# Studying the Interpretability of Brain MR Image Deep Learning Modelling

Data Science Institute
COLUMBIA UNIVERSITY

Justin Kennedy[1], Siyuan Shi[2], Jiook Cha[3]
[1]Data Science, [2]Data Science, [3]Pyschiatry, Columbia University

Data Science Capstone Project
with Professor Jiook Cha

## Overview

In this project, we looked to identify parts of the brain that are relevant towards prediction of specific levels of cognitive impairment, from cognitively normal to diagnosed alzheimer's disease patients. To do this, we applied several machine learning algorithms in aggregation on patient brain MR image sequences to generate region-specific heatmaps of each patient sample. In building off the resulting interpretations, this project has implications on forecasting an individual's neurocognitive development, cognition, and behaviors.
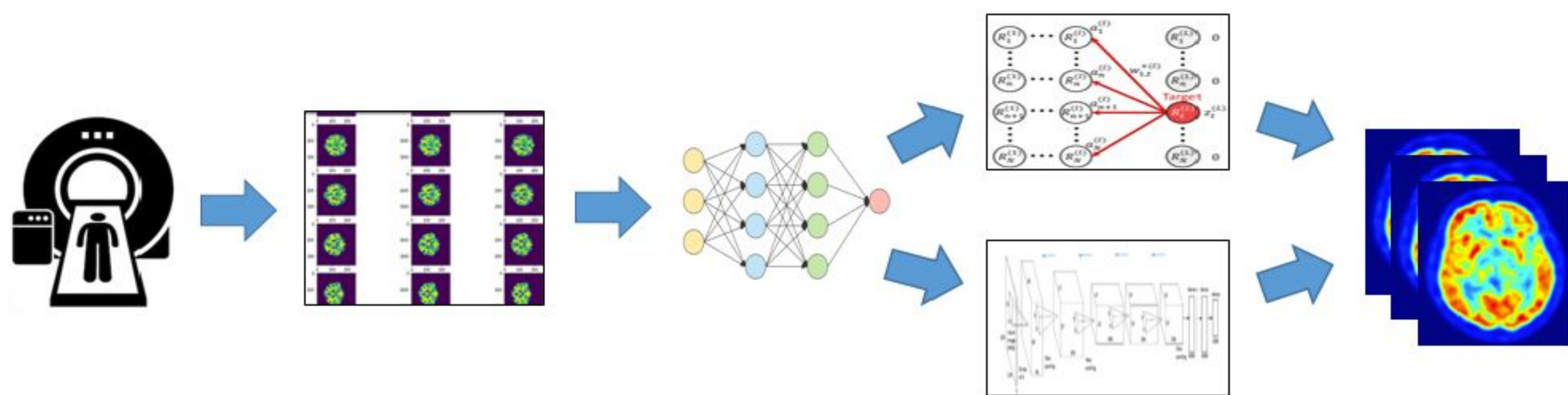


Figure 1. A patient's 3D MRI sequence is passed into a 3D CNN model that predicts the patient's cognitive group. Multiple MLI algorithms are applied to the CNN to produce a 'summary score' 3D pixel-wise contribution heatmap.

## Method

Layer-wise relevance propagation (LRP) works by creating a heatmap of per-pixel contributions to an input image's output label by computing relevances of neurons in each layer in a backward pass. In Guided Backpropagation (GBP), the gradient is backpropagated and all negative gradients are set to 0 to identify the neurons that contribute positively to the predicted group. Grad-CAM(++) works by constructing a weighted sum of feature maps using gradients in the network.
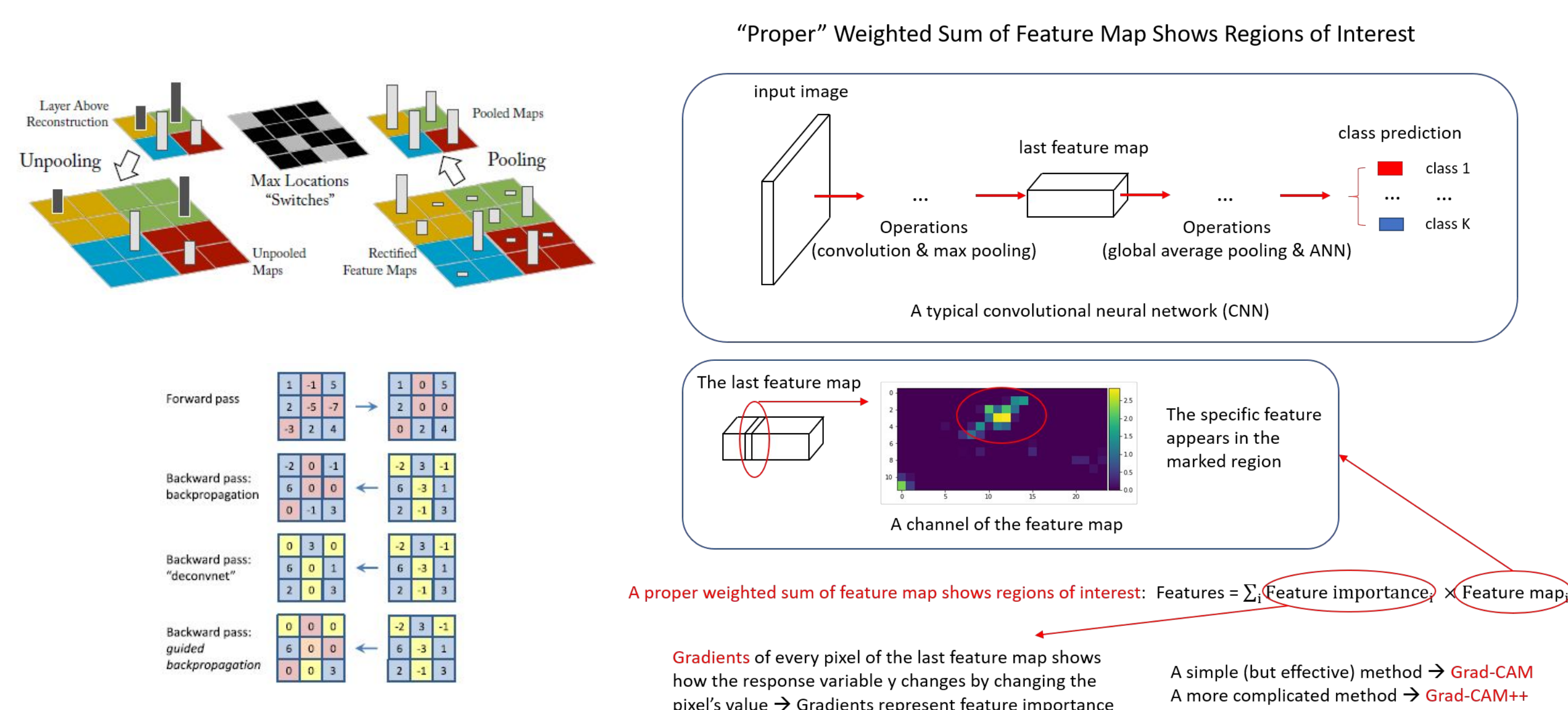


Figure 2. Core ideas of MLI algorithms. Left: guided backpropagation. Right: Grad-CAM (Grad-CAM++).

## Results

A subsample of 133 instances was chosen to test the MLI algorithms on. The model's predictions on the sample are summarized as a confusion matrix in Figure 3.

Here we show the resulting MLI interpretations of one of the instances. The LRP and GBP algorithms give per-pixel details; the colored pixels are important for making the prediction. The Grad-CAM (Grad-CAM++) result focuses on the region-level details; the colored regions contain features that the model thinks belong to the predicted class.
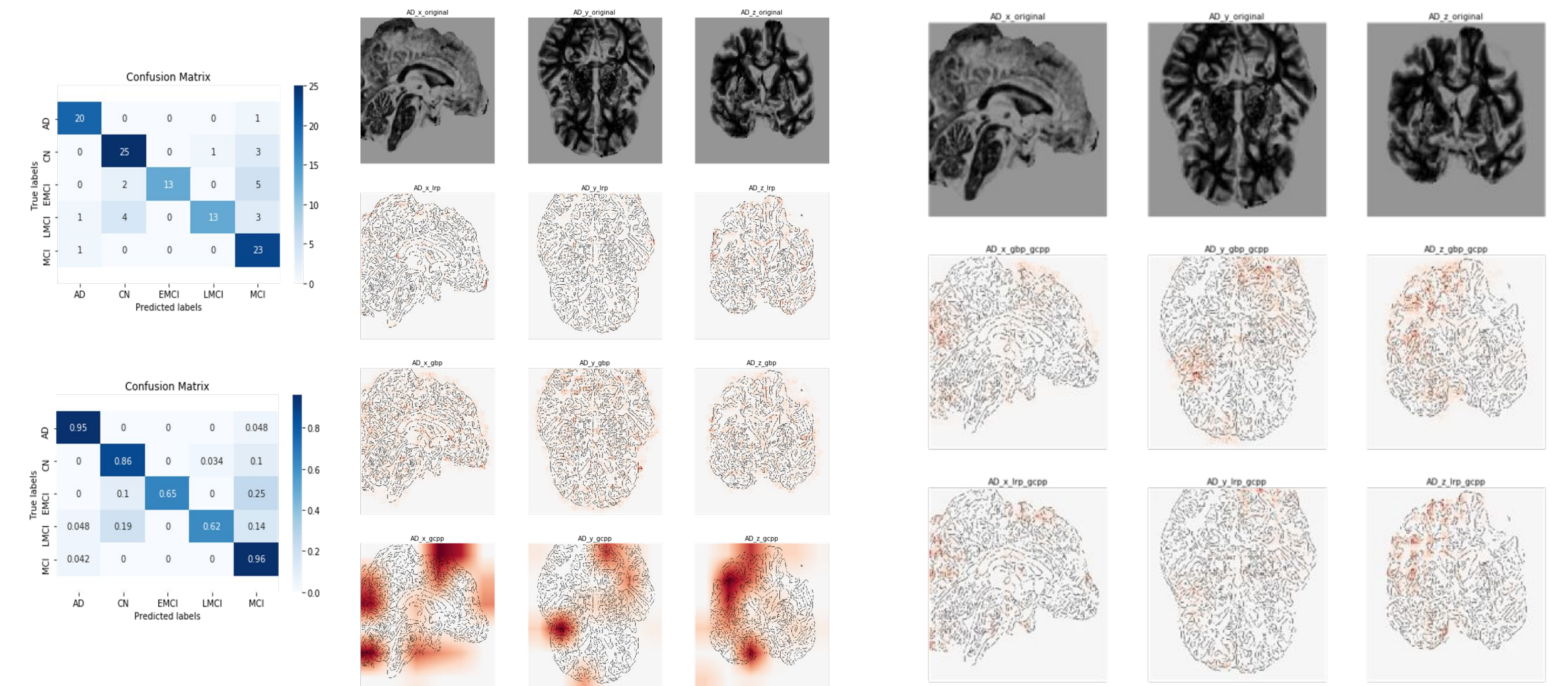


Figure 3. Left: confusion matrix of the subsample. Middle: input image, LRP, GBP, and Grad-CAM++ results. Right: Layer-wise relevance Grad-CAM++ & Guided Grad-CAM++.

## Conclusion

Though difficult to assess the accuracy of the resulting MLI heatmaps, especially given their dependence on the accuracy of the model they're based on, we can observe a consistency between separate MLI methods in certain brain regions being commonly specified as predictive of cognitive group. As the model's accuracy improves, these results could help infer area-specific characteristics related to cognitive conditions.

## Acknowledgments

We would like to thank Seungwook Han and Professor Jiook Cha for providing helpful feedback over the course of this project.

## References

Bach S, Binder A, Montavon G, Klauschen F, Müller K-R, Samek W, "On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation," PLoS ONE 10(7): e0130140, 2015.

"Grad-CAM: Visual Explanations from Deep ...." 21 Mar. 2017, https://arxiv.org/abs/1610.02391.

"Grad-CAM++: Generalized Gradient-Based ...." https://ieeexplore.ieee.org/document/8354201.