

2D HAND POSE ESTIMATION

USING U-NET LIKE ARCHITECTURE



-OISHIK DASGUPTA

-BDA 1st Year

-B2030035

INTRODUCTION

Passive sensing of human hand and limb motion is important for a wide range of applications from human-computer interaction to athletic performance measurement. Human hand pose estimation is a long standing problem in the computer vision and graphics research fields, with a plethora of applications such as:

- machine control,
- augmented and virtual reality.

PROBLEMS

Properties such as:

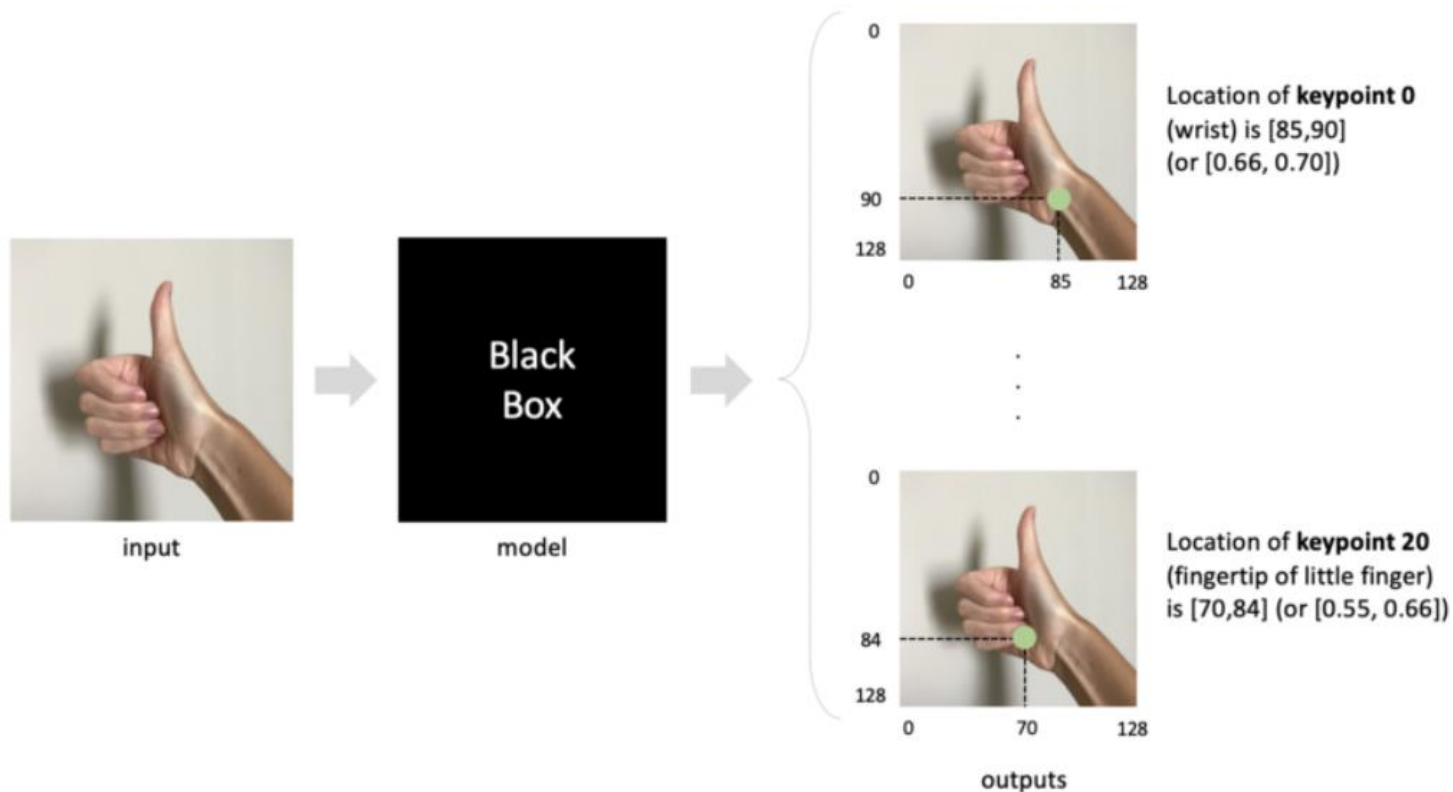
- the hand's morphology,
- occlusions due to interaction with objects,
- appearance diversity due to clothing and jewelry,
- varying lightning conditions and
- different backgrounds,

add extra burden to the nature of the problem.

WHAT EXACTLY ARE WE ESTIMATING IN THE HAND POSE ESTIMATION?

WHAT ABOUT THE SECOND HAND?

A TYPICAL 2D HAND POSE ESTIMATOR LOOKS SOMETHING LIKE THIS:



DATASET DESCRIPTION -A



FreiHAND Dataset:

It consists of **130240** images (of the right hand) of dimensions **224*224*3**.

I have worked with the **first 32560 images**. Other images in the dataset are exactly the same as raw ones but with background augmentation.

DATASET SPLIT:

Train Set: 0-25999 (80%)

Val. Set: 26000-30999 (15%)

Test Set: 30000-32560 (5%)

“training_K.json”:
Array of Camera
Intrinsic Matrix
for 32560 images

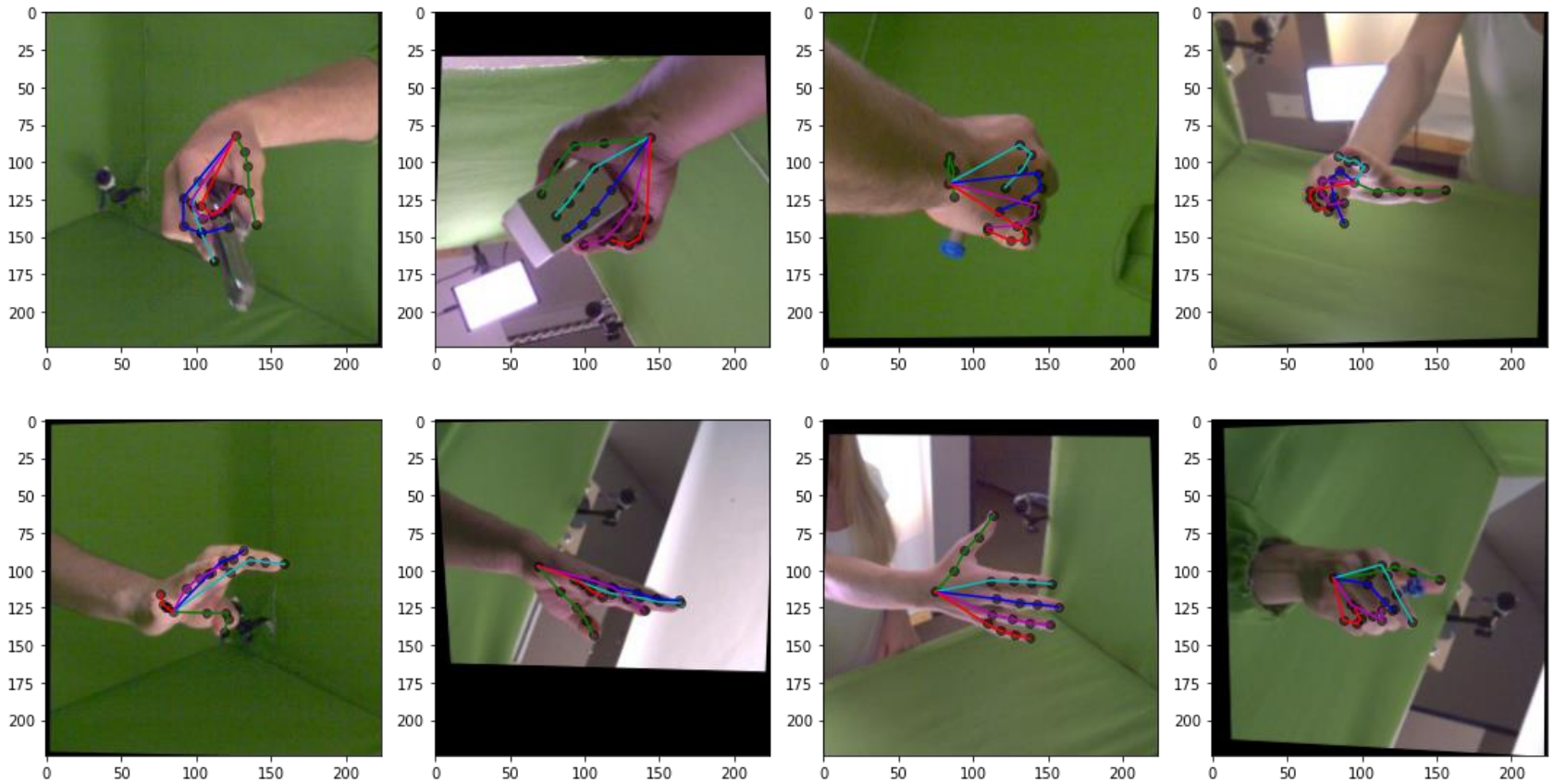
JSON	Raw Data	Headers
Save	Copy	Collapse All Expand All (slow)
▼ 0:		
▼ 0:		
0:	388.9018310596544	
1:	0	
2:	112	
▼ 1:		
0:	0	
1:	388.71231836584275	
2:	112	
▶ 2:	[...]	
▶ 1:	[...]	
▶ 2:	[...]	
▶ 3:	[...]	
▶ 4:	[...]	

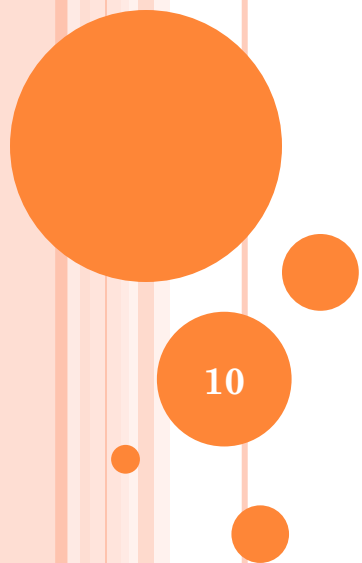
“training_xyz.json”: Array of annotations of
21 keypoints for 32560 images

JSON	Raw Data	Headers
Save	Copy	Collapse All Expand All (slow) 🔍 F
▼ 0:		
▼ 0:		
0:	0.029402047395706177	
1:	-0.027920207008719444	
2:	0.5870807766914368	
▶ 1:	[...]	
▶ 2:	[...]	
▶ 3:	[...]	
▶ 4:	[...]	
▶ 5:	[...]	
▶ 6:	[...]	
▶ 7:	[...]	
▶ 8:	[...]	
▶ 9:	[...]	
▶ 10:	[...]	
▶ 11:	[...]	
▶ 12:	[...]	
▶ 13:	[...]	
▶ 14:	[...]	
▶ 15:	[...]	
▶ 16:	[...]	
▶ 17:	[...]	
▶ 18:	[...]	
▶ 19:	[...]	
▼ 20:		
0:	-0.06662826985120773	
1:	0.07311560213565826	
2:	0.6494264602661133	
▶ 1:	[...]	
▶ 2:	[...]	
▶ 3:	[...]	



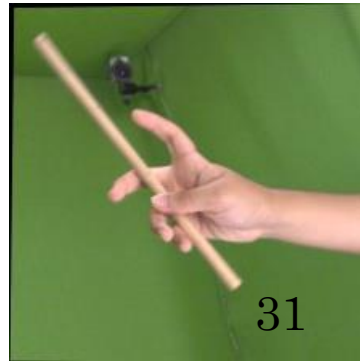
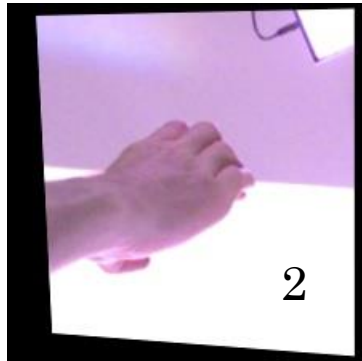
DATASET VISUALIZATION



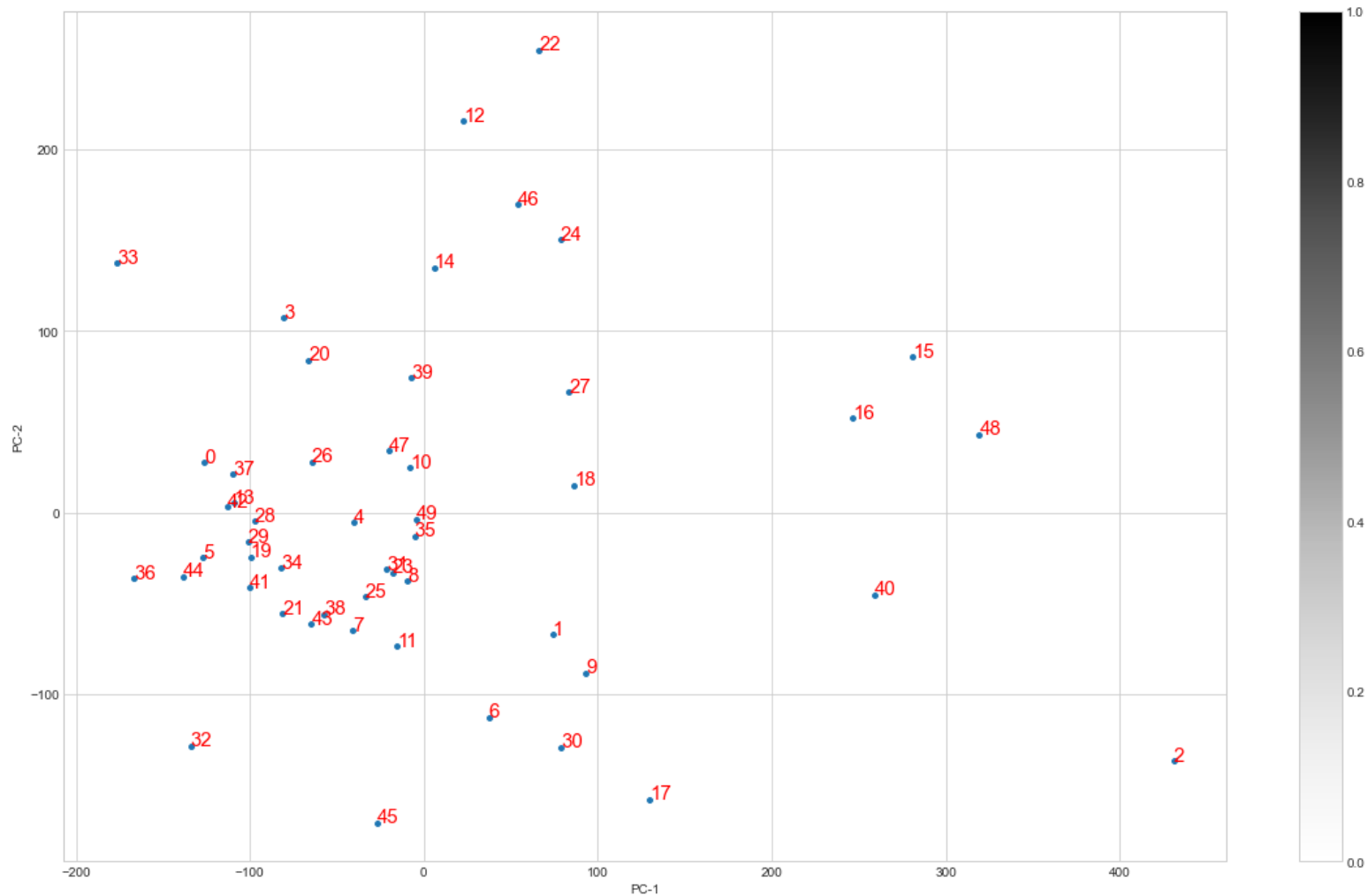


EDA

A FEW SAMPLE IMAGES

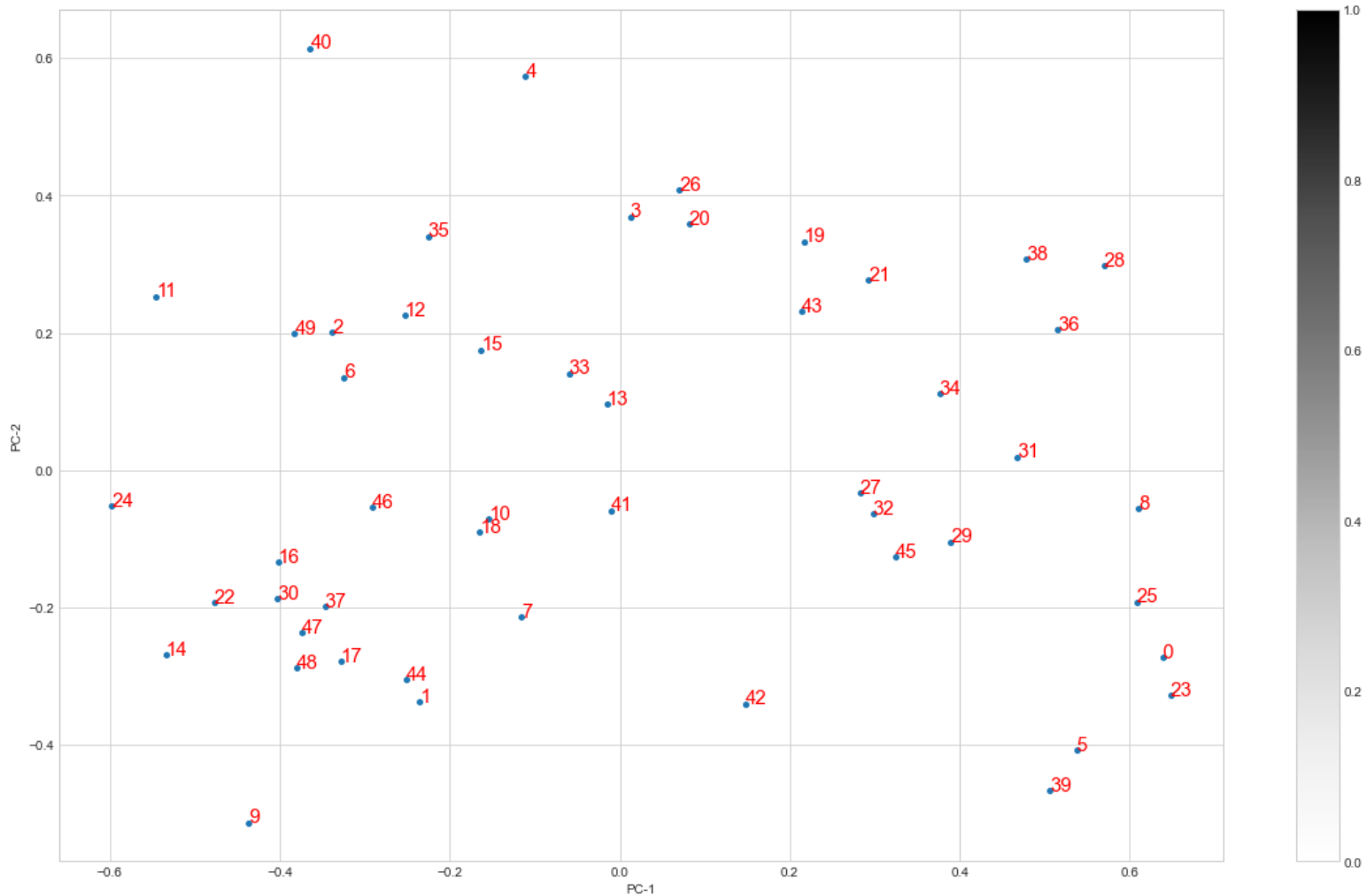


PCA OF 50 SAMPLE IMAGES FROM DATASET



PCA OF 50 SAMPLE IMAGES' KEYPOINTS

FROM DATASET





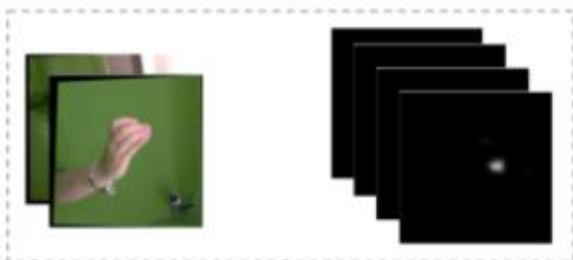
TRAINING WORKFLOW

14



DataLoader Class

- Load several images and their 2D locations
- Preprocess each images
 - Resize to 128x128
 - Min-Max Scale
 - Standard normalize
 - Create heatmaps
- Concatenate images and heatmaps into a single batch



$\text{batch_size} * 3 * 128 * 128$

$\text{batch_size} * 21 * 128 * 128$

Trainer Class

Loop over:

- Train step
 - Get batch of train data from dataloader
 - Give to model, get prediction
 - Calculate train loss
 - Take optimization step
- Validation Step
 - Get batch of validation data
 - Give to model, get predictions
 - Calculate val loss
 - Stop training if val loss does not decrease any more

Batch of images

Model Class



- Input: images
 $\text{batch_size} * 3 * 128 * 128$
- Output: heatmaps
 $\text{batch_size} * 21 * 128 * 128$

Model predictions (heatmaps)

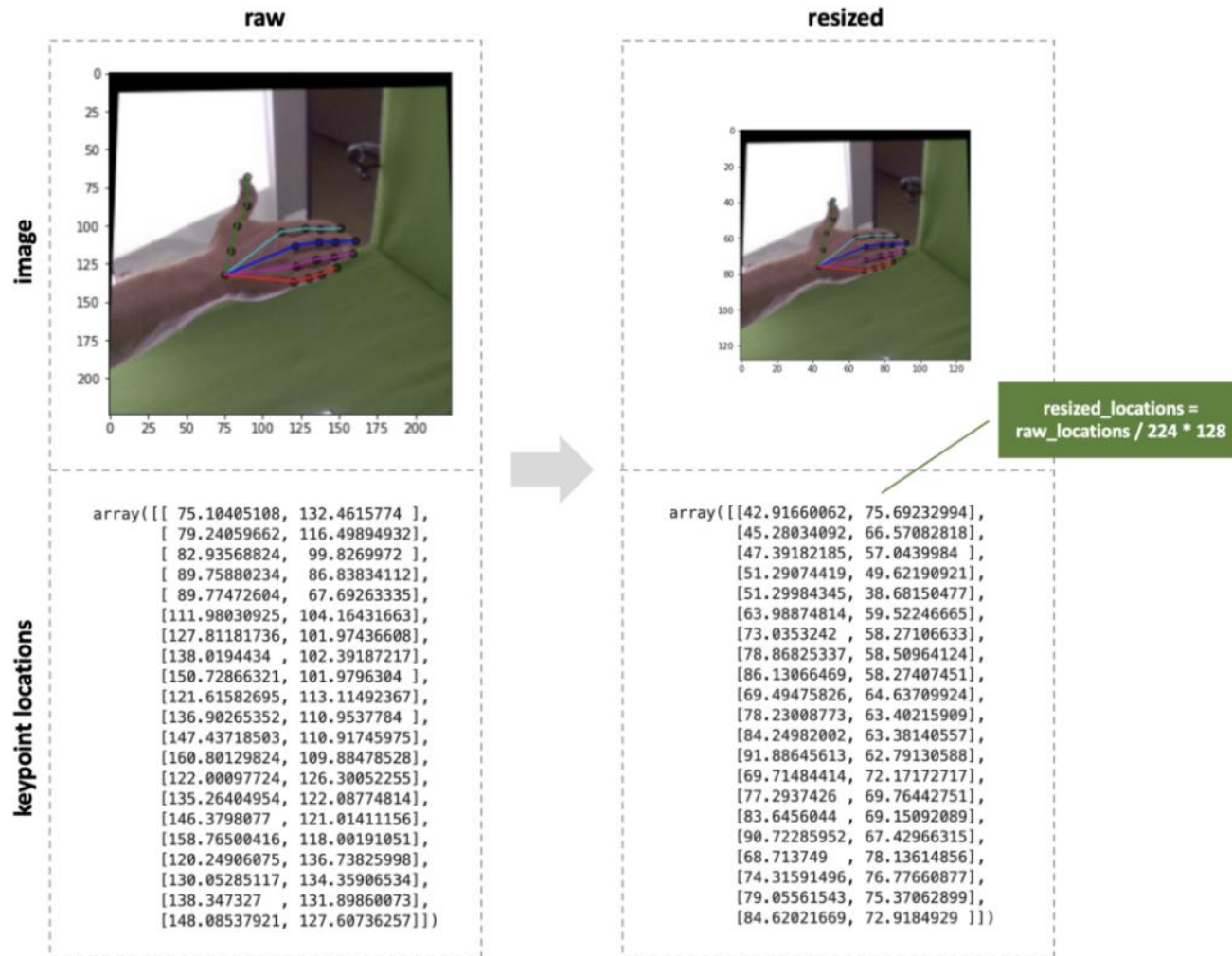
Batch of data

CALCULATION OF 2D LOCATIONS OF 21 KEYPOINTS FROM GIVEN 3D LOCATIONS

Code Snippet:

```
def projectPoints(xyz, K):  
    xyz = np.array(xyz)  
    K = np.array(K)  
    uv = np.matmul(K, xyz.T).T  
    return uv[:, :2] / uv[:, -1:]
```


WHEN RESIZING AN IMAGE, KEYPOINT LOCATIONS SHOULD BE ALSO “RESIZED”

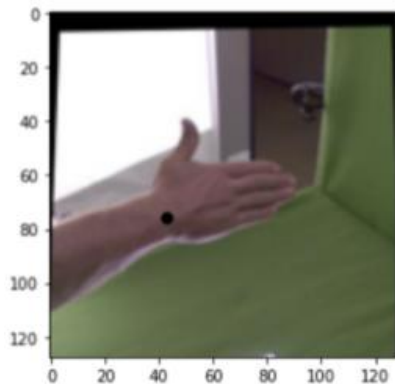


CREATE HEATMAPS

We need to create a separate heatmap for each keypoint, so there will be 21 heatmaps in total.

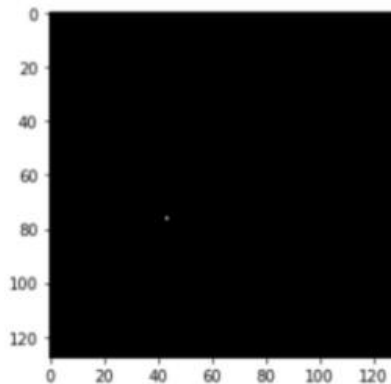
STEP 1

Choose a keypoint.
For instance, a wrist,
It's location on resized
image is (43,76)



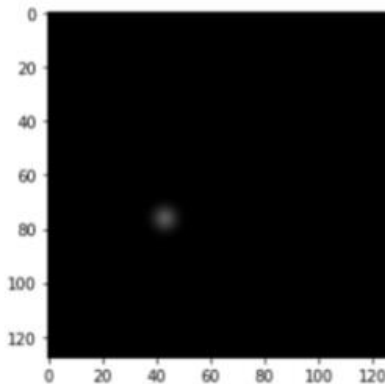
STEP 2

Create array with zeros
of size 128x128
Change entry [43,76] to
be 1.



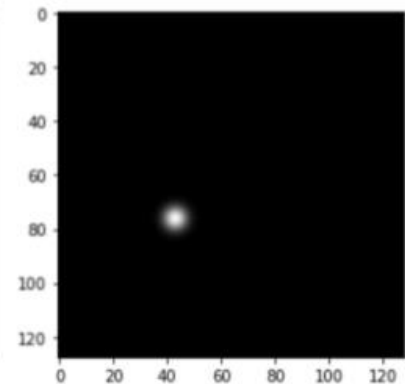
STEP 3

Gaussian Blur



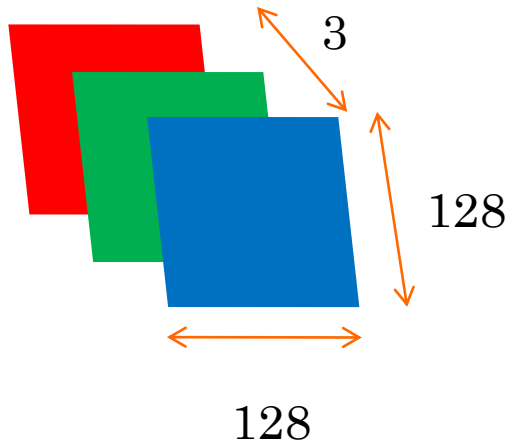
STEP 4

Divide array by its max
value, so array values
are in range [0,1]

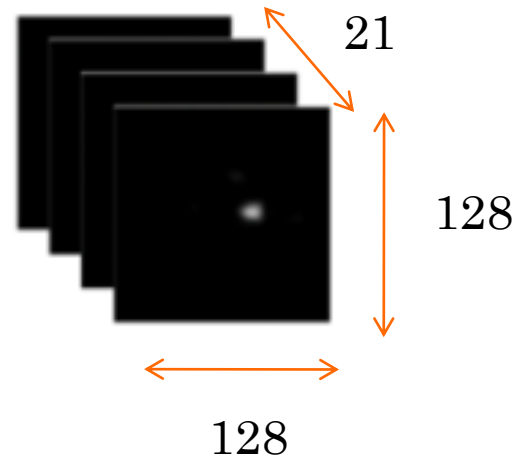


- ❑ X is an image of size $3 \times 128 \times 128$.
- ❑ Y is an array of size $21 \times 128 \times 128$, that contains 21 stacked heatmaps.

X:



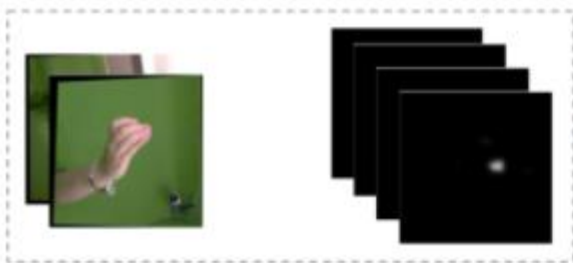
Y:





DataLoader Class

- Load several images and their 2D locations
- Preprocess each images
 - Resize to 128x128
 - Min-Max Scale
 - Standard normalize
 - Create heatmaps
- Concatenate images and heatmaps into a single batch



$\text{batch_size} * 3 * 128 * 128$

$\text{batch_size} * 21 * 128 * 128$

Trainer Class

Loop over:

- Train step
 - Get batch of train data from dataloader
 - Give to model, get prediction
 - Calculate train loss
 - Take optimization step
- Validation Step
 - Get batch of validation data
 - Give to model, get predictions
 - Calculate val loss
 - Stop training if val loss does not decrease any more

Batch of images

Model Class



- Input: images
 $\text{batch_size} * 3 * 128 * 128$
- Output: heatmaps
 $\text{batch_size} * 21 * 128 * 128$

Model predictions (heatmaps)

Batch of data

$\text{batch_size}=48$

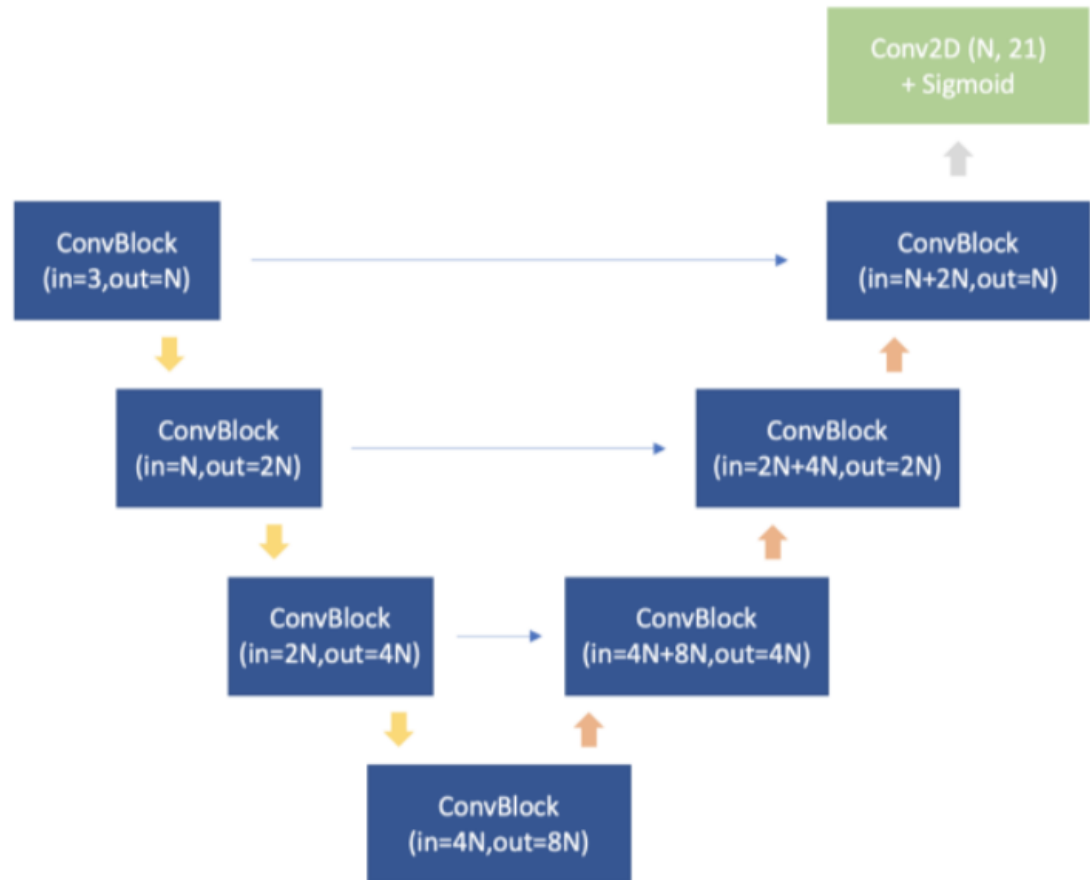
MODEL

Model Parameters

- N Init Neurons = 16
- ↓ MaxPool (size=2)
- ↑ UpSample (size=2)
- Skip connection (concatenate to layer input)

ConvBlock (in, out)

BatchNorm2D
Conv2D(in, out, size=3, padding=1)
ReLU
BatchNorm2D
Conv2D(out, out, size=3, padding=1)
ReLU



LOSS

WHY IoU LOSS?

LOSS EXPRESSION:

$$\mathbf{I} = \sum_i y_i * t_i$$

$$\mathbf{U} = \sum_i y_i * y_i + \sum_i t_i * t_i - \sum_i y_i * t_i$$

$$\mathbf{IoU} = \frac{\mathbf{I}}{\mathbf{U}}$$

$$\mathbf{L_{IoU}} = 1 - \mathbf{IoU}$$

y_i – predicted values

t_i – target values for a pixel in a heatmap.

Loss is calculated for each heatmap separately, then averaged among the images in the batch.

LOSS CALCULATION:

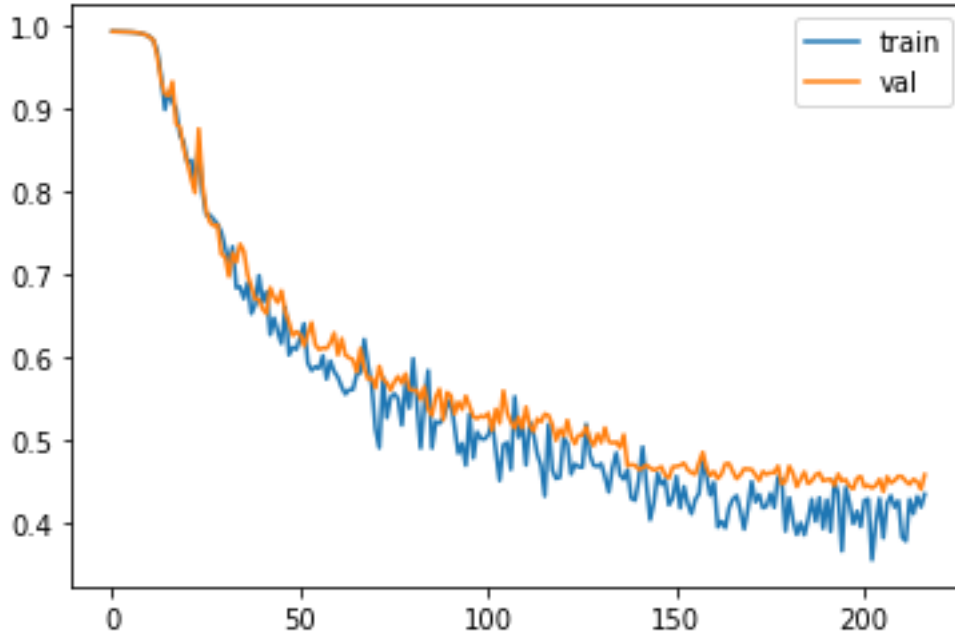
```
class IoULoss(nn.Module):
    """
    Intersection over Union Loss.
    IoU = Area of Overlap/Area of Union
    IoU loss is modified to use for heatmaps.
    """

    def __init__(self):
        super(IoULoss, self).__init__()
        self.EPSILON = 1e-6

    def _op_sum(self, x):
        return x.sum(-1).sum(-1)

    def forward(self, y_pred, y_true):
        inter = self._op_sum(y_true * y_pred)
        union = (
            self._op_sum(y_true ** 2)
            + self._op_sum(y_pred ** 2)
            - self._op_sum(y_true * y_pred)
        )
        iou = (inter + self.EPSILON) / (union + self.EPSILON)
        iou = torch.mean(iou)
        return 1 - iou
```

TRAINING RESULTS:



Training Loss v/s
Validation Loss with
progress in epochs

Epoch: 215/1000, Train Loss=0.4323002696, Val Loss=0.4503314694
Epoch: 216/1000, Train Loss=0.4195953608, Val Loss=0.4410485086
Epoch: 217/1000, Train Loss=0.4350762367, Val Loss=0.458478973

EVALUATION

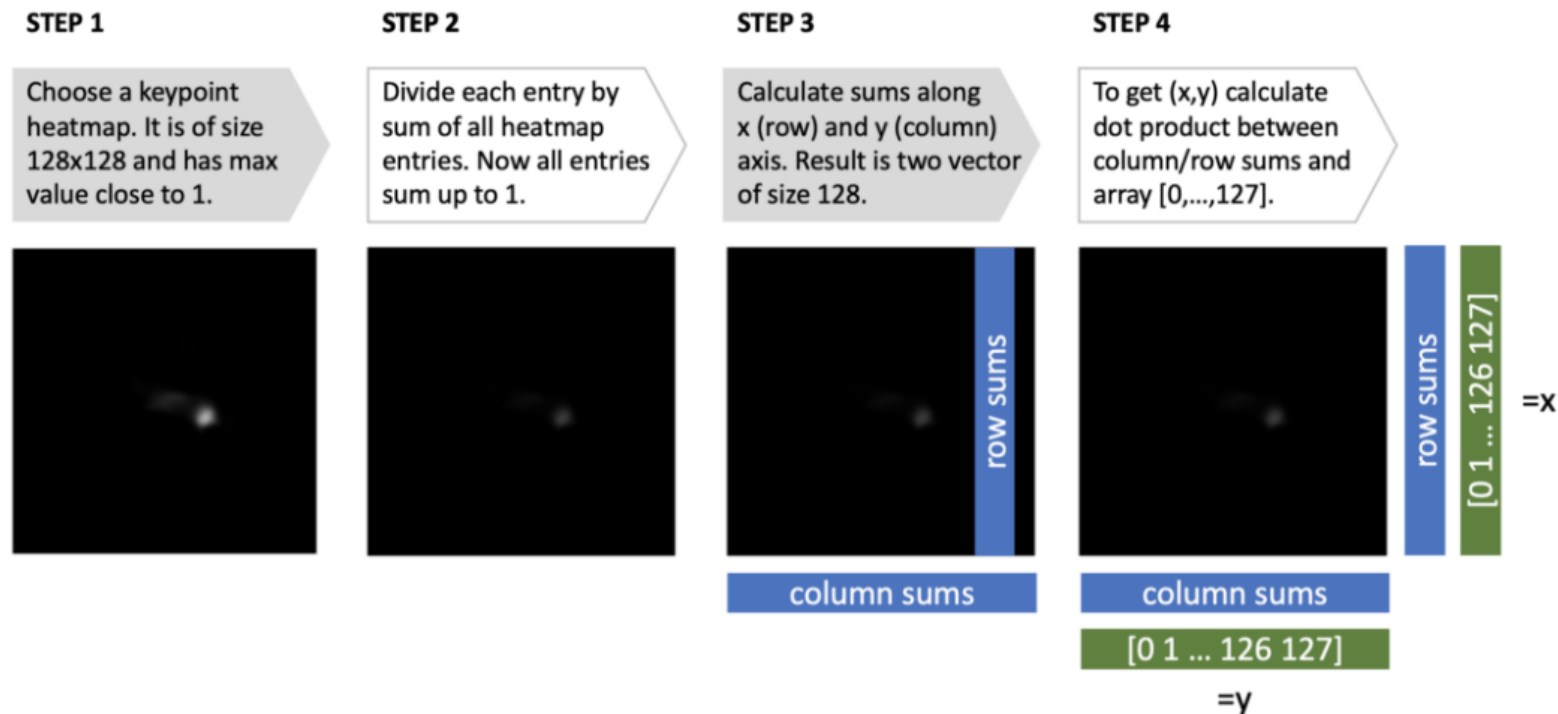
27

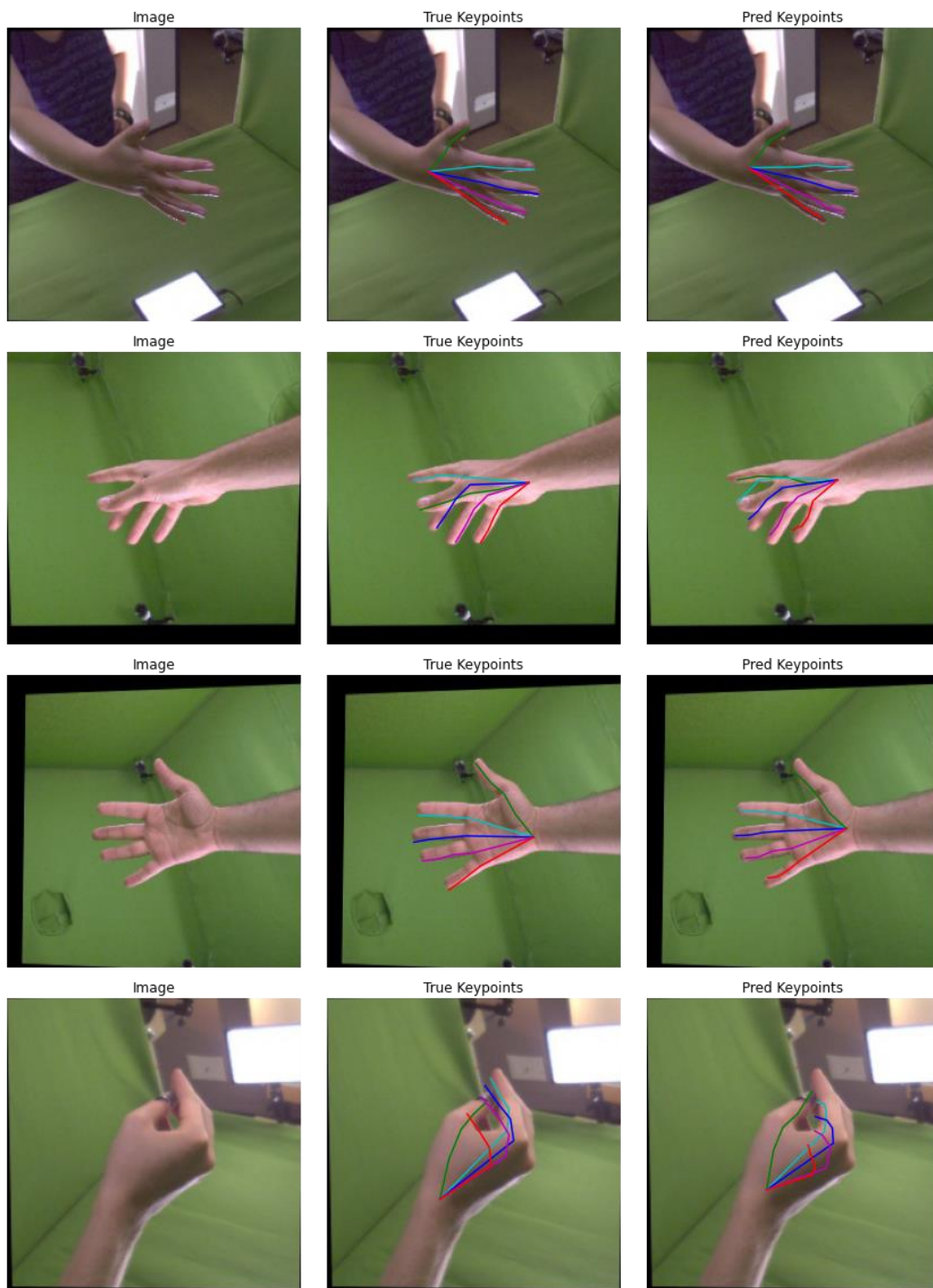
HEATMAP TO COORDINATES

There are 2 options:

➤ Most simply, we may just find a pixel with the largest value in the heatmap. (x,y) location of this pixel is the keypoint location.

➤ But a more robust way would be to calculate the average among all heatmap values.





EVALUATION METRIC:

- (1) Mean error for each joint:

Average error per keypoint: 4.2% from image size

Average error per keypoint: 5 pixels for image 128x128

Average error per keypoint: 10 pixels for image 224x224

- (2) Success rate:

- The proportion of test frames whose average error falls below a threshold

SUCCESS RATE FOR DIFFERENT THRESHOLD VALUES

