

# Visualization of neuronal activations under class-specific perturbations

Prepared by : Sinem Bilge Güler  
Supervised by: Annika Mikliss

# Tasks

## Task 1

Systematically generate class-specific transformations

---

## Task 2

Systematically investigate how the output changes

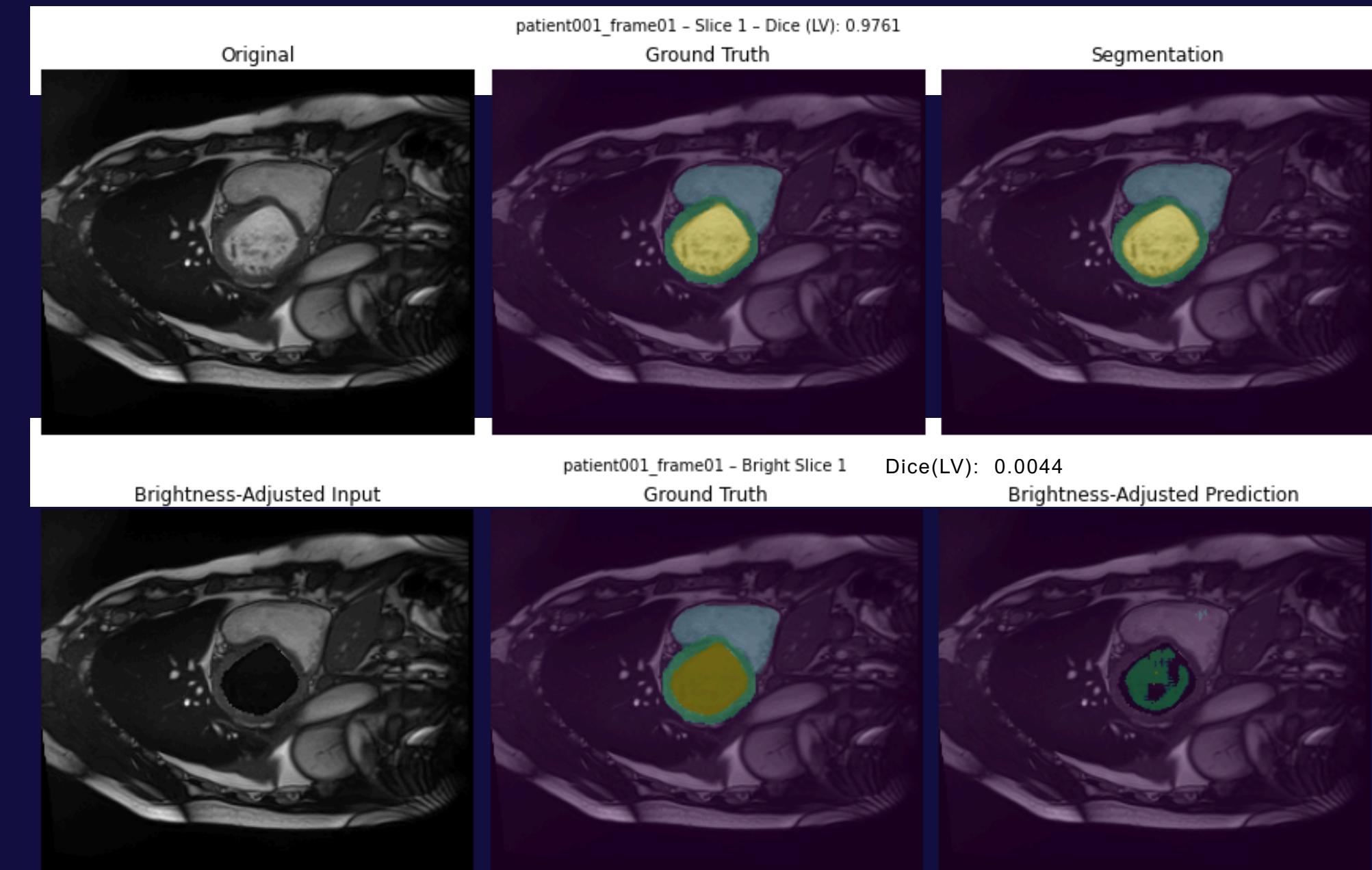
---

## Task 3

Investigate and visualize which and how important feature maps

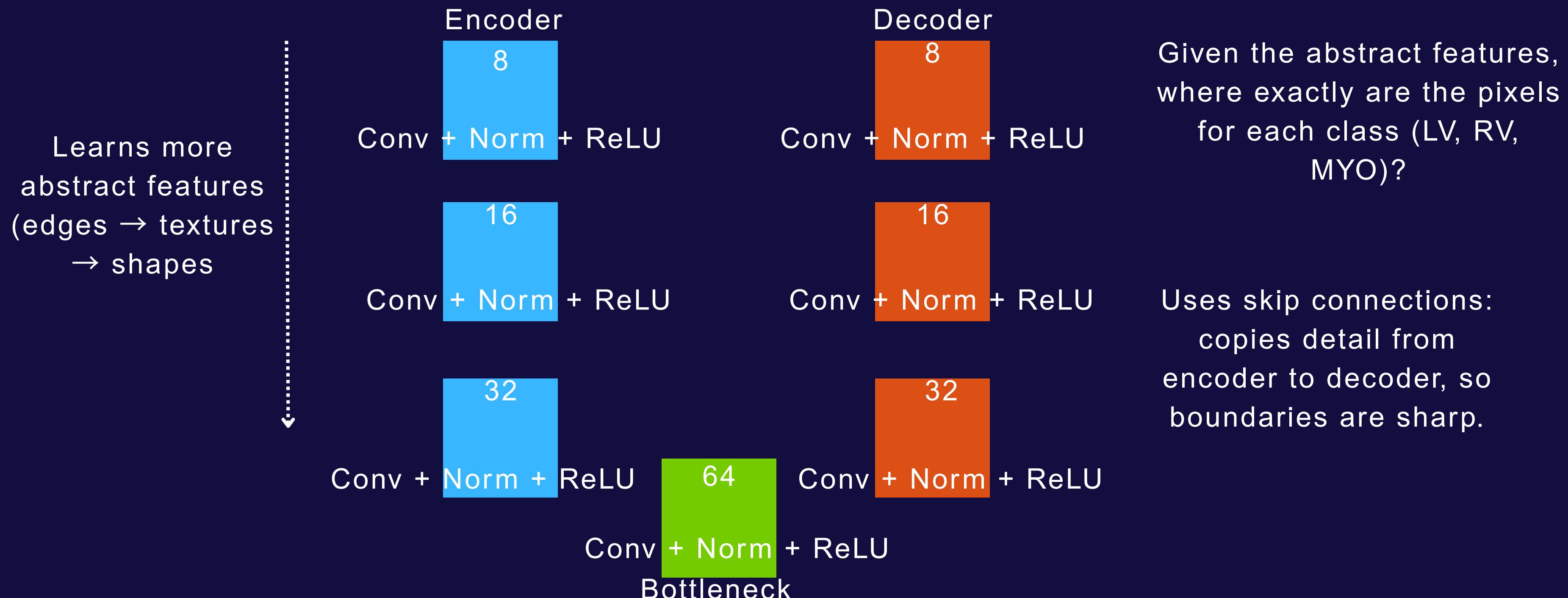
- DL models are often referred to as black box models because it is not clear which strategies the models learn to solve a task.

- Problem: How do class-specific transformations influence the activations inside of a DL segmentation model?



# Unet Model

model = UNetHeart\_8\_4\_4()



# Data Processing Steps



## 10 times 3D Images

10 2D images for each 10 patient images

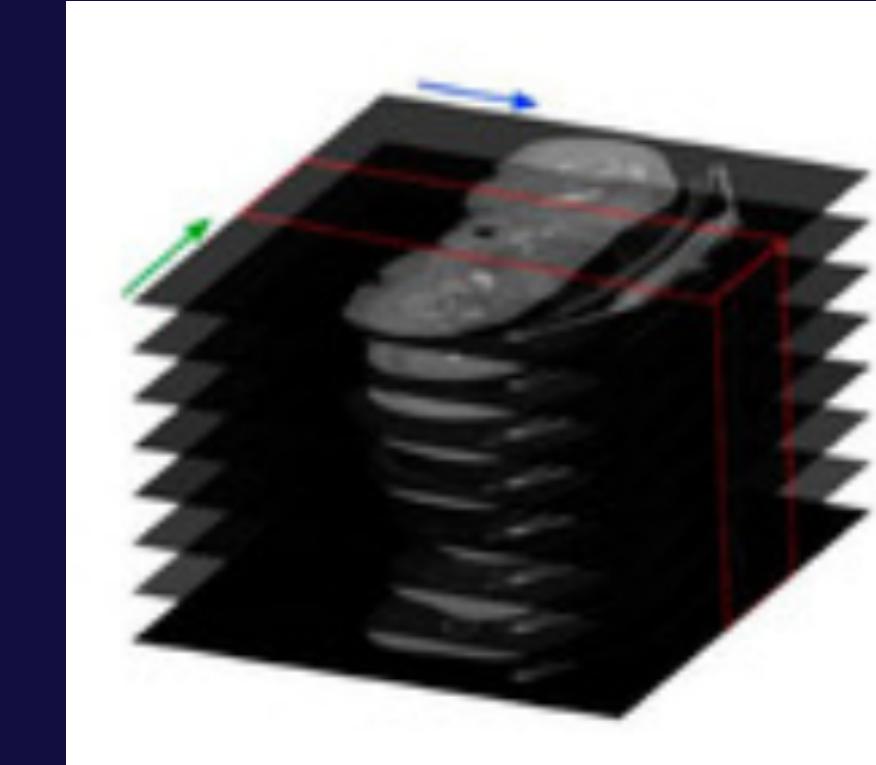
4-64 channels/feature maps

Filter out images without the specified class

LV:98

RV:79

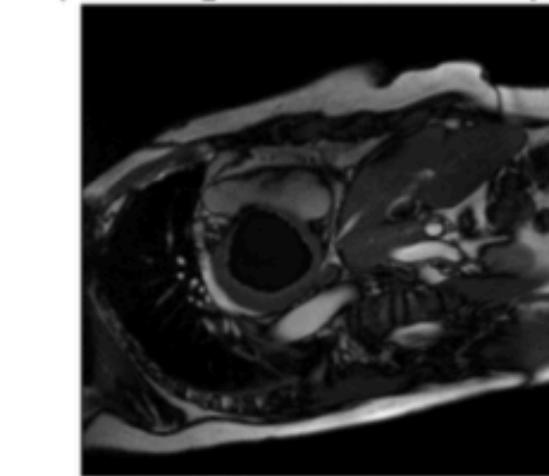
MYO:98



(Bertolini, Rossoni, & Colombo, 2021)

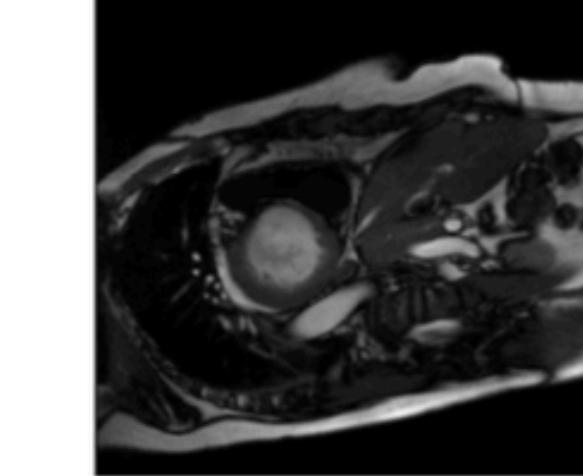
## Brightness Adjustments Dice Score Calculations

patient004\_frame15 - Slice 0 (LV Adjusted)



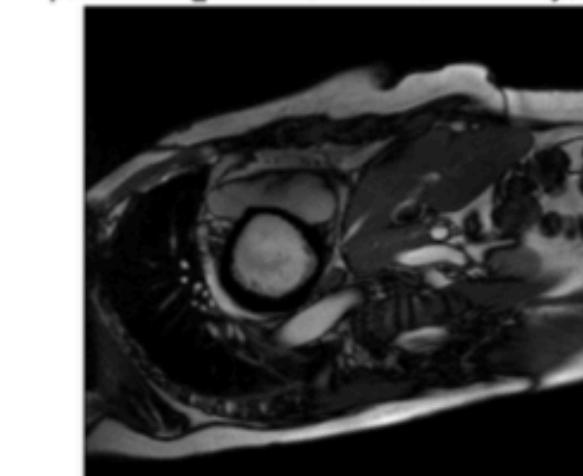
(a) LV Brightness Adjusted

patient004\_frame15 - Slice 0 (RV Adjusted)



(b) RV Brightness Adjusted

patient004\_frame15 - Slice 0 (MYO Adjusted)



(c) MYO Brightness Adjusted

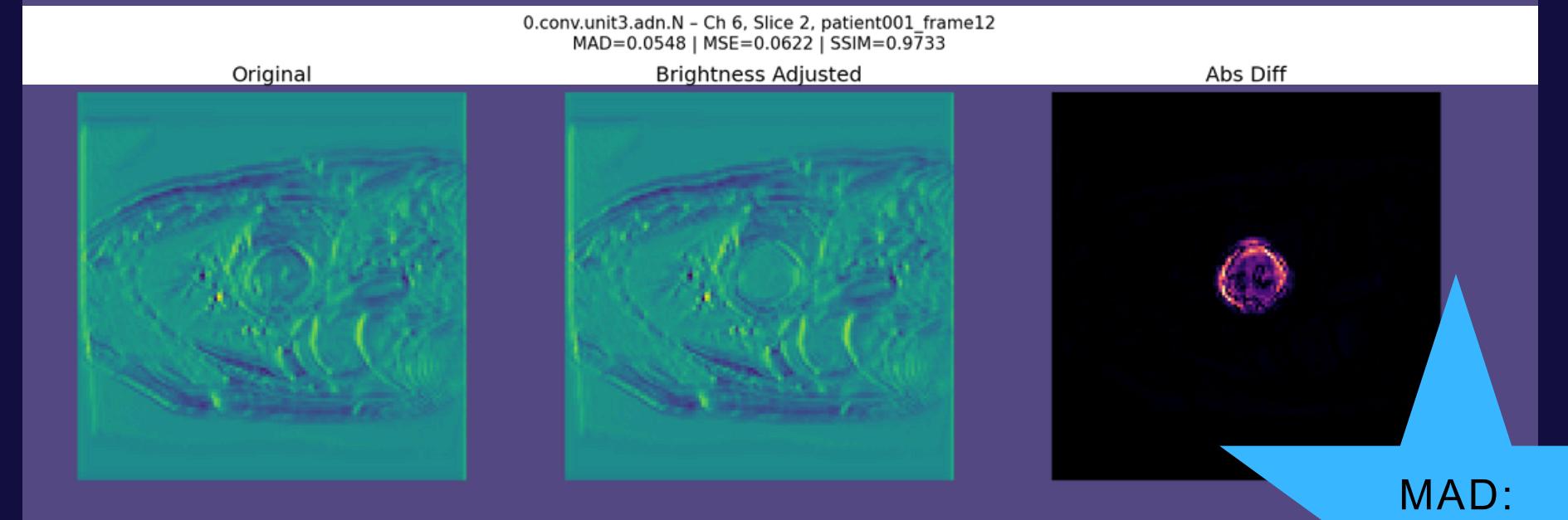
# Hooks on adn.N why?

MAD

21 Hooks to different Layers

L1	0.conv.unit0.adn.N
L2	0.conv.unit1.adn.N
L3	0.conv.unit2.adn.N
L4	0.conv.unit3.adn.N
L5	1.submodule.0.conv.unit0.adn.N
L6	1.submodule.0.conv.unit1.adn.N
L7	1.submodule.0.conv.unit2.adn.N
L8	1.submodule.0.conv.unit3.adn.N
L9	1.submodule.1.submodule.0.conv.unit0.adn.N
L10	1.submodule.1.submodule.0.conv.unit1.adn.N
L11	1.submodule.1.submodule.0.conv.unit2.adn.N
L12	1.submodule.1.submodule.0.conv.unit3.adn.N
L13	1.submodule.1.submodule.1.submodule.0.conv.unit0.adn.N
L14	1.submodule.1.submodule.1.submodule.0.conv.unit1.adn.N
L15	1.submodule.1.submodule.1.submodule.0.conv.unit2.adn.N
L16	1.submodule.1.submodule.1.submodule.0.conv.unit3.adn.N
L17	1.submodule.1.submodule.2.0.adn.N
L18	1.submodule.1.submodule.2.1.conv.unit0.adn.N
L19	1.submodule.2.0.adn.N
L20	1.submodule.2.1.conv.unit0.adn.N
L21	2.0.adn.N

◆ MAD – Mean Absolute Difference  
Measures the average absolute pixel-wise difference between two images.





# Case Categorization of LV

Empty Prediction – <35 total foreground pixels

Only LV – LV present, MYO & RV absent

Only RV – RV >100px, LV <10px, MYO ≤20px

LV → MYO Misclassification – ≥10% LV misclassified as MYO

**Extreme LV → MYO Misclassification** – >50% LV misclassified as MYO

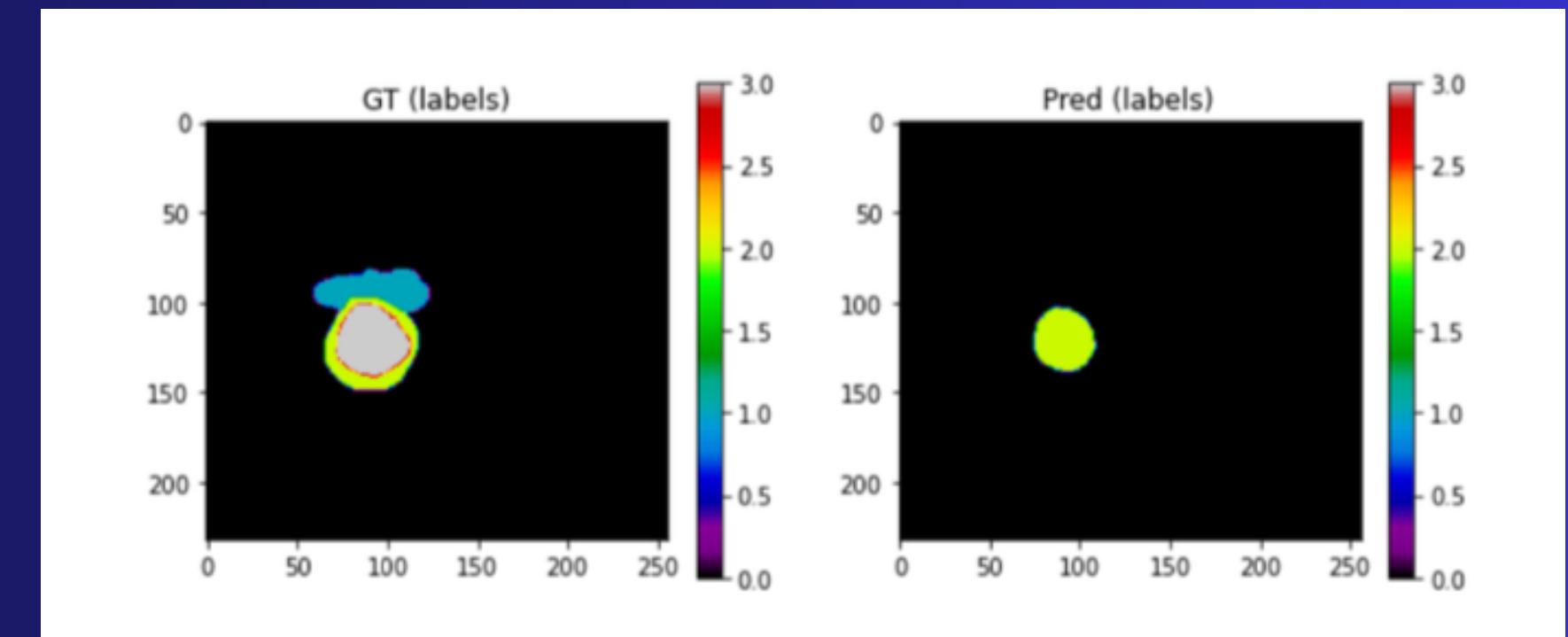


**LV + MYO** – Both inside LV mask, ≥70% MYO within LV



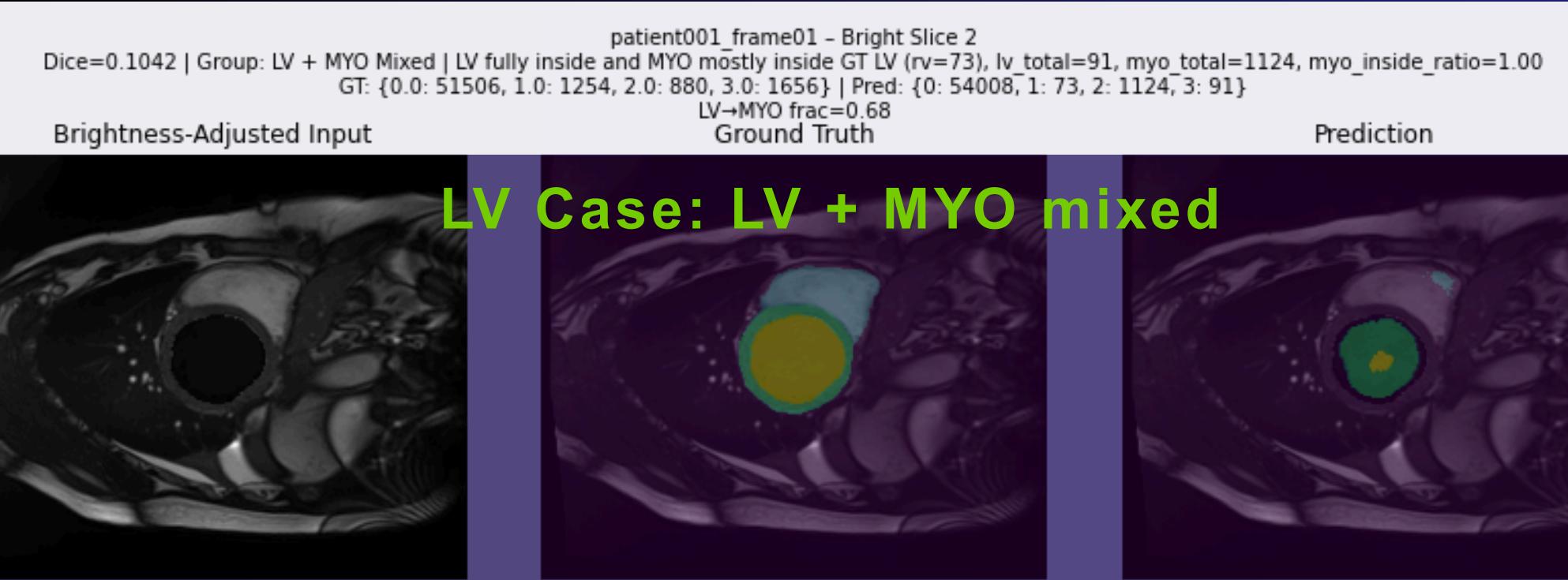
Every Class – LV, RV, and MYO all appear inside LV mask

Error – fallback if thresholds not met

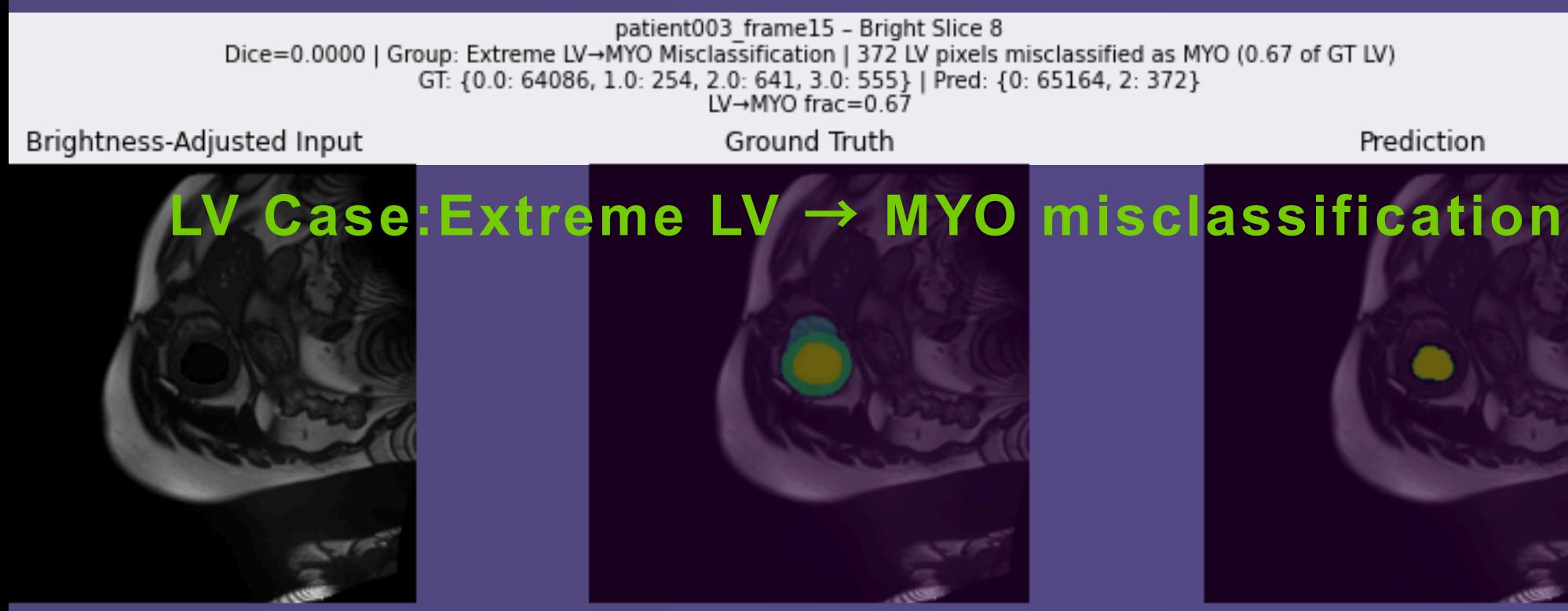


# LV Selected Cases

# LV Not Selected Cases

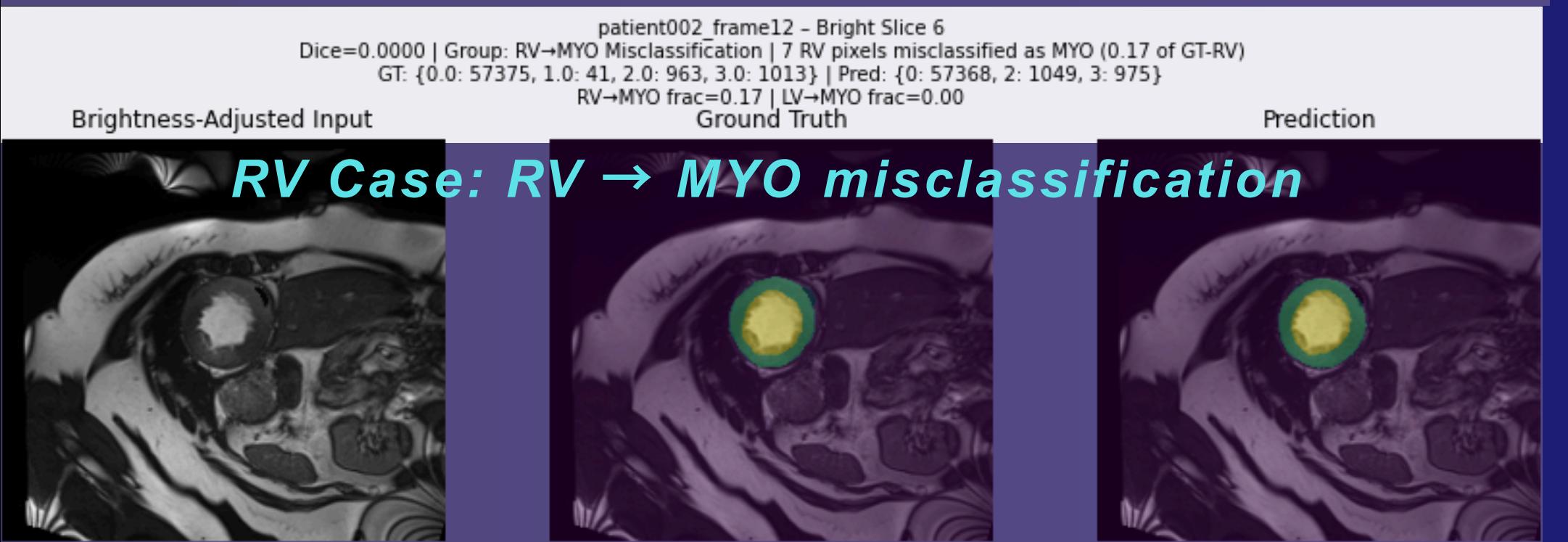
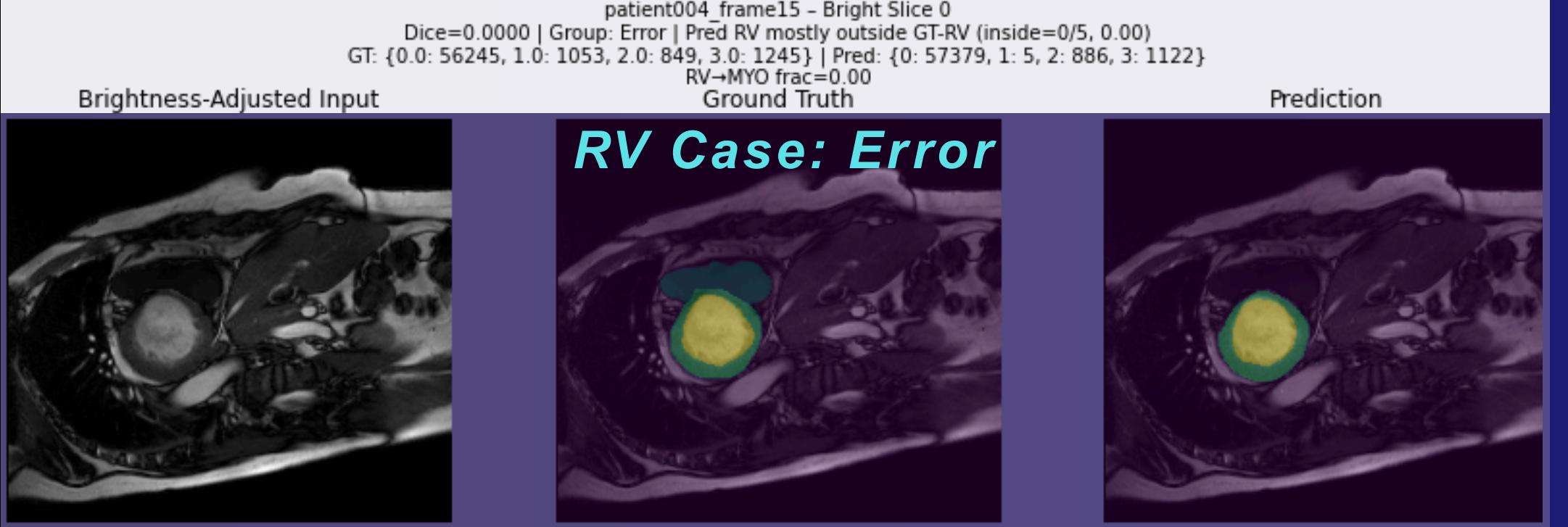


**LV Case: Empty Prediction**



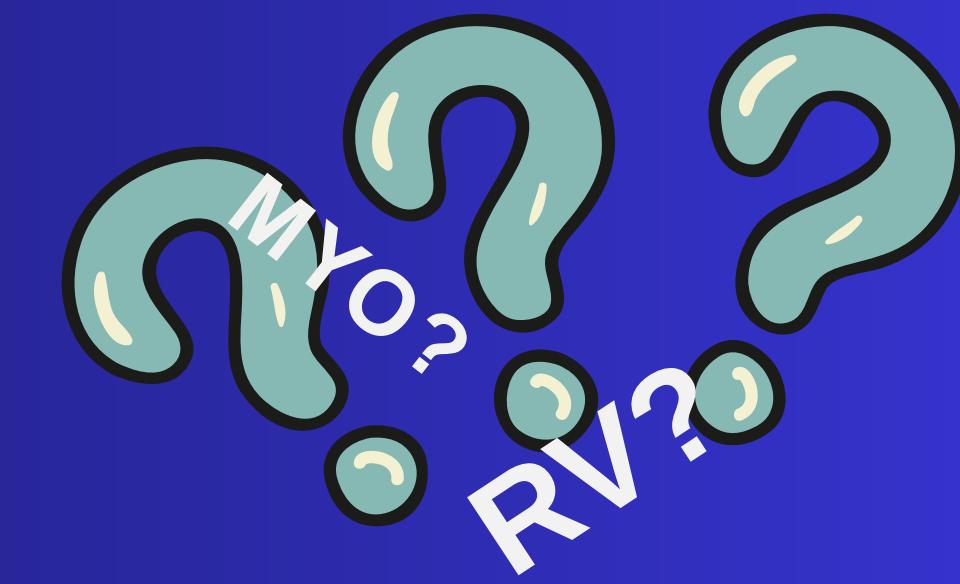
**LV Case: Error**

# Where are the Cases for RV and MYO?

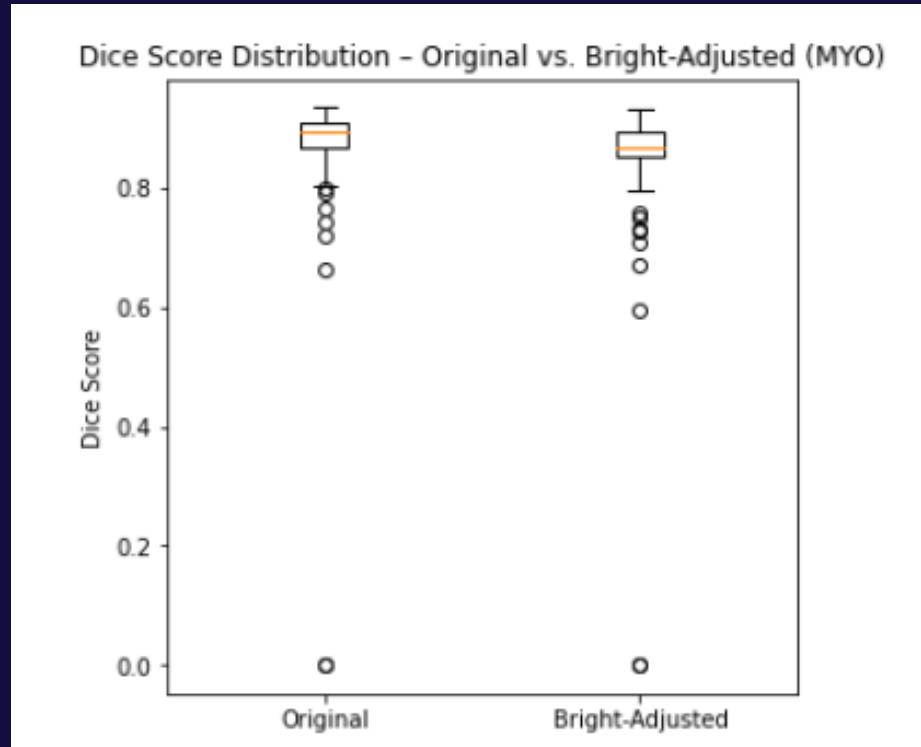
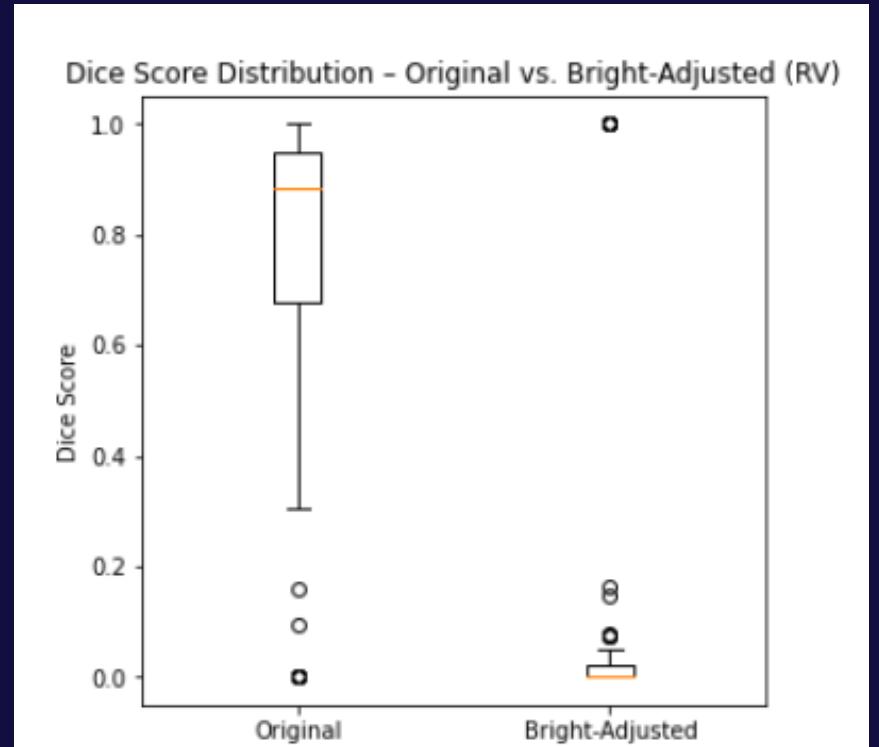
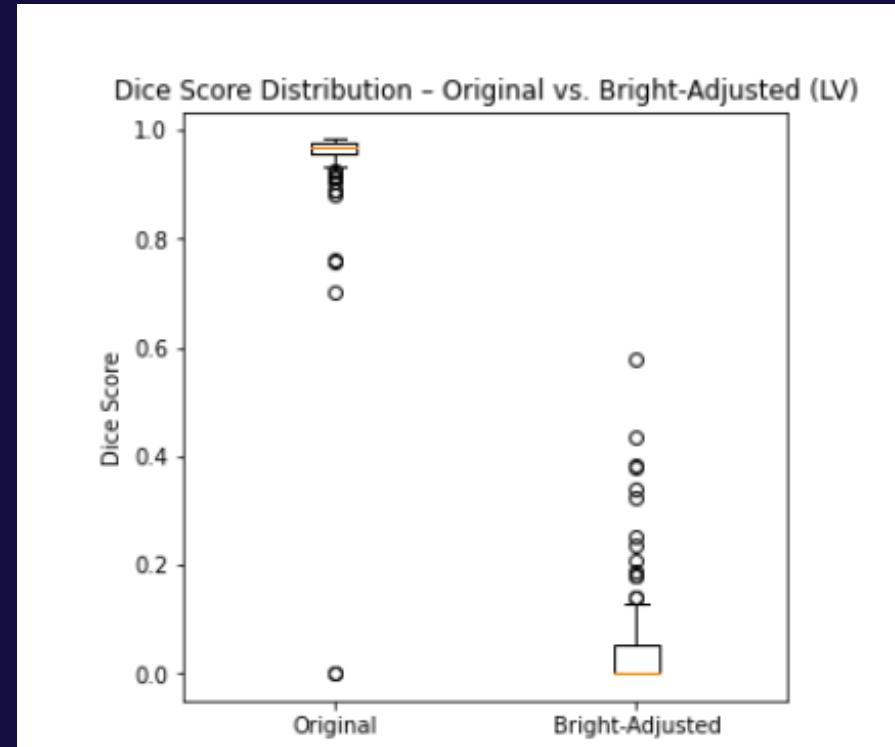


7 pixels misclassified as MYO but it looks almost same as the case Error

We see every classes in MYO segmentation results



# Dice Score Distribution - Original vs Adjusted



## LV Dice Score Change

- Original: 0.94 → Adjusted: 0.05
- Drop  $\approx -0.88$
- Almost complete segmentation failure

## RV Dice Score Change

- Original: 0.73 → Adjusted: 0.12
- Drop  $\approx -0.61$
- Strong disruption, but not total failure

## MYO Dice Score Change

- Original: 0.86 → Adjusted: 0.84
- Drop  $\approx -0.02$
- Largely robust, only few outlier failures

LV and RV are highly sensitive, MYO is comparatively robust against darkening.

# Calculation of the Top Feature Maps

## Goal

Identify model activations sensitive to brightness adjustments.

## Class-Specific Analysis

- LV: compare with case slices (e.g., LV $\rightarrow$ MYO misclassification)
- MYO & RV: compare frequencies across all three classes

## Method

- Under each layer, for each patients and their slices, select Top-3 MAD channels
  - Top-1 = too noisy
  - Top-5+ = too diluted
  - Top-3 = balance
- Count how often each channel appears overall  $\rightarrow$  Top-7 most frequent per layer

## Insights

- General distribution vs case-specific distribution
- Shows which channels are:
  - Consistently responsible for class detection
  - Most sensitive to brightness adjustments

```
for each layer:  
    for each patient, each slice:  
        compute MAD(channel) for all channels  
        select Top-3 channels with highest MAD  
        record these channels  
  
for each layer:  
    count frequency of channels across all patients & slices  
    select Top-7 most frequent channels  
  
for each case (e.g., LV+MYO, Extreme LV $\rightarrow$ MYO):  
    repeat Top-3 selection restricted to case slices  
    compute Top-7 frequencies within this subset
```

# Evaluation of the Top Channel Maps

## LV Frequency Map

L7

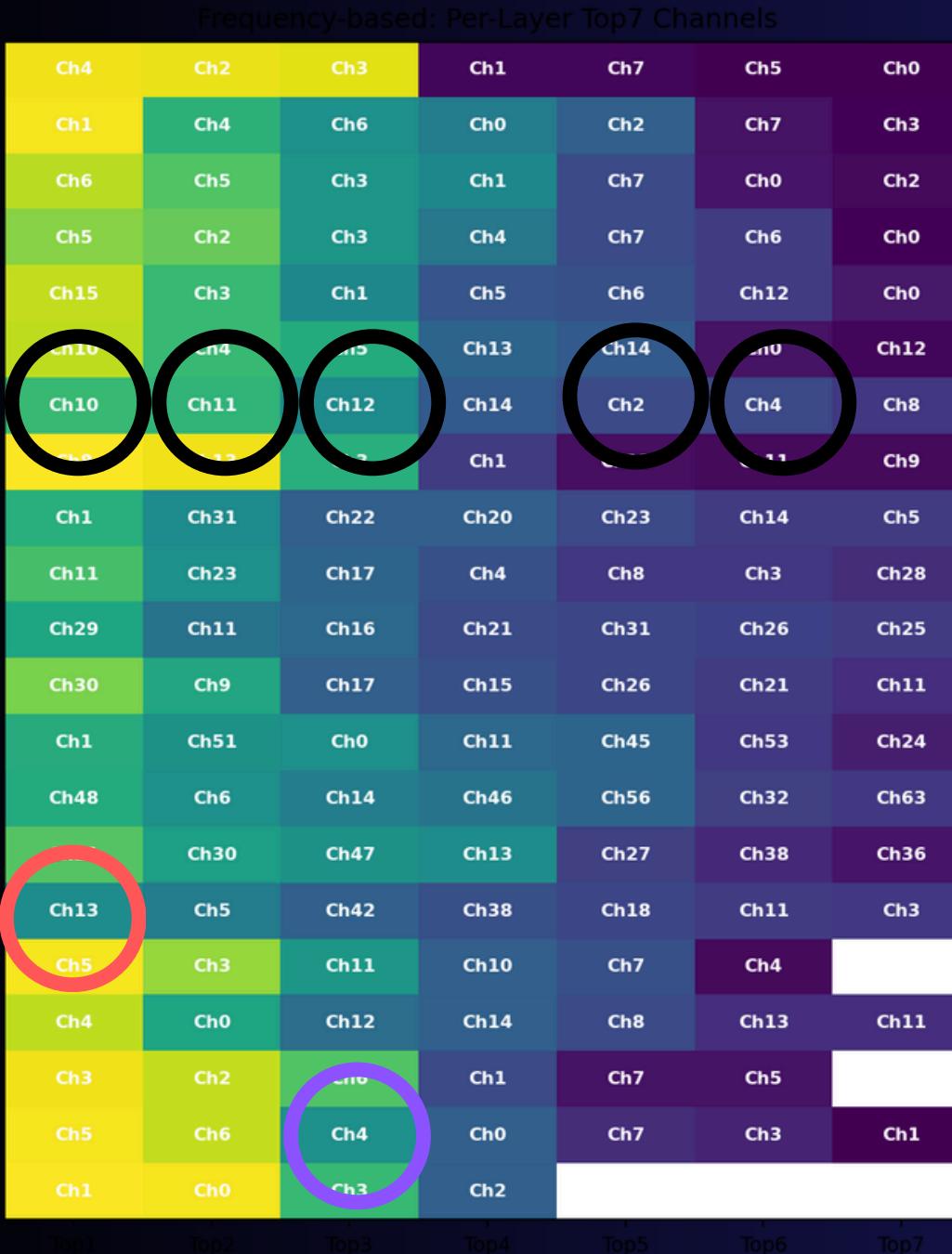
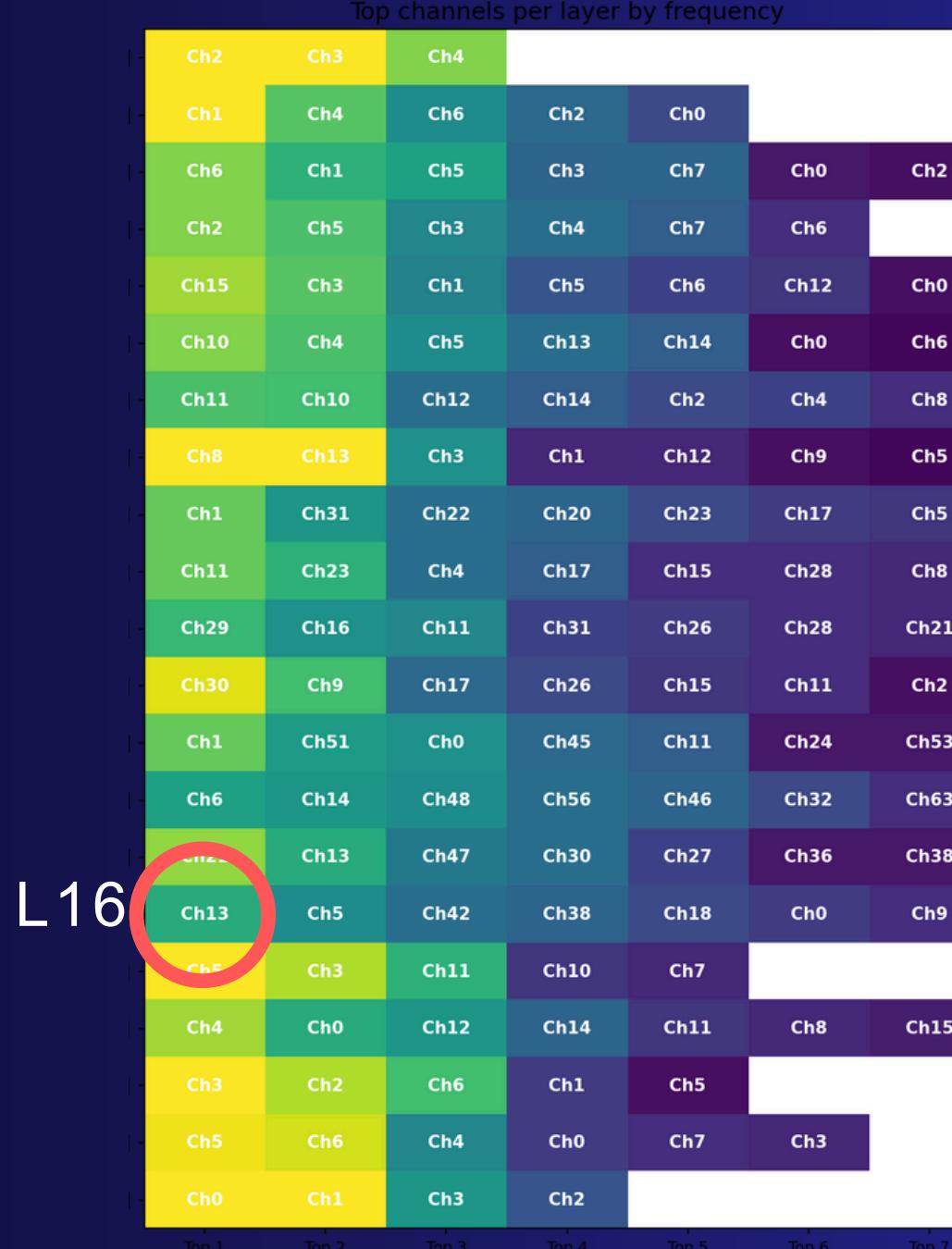


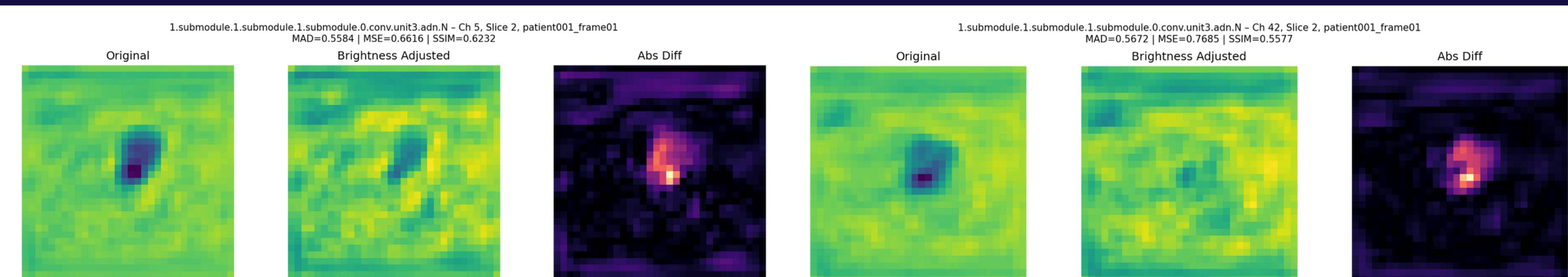
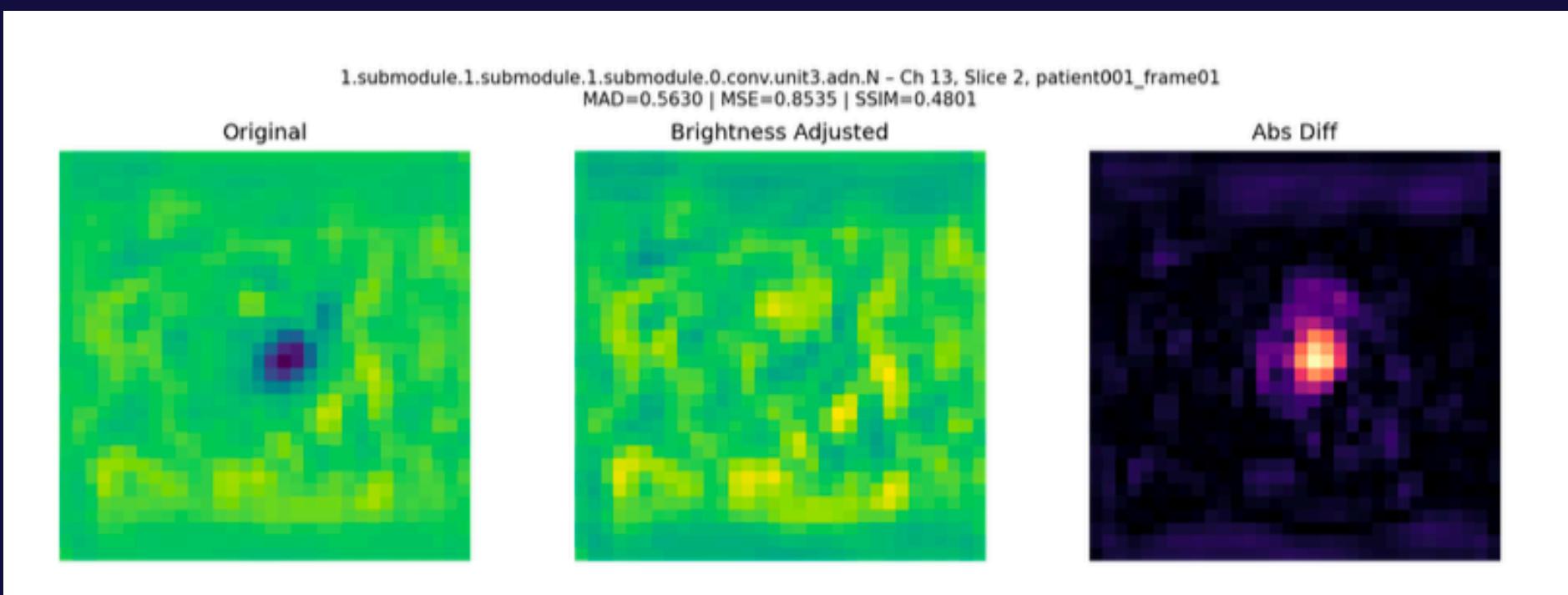
Figure 1

## LV+MYO case



# LV-Brightness Sensitive LV class dependent

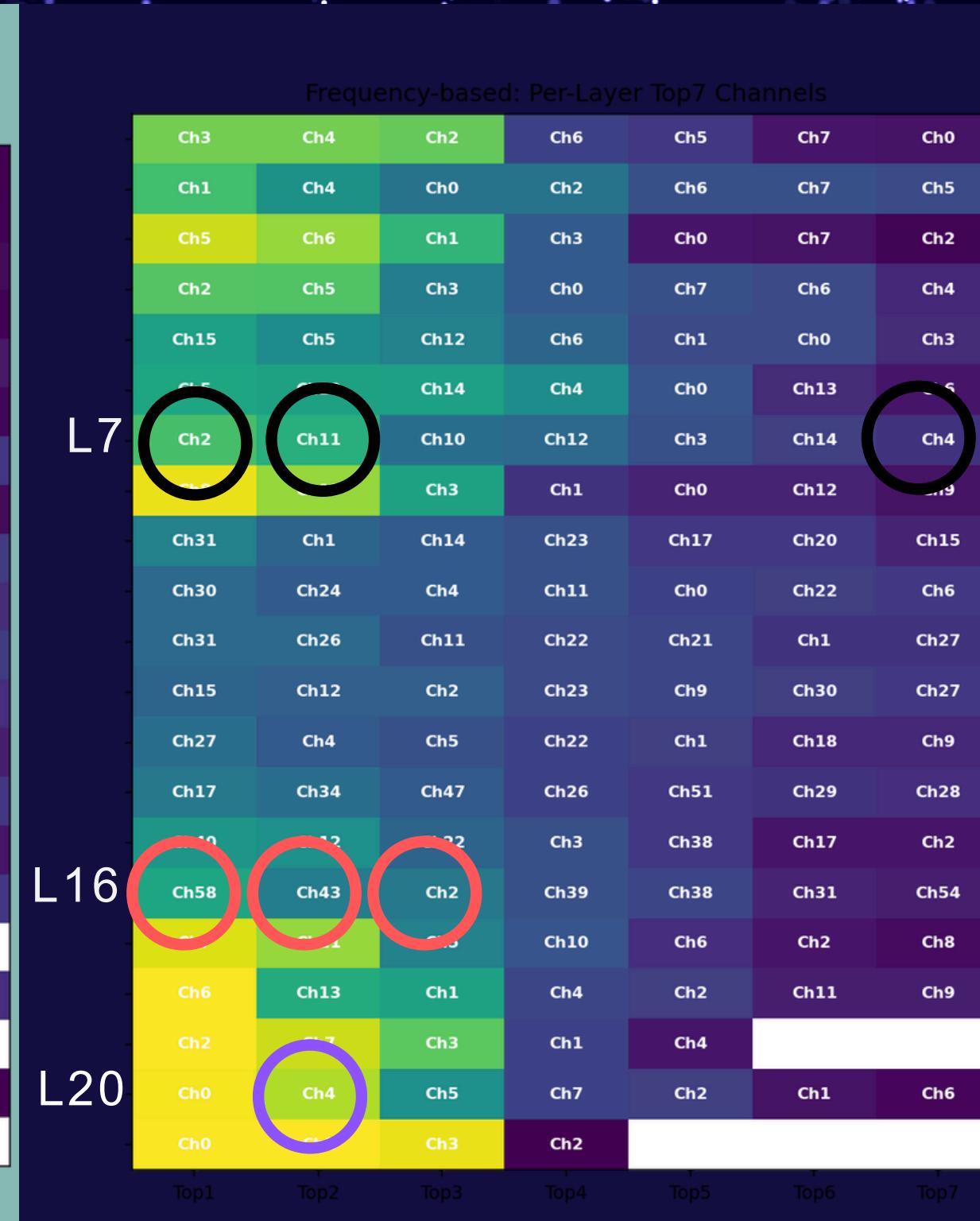
## L16-Ch13, L16-Ch42, and L16-Ch05



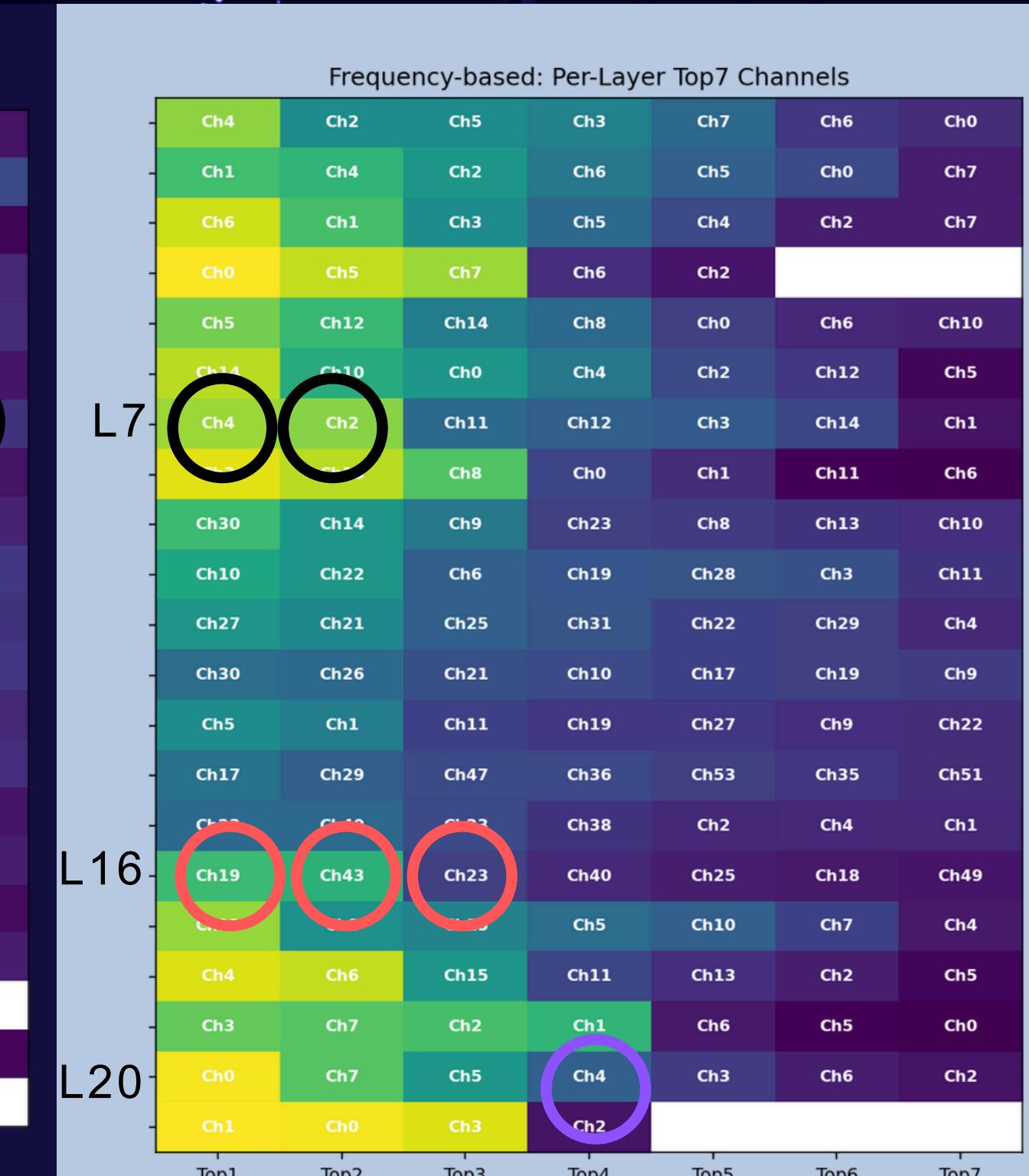
# LV Frequency Map



# RV Frequency Map

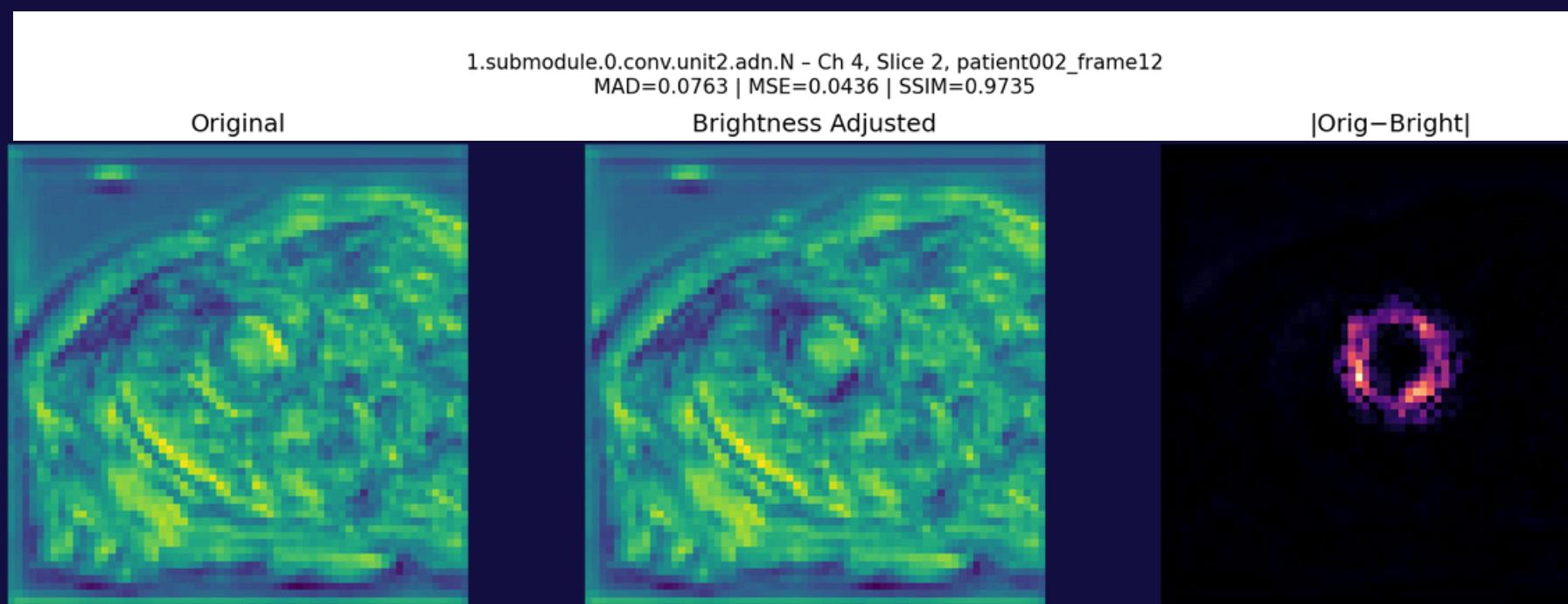
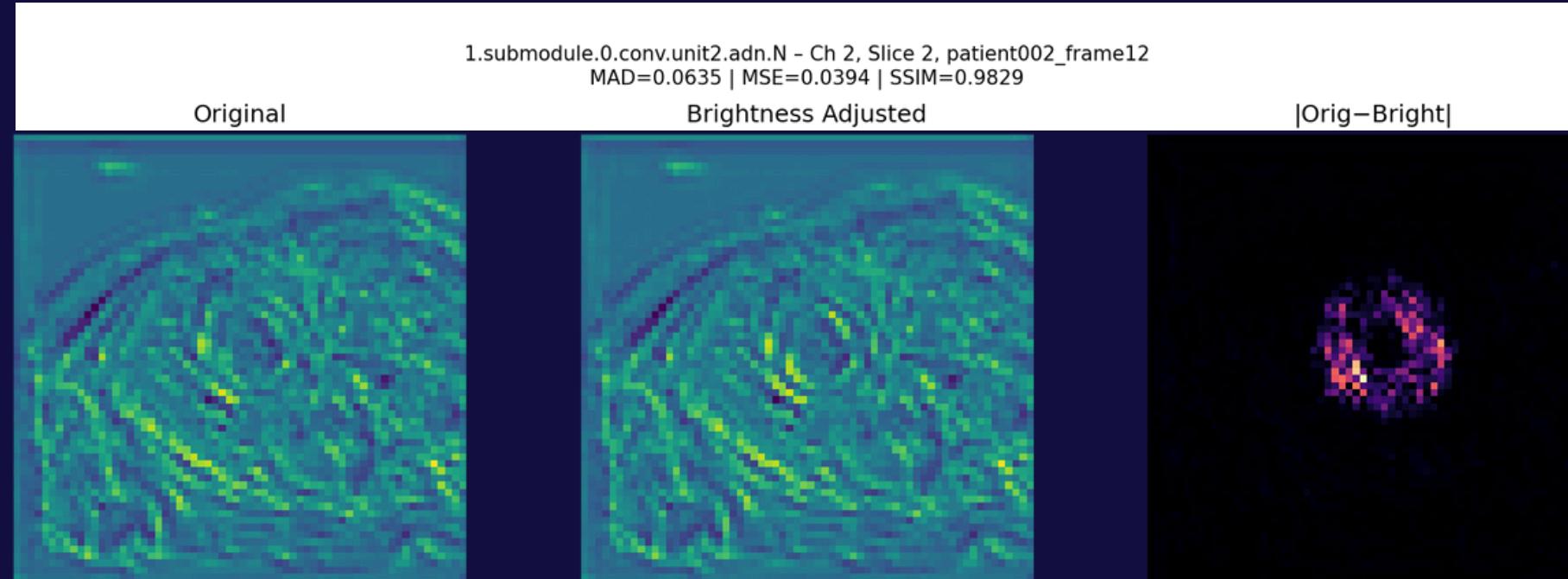


# MYO Frequency Map



# MYO-Brightness Sensitive MYO class dependent

## L7-Ch04



Thank you for  
listening!



*Why do we have brightness  
specific channels just for a  
particular patient group  
(patient1 frame1 and frame12)?*

*Some layers have  
different MAD value  
ranges what does this  
indicate?*

# References

- Bertolini, M., Rossoni, M., & Colombo, G. (2021). Operative Workflow from CT to 3D Printing of the Heart: Opportunities and Challenges. *Bioengineering*, *8*(10), 130.
- Bertolini, M., Rossoni, M., & Colombo, G. (2021). Operative Workflow from CT to 3D Printing of the Heart: Opportunities and Challenges. Bioengineering, 8(10), 130.