# Abstract

In the technological era of deep learning, neural network architecture has become incredibly well-known for its ability to look for patterns in data and use those patterns. However, the way to discover these patterns is still unknown to most people. Many people are still unable to comprehend the mechanisms underlying these discoveries.

Our project aims to shed light into the decision-making process of the pretrained models showing the transparency and interpretability of AI systems using Gradient-weighted class activation mapping (Grad-Cam). Grad-CAM highlights important regions in input data for specific class predictions, offering visual explanations for the model's decisions.

We will start with implementation of Grad-Cam on simple CNN models and visualize layer specific embeddings to analyze the model's behavior and outcome, subsequently the framework will then be modified to fit and further extended with pretrained models and identify the insights in data that influence the model's predictions,enabling a more thorough comprehension of the model's decision-making procedure.

Overall, deep learning models are great at recognizing patterns, but their lack of transparency makes it difficult to trust them completely. We aim to produce a framework for Explainable AI using Grad-CAM and layer embeddings. By these techniques, we can uncover how pre-trained models make decisions, which in turn promotes trust and interpretability.

# Introduction and Motivation

The deep neural networks have revolutionized various domains with their remarkable performance. However, the nature of these neural networks, "black-box" often raises concerns about transparency and interpretability, hindering trust and broader adoption in critical applications.In order to understand and address this, explainable AI has become crucial, aiming to unveil the inner workings of these complex models.

This research mainly focuses on exploring and explaining the behavior of deep neural networks through two key methods: visualization of layer-specific embeddings and Gradient-weighted Class Activation Mapping (Grad-CAM). Layer-specific embeddings reveal the learned representations at different network layers, providing insights into the model's information processing. On the other hand, Grad-CAM highlights important regions in input data for specific class predictions, offering visual explanations for the model's decisions.

## • Background

The ability of deep learning models to extract the features from the raw data has made significant changes in computer vision. Two important papers "Learning Deep Features for Discriminative Localization" and "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization" have played important role in enhancing the transparency and interpretability of CNNs through innovative techniques, Class activation mapping (CAM) and Gradient-Weighted class activation mapping (Grad-CAM).

## • Research Gap Identification

Current research does not fully grasp the literature on how complex Convolutional Neural Networks (CNNs) make decisions, especially when used with pretrained models. Because lack of understanding makes it difficult to trust and use these models in real-world situations. Furthermore, there is also a lack of standardized methods and evaluations for interpreting these models, which adds to the challenge. These issues are significant because they prevent widespread use of AI technologies due to concerns about reliability and bias.we hope to make AI systems more understandable and trustworthy. We are also focusing on practical applications and developing a framework to address these gaps in knowledge.

## • Research Objectives or Hypotheses

Our research objectives focus on implementing Grad-CAM on pre-trained CNN models to provide visual explanations for predictions and evaluating its effectiveness across diverse datasets. Moreover, we also aim to establish a standardized framework for evaluating model interpretability and to uncover how Grad-CAM influences user trust in AI reliability. Additionally, the real-world concerns related to AI bias and transparency is also a key objective. Altogether, these objectives aim to improve model interpretability and trustworthiness, filling the identified gaps in the current literature.

- **Significance of the Study**

This research is essential as it advances the theoretical understanding of AI as well as its practical applications. Our method provides practical insights into complex CNN decision-making along with visual explanations that can improve trust and understanding of AI systems. Our standardized framework for model interpretability evaluation, and our research contributes valuable insights for future studies in Explainable AI, informing and inspiring further research and applications. In a nutshell our research attempts to fill in the current gaps, enhance AI trust, and have a long-lasting effect on the field.

# Related Work

A significant breakthrough in bridging this gap came with the introduction of Class Activation Mapping (CAM) by Zhou et al. (2016) [1]. Zhou et al. introduced CAM for CNNs with global average pooling. This enabled the pre-trained CNN to perform object localization. The class activation mapping allowed them to predict the score of the given image and highlight the data features detected by CNN.

Expanding on this Selvaraju et al. (2017) [2] introduced Grad-CAM to make CNN models more transparent by producing a visual explanation. They also combined Grad-Cam with existing fine-grained visualization to create class discriminative visualization called Guided grad cam and applied it to image classification, and image captioning including ResNet-Based architecture further improving deep neural network interpretability.

# Data Material and Preprocessing

- **Data collections:** We will gather the necessary data for our study, consisting of pre-trained models and testing datasets. These models can be sourced from platforms such as huggingfaces.co, Kaggle, or pre-built models within PyTorch libraries. Similarly, testing data will be obtained from similar sources, ensuring compatibility with the models.
- **Data description:** Our study will primarily utilize image data, typically in JPEG or PNG formats. Additionally, specific model requirements will guide data preprocessing, ensuring alignment with model input specifications.
- **Data Preprocessing:** Raw image data will be initially represented as pixel values. To prepare this data for model consumption, we'll employ PyTorch to convert these values into tensors with model-specific dimensions. These processed tensors will then serve as inputs for the pre-trained models.
- **Software and tools:** The entire project will be conducted using Python with a Jupyter Notebook environment. Essential libraries such as PyTorch, Grad-CAM, torchvision, matplotlib, numpy, and pandas will be utilized to facilitate various aspects of the project.

# Methodology

We set out to create a framework that will help us to understand how pre-trained AL models make decisions. We will focus on CNNs and use Grad-CAM to visualize what parts of an image are important for the model's predictions.

Our work will be carried out in three parts:

**1. Building a simple Framework:**

- First, we will create a basic framework. It will take a pre-trained model and some images as input.
- The pre-trained model can be downloaded from Hugging Face, kaggle or can be imported from Python libraries and images can be downloaded from anywhere.
- The framework should highlight to us the regions and important features of the data on which the model has made certain predictions.

**2.Improving the Framework:**

- Next improvements will be made to the framework to make it better. It should be able to work with many other pre-trained models.

**3.Testing the Framework:**

- Finally, we will put the framework through the paces and test it with various pre-trained models and different sets of images to make sure it does a good job of showing us what's important in the data.

# References

1. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A. and Torralba, A., 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2921-2929).
2. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D. and Batra, D., 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision* (pp. 618-626).