

Semantic Segmentation of RGB-D Images

- The DeResUNet
- BatchNorm2D & PixelShuffle
- Albulmentations
- PyTorch Lightning

Semantic Segmentation

- “*Semantic*” refers to the meaning or context associated to objects present in an image
- “*Segmentation*” is the process of partitioning an image into segments, each grouping pixels sharing common features
- In our Case:

$$f : I_x(u, v) \rightarrow c \in C^N$$

$$\|\phi(I_x(u1, v1)) - \phi(I_y(u2, v2))\| < \epsilon \Rightarrow f(I_x(u1, v1)) = f(I_y(u2, v2)) = c_i$$

where $f \hat{=}$ segmentation function and $\phi(I_x) \hat{=}$ repr. of I_x in feature space

The SunRGBD Dataset



Intel Realsense



Asus Xtion



Kinect v1



Kinect v2

Color



Raw depth



Improved depth



10.335 RGB-D



6.4 GB



Toolbox
provided



Annotations for
2D/3D
146,617 2D polygons
58,657 3D bounding
boxes

Scene Classification



Semantic Segmentation



Room Layout



Detection and Pose



Total Scene Understanding

Augmentation

Albumentations

vs torchvision.transforms

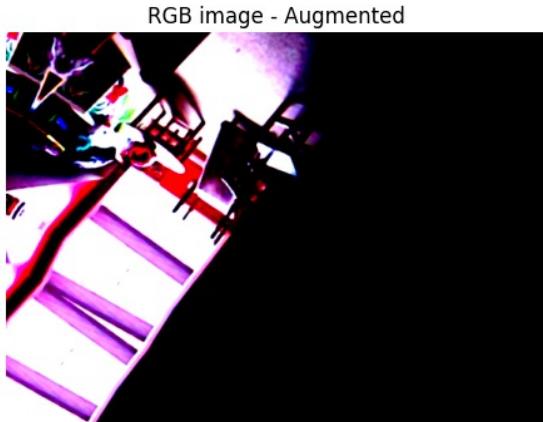
- ⊕ Performance (up to 23x faster)^[A1]
- ⊕ More variety – Extreme Augmentations
- ⊕ Usage with any DL framework
- ⊕ Easy transforms for Segmentation and Detection Tasks

Optimized for NumPy and OpenCV

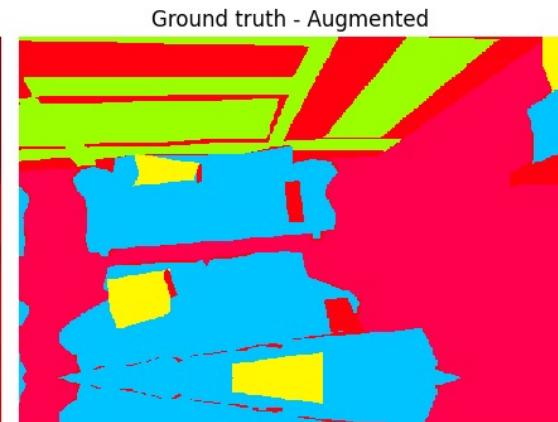
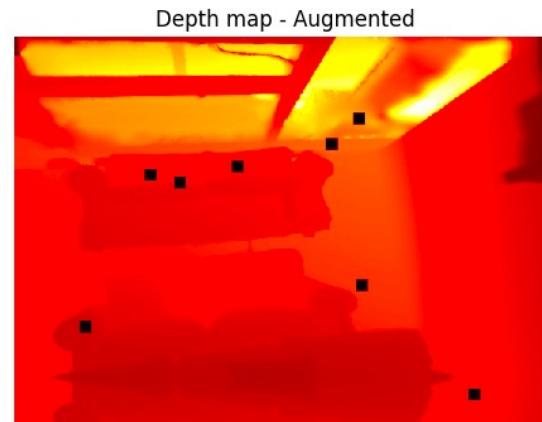
```
A.RGBShift(  
    r_shift_limit=20, g_shift_limit=20, b_shift_limit=20, p=0.95  
)  
,  
A.OneOf(  
    [  
        A.Blur(blur_limit=3, p=0.5),  
        A.ColorJitter(p=0.8),  
    ],  
    p=0.8,  
)  
,  
A.CLAHE(p=0.5),  
A.RandomBrightnessContrast(p=0.6),  
A.GaussianBlur(p=0.4),  
A.RandomGamma(p=0.6),  
]  
)  
self.tf_augment = A.Compose(  
    [  
        A.RandomCrop(  
            width=int(self.width * 0.5), height=int(self.height * 0.5), p=0.45  
)  
,  
        A.HorizontalFlip(p=0.5),  
        A.VerticalFlip(p=0.2),  
        A.RandomScale(scale_limit=0.2, p=0.8),  
        A.ElasticTransform(alpha=32, sigma=50, alpha_affine=120 * 0.03, p=0.6),  
        A.OpticalDistortion(p=0.3),  
        A.GridDistortion(num_steps=2, distort_limit=0.2, p=0.4),  
        A.ShiftScaleRotate(  
            shift_limit=0.0625, scale_limit=0.1, rotate_limit=45, p=0.6  
)  
,  
        A.RandomRotate90(p=0.25),  
        A.CoarseDropout(max_holes=8, max_height=16, max_width=16, p=0.3),  
    ],  
)  
self.tf_norm_and_resize = A.Compose(  
    [  
        A.Resize(  
            width=256, height=256, p=1.0),  
        A.Normalize(mean=[0.485, 0.45, 0.406], std=[0.229, 0.224, 0.225], p=1.0),  
    ]  
)
```

Augmentation

Sample 4637 - normalized: True - split_type: train



Sample 5443 - normalized: True - split_type: train

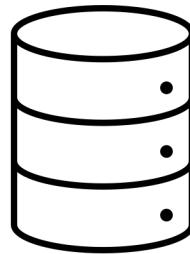


The Network

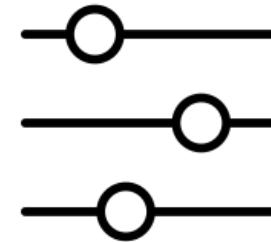
*DeResUNet*¹



Layers: 399
(not)



Size: 777.36 MB



Params: 203,780,570

¹ *DualencoderResidualUNet*

The Network

- 2 x parallel ResNet50 ([ResNet50_Weights.DEFAULT](#))

\longleftrightarrow ResNet4Channel

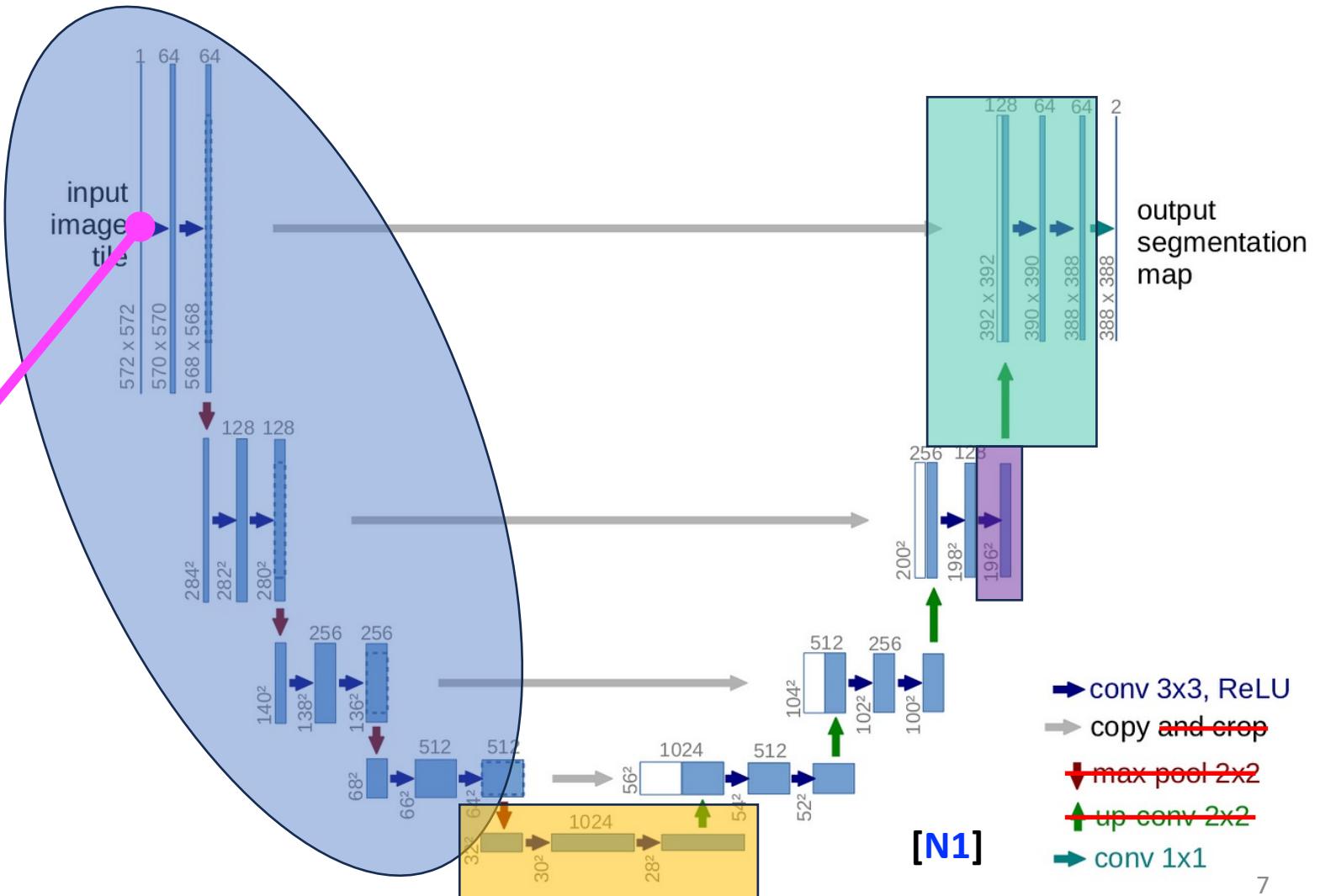
- Bottleneck (latent space)

- UpConvPxShuffle¹

made out of:

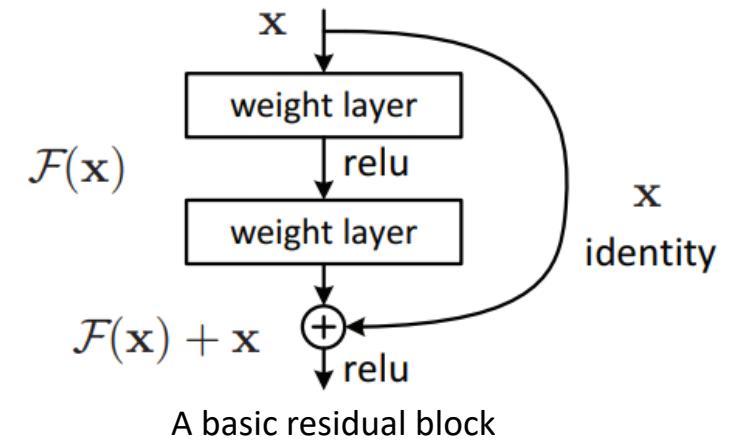
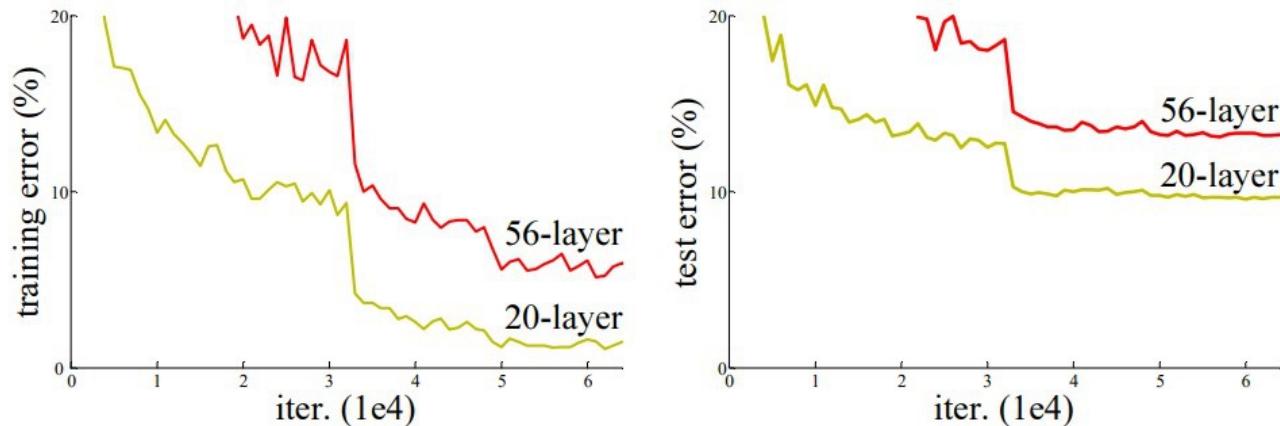
- 2x ConvNormRelu

$$T_{in} \in \mathbb{R}^{B \times 4 \times H \times W}$$



The Network

A Residual Block

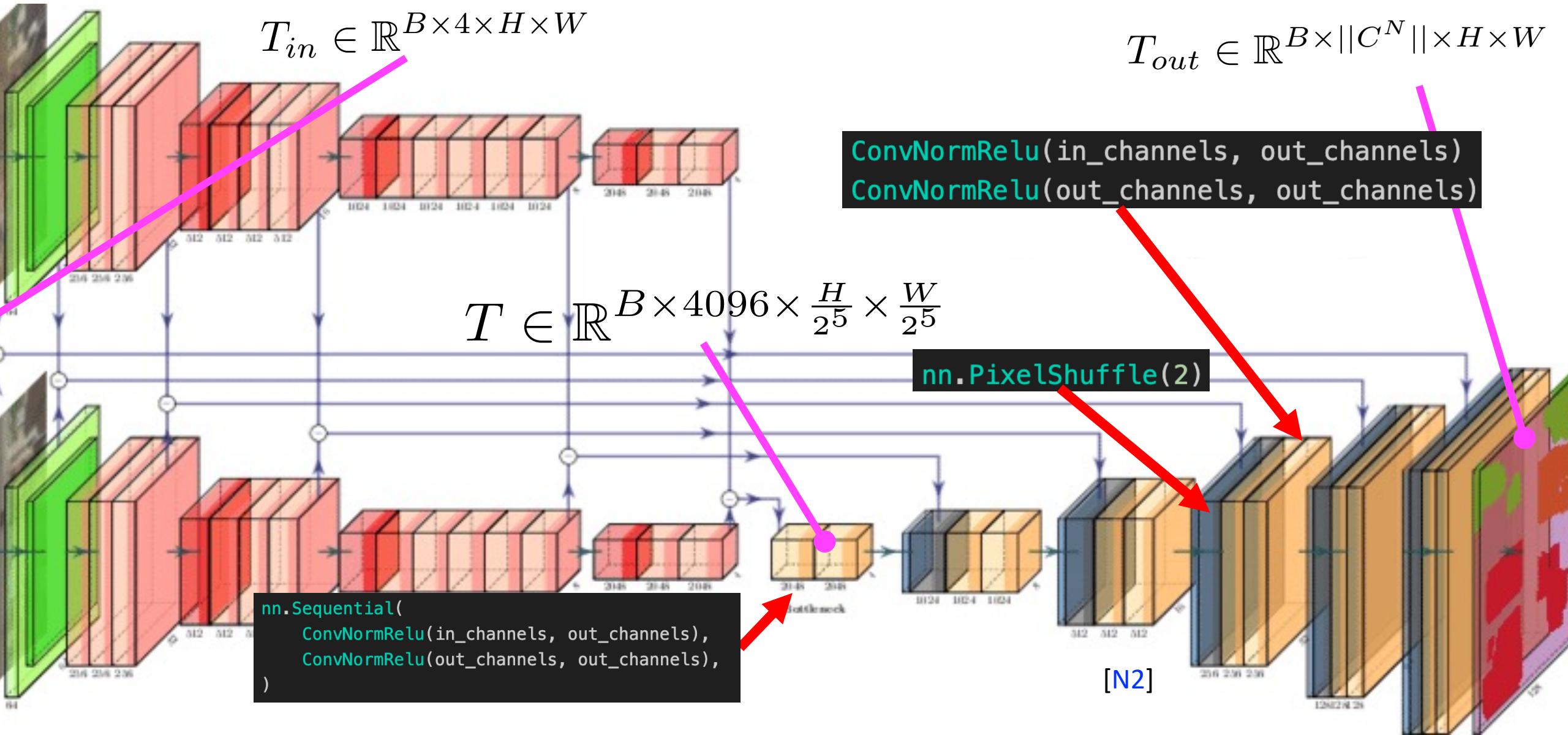


- Deeper networks -> vanishing gradient problem
- “Skip connections” that perform identity mapping
- Blocks can be stacked

Variants of Residual Blocks:

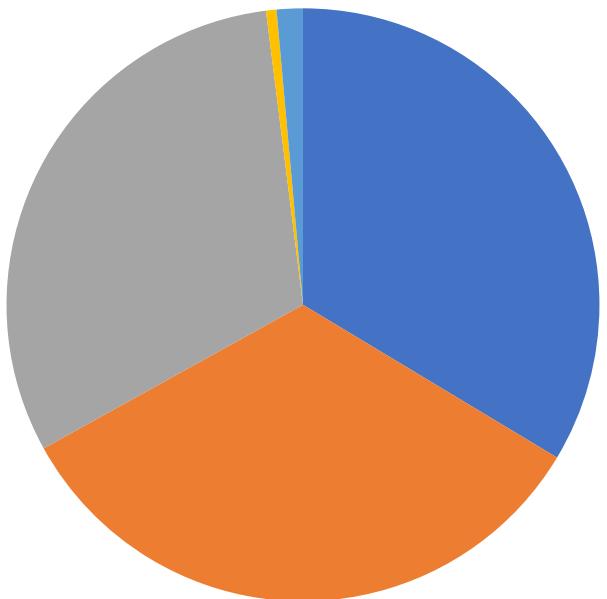
- Basic Block
- Bottleneck Block
- Pre-activation Block
- Transformer Block

The Network



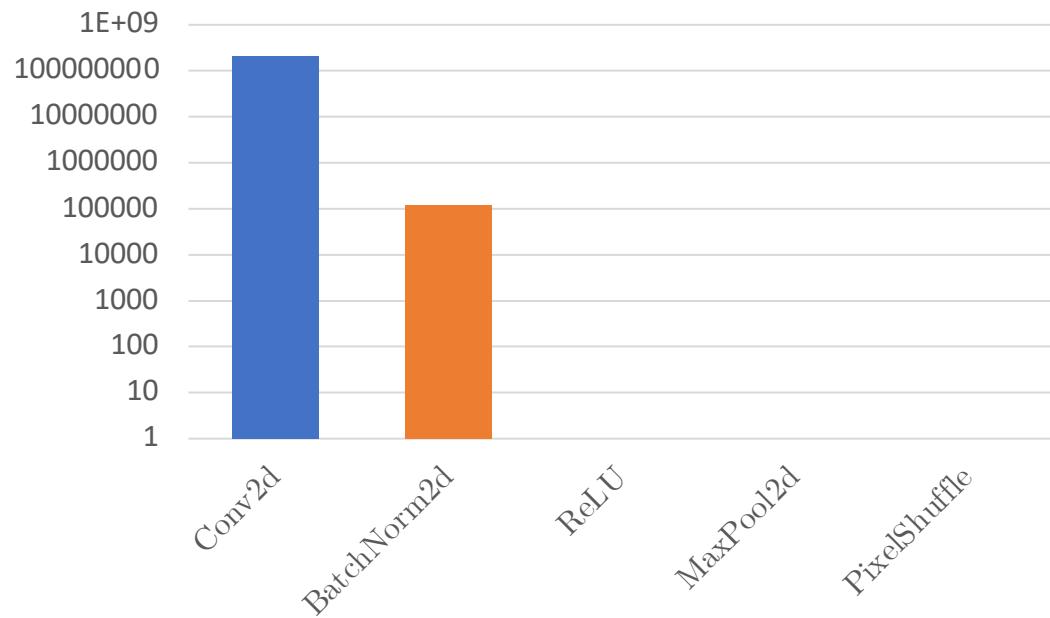
The Network

Number of Layers



- Conv2d
- BatchNorm2d
- ReLU
- MaxPool2d
- PixelShuffle

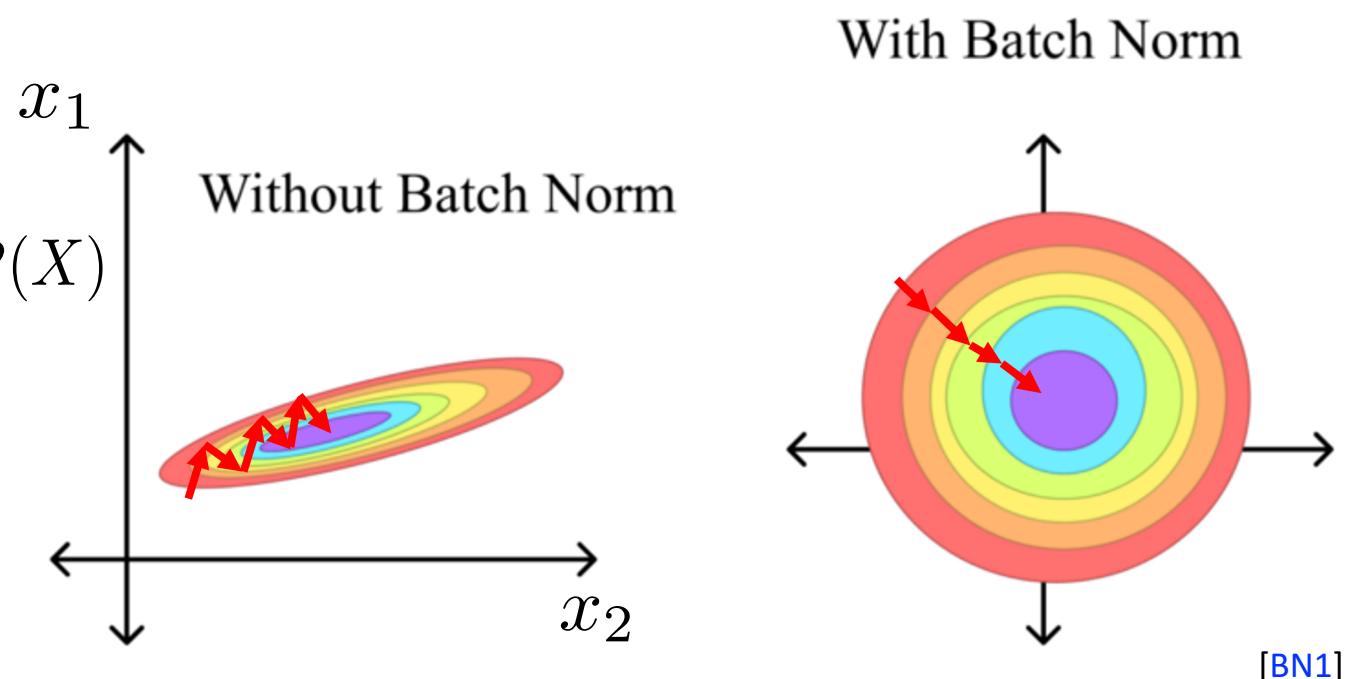
Number of Parameters



- Conv2d
- BatchNorm2d
- ReLU
- MaxPool2d
- PixelShuffle

BatchNorm2D

- Reduces *internal covariate shift*
- Regularization:
$$\mu_x = \mathbb{E}[X_{MB}]$$
$$\sigma^2 = \mathbb{E}[(X_{MB} - \mu)^2]$$
- Allows higher learning rates
 \Rightarrow faster convergence!
- Decreases bias towards *initial weights*



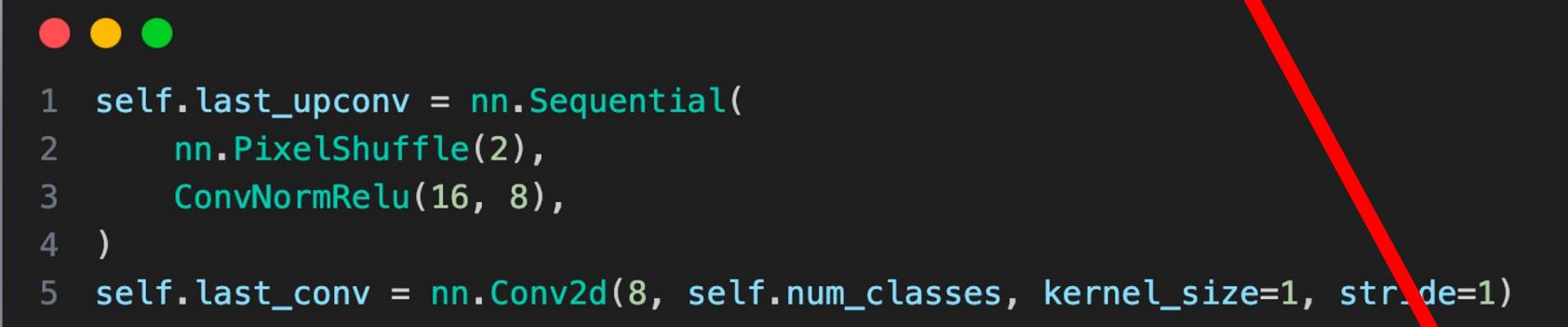
$$\hat{X}_{MB} = \gamma \left(\frac{X_{MB} - \mu_{MB}}{\sqrt{\sigma_{MB}^2 + \epsilon}} \right) + \beta$$

Learnable params

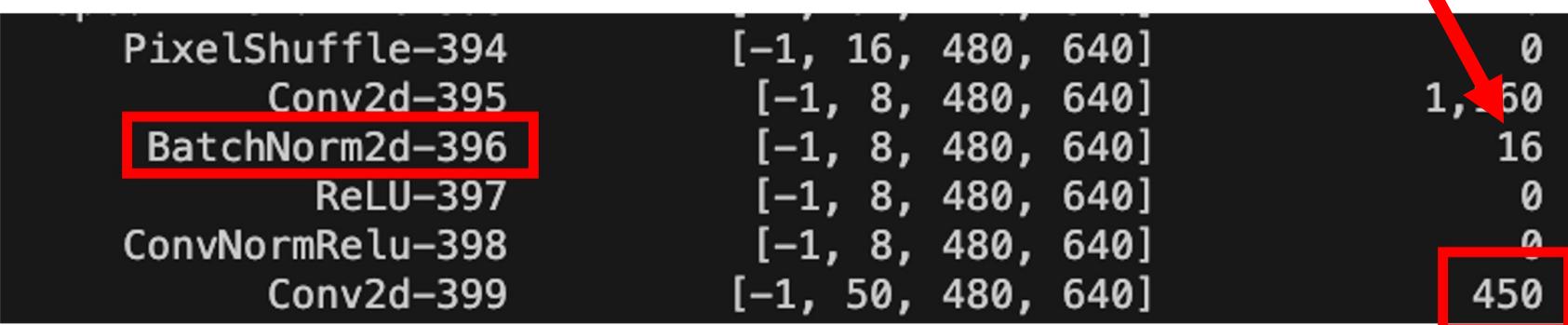
BatchNorm2D

- Last 1x1 Conv2D layer probably can't deal with the introduced *stochasticity!*

γ, β



```
1 self.last_upconv = nn.Sequential(  
2     nn.PixelShuffle(2),  
3     ConvNormRelu(16, 8),  
4 )  
5 self.last_conv = nn.Conv2d(8, self.num_classes, kernel_size=1, stride=1)
```



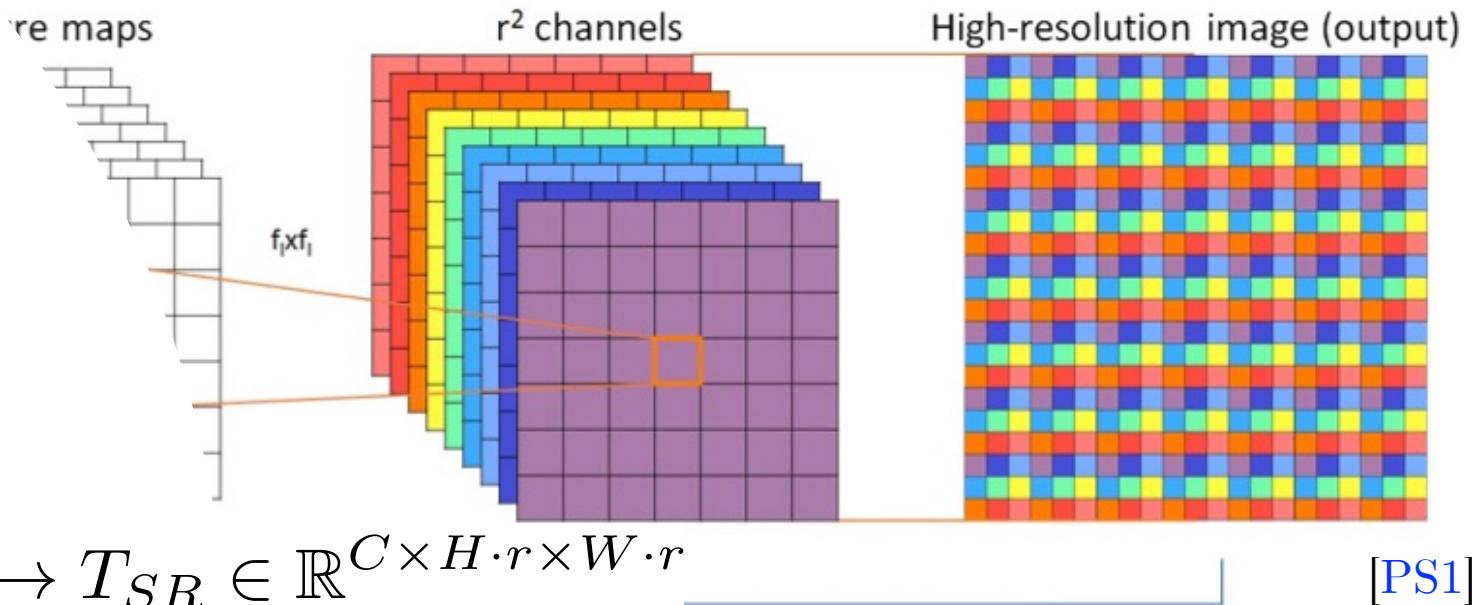
Layer	Output Shape	Count
PixelShuffle-394	[-1, 16, 480, 640]	0
Conv2d-395	[-1, 8, 480, 640]	1,150
BatchNorm2d-396	[-1, 8, 480, 640]	16
ReLU-397	[-1, 8, 480, 640]	0
ConvNormRelu-398	[-1, 8, 480, 640]	0
Conv2d-399	[-1, 50, 480, 640]	450

PixelShuffle

\mathcal{P} → periodic shuffling operator

LR → low resolution

SR → super resolution

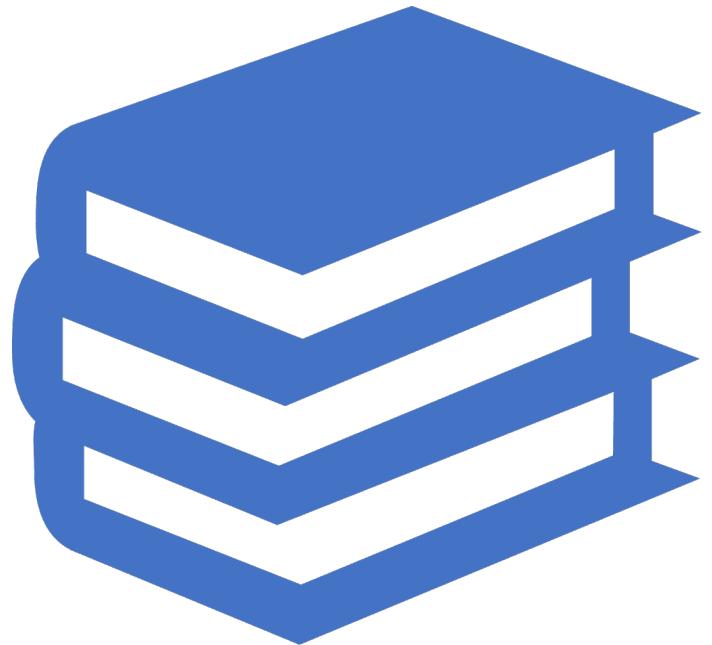


[PS1]

- Performs a “sub-pixel convolution” ($\text{stride} = 1/r$) by re-arranging its input pixels

- ⊕ efficient
- ⊕ no learnable params
- ⊕ less artefacts (than ConvTranspose)

```
class ConvPxShuffle(nn.Module):  
    def __init__(self, in_channels, out_channels, in_channels_up=None):  
        super().__init__()  
        if in_channels_up is not None:  
            in_channels = in_channels + in_channels_up  
  
        self.upscale = nn.PixelShuffle(2)  
        self.conv1 = ConvNormRelu(in_channels, out_channels)  
        self.conv2 = ConvNormRelu(out_channels, out_channels)  
  
    def forward(self, x, x_cat):  
        x = self.upscale(x)  
        x = torch.cat([x, x_cat], dim=1)  
        x = self.conv1(x)  
        x = self.conv2(x)
```



Training

- Batch-size: 6
- Initial Learning Rate: 1e-4
 - With `ReduceLROnPlateau`
- Epochs: 50
 - With `EarlyStopping`
- Input-size: (4, 240, 320)
- Cost function: `CrossEntropyLoss`

Examples

Segmentation of a random sample

Predicted segmentation



Ground truth



Examples

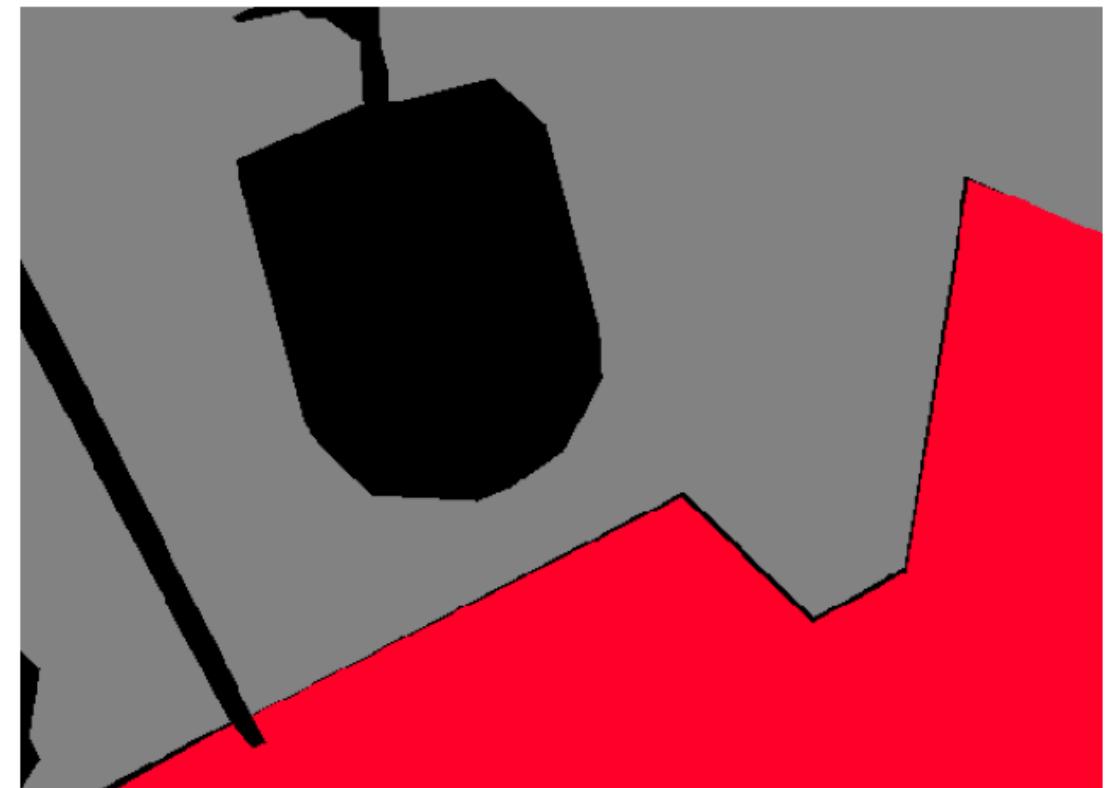


Segmentation of a random sample

Predicted segmentation

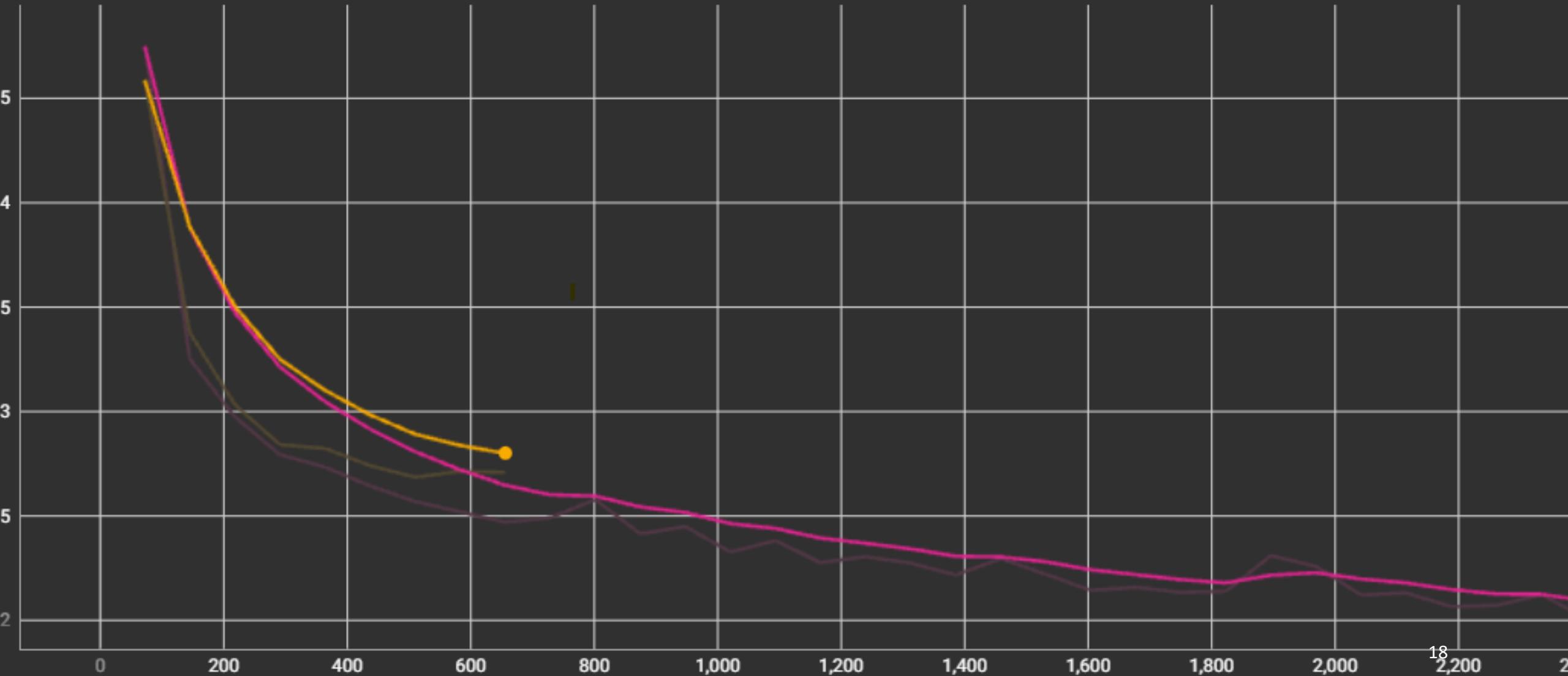


Ground truth



Run	Smoothed Value	Step	Time	Relative
version_0	6.824	6.824	3'	5/12/23, 8:42 PM 0
version_2	2.084	2.055	2.481	5/12/23, 7:56 PM 2
version_3	2.801	2.711	656	5/12/23, 8:34 PM 32.13 min

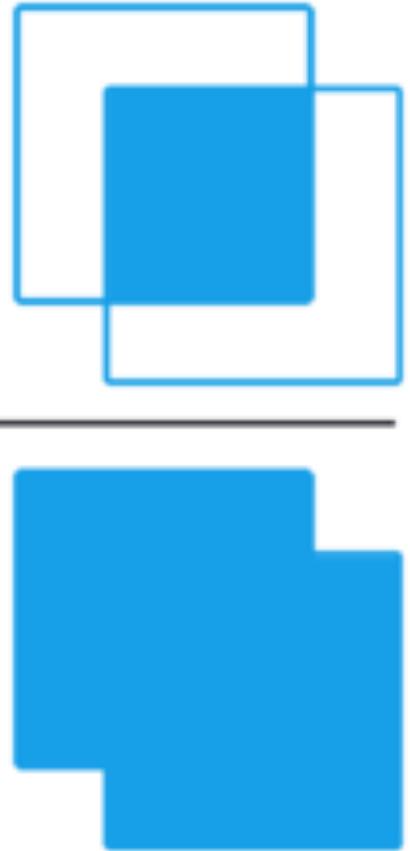
Performance – Loss values



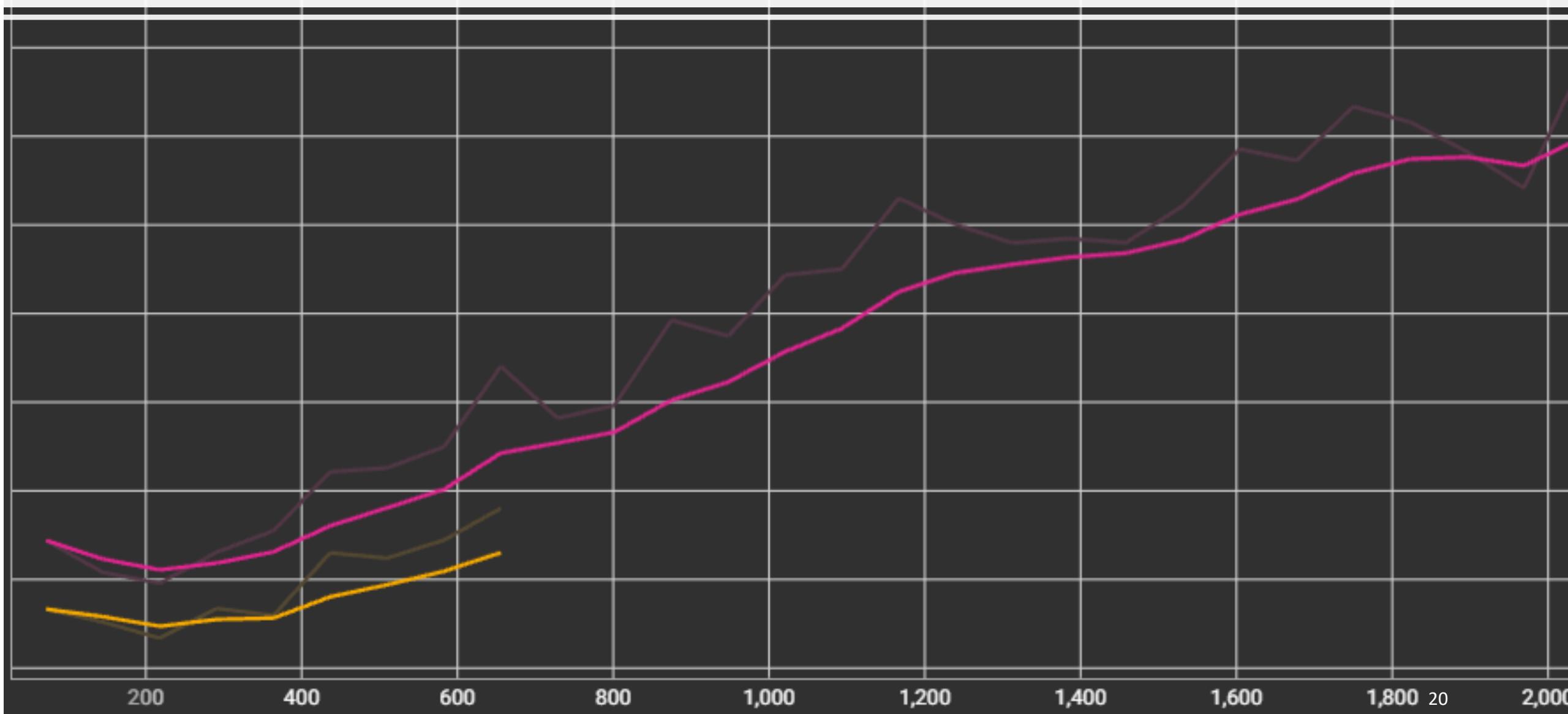
Performance - IoU

-
- Metric for segmentation & object detection tasks
 - Can be used as regression loss
 - Plateau for non-overlapping areas
 - Addressed by GIoU

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



Performance – IoU values





PyTorch Lightning

[#logging](#)

[#easyScaling](#)

[#checkpoints](#)

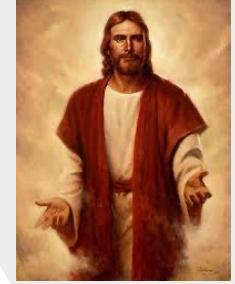
[#metrics](#)

[#earlyStopping](#)

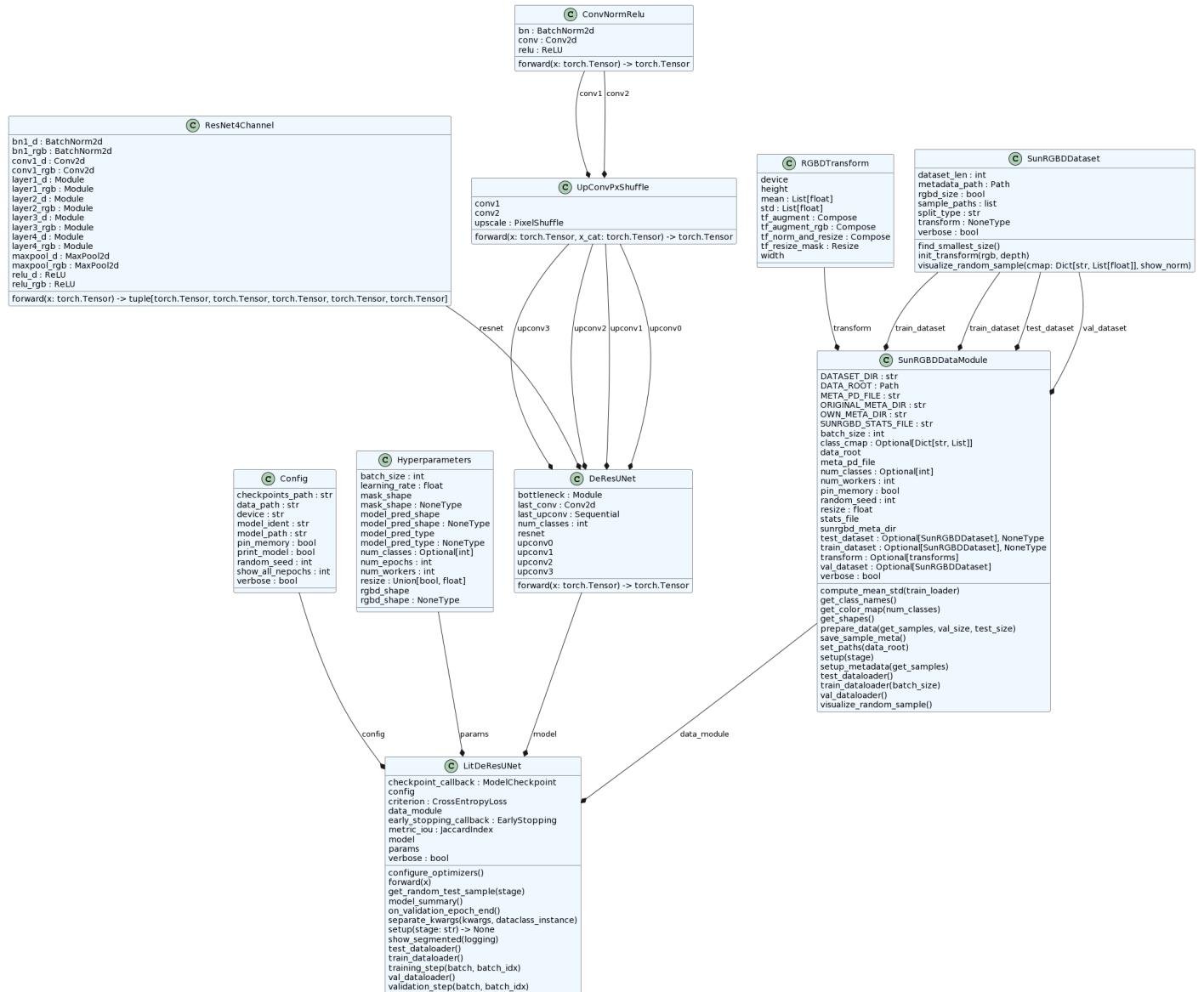
[#visualization](#)

[#noBoilerplate](#)

- [PyTorch Lightning](#)
- [PyTorch Lightning Tutorials by Aladdin Persson](#)



Full Project UML



- [S1] https://pbs.twimg.com/media/Fv1EjNuakAE44z_.jpg:large
- [N1] <https://lmb.informatik.uni-freiburg.de/people/ronneber/u-net/>
- [N2] https://www.researchgate.net/publication/346089961_Assessing_out-of-domain_generalization_for_robust_building_damage_detection
- [PS1] <https://arxiv.org/pdf/1609.05158v2.pdf>
- [BN1] <https://www.linkedin.com/pulse/ways-improve-your-deep-learning-model-batch-adam-albuquerque-lima/>
- [A1] <https://towardsdatascience.com/getting-started-with-albumentation-winning-deep-learning-image-augmentation-technique-in-pytorch-47aaba0ee3f8>

