

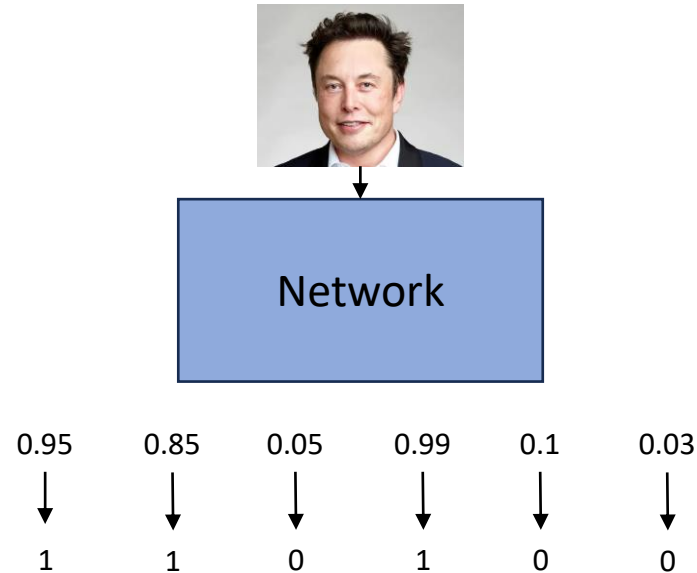
# Neural Field Classifiers via Target Encoding and Classification Loss

Anonymous authors

Efficient Learning Lab@POSTECH

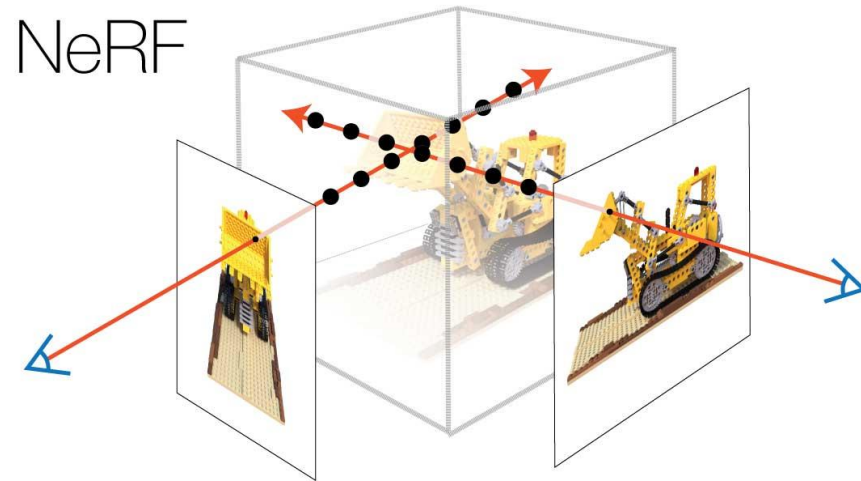
Junwon Seo

# Recap: LEARNING LABEL ENCODINGS FOR DEEP REGRESSION



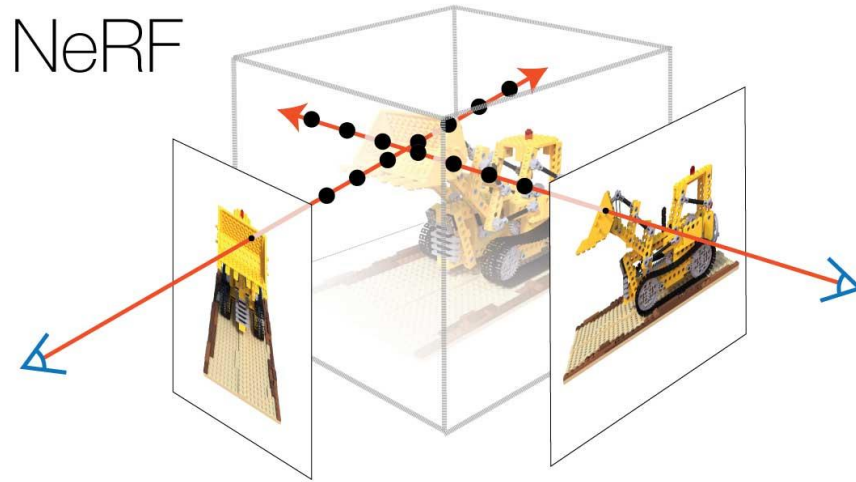
- Example
  - The age of 52 would be 110100 using binary conversion as the encoder.  
(Training phase)
  - 110100 would be converted to real-value prediction using a decoding function (decimal convert)  
(Inference phase)

# Neural Field



- Neural field methods emerge as promising methods for parameterizing **a field**
  - that has a **target value** for **each point in space and time**
- Try to predict some coordinate-based continuous target values
  - RGB for Neural Radiance Field

# NeRF Basics



$$\hat{\mathbf{C}}(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(t) \mathbf{c}(t) dt$$

$$\text{where } T(t) = \exp\left(-\int_{t_n}^t \sigma(s) ds\right)$$

$$L(\Theta) = \frac{1}{\|\mathcal{R}\|} \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|_2^2$$

- NeRF **regresses** from a single 5D representation  $(x, y, z, \theta, \phi)$  to  $(\sigma, c)$ 
  - 3D coordinates  $x = (x, y, z)$  / viewing directions  $d = (\theta, \phi)$
  - Volume density  $\sigma$  and view-dependent color  $c = (r, g, b)$

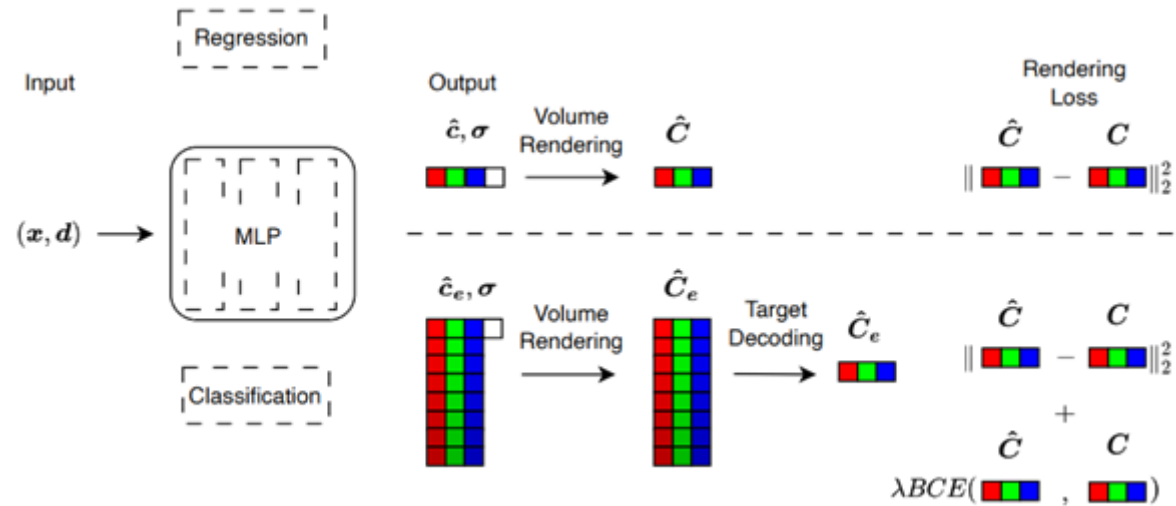
# Motivation

Are **regression formulations** really better than **classification formulations** for neural field methods?

# Motivation

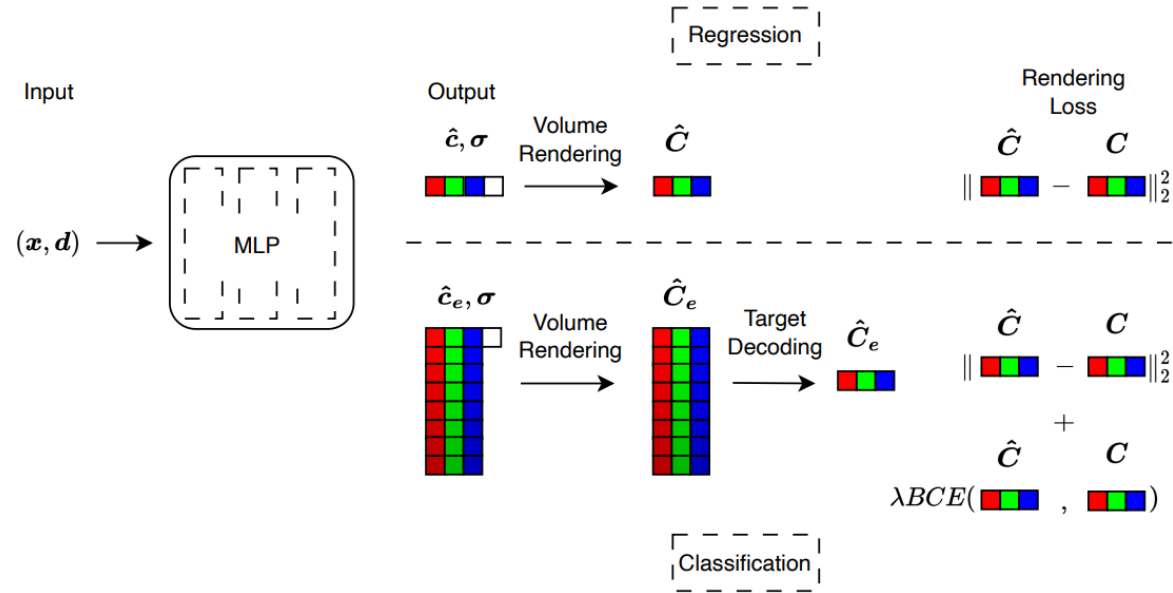
- People naturally formulate neural fields as regression models
  - Because the targets are continuous values
- The authors claimed that there exist overlooked pitfalls
  - Each data point has its own ground-truth label in classical supervised learning methods
  - NeRF output N points' predictions per pixel.
  - Supervision signals for NeRF are obviously very weak and insufficient

# Methodology – Target Encoding and Decoding



- Naïve target encoding rule is one-hot encoding
  - Directly use  $y$  as the class label (i.e. 256-class classification)
    - The number of logits increases to 768 from (computational cost  $\uparrow$ )
    - Ignore the relevant information carried by the classes
      - Suppose ground-truth label of a sample is 0
      - A model predicts 1 and another model predicts 255
        - Their loss will be equally high

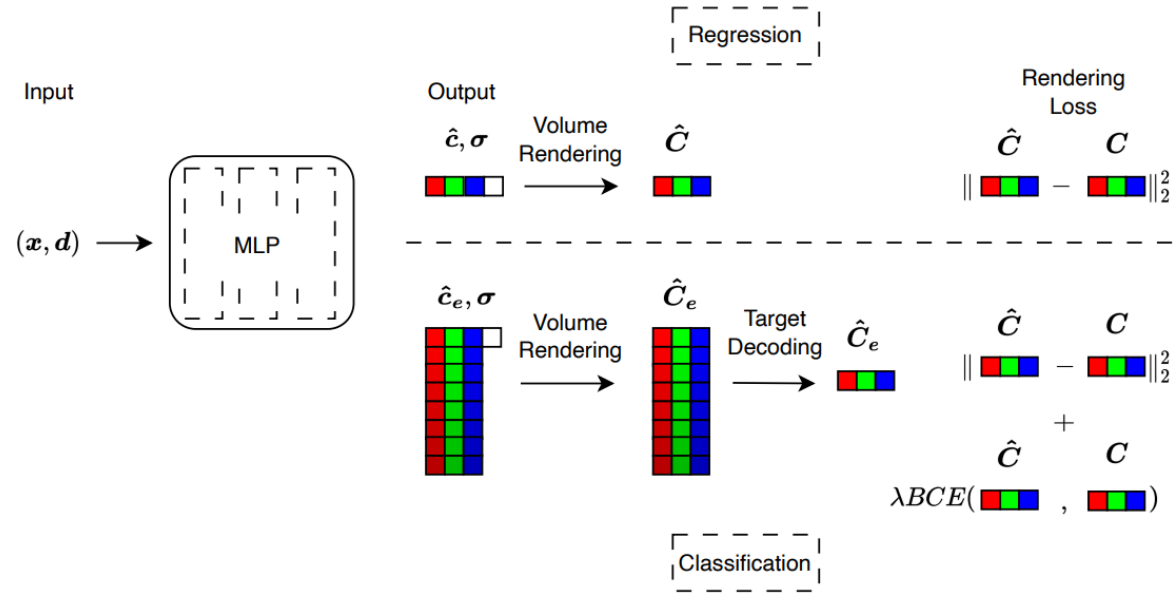
# Methodology – Target Encoding and Decoding



- Fortunately, the authors discover that the classical binary-number system can work well
  - For example,  $y = \text{BinaryEncod}(203) = [1, 1, 0, 0, 1, 0, 1, 1]$



# Methodology – Classification Loss

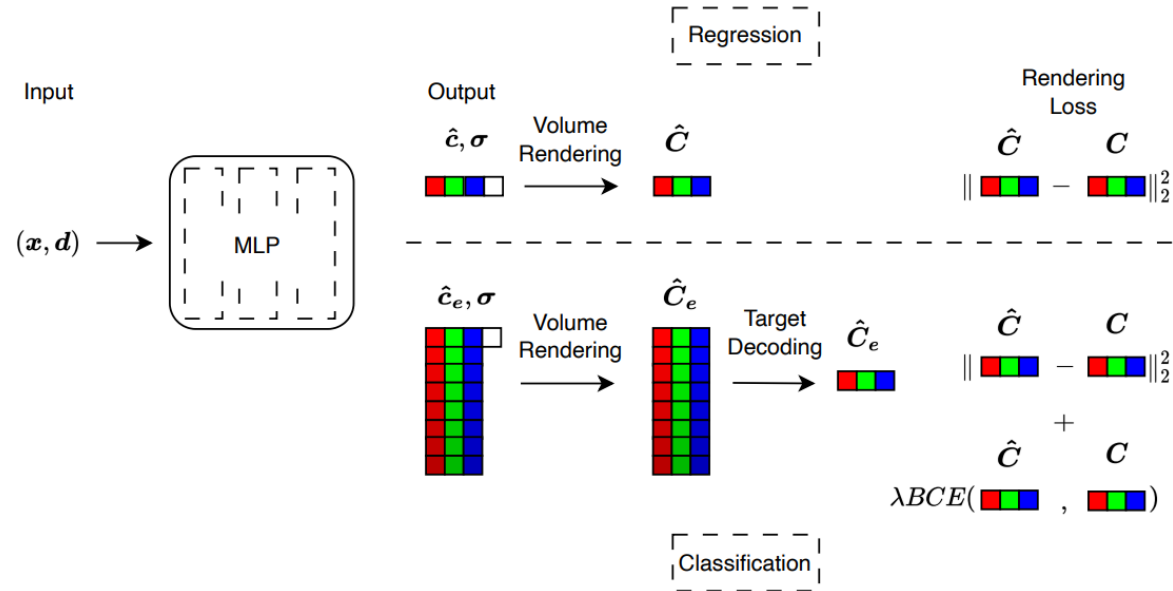


- **Classification Loss**

- $2^{j-1}$  assigns a higher weight to the class with a higher place value.

$$l_b(\hat{y}, y) = \frac{1}{255} \sum_{j=1}^8 2^{j-1} \text{BCE}(\hat{y}^{(j)}, y^{(j)})$$

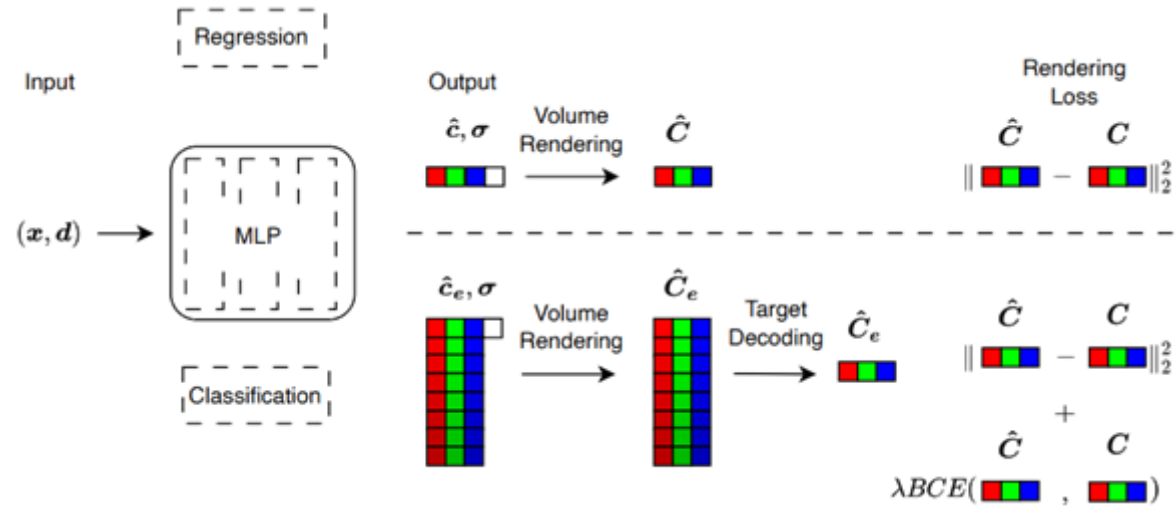
# Methodology – Classification Loss



- Alternative choice is..
  - Decode the predict probability  $\hat{y}$  back into a continuous value

$$l_c(\hat{C}, C) = BCE(\hat{C}, C) = BCE\left(\frac{1}{255} \text{BinaryDecod}(\hat{y}), C\right)$$

# Methodology – Classification Loss



- Predicted probability  $\hat{C}$  (or  $\hat{y}$ ) of NFC does not strictly lie in  $(0, 1)$

$$l_c(\hat{C}, C) = \text{BCE}(\min(\hat{C}, 1 - \epsilon), C)$$

$$L_{\text{NFC}}(\hat{C}, C) = \|\hat{C} - C\|_2^2 + \lambda \text{BCE}(\min(\hat{C}, 1 - \epsilon), C)$$

# Experiments Settings

- Neural Field backbones
  - DVGO, vanilla NeRF, D-NeRF, NeuS
- Dataset
  - Static: Replica Dataset, Tanks and Temples Advanced
  - Dynamic: Lego and Hook
  - Challenging: Tanks and Temples Advanced
- Hyperparameters
  - Keep all hyperparameters same for them

# Results

- Static Scene



# Results

- Static Scene

Table 1: Quantitative results of DVGO on Replica Dataset.

Scene	Mode	PSNR( $\uparrow$ )	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )
Scene 1	Regression	13.03	0.508	0.726
	Classification	<b>34.63</b>	<b>0.934</b>	<b>0.0582</b>
Scene 2	Regression	14.81	0.654	0.640
	Classification	<b>33.82</b>	<b>0.942</b>	<b>0.0660</b>
Scene 3	Regression	15.66	0.661	0.634
	Classification	<b>34.04</b>	<b>0.965</b>	<b>0.0451</b>
Scene 4	Regression	18.17	0.696	0.546
	Classification	<b>36.52</b>	<b>0.977</b>	<b>0.0311</b>
Scene 5	Regression	15.17	0.640	0.504
	Classification	<b>35.93</b>	<b>0.974</b>	<b>0.0576</b>
Scene 6	Regression	21.33	0.854	0.254
	Classification	<b>29.75</b>	<b>0.941</b>	<b>0.0994</b>
Scene 7	Regression	22.54	0.865	0.231
	Classification	<b>34.77</b>	<b>0.966</b>	<b>0.0432</b>
Scene 8	Regression	15.89	0.724	0.519
	Classification	<b>33.40</b>	<b>0.952</b>	<b>0.0775</b>
Mean	Regression	17.08	0.700	0.507
	Classification	<b>34.11</b>	<b>0.956</b>	<b>0.0598</b>

Table 2: Quantitative results of DVGO, (vanilla) NeRF, NeuS on T&T Dataset. The mean PSNR, SSIM, LPIPS are computed over four scenes of T&T.

Model	Mode	PSNR( $\uparrow$ )	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )
DVGO	Regression	22.41	0.776	0.236
	Classification	<b>23.18</b>	<b>0.810</b>	<b>0.178</b>
NeRF	Regression	22.16	0.679	0.382
	Classification	<b>22.68</b>	<b>0.716</b>	<b>0.315</b>
NeuS	Regression	19.97	0.620	0.413
	Classification	<b>21.67</b>	<b>0.679</b>	<b>0.317</b>

# Results

- Dynamic Scene
  - requires the capability of neural fields to model time domain.

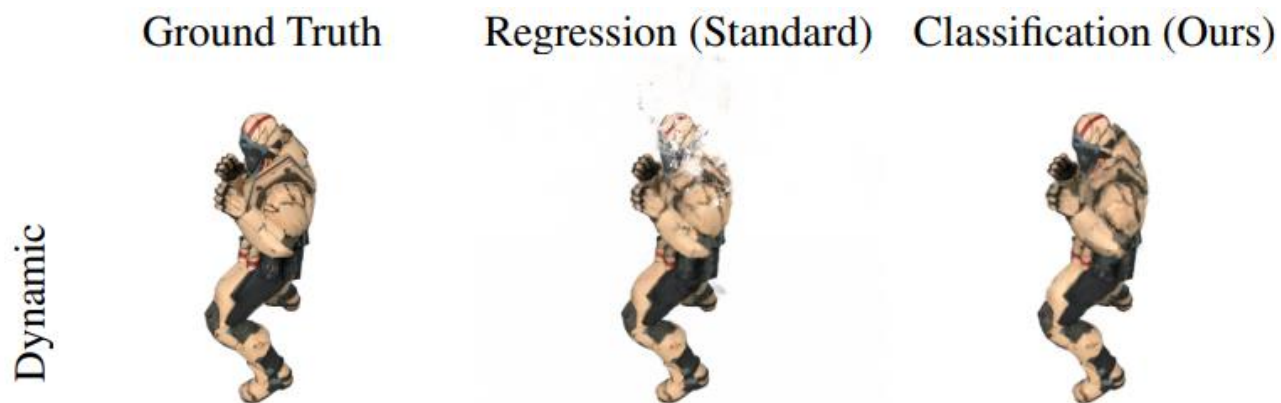


Table 3: Quantitative results of D-NeRF on dynamic scenes, LEGO and Hook.

Scene	Mode	PSNR(↑)	SSIM(↑)	LPIPS(↓)
Lego	Regression	21.64	0.839	0.165
	Classification	<b>23.11</b>	<b>0.886</b>	<b>0.121</b>
Hook	Regression	29.25	0.965	0.118
	Classification	<b>29.45</b>	<b>0.967</b>	<b>0.0392</b>

# Results



- Challenging Scenes
  - Sparse: Fewer data for an object
  - Corrupted: Digital images often contain Gaussian noise



Table 4: Quantitative results of neural rendering with sparse training images.

Data Size	Mode	PSNR( $\uparrow$ )	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )
20%	Regression	14.87	0.530	0.580
	Classification	<b>19.38</b>	<b>0.629</b>	<b>0.395</b>
40%	Regression	18.76	0.637	0.426
	Classification	<b>21.73</b>	<b>0.711</b>	<b>0.340</b>
60%	Regression	21.02	0.682	0.394
	Classification	<b>22.27</b>	<b>0.728</b>	<b>0.329</b>
80%	Regression	21.72	0.698	0.386
	Classification	<b>22.46</b>	<b>0.734</b>	<b>0.322</b>



# Results



- Challenging Scenes
  - Sparse: Fewer data for an object
  - Corrupted: Digital images often contain Gaussian noise



Table 5: Quantitative results of neural rendering with corrupted training images.

Noise Scale	Mode	PSNR( $\uparrow$ )	SSIM( $\uparrow$ )	LPIPS( $\downarrow$ )
0.2	Regression	21.97	0.694	0.406
	Classification	<b>22.33</b>	<b>0.719</b>	<b>0.353</b>
0.4	Regression	21.33	0.663	0.451
	Classification	<b>22.08</b>	<b>0.692</b>	<b>0.391</b>
0.6	Regression	19.67	0.615	0.512
	Classification	<b>21.45</b>	<b>0.662</b>	<b>0.429</b>

# Ablation Study

$$L_{\text{NFC}}(\hat{C}, C) = \|\hat{C} - C\|_2^2 + \lambda \text{BCE}(\min(\hat{C}, 1 - \epsilon), C)$$

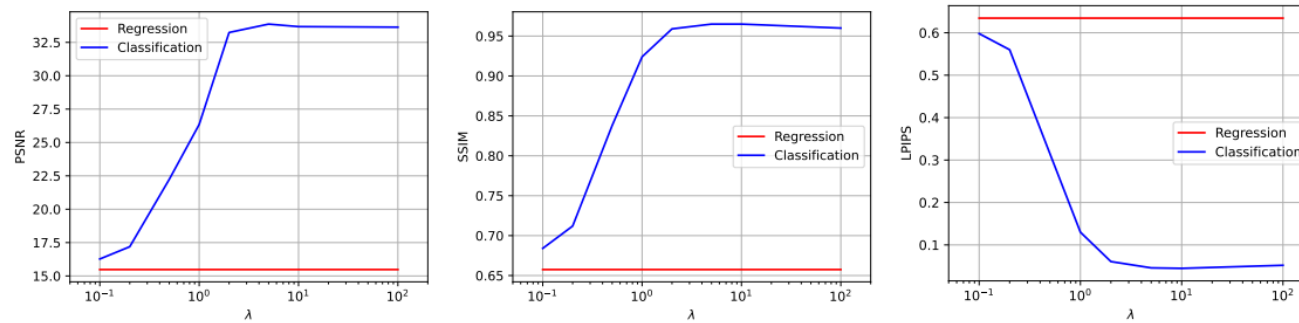


Figure 6: The curves of PSNR, SSIM, and LPIPS with respect to the hyperparameter  $\lambda$ . NFC is robust to a wide range of  $\lambda$ . Model: DVGO. Dataset: Replica Scene 3.

- Robustness to the hyperparameter
  - The quantitative results shows that a wide value range of  $\lambda$  can enhance the performance.

# Limitation

- Target Encoding
  - In channel-wise classification loss
    - Is distance(relevant information) of binary-number system maintained?

$$l_c(\hat{C}, C) = \text{BCE}(\hat{C}, C) = \text{BCE}\left(\frac{1}{255} \text{BinaryDecod}(\hat{y}), C\right)$$

1111  
↕  
0001  
↕  
1000

- Ignore the properties of the Neural Field.
  - Continuous!

Q & A