# Global Pooling, More than Meets the Eye: Position Information is Encoded Channel-Wise in CNNs

Md Amirul Islam[1,6]    Matthew Kowal[2,6]    Sen Jia[4]    Kosta Derpanis[2,5,6]    Neil Bruce[3,6]

# *Motivation*

**Islam et. al (ICLR 2020)**



**Kayhan et. al (CVPR 2020)**



**Islam et. al (arXiv 2021)**

## CNNs encode absolute position information

# *Motivation*

How does a CNN contain positional information in the representations after a Global Average Pooling (GAP) layer?

# *Hypothesis*

**CNNs encode absolute position information along the ordering of the channel dimension**

# Learning Position with a GAPNet

# *Learning Position with a GAPNet*

# *Learning Position with a GAPNet*

# *Learning Position with ShuffleNet*

# *Evaluation of Channel-wise Position Encoding*

| Network | | Loc. Classification | | Image Classification | |
|---------|--------|------|------|------|------|
| | | 3x3 | 7x7 | 3x3 | 7x7 |
| **Res18** | GAPNet | 100 | 100 | 82.6 | 82.1 |
| | *PermuteNet* | 78.8 | 21.4 | 82.1 | 69.9 |

*Results are on CIFAR-10 dataset*  9

# *Applicability of Channel-wise Positional Encoding*

**Learning Translation Invariant Representation**

**Attacking Position Encoding Channels**

# *Learning Translation Invariant Representations*

# Results: Translation Invariance

| Network | CIFAR-10 | | CIFAR-100 | |
|---|---|---|---|---|
| | Top-1 Acc. | Consistency | Top-1 Acc. | Consistency |
| ResNet-18 | 93.1 | 90.8 | 72.6 | 70.1 |
| Blurpool | 92.5 | 92.5 | 72.4 | 78.2 |
| AugShift (Ours) | 92.1 | 94.8 | 72.6 | 85.6 |

# *Results: Translation Invariance*

| Network | CIFAR-10 | | CIFAR-100 | |
|---|---|---|---|---|
| | *Top-1 Acc.* | *Consistency* | *Top-1 Acc.* | *Consistency* |
| *ResNet-18* | 93.1 | 90.8 | 72.6 | 70.1 |
| *Blurpool* | 92.5 | 92.5 | 72.4 | 78.2 |
| *AugShift (Ours)* | 92.1 | **94.8** | 72.6 | **85.6** |

# *Attacking the Position Encoding Channels*

1) **Identify the position-specific channels**

2) **Target the position-specific channels**

# *Identifying the Overall Position Encoding Channels*



Image

Flipped Image

$$\hat{z} = \mathrm{argsort}_{j \in C} \left[ \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} |\Delta z_i| \right]$$

15

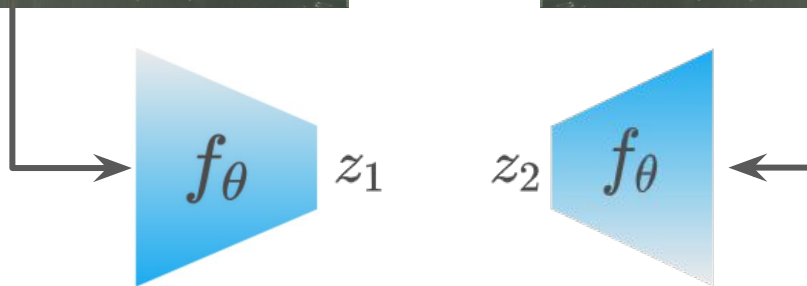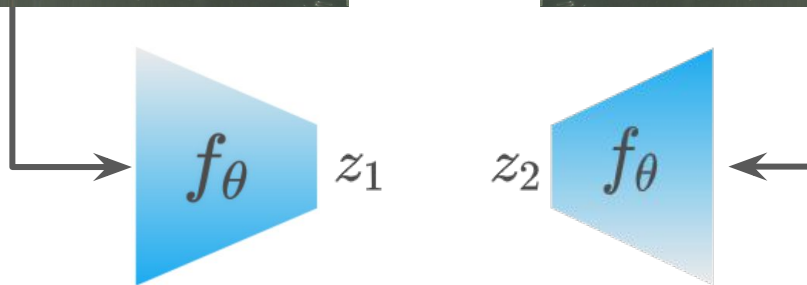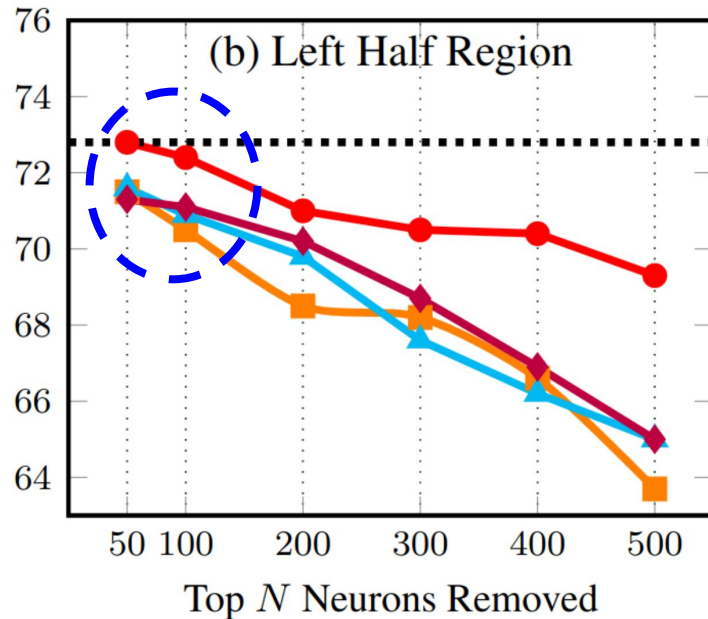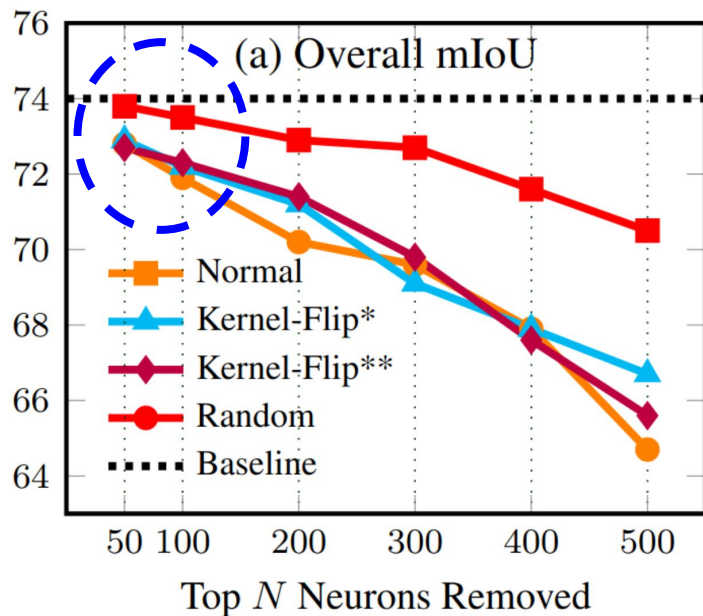# *Region-Specific Position Encoding Channels*
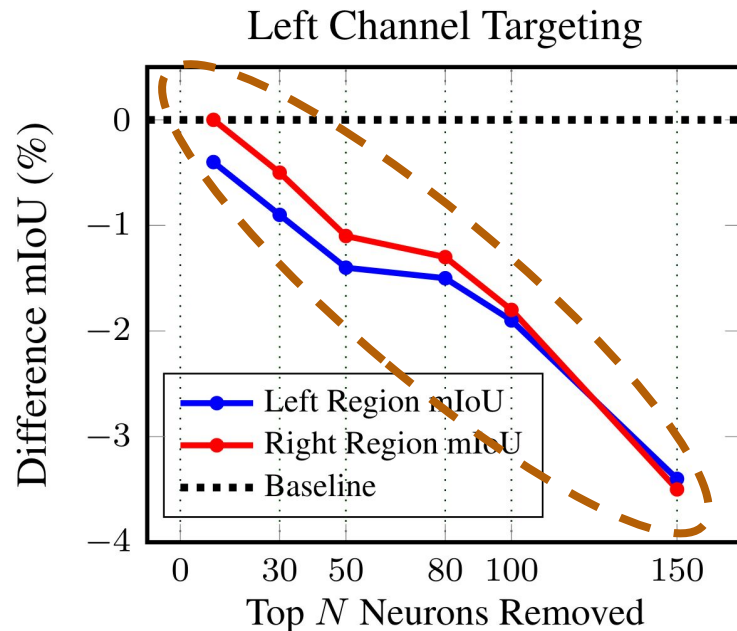


Image

Flipped Image
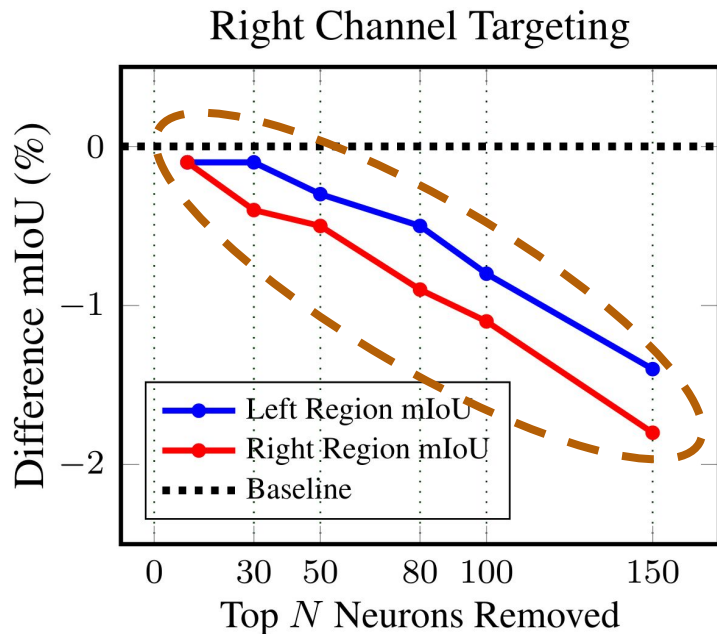
$$\hat{z}^l = \texttt{argsort}_{j \in C}\left[\frac{1}{|\mathcal{D}|}\sum_{i=1}^{|\mathcal{D}|}\Delta z_i\right]$$

16

# *Targeting Position-specific Channels*

*Evaluated on Cityscapes pre-trained DeepLabv3-ResNet50 model*

# *Targeting Region-specific Channels*



*Evaluated on Cityscapes pre-trained DeepLabv3-ResNet50 model*

# *Take Away*

- **Position information is encoded based on the *ordering* of the channels while semantic information is largely not.**

- **Introduced a simple data augmentation strategy to improve translation invariance of CNNs.**

- **Introduced an intuitive technique to identify the position-specific neurons in a network's latent representation.**

# *Thanks for Listening*

Code is available at: *https://github.com/islamamirul/PermuteNet*