

Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles

Bhaskar Joshi (2019111002)





Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles

- **Outline**
 - Feature Learning by Solving Jigsaw Puzzles
 - Context Free Network (CFN): Network Architecture
 - Experimentation

Brief Introduction of Paper

- Paper Objective:

The objective of this is to study the problem of image representation learning without human annotation. The paper aims to build a CNN that could be trained for solving Jigsaw puzzle without the help of manual labelling. It also solves object classification and detection task.

The paper talks about new self supervised task i.e. the jigsaw reassembly problem that is able to perform well when transferred to detection and classification tasks.

Brief Introduction of Paper

- The Context Free Network (CFN), is designed to take picture tiles as input and output the correct spatial arrangement. This way, useful feature representation is learned and applied to a variety of transfer learning benchmarks (downstream tasks).
- Conceptual Idea:

Two different-colored cars and two dogs with various fur patterns. Because they are invariant to shared patterns, the characteristics learned to solve puzzles in one (car/dog) image will also apply to the second (car/dog) image.



Fig. 2: Most of the shape of these 2 pairs of images is the same (two separate instances within the same categories). However, some low-level statistics are different (color and texture). The Jigsaw puzzle solver learns to ignore such statistics when they do not help the localization of parts.

Brief Introduction of Paper

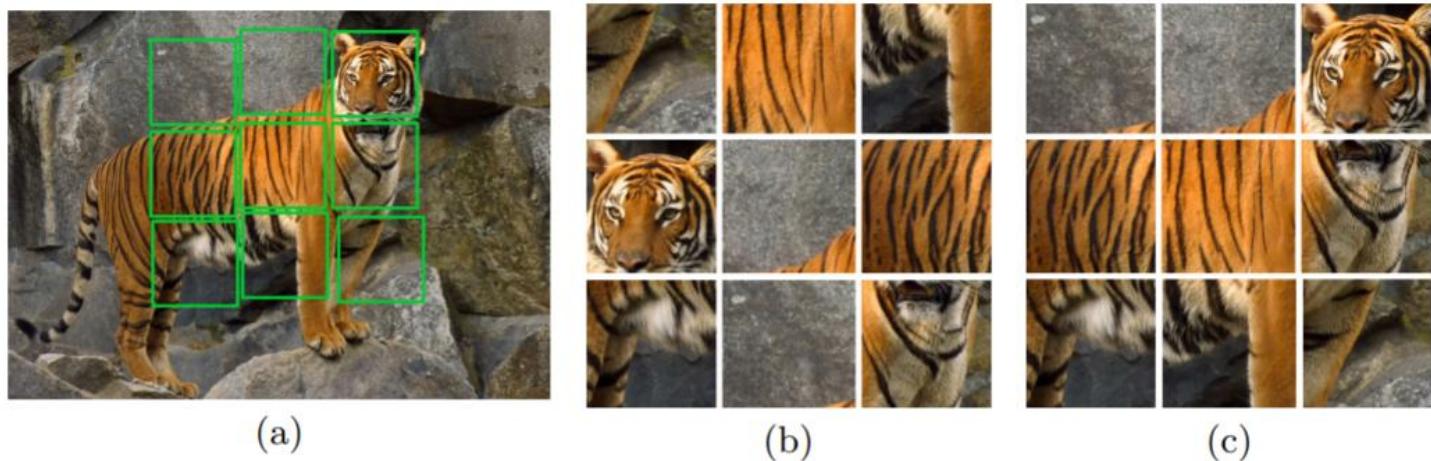


Fig. 1: Learning image representations by solving Jigsaw puzzles. (a) The image from which the tiles (marked with green lines) are extracted. (b) A puzzle obtained by shuffling the tiles. Some tiles might be directly identifiable as object parts, but others are ambiguous (*e.g.*, have similar patterns) and their identification is much more reliable when all tiles are jointly evaluated. In contrast, with reference to (c), determining the relative position between the central tile and the top two tiles from the left can be very challenging [10].

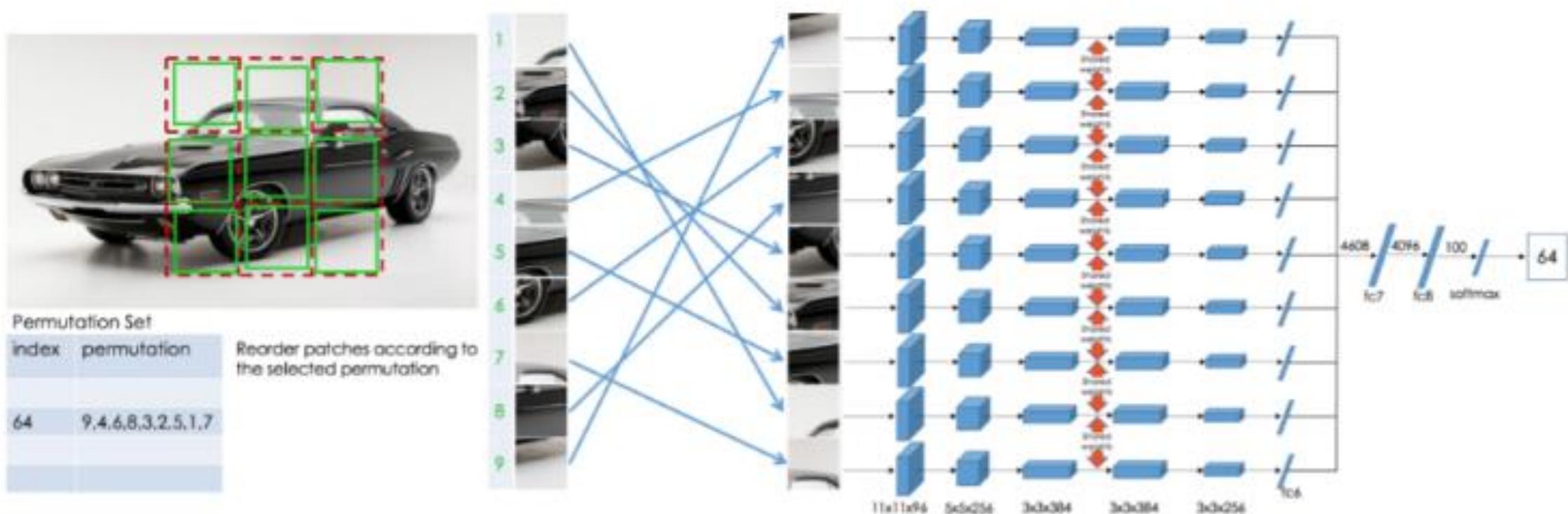
Brief Introduction of Paper **

- Brief description of how authors arrived at convolutional architecture that is capable to solve jigsaw puzzle is to stack the tiles of the puzzle along the channels (i.e., the input data would have $9 \times 3 = 27$ channels) and input these channels into a CNN to solve the Jigsaw puzzles. (**Naïve Stacked Patches NOT Working**)
- The problem with this design is that the network prefers to identify correlations between low-level texture statistics across tiles rather than between the high-level primitives. A CNN with only low-level features learnt is NOT what we want.

*

*(Can be skipped)

Paper Implementation Description





Paper Implementation Description

- The given architecture is called context-free network (CFN) because the data flow of each patch is explicitly separated until the fully connected layer and context is handled only in the last fully connected layers.
- GitHub Code Link: <https://github.com/BhaskarJoshi-01/JigSawPuzzle>



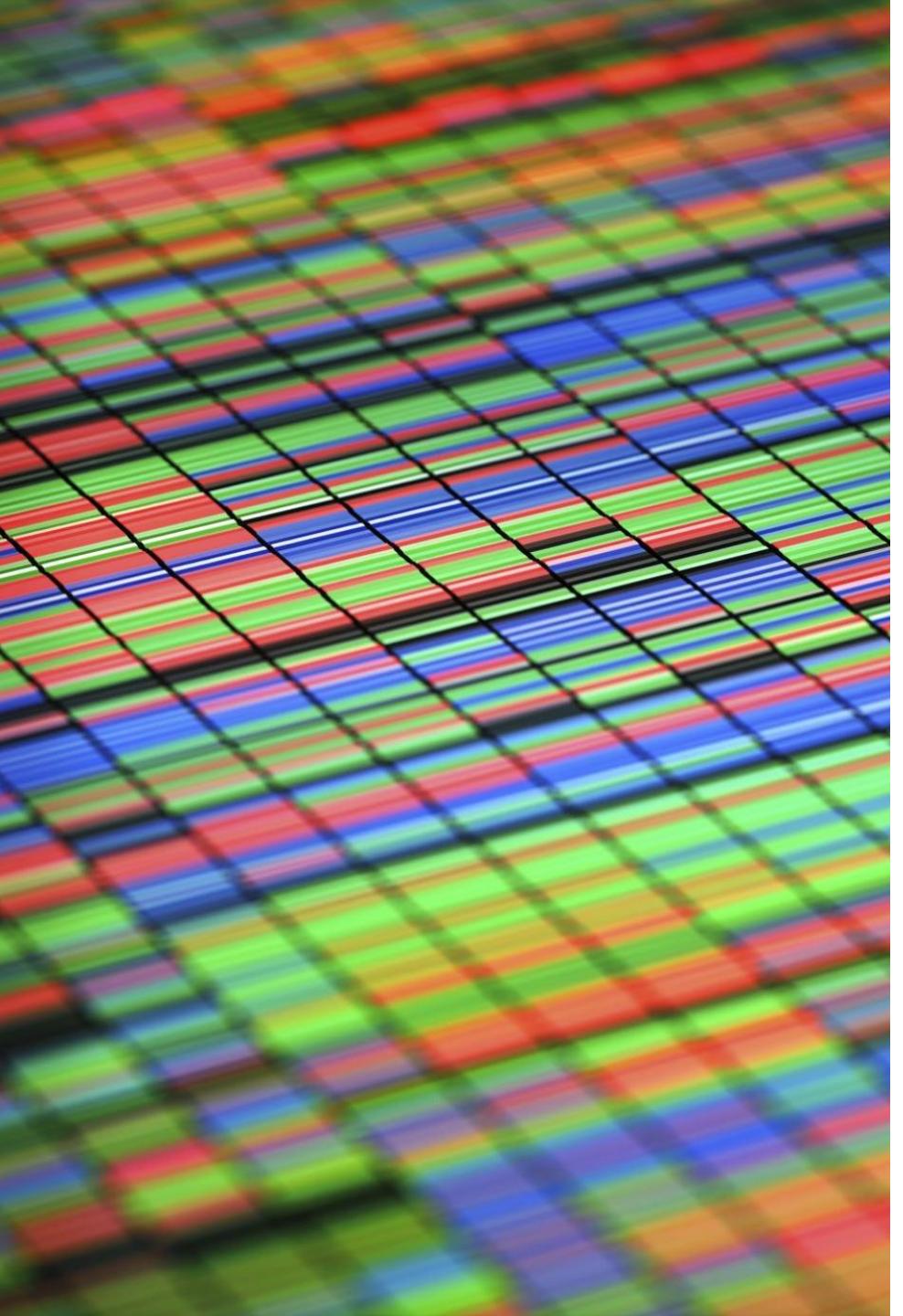
Network Framework Description

First, the Jigsaw puzzle is permuted before being fed into an AlexNet-like CNN.

Then, each patch/puzzle is processed by the CNN. In this case, puzzles in the network's early layers will NOT communicate with one another.

The weights are the same for each branch.

There are numerous permutation options. The goal is to correctly predict the index of the chosen permutation (technically, we define as output a probability vector with 1 at the 64-th location and 0 elsewhere).



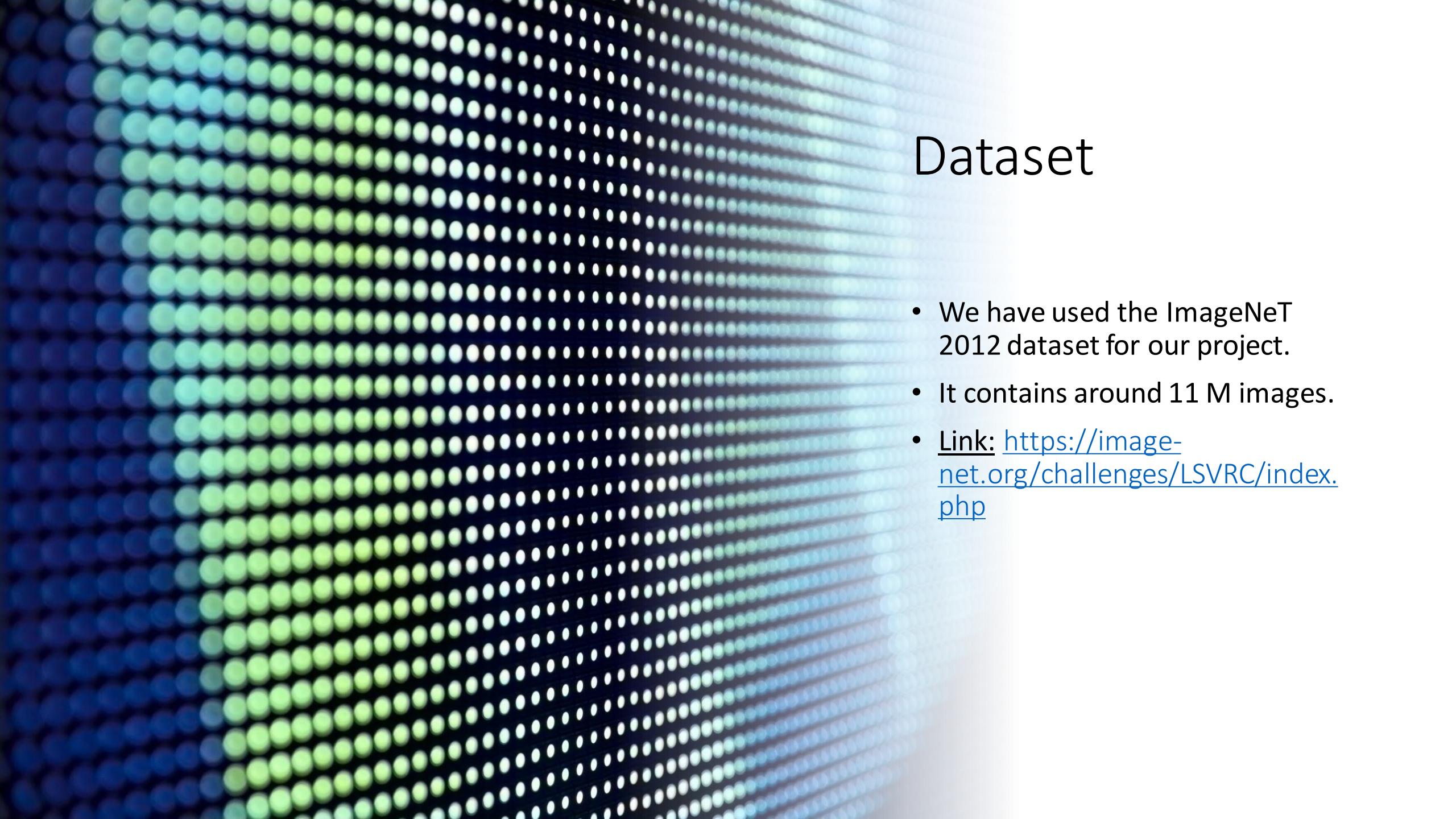
Network Training Description

- During training, each input image is resized until either the height or the width matches 256 pixels and preserve the original aspect ratio. Then, a random region is cropped from the resized image of size 225×225 and is split into a 3×3 grid of 75×75 pixels tiles.
- This network can thus be used interchangeably for different tasks including detection and classification.
- This may sound simple but the paper presenters took 2.5 days to train on TitanX GPU. We trained on Ada by submitting it as batch job for 4 days.



Phases of our project work

- Dataset loading : Trying out on small miniset and then on 170Gbs of ImageNet dataset. Yeah manipulating data was big problem!!
- Creating the framework: We created a CFN model, and added layers of convolution, ReLU, etc.
- Training: Setting the hyperparameters like learning rate, loss function, etc.

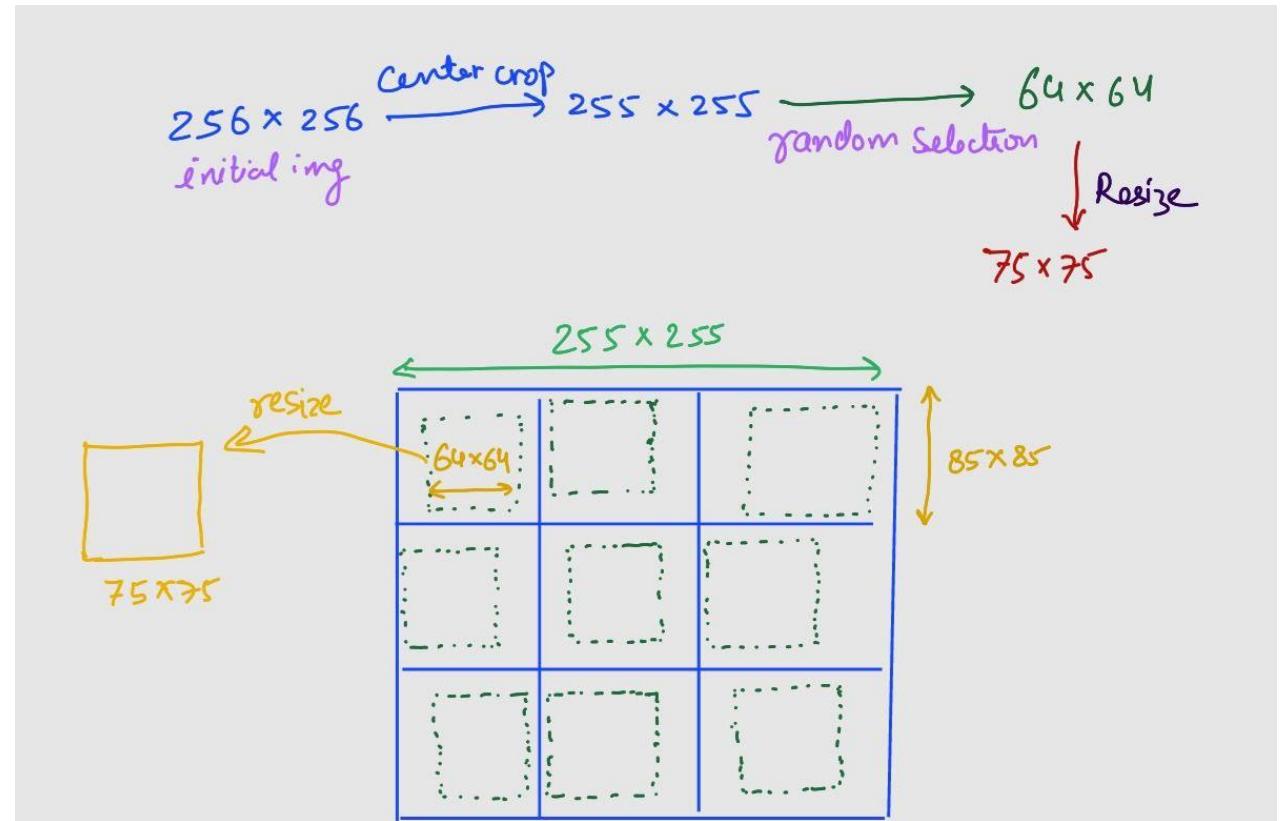
The background of the slide features a dense, abstract pattern of small, glowing dots arranged in horizontal rows. The dots transition through various colors, including shades of blue, green, and white, creating a digital or futuristic aesthetic.

Dataset

- We have used the ImageNeT 2012 dataset for our project.
- It contains around 11 M images.
- Link: <https://image-net.org/challenges/LSVRC/index.php>

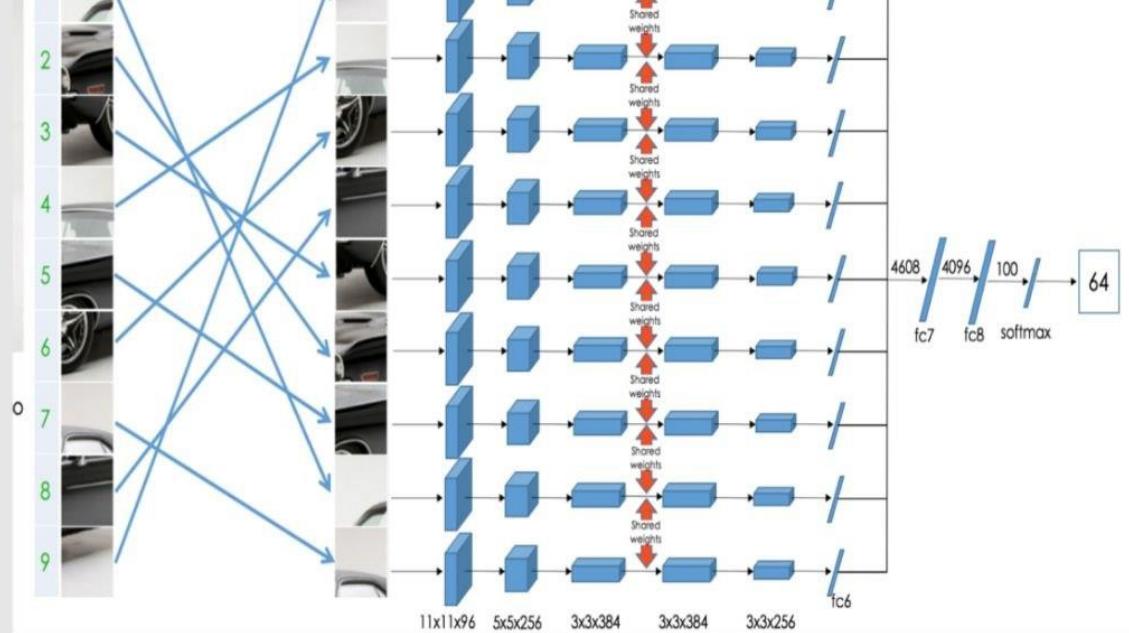
Generating Permutations

- We generated all $9!$ possible permutations (only done once). Then we picked the best 1000 permutations based on maximizing the metric of Hamming distance and give each one a unique index.
- Now, we permute the image randomly based on one of these 1000 permutations.



Framework

- We have implemented a Context Free Network (CFN) which is like an AlexNet CNN.
- Then we pass each patch/puzzle through our CFN.
- These patches will not communicate with each other within the network layers until the last 3 fully connected layers. This is done to ensure that patches do not learn from each other, preventing the model from learning patterns in the images (low-level feature).
- Each branch shares the same weights.
- There are a lot of permutation choices. The task is to predict the index of the chosen permutation (technically, we define as output a probability vector of size 1000 with 1 at the i^{th} location and 0 elsewhere, where i is the index of the vector predicted).

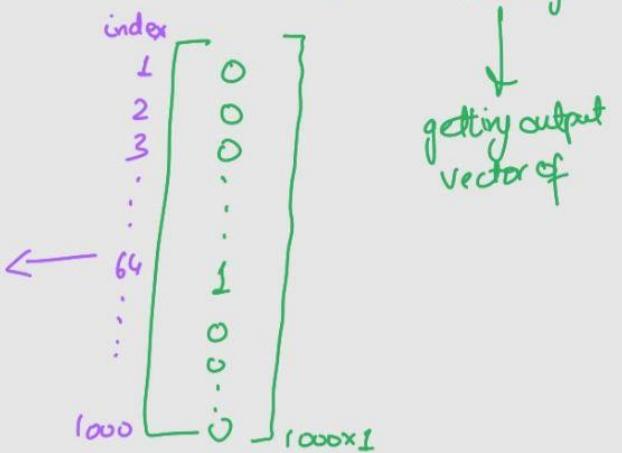


Picking vector from
1000 permutation
set & they have
unique index

Alexnet like network

Last 3
layers
to get
Representation
Learning

1 at index 64 means
the vector that was
picked corresponds to
order of 75th index



JigsawNetwork.py

1

Define a network
(add layers
Convolution, ReLU,
max-pool).

2

Load the model.

3

Weight Initialization
(Xavier
normalization).

4

Forward propagation,
load/create
Checkpoints (saving
various states of
model).

Training

Load train/test/validation data, define batch size = 128.

Learning rate (α) = 0.001

Loss function = Cross-Entropy Loss

Optimizer = Adam

No. of epochs = 70

Experiments/Analysis

- Since, the validation loss was not decreasing during the validation phase, we tried out a few different experiments to try and figure out the problem:

Hyperparameters

Learning Rate (α)

We experimented with various learning rates to see if something was working xD.

The learning rates we tried were: {1e-3,1e-4,1e-2,1e-1}

None of these values for the learning rate seemed to reduce the training loss.

Something fishy seems to be happening in the code

Also, batch sizes of 128 and 64 were tried and none seemed to show any significant positive result.

Dataset Change: Checking Data loader part

MNIST

- The first dataset we tried was the MNIST dataset.
- However, the images from this dataset were stored in .csv format, so we had to read data of a particular buffer length from the file, and then manually reshape them, and then finally convert to .png format.
- We divided the dataset into about ~ 40000 images for the training set and ~ 10000 images for the validation set.



Dataset Change : Checking Data loader part

Smaller ImageNet dataset

- We also tried training our model on a single class dataset of Dog's images obtained from Image NET consisting of ~ 16000 images for the training set, and ~ 2000 images for the validation set.
- Permutation and reordering of the images seemed to work well so it ensures that Dataset loader section was correct, and we need to verify the network.

```
dataset_mnist.py - cv [SSH: gnode] - Visual Studio Code

File Edit Selection View Go Run Terminal Help

dataset.py M dataset_mnist.py 1 original.jpeg mnist.py train.py TrainingUtils.py JigsawTrain.py S, M JigsawNet D v

EXPLORER OPEN EDITORS
dataset.py JigSawCustom M
dataset_mnist.py JigSawCustom... 1
original.jpeg JigSawCustom
mnist.py JigSawCustom
train.py JigSawCustom
TrainingUtils.py JigSawPuzzlePyton...
JigsawTrain.py JigSawPuzz... S, M
JigsawNetwork.py JigSawCu... U

CV [SSH: GNODE]
> mnist
cropped.jpeg

dataset_mnist.py I
dataset.py M
JigsawNetwork.py U
Layers.py U
mnist.ipynb
mnist.py U
orig.jpeg
original.jpeg
original0.jpeg
original1.jpeg
original2.jpeg
original3.jpeg
original4.jpeg
original5.jpeg
original6.jpeg
original7.jpeg
original8.jpeg
perm.jpeg
permutations_1000.npy

> OUTLINE
> TIMELINE

dataset.py M dataset_mnist.py 1 original.jpeg mnist.py train.py TrainingUtils.py JigsawTrain.py S, M JigsawNet D v

84
85     @staticmethod
86     def rgb_jittering(im):
87         im = np.array(im, 'int32')
88         for ch in range(3):
89             im[:, :, ch] += np.random.randint(-2, 2)
90         im[im > 255] = 255
91         im[im < 0] = 0
92         return im.astype('uint8')

93
94
95 import matplotlib.pyplot as plt
96 if __name__ == '__main__':
97     data_path = 'data/ILSVRC2012_img_train'
98     dataset = ImageData(data_path)
99     permuted,order,original = dataset[0]
100    print("Dataset Permutation : ",dataset.permutations[order])
101    transform=T.ToPILImage()
102    f, ax = plt.subplots(3, 3)
103    for i,tile in enumerate(original):
104        tile=transform(tile)
105        ax[i//3][i%3].imshow(tile)
106        # tile.save(f"original{i}.jpeg")
107    plt.savefig("orig.jpeg")

108
109
110    f, ax = plt.subplots(3, 3)
111    for i,tile in enumerate(permuted):
112        tile=transform(tile)
113        ax[i//3][i%3].imshow(tile)
114        # tile.save(f"permuted{i}.jpeg")
115    plt.savefig("perm.jpeg")
```

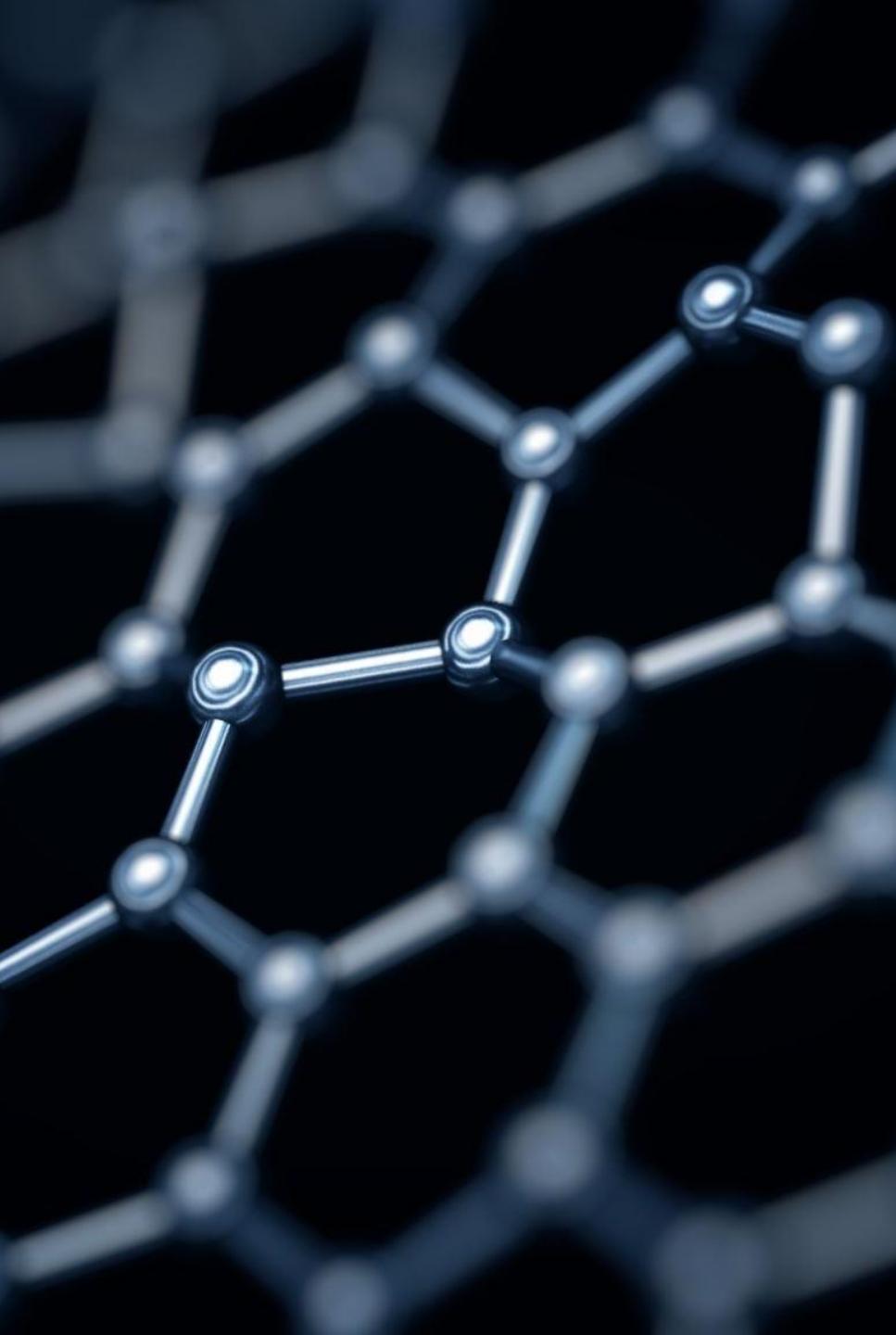
LRN (Local Response Normalization) Layer: Changes to Network

The reason we may want to have normalization layers in our CNN is that we want to have some kind of lateral inhibition scheme, where an excited neuron subdues its neighbors.

This results in a significant peak so that we have a form of local maxima.

ReLU neurons have unbounded activations and we need LRN to normalize that. We want to detect high frequency features with a large response. If we normalize around the local neighborhood of the excited neuron, it becomes even more sensitive as compared to its neighbors.

At the same time, it will dampen the responses that are uniformly large in any given local neighborhood, and boost the neurons with relatively larger activations.



LRN (Local Response Normalization) Layer

LRN can be performed within the same channel or across multiple channels.

In our code, we tried Response Normalization within the same channel.

Transfer Learning

- We tested model by importing AlexNet model from Pytorch as it was already per-trained on the existing ImageNet dataset.
- The next thing we tried was to freeze the weights of AlexNet model and then try training.
- But this also didn't give us any useful output.

Trying out GitHub Code

- Next thing was to find some online implementation of the paper.
- Their code was in python 2.7 and there was some conflict in version of pytorch and tensorflow while making conda environment.
- There was an issue with dataloader part of github code as it was choosing some custom and existing dataset images so we needed to tweak it accordingly.
- Next, we trained out on small dataset of github's implemendation and found that their code was not converging and so we trained it for 4 days on Imagenet dataset.

Trying out GitHub Code

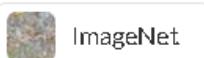
Code

Edit

- | | | |
|--|----------------------|---|
|  facebookresearch/vissl | ★ 2,599 |  |
| <small>↳ Quickstart in  Colab</small> | | PyTorch |
|  bbrattoli/JigsawPuzzlePytorch | ★ 132 |  |
|  gitlimlab/Representation-Learning-b... | ★ 107 |  |
|  Confusezius/selfsupervised_learning | ★ 9 |  |
|  virtualgraham/sc_patch | ★ 9 |  |
|  ferrannoguera/MLMI-Transfromers | ★ 2 |  |
|  SharadGitHub/Self-Supervised-Autoen... | ★ 1 |  |

[Collapse 7 implementations](#)

Datasets



Author's Train Loss

```
Using GPU 0
Process number: 36256
Images: train 16464, validation 2058
Start training: lr 0.000100, batch size 128, classes 1000
Checkpoint: checkpoints/
Learning Rate 0.000100
[ 1/70]    0) [batch load 4.356sec, net 0.28sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 1/70]   20) [batch load 0.227sec, net 0.03sec], LR 0.00010, Loss: 6.906, Accuracy 0.78%
[ 1/70]   40) [batch load 0.130sec, net 0.02sec], LR 0.00010, Loss: 6.908, Accuracy 0.00%
[ 1/70]   60) [batch load 0.094sec, net 0.02sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 1/70]   80) [batch load 0.075sec, net 0.02sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 1/70]  100) [batch load 0.020sec, net 0.02sec], LR 0.00010, Loss: 6.908, Accuracy 0.78%
[ 1/70]  120) [batch load 0.017sec, net 0.01sec], LR 0.00010, Loss: 6.906, Accuracy 0.00%
Learning Rate 0.000100
[ 2/70]  140) [batch load 0.043sec, net 0.02sec], LR 0.00010, Loss: 6.906, Accuracy 0.00%
[ 2/70]  160) [batch load 0.042sec, net 0.01sec], LR 0.00010, Loss: 6.906, Accuracy 0.00%
[ 2/70]  180) [batch load 0.043sec, net 0.02sec], LR 0.00010, Loss: 6.906, Accuracy 0.78%
[ 2/70]  200) [batch load 0.044sec, net 0.02sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 2/70]  220) [batch load 0.047sec, net 0.02sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 2/70]  240) [batch load 0.017sec, net 0.01sec], LR 0.00010, Loss: 6.907, Accuracy 0.78%
Learning Rate 0.000100
[ 3/70]  260) [batch load 0.045sec, net 0.01sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 3/70]  280) [batch load 0.045sec, net 0.02sec], LR 0.00010, Loss: 6.905, Accuracy 0.00%
[ 3/70]  300) [batch load 0.046sec, net 0.02sec], LR 0.00010, Loss: 6.906, Accuracy 1.56%
[ 3/70]  320) [batch load 0.045sec, net 0.02sec], LR 0.00010, Loss: 6.906, Accuracy 0.00%
[ 3/70]  340) [batch load 0.047sec, net 0.02sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
[ 3/70]  360) [batch load 0.020sec, net 0.01sec], LR 0.00010, Loss: 6.906, Accuracy 0.78%
[ 3/70]  380) [batch load 0.016sec, net 0.01sec], LR 0.00010, Loss: 6.907, Accuracy 0.00%
```

Learning Rate 0.000100

```
[70/70] 8960) [batch load 0.049sec, net 0.01sec], LR 0.00010, Loss: 6.900, Accuracy 0.00%
[70/70] 8980) [batch load 0.048sec, net 0.01sec], LR 0.00010, Loss: 6.901, Accuracy 0.00%
[70/70] 9000) [batch load 0.050sec, net 0.01sec], LR 0.00010, Loss: 6.902, Accuracy 0.78%
[70/70] 9020) [batch load 0.049sec, net 0.01sec], LR 0.00010, Loss: 6.902, Accuracy 0.00%
[70/70] 9040) [batch load 0.053sec, net 0.01sec], LR 0.00010, Loss: 6.903, Accuracy 0.78%
[70/70] 9060) [batch load 0.025sec, net 0.01sec], LR 0.00010, Loss: 6.902, Accuracy 0.00%
[70/70] 9080) [batch load 0.023sec, net 0.01sec], LR 0.00010, Loss: 6.904, Accuracy 0.78%
```

Author's Validation Loss

logs_val_0.0001_.txt

```
1 Using GPU 0
2 Process number: 1409
3 Images: train 16464, validation 2058
4 ('Starting from: ', 'jps_007_001000.pth.t
5 Evaluating network.....
6 TESTING: 0), Accuracy 0.05%
7 |
```

```
train_logs_custom_0.001.txt
1 Using device: cuda
2 Iteration: 20, Epoch: 1, Loss: 6.9091
3 Iteration: 40, Epoch: 1, Loss: 6.9132
4 Iteration: 60, Epoch: 1, Loss: 6.9106
5 Iteration: 80, Epoch: 1, Loss: 6.9055
6 Iteration: 100, Epoch: 1, Loss: 6.9062
7 Iteration: 120, Epoch: 1, Loss: 6.9080
8 Iteration: 140, Epoch: 1, Loss: 6.9078
9 Iteration: 160, Epoch: 1, Loss: 6.9151
10 Iteration: 180, Epoch: 1, Loss: 6.9023
11 Iteration: 200, Epoch: 1, Loss: 6.9165
12 Iteration: 220, Epoch: 1, Loss: 6.9078
13 Iteration: 240, Epoch: 1, Loss: 6.9053
14 Epoch [1/70], Loss: 6.9049
15 Accuracy of the network on the 2058 test images: 0.04295432458697766
16
17 Iteration: 20, Epoch: 2, Loss: 6.9083
18 Iteration: 40, Epoch: 2, Loss: 6.9096
19 Iteration: 60, Epoch: 2, Loss: 6.9020
20 Iteration: 80, Epoch: 2, Loss: 6.9048
21 Iteration: 100, Epoch: 2, Loss: 6.9097
22 Iteration: 120, Epoch: 2, Loss: 6.9052
23 Iteration: 140, Epoch: 2, Loss: 6.9075
24 Iteration: 160, Epoch: 2, Loss: 6.9065
25 Iteration: 180, Epoch: 2, Loss: 6.9075
26 Iteration: 200, Epoch: 2, Loss: 6.9022
27 Iteration: 220, Epoch: 2, Loss: 6.9112
28 Iteration: 240, Epoch: 2, Loss: 6.9041
29 Epoch [2/70], Loss: 6.9233
30 Accuracy of the network on the 2058 test images: 0.04859086491739553
31
```

Our Train and Validation Loss (initial)

- implemented without the AlexNet pre-trained model

```
Iteration: 20, Epoch: 69, Loss: 6.9064
Iteration: 40, Epoch: 69, Loss: 6.9058
Iteration: 60, Epoch: 69, Loss: 6.9080
Iteration: 80, Epoch: 69, Loss: 6.9108
Iteration: 100, Epoch: 69, Loss: 6.9101
Iteration: 120, Epoch: 69, Loss: 6.9065
Iteration: 140, Epoch: 69, Loss: 6.9168
Iteration: 160, Epoch: 69, Loss: 6.9159
Iteration: 180, Epoch: 69, Loss: 6.9064
Iteration: 200, Epoch: 69, Loss: 6.9084
Iteration: 220, Epoch: 69, Loss: 6.9064
Iteration: 240, Epoch: 69, Loss: 6.9029
Epoch [69/70], Loss: 6.8996
Accuracy of the network on the 2058 test images: 0.0457725947521866
```

```
Iteration: 20, Epoch: 70, Loss: 6.9059
Iteration: 40, Epoch: 70, Loss: 6.9130
Iteration: 60, Epoch: 70, Loss: 6.9089
Iteration: 80, Epoch: 70, Loss: 6.9034
Iteration: 100, Epoch: 70, Loss: 6.9100
Iteration: 120, Epoch: 70, Loss: 6.9035
Iteration: 140, Epoch: 70, Loss: 6.9083
Iteration: 160, Epoch: 70, Loss: 6.9031
Iteration: 180, Epoch: 70, Loss: 6.9063
Iteration: 200, Epoch: 70, Loss: 6.9087
Iteration: 220, Epoch: 70, Loss: 6.9053
Iteration: 240, Epoch: 70, Loss: 6.9061
Epoch [70/70], Loss: 6.8889
Accuracy of the network on the 2058 test images: 0.04859086491739553
```

Our Train and Validation Loss (initial)

```
train_logs_alex_0.0001.txt
```

```
1 Using device: cpu
2 Iteration: 20, Epoch: 1, Loss: 6.9102
3 Iteration: 40, Epoch: 1, Loss: 6.9109
4 Iteration: 60, Epoch: 1, Loss: 6.9085
5 Iteration: 80, Epoch: 1, Loss: 6.9053
6 Iteration: 100, Epoch: 1, Loss: 6.9065
7 Iteration: 120, Epoch: 1, Loss: 6.9129
8 Iteration: 140, Epoch: 1, Loss: 6.9091
9 Iteration: 160, Epoch: 1, Loss: 6.9039
10 Iteration: 180, Epoch: 1, Loss: 6.8944
11 Iteration: 200, Epoch: 1, Loss: 6.9036
12 Iteration: 220, Epoch: 1, Loss: 6.9057
13 Iteration: 240, Epoch: 1, Loss: 6.9091
14 Epoch [1/70], Loss: 6.9199
15 Accuracy of the network on the 2058 test images: 0.1457725947521866
16
```

AlexNet Training Loss and Accuracy

```
Iteration: 20, Epoch: 2, Loss: 6.9015
Iteration: 40, Epoch: 2, Loss: 6.9072
Iteration: 60, Epoch: 2, Loss: 6.9112
Iteration: 80, Epoch: 2, Loss: 6.9063
Iteration: 100, Epoch: 2, Loss: 6.9052
Iteration: 120, Epoch: 2, Loss: 6.9092
Iteration: 140, Epoch: 2, Loss: 6.9067
Iteration: 160, Epoch: 2, Loss: 6.9085
Iteration: 180, Epoch: 2, Loss: 6.9060
Iteration: 200, Epoch: 2, Loss: 6.9075
Iteration: 220, Epoch: 2, Loss: 6.9068
Iteration: 240, Epoch: 2, Loss: 6.9068
Epoch [2/70], Loss: 6.9052
Accuracy of the network on the 2058 test images: 0.04859086491739553
```

```
Iteration: 20, Epoch: 3, Loss: 6.9113
Iteration: 40, Epoch: 3, Loss: 6.9080
Iteration: 60, Epoch: 3, Loss: 6.9068
Iteration: 80, Epoch: 3, Loss: 6.9065
Iteration: 100, Epoch: 3, Loss: 6.9041
Iteration: 120, Epoch: 3, Loss: 6.9080
Iteration: 140, Epoch: 3, Loss: 6.9088
Iteration: 160, Epoch: 3, Loss: 6.9088
Iteration: 180, Epoch: 3, Loss: 6.9062
Iteration: 200, Epoch: 3, Loss: 6.9054
Iteration: 220, Epoch: 3, Loss: 6.9027
Iteration: 240, Epoch: 3, Loss: 6.9047
Epoch [3/70], Loss: 6.9040
Accuracy of the network on the 2058 test images: 0.09718172983479106
```

AlexNet Training Loss and Accuracy (with frozen weights)

Comparisons

Github's Code

Training Accuracy → 0.78%

Validation Accuracy → 0.05%

Loss start → 6.907

Loss end → 6.904

Comparisons

Our Model's [Initial]

Training Accuracy → 1.56%

Validation Accuracy → 0.048%

Loss start → 6.9091

Loss end → 6.900

Comparisons

Using Modified AlexNet (without weight freezing)

Training Accuracy → 1.56 %.

Validation Accuracy → 0.145 %.

Loss start → 6.9102

Loss end → 6.9091

Comparisons

Using Modified AlexNet (with weight freezing)

Training Accuracy → 1.56%.

Validation Accuracy → 0.097%.

Loss start → 6.9015

Loss end → 6.9017

Avoiding Shortcuts

- Feature association with position: The CFN may accidentally learn to associate features with absolute positions if only 1 Jigsaw puzzle is generated per image.

$$p(S|A_1, A_2, \dots, A_9) = p(S|F_1, F_2, \dots, F_9) \prod_{i=1}^9 p(F_i|A_i) \quad (1)$$

- This shortcut can be avoided by feeding the CFN multiple Jigsaw puzzles of the same image, making sure that tiles are shuffled as much as possible and have adequately large average Hamming distance, ensuring that the same tile can be assigned to multiple positions within the grid and avoiding the correlation of an absolute position.

Algorithm 1. Generation of the *maximal* Hamming distance permutation set

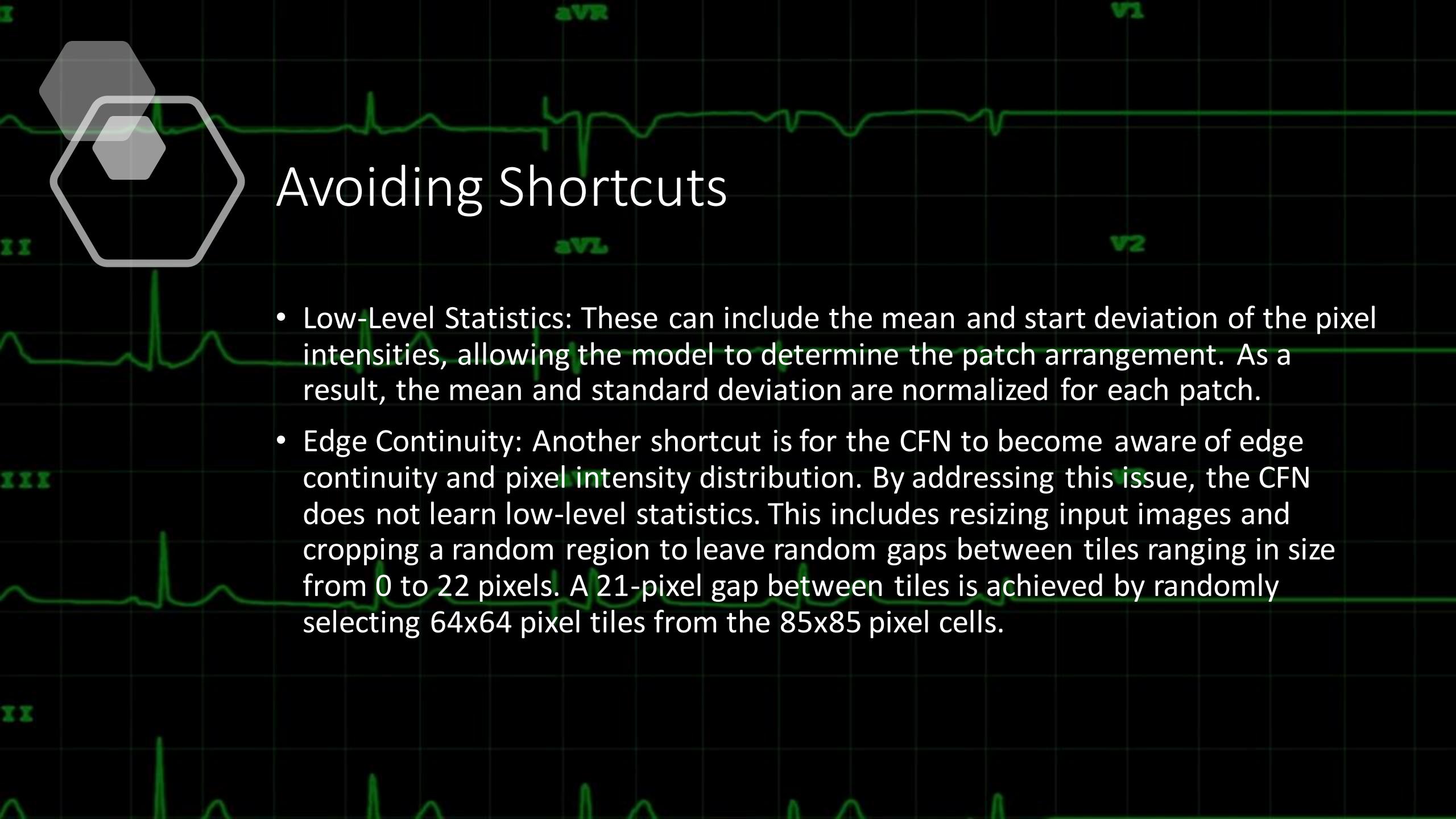
Input: N

\\ number of permutations

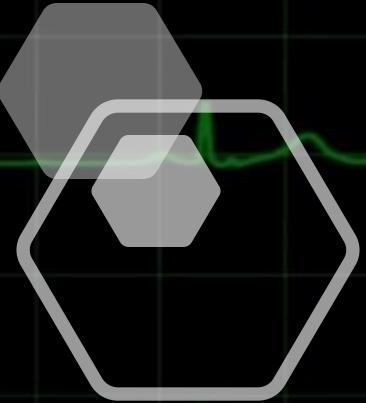
Output: P

\\ maximal permutation set

1: $\bar{P} \leftarrow$ all permutations $[\bar{P}_1, \dots, \bar{P}_{9!}]$ \\ \bar{P} is a $9 \times 9!$ matrix2: $P \leftarrow \emptyset$ 3: $j \sim \mathcal{U}[1, 9!]$ \\ uniform sample out of $9!$ permutations4: $i \leftarrow 1$ 5: **repeat**6: $P \leftarrow [P \ \bar{P}_j]$ \\ add permutation \bar{P}_j to P 7: $\bar{P} \leftarrow [\bar{P}_1, \dots, \bar{P}_{j-1}, \bar{P}_{j+1}, \dots]$ \\ remove \bar{P}_j from \bar{P} 8: $D \leftarrow \text{Hamming}(P, P')$ \\ D is an $i \times (9! - i)$ matrix9: $\bar{D} \leftarrow \mathbf{1}^T D$ \\ \bar{D} is a $1 \times (9! - i)$ row vector10: $j \leftarrow \arg \max_k \bar{D}_k$ \\ \bar{D}_k denotes the k -th entry of \bar{D} 11: $i \leftarrow i + 1$ 12: **until** $i \leq N$



Avoiding Shortcuts



- Low-Level Statistics: These can include the mean and start deviation of the pixel intensities, allowing the model to determine the patch arrangement. As a result, the mean and standard deviation are normalized for each patch.
- Edge Continuity: Another shortcut is for the CFN to become aware of edge continuity and pixel intensity distribution. By addressing this issue, the CFN does not learn low-level statistics. This includes resizing input images and cropping a random region to leave random gaps between tiles ranging in size from 0 to 22 pixels. A 21-pixel gap between tiles is achieved by randomly selecting 64x64 pixel tiles from the 85x85 pixel cells.

Avoiding Shortcuts

- Chromic Aberration: A third shortcut can be attributed to chromatic aberration, which is a relative spatial shift between colour channels that increases from the center images to the borders, assisting the network in estimating the tiles. This is resolved by cropping and resizing the original image's central square, training with both colour and grayscale images, jittering colour channels, and using greyscale images.

Some References

- <https://arxiv.org/pdf/1603.09246v3.pdf>
- <https://paperswithcode.com/paper/unsupervised-learning-of-visual-1#code>
- <https://towardsdatascience.com/difference-between-local-response-normalization-and-batch-normalization-272308c034ac>
- <https://towardsdatascience.com/transfer-learning-using-pre-trained-alexnet-model-and-fashion-mnist-43898c2966fb>

The background of the slide is a photograph of a long bridge spanning across a body of water. The bridge has multiple lanes of traffic, including several trucks and cars. The water below is a deep teal color.

Thank You !

We would like to thank Samartha Sir for being available for meet without prescheduling and providing idea about how to try to improve our paper implementation.