|}{\rho}\right) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}),$$

$\rho$  is parameter controlling the decay curve.

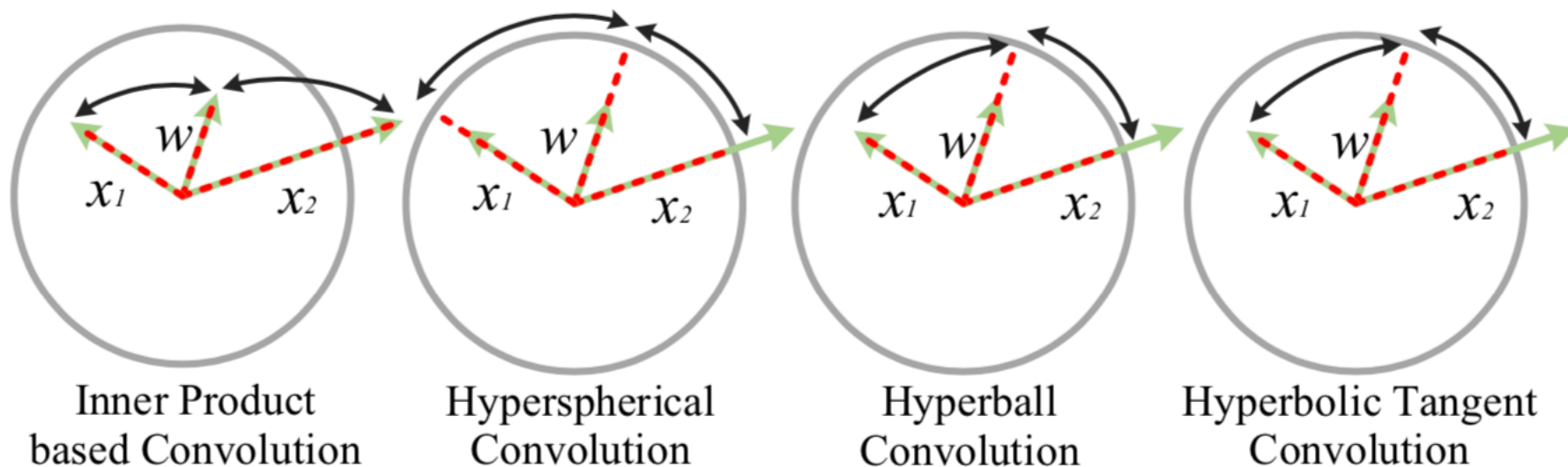
The TanhConv can be viewed as a smooth version of BallConv.  
(more convergence gain due to its smoothness)



hyperbolic tangent

**More convergence gain**

# Bounded Decoupled Operators



# Unbounded Decoupled Operators

## Linear Convolution

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = \alpha \|\boldsymbol{x}\| \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

$\alpha$  controls the output scale.

LinearConv projects the  $w$  to a hypersphere and has a parameter to control the slope

# Unbounded Decoupled Operators

**Segmented Convolution.**

$$f_d(\mathbf{w}, \mathbf{x}) = \begin{cases} \alpha \|\mathbf{x}\| \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}), & 0 \leq \|\mathbf{x}\| \leq \rho \\ (\beta \|\mathbf{x}\| + \alpha\rho - \beta\rho) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}), & \rho < \|\mathbf{x}\| \end{cases} ,$$

$\alpha$  controls the slope when  $\|\mathbf{x}\| \leq \rho$  ,  $\beta$  controls the slope when  $\|\mathbf{x}\| \geq \rho$

Both LinearConv and BallConv are special cases of SegConv.



# Unbounded Decoupled Operators

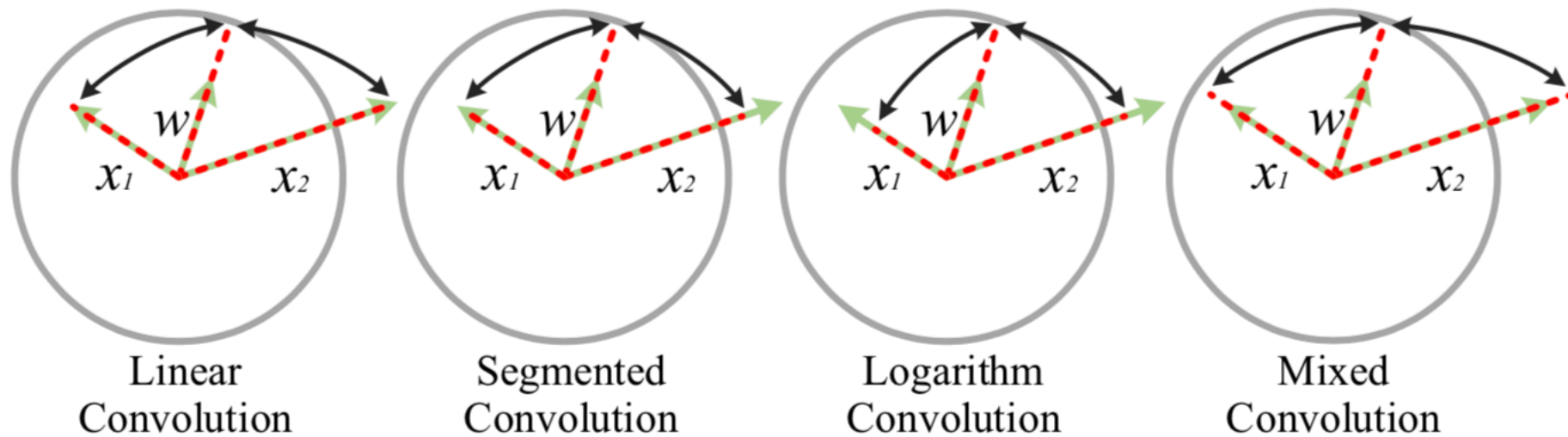
**Logarithm Convolution.**

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = \alpha \log(1 + \|\boldsymbol{x}\|) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

**Mixed Convolution.**

$$f_d(\boldsymbol{w}, \boldsymbol{x}) = (\alpha \|\boldsymbol{x}\| + \beta \log(1 + \|\boldsymbol{x}\|)) \cdot g(\theta_{(\boldsymbol{w}, \boldsymbol{x})}),$$

# Unbounded Decoupled Operators



# 3 Properties of Decoupled Operators

## Operator Radius.

describe the gradient change point of the magnitude function.

**BallConv** 
$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \cdot \frac{\min(\|\mathbf{x}\|, \rho)}{\rho} \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}),$$

**SegConv** 
$$f_d(\mathbf{w}, \mathbf{x}) = \begin{cases} \alpha \|\mathbf{x}\| \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}), & 0 \leq \|\mathbf{x}\| \leq \rho \\ (\beta \|\mathbf{x}\| + \alpha\rho - \beta\rho) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}), & \rho < \|\mathbf{x}\| \end{cases},$$

# 3 Properties of Decoupled Operators

## Boundedness

affect network's convergence speed and robustness.

- Lead to **better problem conditioning** in training a deep network via stochastic gradient descent.
- Make the **variance of the output small** and **partially address the internal covariate shift problem**.
- **Constrain the Lipschitz constant of a neural network**, making the entire network more smooth.

(Lipschitz constant of a neural network is shown to be closely related to its robustness against adversarial perturbation)

# Design of the Angular Function

Linear angular activation

$$g(\theta_{(\mathbf{w}, \mathbf{x})}) = -\frac{2}{\pi} \theta_{(\mathbf{w}, \mathbf{x})} + 1,$$

Cosine angular activation

$$g(\theta_{(\mathbf{w}, \mathbf{x})}) = \cos(\theta_{(\mathbf{w}, \mathbf{x})}),$$

Sigmoid angular activation

$$g(\theta_{(\mathbf{w}, \mathbf{x})}) = \frac{1 + \exp(-\frac{\pi}{2k})}{1 - \exp(-\frac{\pi}{2k})} \cdot \frac{1 - \exp(\frac{\theta_{(\mathbf{w}, \mathbf{x})}}{k} - \frac{\pi}{2k})}{1 + \exp(\frac{\theta_{(\mathbf{w}, \mathbf{x})}}{k} - \frac{\pi}{2k})},$$

Square cosine angular activation

$$g(\theta_{(\mathbf{w}, \mathbf{x})}) = \text{sign}(\cos(\theta)) \cdot \cos^2(\theta),$$

# Design of the Angular Function

**Linear angular activation**  $g(\theta_{(w,x)}) = -\frac{2}{\pi}\theta_{(w,x)} + 1,$

**Sigmoid angular activation**  $g(\theta_{(w,x)}) = \frac{1 + \exp(-\frac{\pi}{2k})}{1 - \exp(-\frac{\pi}{2k})} \cdot \frac{1 - \exp(\frac{\theta_{(w,x)}}{k} - \frac{\pi}{2k})}{1 + \exp(\frac{\theta_{(w,x)}}{k} - \frac{\pi}{2k})},$

**Cosine angular activation**  $g(\theta_{(w,x)}) = \cos(\theta_{(w,x)}),$

**Square cosine angular activation**  $g(\theta_{(w,x)}) = \text{sign}(\cos(\theta)) \cdot \cos^2(\theta),$

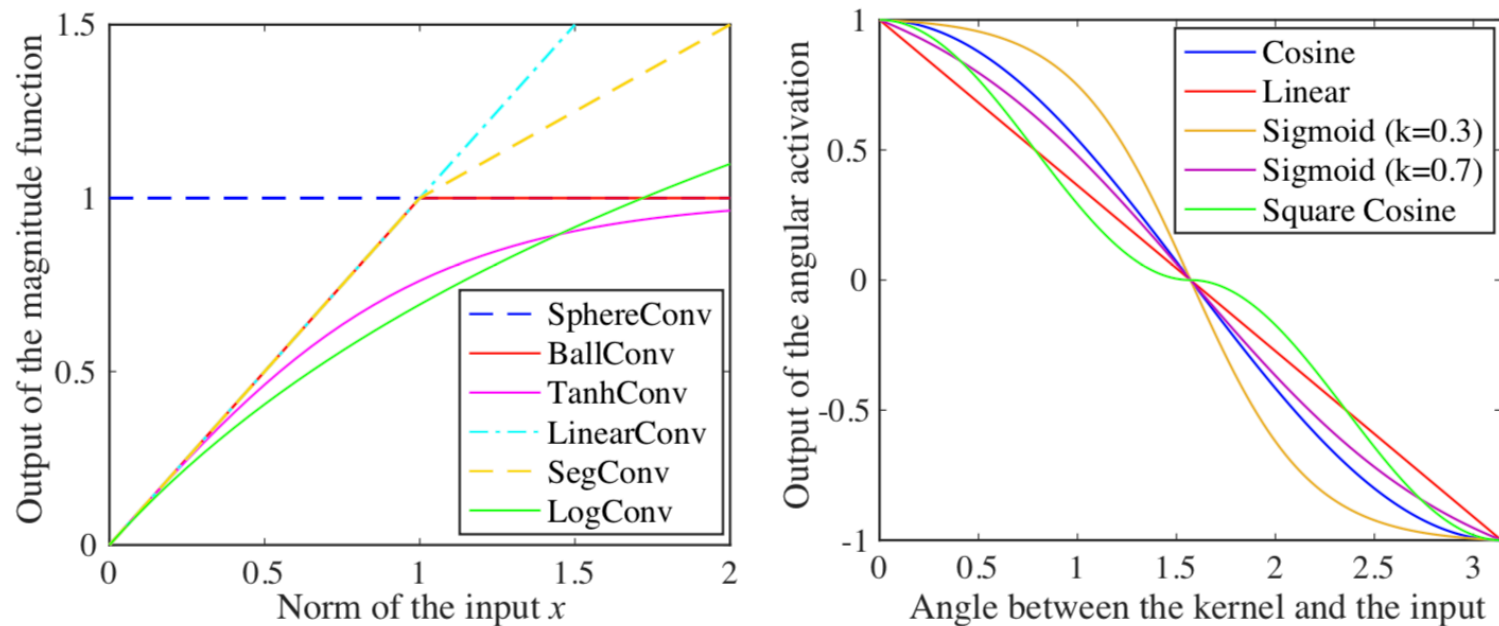


Figure 3: Magnitude function ( $\rho = 1$ ) and angular activation function.

# Weighted Decoupled Operators

Although empirically we find that the **standard decoupled operators work better** than the weighted ones in most cases, we still consider **weighted decoupled operators**, in order to improve the **operator's flexibility**

**Linearly Weighted Decoupled Operator.**

$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}), \longrightarrow f_d(\mathbf{w}, \mathbf{x}) = \alpha \|\mathbf{w}\| \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}).$$

**Nonlinearly Weighted Decoupled Operator.**

$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \tanh\left(\frac{\|\mathbf{x}\|}{\rho}\right) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}), \longrightarrow f_d(\mathbf{w}, \mathbf{x}) = \alpha \tanh\left(\frac{1}{\rho} \|\mathbf{x}\| \cdot \|\mathbf{w}\|\right) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}).$$
$$f_d(\mathbf{w}, \mathbf{x}) = \alpha \tanh\left(\frac{1}{\rho} \|\mathbf{w}\|\right) \cdot \tanh\left(\frac{1}{\rho} \|\mathbf{x}\|\right) \cdot g(\theta_{(\mathbf{w}, \mathbf{x})}).$$

# Learnable Decoupled Operators

By making hyperparameters learnable,

Representational power  $\leftarrow$  Trade Off  $\rightarrow$  Prohibitive amount of training data

In order to achieve an effective trade-off,

use only one learning parameter :  $\rho$  (operator radius)



# Improving the Optimization

# Weight Projection

Large  $\|w\|$  can make the gradients very small.

$$\frac{\partial}{\partial \mathbf{w}} (\hat{\mathbf{w}}^\top \hat{\mathbf{x}}) = \frac{\hat{\mathbf{x}} - \hat{\mathbf{w}}^\top \hat{\mathbf{x}} \cdot \hat{\mathbf{w}}}{\|\mathbf{w}\|}, \quad \text{where } \hat{\mathbf{w}} = \mathbf{w} / \|\mathbf{w}\| \text{ and } \hat{\mathbf{x}} = \mathbf{x} / \|\mathbf{x}\|.$$

Weight projection performs  $\mathbf{w} \leftarrow s \cdot \hat{\mathbf{w}}$

$s$  is positive constant which controls the norm of the gradient (use  $s = 1$  in experiments)

# Weight Gradients

We could simply multiply  $\|w\|$ .

$$\Delta \boldsymbol{w} = \|\boldsymbol{w}\| \cdot \frac{\partial}{\partial \boldsymbol{w}} (\hat{\boldsymbol{w}}^\top \hat{\boldsymbol{x}}) = \hat{\boldsymbol{x}} - \hat{\boldsymbol{w}}^\top \hat{\boldsymbol{x}} \cdot \hat{\boldsymbol{w}},$$

# Experiments

# CIFAR-10 & CIFAR-100

## Weighted Decoupled Operators & Optimization Tricks

Method	Error
Linearly Weighted Decoupled Operator	22.95
Nonlinearly Weighted Decoupled Operator (Eq. (15))	23.03
Nonlinearly Weighted Decoupled Operator (Eq. (16))	23.38
Decoupled Operator (Standard Gradients)	23.09
Decoupled Operator (Weight Projection)	<b>21.17</b>
Decoupled Operator (Weighted Gradients)	21.45

Weighted decoupled operators do not show obvious advantages.

# CIFAR-10 & CIFAR-100

## Weighted Decoupled Operators & Optimization Tricks

Method	Error
Linearly Weighted Decoupled Operator	22.95
Nonlinearly Weighted Decoupled Operator (Eq. (15))	23.03
Nonlinearly Weighted Decoupled Operator (Eq. (16))	23.38
Decoupled Operator (Standard Gradients)	23.09
Decoupled Operator (Weight Projection)	<b>21.17</b>
Decoupled Operator (Weighted Gradients)	21.45

Weight projection and weighted gradients work much better than the competing methods.

# CIFAR-10 & CIFAR-100

Learning without Batch Normalization

Method	Linear	Cosine	Sq. Cosine
CNN Baseline	-	35.30	-
LinearConv	33.39	31.76	N/C
TanhConv	<b>32.88</b>	31.88	<b>34.26</b>
SegConv	34.69	<b>30.34</b>	N/C

Ours can perform much better than the original convolution even without BN

# CIFAR-10 & CIFAR-100

Learning without ReLU.

Method	Cosine w/o ReLU	Sq. Cosine w/o ReLU	Cosine w/ ReLU	Sq. Cosine w/ ReLU
Baseline	58.24	-	26.01	-
SphereConv	33.31	25.90	26.00	26.97
BallConv	<b>31.81</b>	25.43	25.18	26.48
TanhConv	32.27	25.27	25.15	26.94
LinearConv	36.49	24.36	<b>24.81</b>	25.14
SegConv	33.57	<b>24.29</b>	24.96	<b>25.04</b>
LogConv	33.62	24.91	25.17	25.85
MixConv	33.46	24.93	25.27	25.77

Table 3: Testing error rate (%) of plain CNN-9 on CIFAR-100. Note that, BN is used in all compared models. Baseline is the original plain CNN-9.

**Ours can perform much better than the original convolution even without BN**



# CIFAR-10 & CIFAR-100

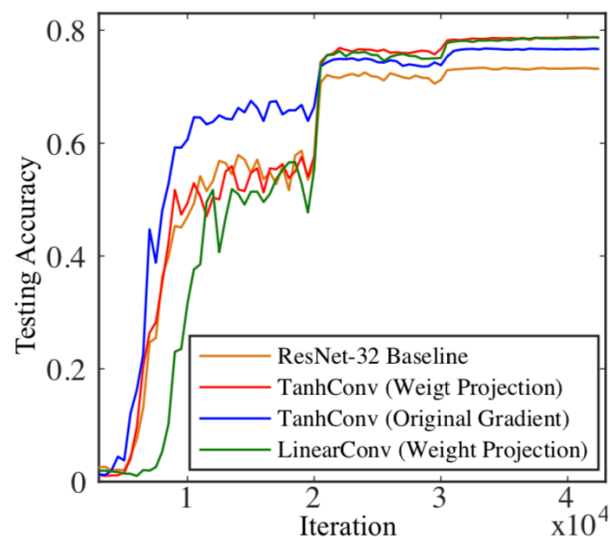
Comparison among Different Decoupled Operators.

Method	Linear	Cosine	Sq. Cosine
ResNet Baseline	-	26.69	-
SphereConv	21.79	21.44	24.40
BallConv	21.44	21.12	24.31
TanhConv	21.6	21.17	24.77
LinearConv	21.09	22.17	21.31
SegConv	<b>20.86</b>	<b>20.91</b>	<b>20.88</b>
LogConv	21.84	21.08	22.86
MixConv	21.02	21.28	21.81

Table 4: Testing error rate (%) of ResNet-32 on CIFAR-100.

# CIFAR-10 & CIFAR-100

Convergence



Acc.

Method	CIFAR-10	CIFAR-100
ResNet-110-original [5]	6.61	25.16
ResNet-1001 [6]	4.92	<b>22.71</b>
ResNet-1001 (64 mini-batch size) [6]	<b>4.64</b>	-
DCNet-32 (TanhConv + Cosine)	<b>4.75</b>	21.12
DCNet-32 (LinearConv + Sq. Cos.)	5.34	<b>20.23</b>

Table 5: Comparison to the state-of-the-art on CIFAR-10 and CIFAR-100.

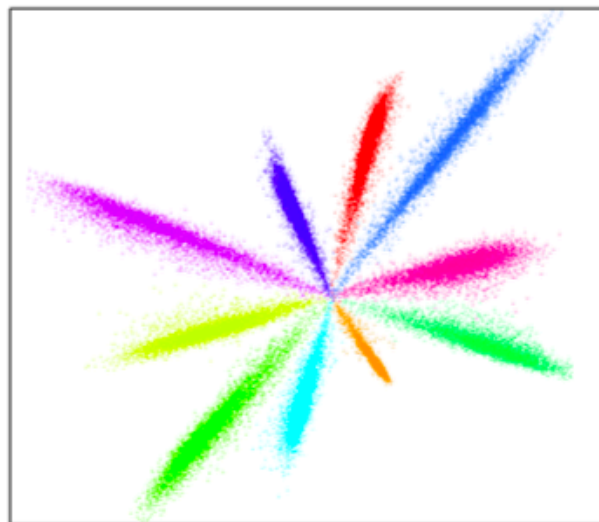
# Adversarial attacks

Attack	Target models			
	Baseline	SphereConv	BallConv	TanhConv
None	85.35	88.58	91.13	<b>91.45</b>
FGSM	50.90	<b>56.71</b>	49.50	50.61
BIM	36.22	<b>43.10</b>	27.48	29.06
None	83.70	87.41	87.47	<b>87.54</b>
FGSM	77.57	76.29	78.67	<b>80.38</b>
BIM	78.55	77.79	80.59	<b>82.47</b>

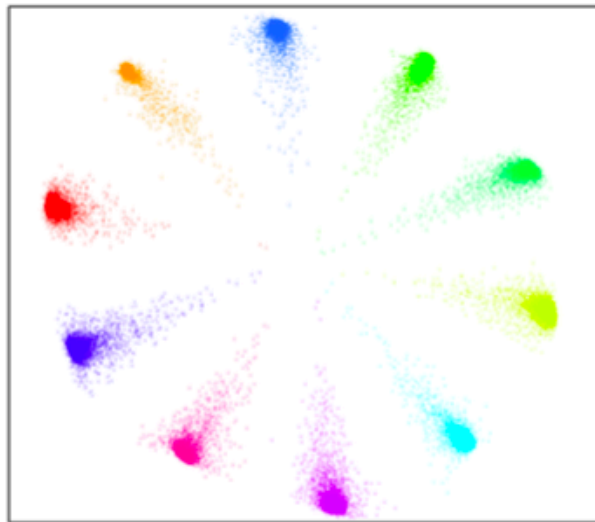
Table 8: Black-box attacks on CIFAR-10. Performance is measured in accuracy (%). The first three rows are results of naturally trained models, and the last three rows are results of adversarially trained models.

Attack	Target models			
	Baseline	SphereConv	BallConv	TanhConv
None	85.35	88.58	91.13	<b>91.45</b>
FGSM	18.82	43.64	50.47	<b>52.60</b>
BIM	8.67	8.89	7.74	<b>10.18</b>
None	83.70	87.41	87.47	<b>87.54</b>
FGSM	78.96	<b>85.98</b>	82.20	81.46
BIM	7.96	<b>35.07</b>	17.38	19.86

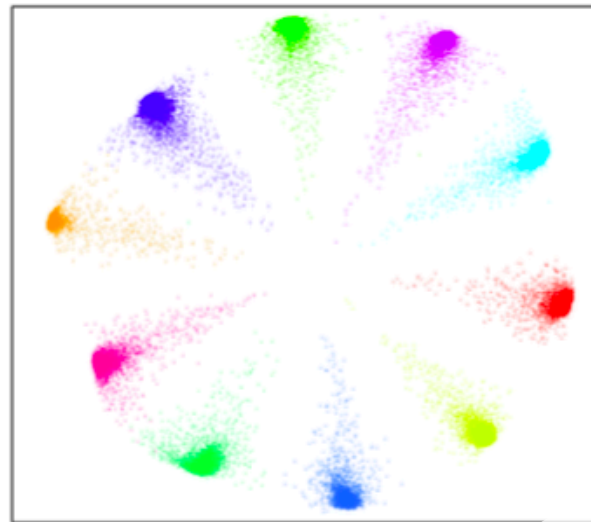
Table 7: White-box attacks on CIFAR-10. Performance is measured in accuracy (%). The first three rows are results of naturally trained models, and the last three rows are results of adversarially trained models.



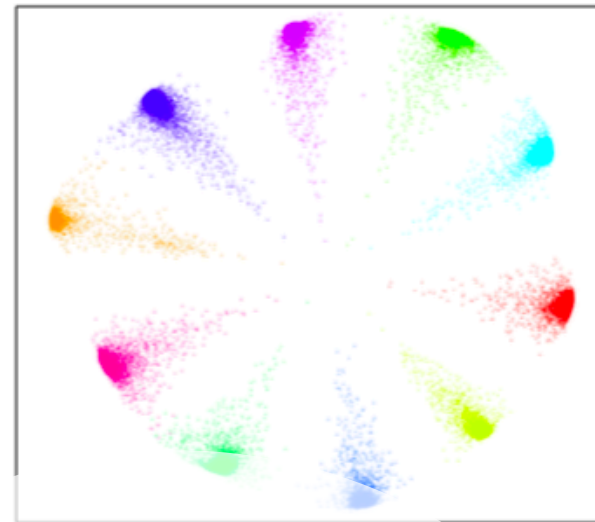
Baseline CNN (Original Convolution)



DCNet (SphereConv)

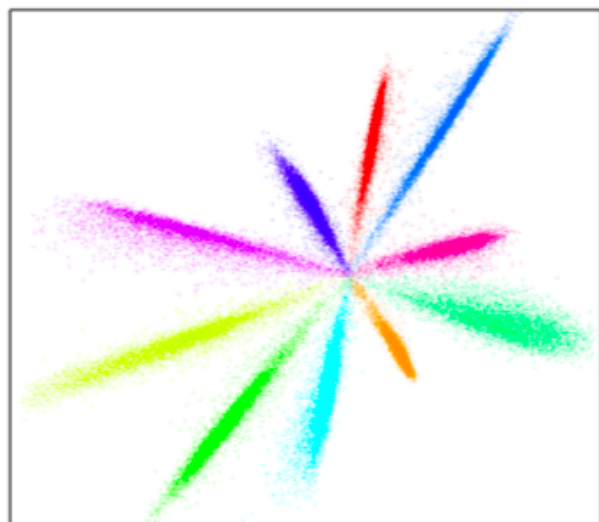


DCNet (BallConv)

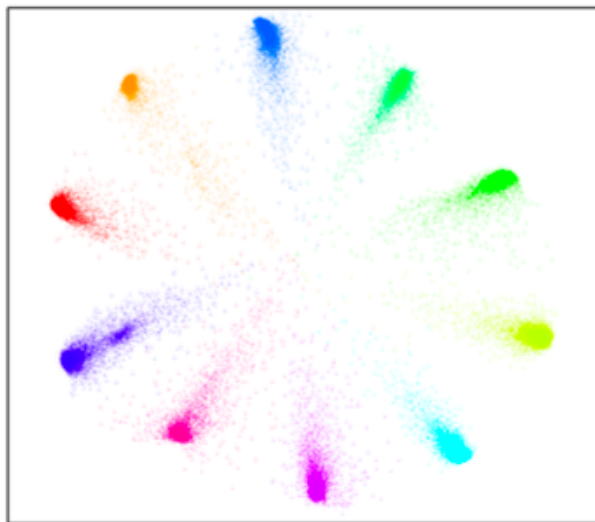


DCNet (TanhConv)

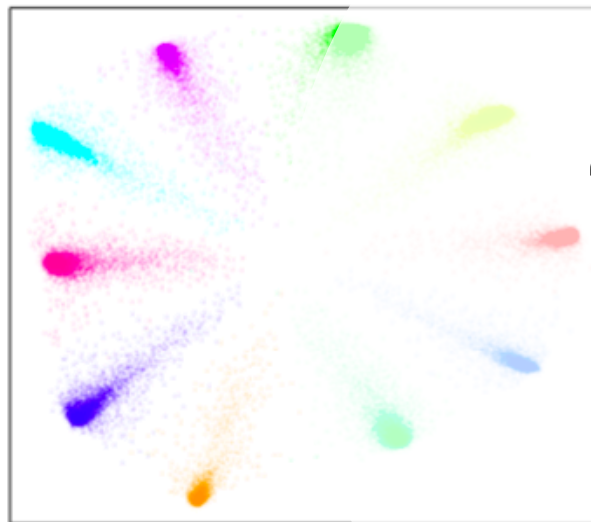
Figure 6: 2D feature visualization on MNIST dataset with natural training.



Baseline CNN (Original Convolution)



DCNet (SphereConv)



DCNet (BallConv)



DCNet (TanhConv)

**2D Plots  
(MNIST)**

Figure 7: 2D feature visualization on MNIST dataset with adversarial training.

**Thanks.**