A decorative background featuring a network diagram with nodes and connecting lines, primarily in light gray and blue, located in the top-left and bottom-right corners.

Convolutional Neural Network Architectures

Hello!

I am Wesley Osborne

Head
Research
Data
Scientist



A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines, with some nodes highlighted in blue.

1.

Problem Statement

What are we researching?

Problem Statement

Osborne Research Labs is a recently started research lab. Our lab is conducting research on computer vision systems for drones and robotics. Given we are just entering the field and we have a basic understanding of convolutional neural networks (CNNs), we wanted to research popular complex CNN architectures. The goal of this project is to learn about the pros and cons of three popular CNN architectures and see how they perform on an indoor scene dataset. We're hoping for accuracy scores above the baseline of 9%.

A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. The nodes are represented by small circles, some of which are larger and have concentric circles, suggesting different levels of connectivity or importance. The lines are thin and gray, creating a mesh-like structure.

2. **About the Data**

Data Collection

- ① Used MIT Indoor Scene Dataset
 - 15,000 images across 67 classes
- ① Used 11 Classes
 - Bathroom, Bedroom, Pantry, Corridor, Staircase, Gym, Living Room, Kitchen, Office, Closet, and elevator
- ① Scraped ~500 new images for each class from google
- ① Total dataset of approx. 7,000 images

A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines, with some nodes highlighted in blue.

3.

Basic Building Blocks of CNNs

CNN Basics

◎ Convolutional Layers (Conv)

◎ Pooling Layer (POOL)

◎ Fully Connected Layer (FC)

a.k.a Dense Layer

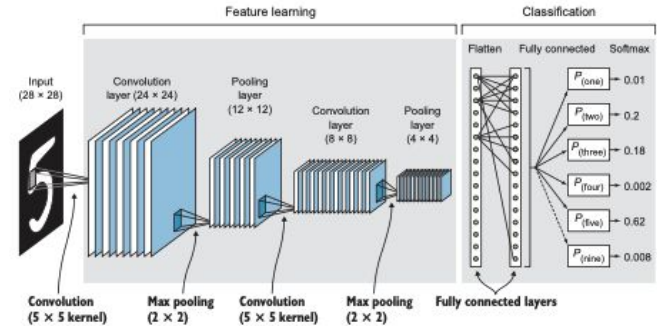


Figure 3.12 The basic components of convolutional networks are convolutional layers and pooling layers to perform feature extraction, and fully connected layers for classification.

A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines, with some nodes highlighted in blue and others in grey.

4. **CNN Architectures**

VGGNet

- Developed in 2014 by The Visual Geometry Group at Oxford University (Keren Simonyan and Andrew Zisserman)
- Most popular configuration is VGG16
- Popular because it's simple to understand and has a uniform architecture
- Inspired by AlexNet and LeNet
- Top-5 error rate of 8.1% on ImageNet

VGGNet Architecture

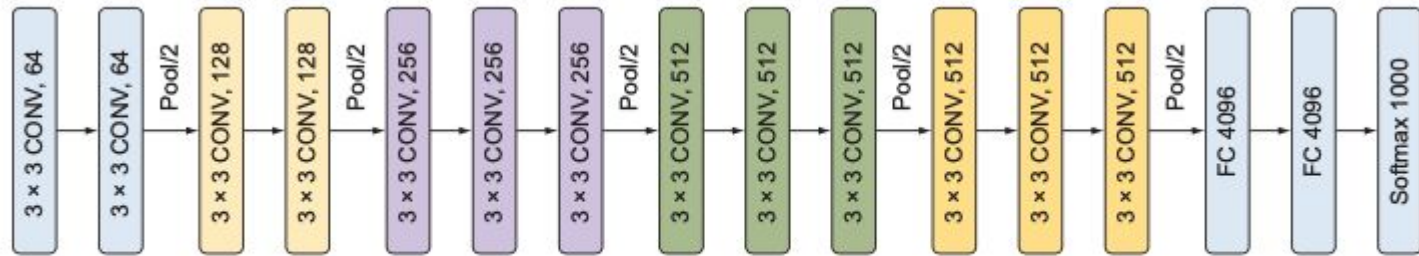
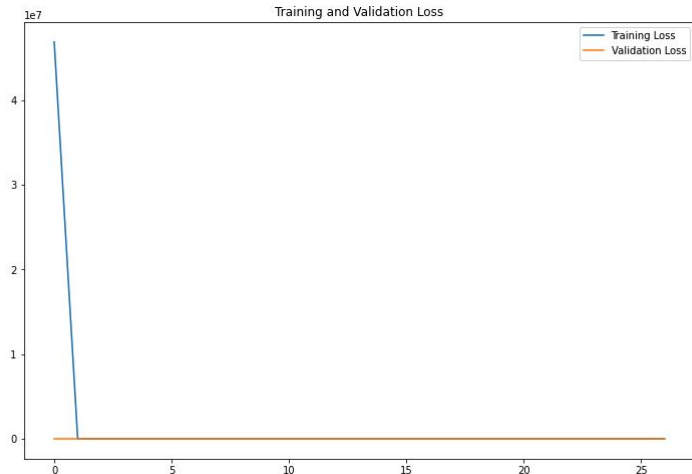
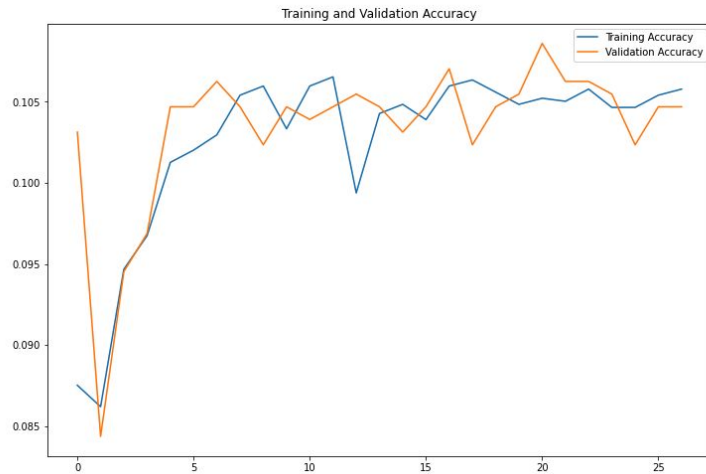


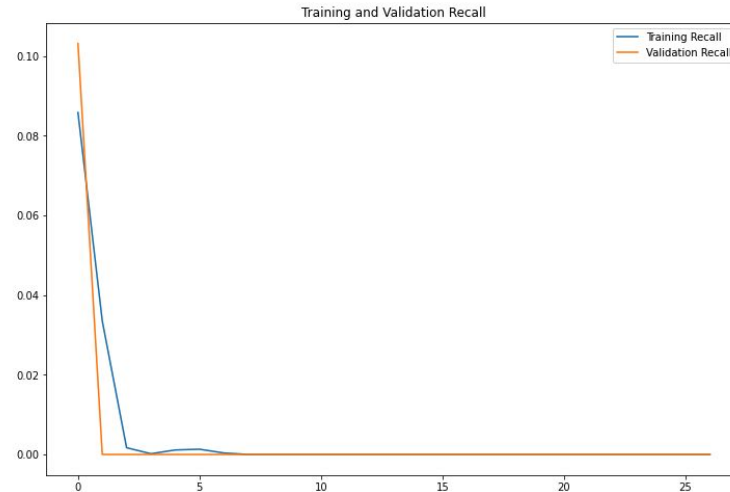
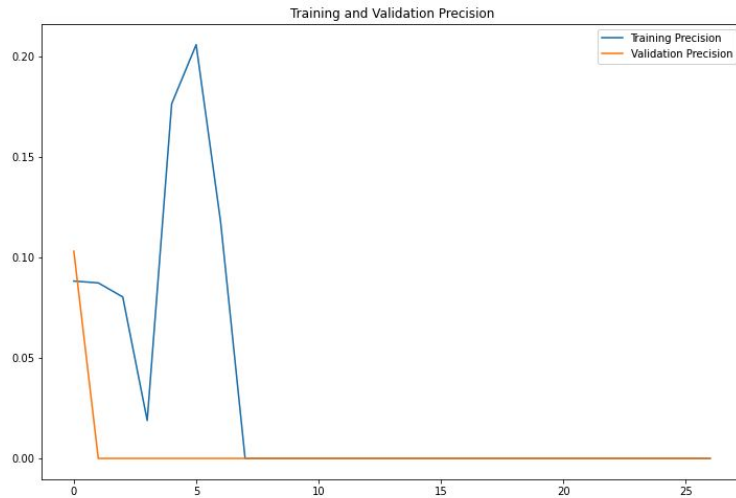
Figure 5.8 VGGNet-16 architecture

VGG16 Performance on Indoor Scenes

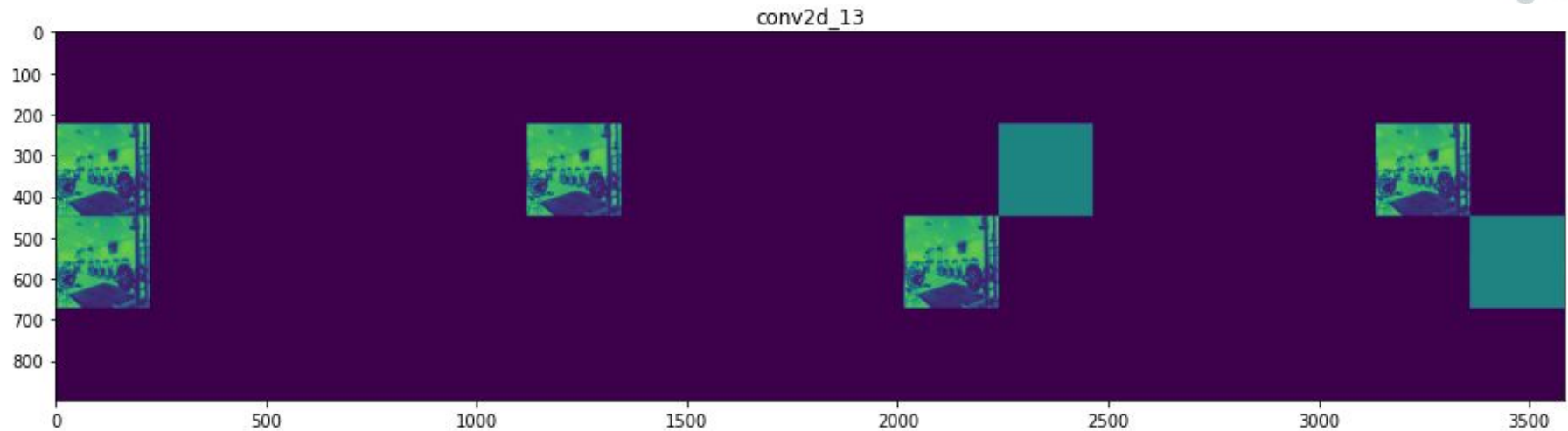
🎯 Top accuracy of around 10%



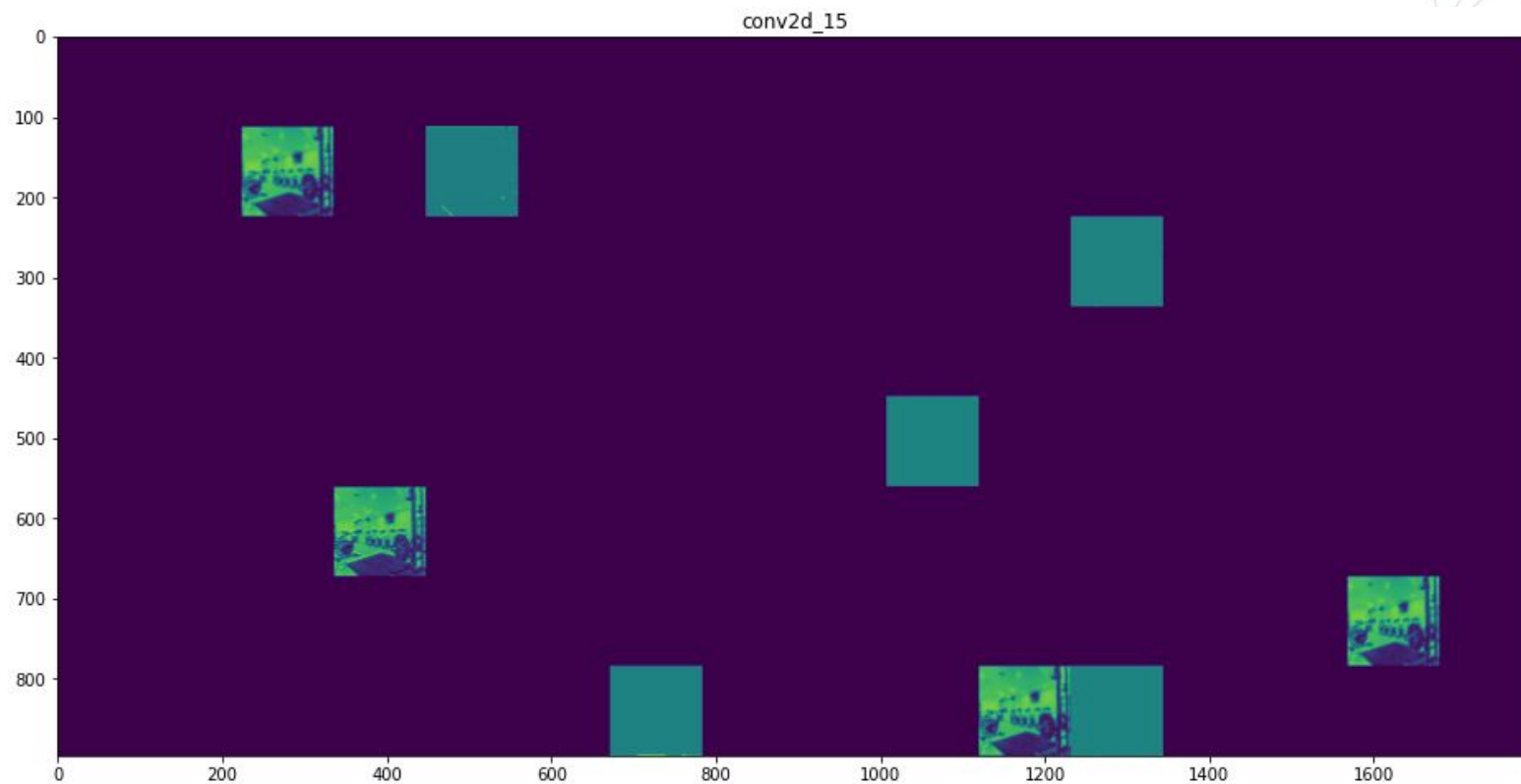
VGG16 Performance on Indoor Scenes



VGG16 Performance on Indoor Scenes



VGG16 Performance on Indoor Scenes



Inception and GoogleNet

- Developed in 2014 by a group of researchers at Google
- GoogleNet is a configuration of the Inception Network
- Created a deeper network than VGGNet while reducing the number of parameters.
- Introduced the Inception Module
- Top-5 error rate of 6.67% on ImageNet

Inception Module

- ◎ Researchers at Google decided to implement different sized convolutional layers and pooling layer all together in one block.
- ◎ The network architecture is developed stacking series of inception modules together

Inception Module

Inception module with dimensionality reduction

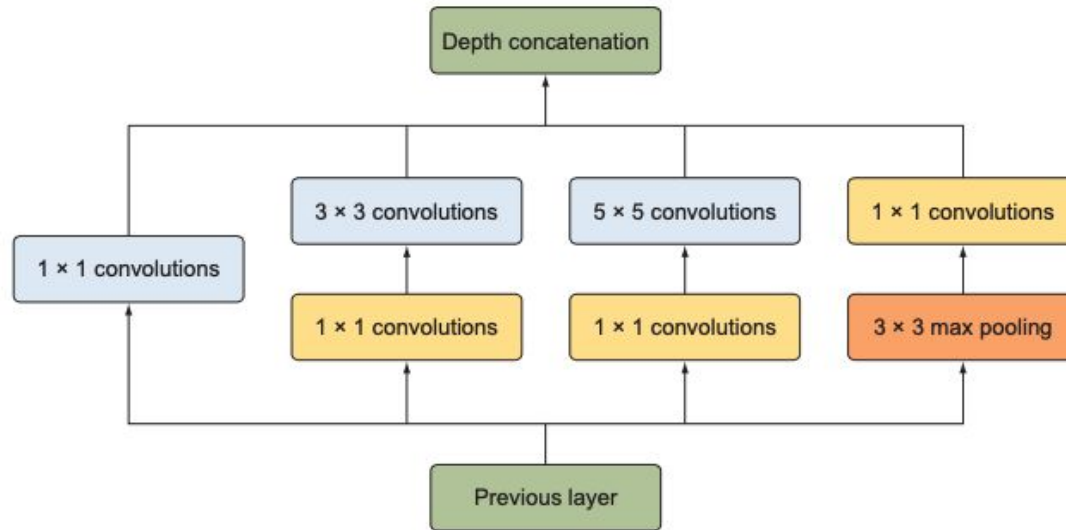


Figure 5.13 Building an inception module with dimensionality reduction

Inception Architecture

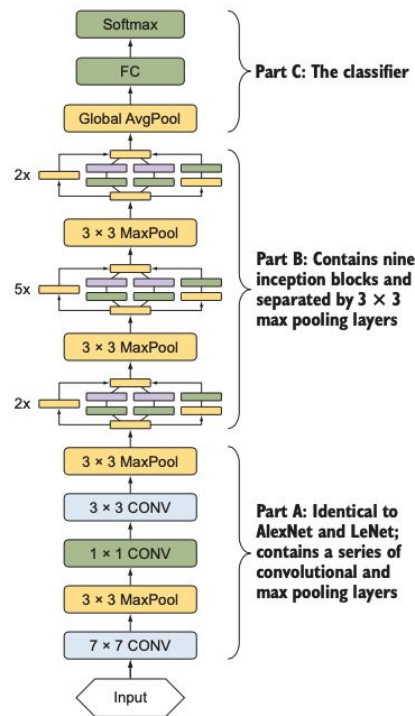
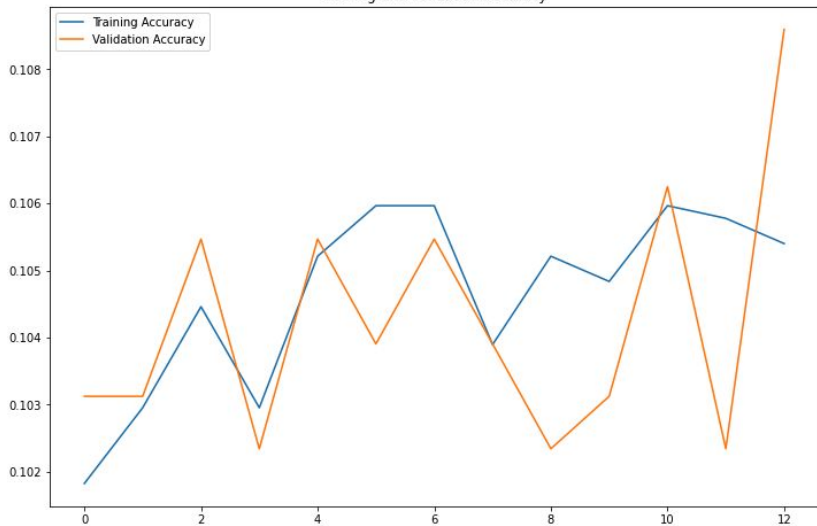


Figure 5.15 The full GoogLeNet model consists of three parts: the first part has the classical CNN architecture like AlexNet and LeNet, the second part is a stack of Inceptions modules and pooling layers, and the third part is the traditional fully connected classifiers.

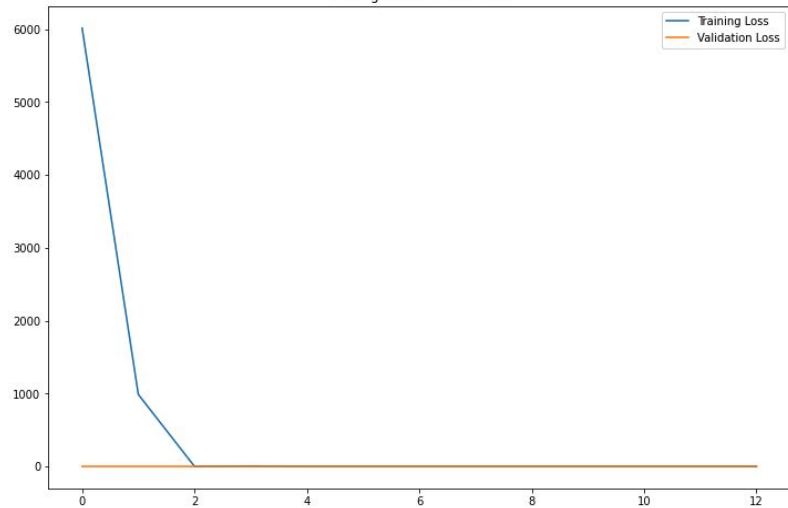
GoogleNet Performance on Indoor Scenes

🎯 Top accuracy of around 10.8%

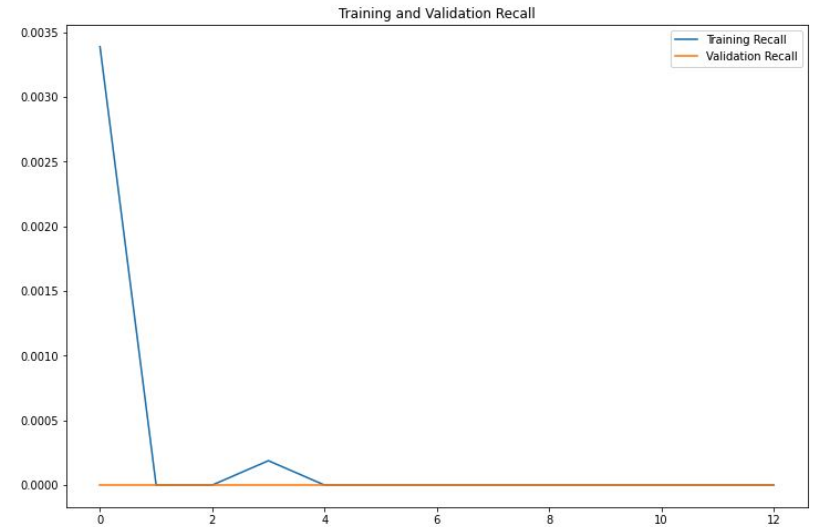
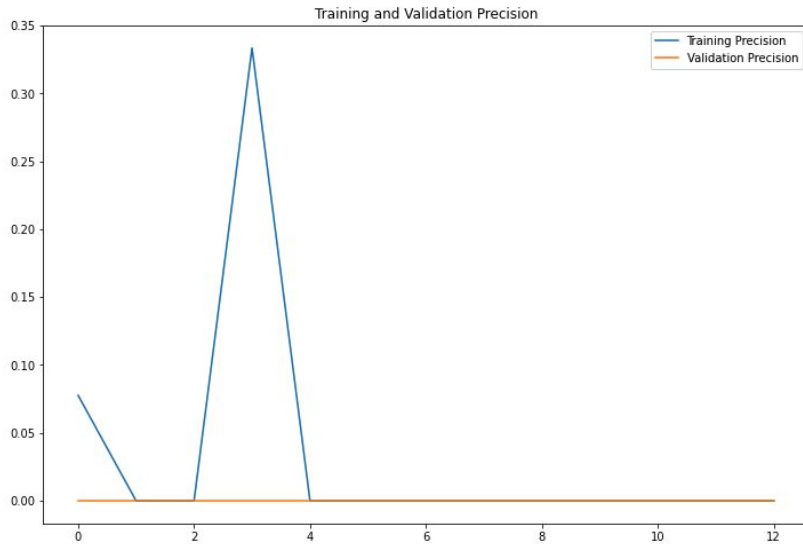
Training and Validation Accuracy



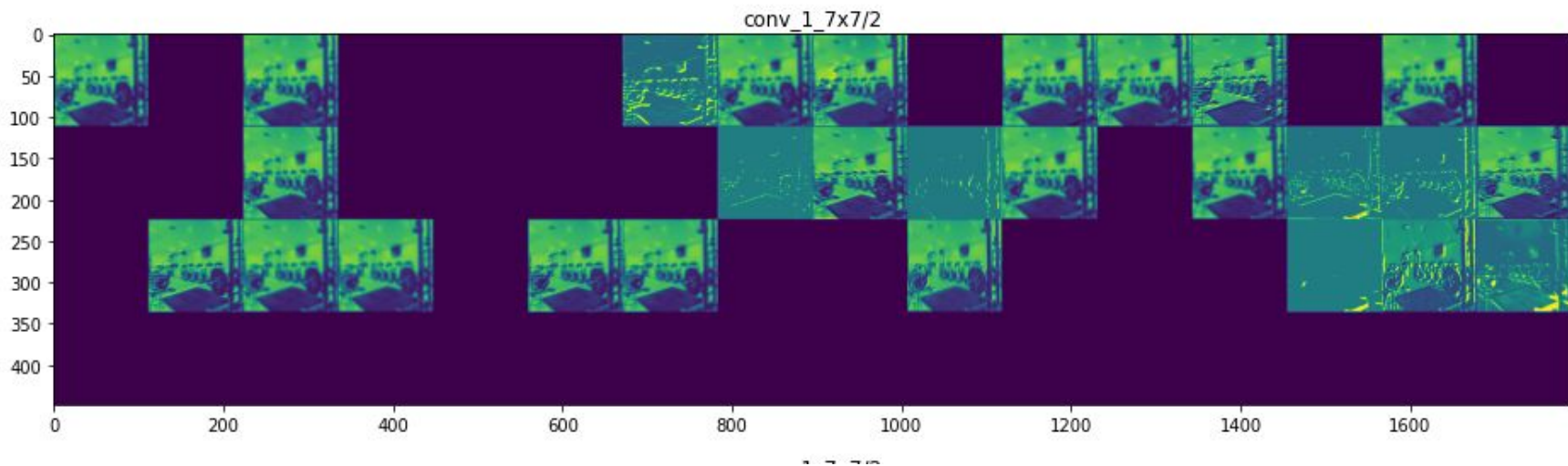
Training and Validation Loss



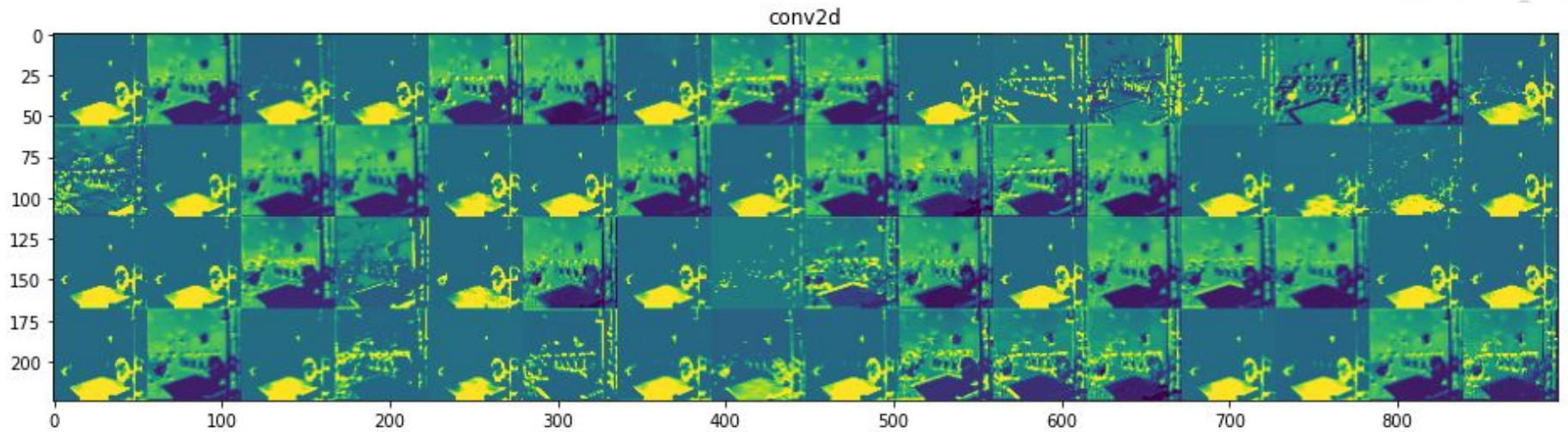
GoogleNet Performance on Indoor Scenes



GoogLeNet Performance on Indoor Scenes



GoogleNet Performance on Indoor Scenes



ResNet

- ◎ Residual Neural Network (ResNet) developed in 2015 by a group of researchers from the Microsoft Research team.
- ◎ Introduced the Residual Module with skip connections
- ◎ Uses Batch Normalization heavily for hidden layers
- ◎ Able to achieve really deep neural networks with 50, 101, and 152 weight layers with lower complexity than smaller networks like VGGNet.

Top-5 error rate of 3.57% on ImageNet

ResNet Skip Connections

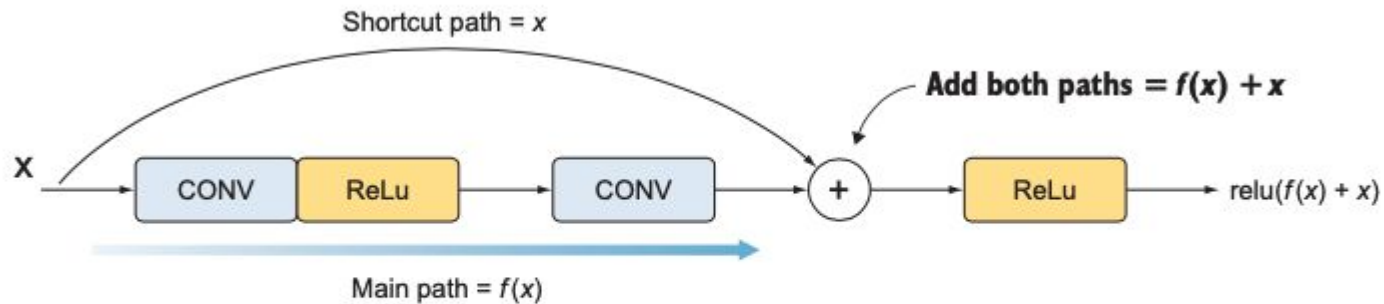


Figure 5.20 Adding the paths and applying the ReLU activation function to solve the vanishing gradient problem that usually comes with very deep networks

Residual Block

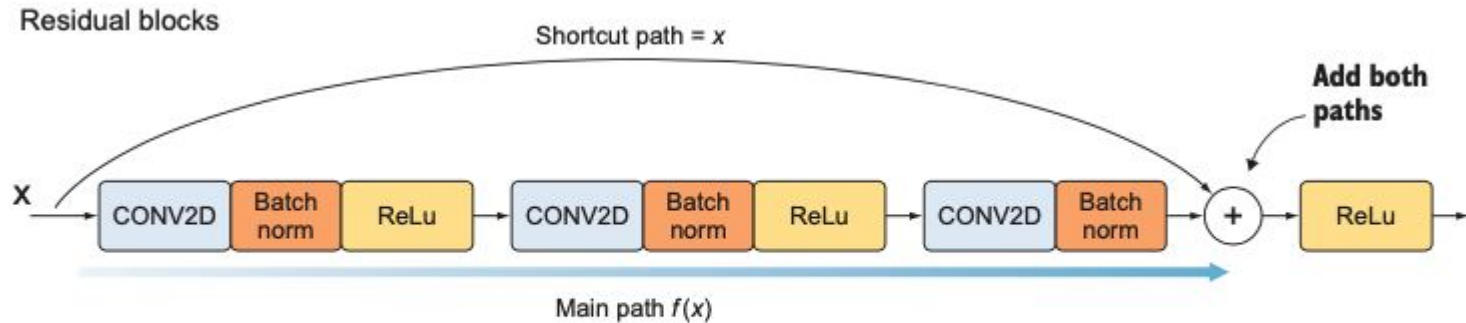


Figure 5.22 The output of the main path is added to the input value through the shortcut before they are fed to the ReLU function.

Residual Block

Bottleneck residual block with reduce shortcut

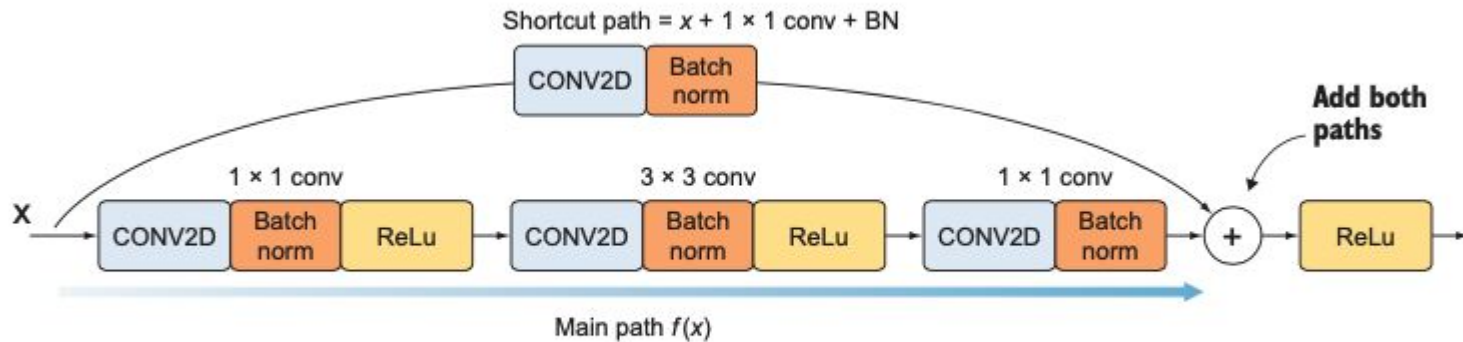


Figure 5.23 To reduce the input dimensionality, we add a bottleneck layer (1×1 convolutional layer + batch normalization) to the shortcut path. This is called the *reduce shortcut*.

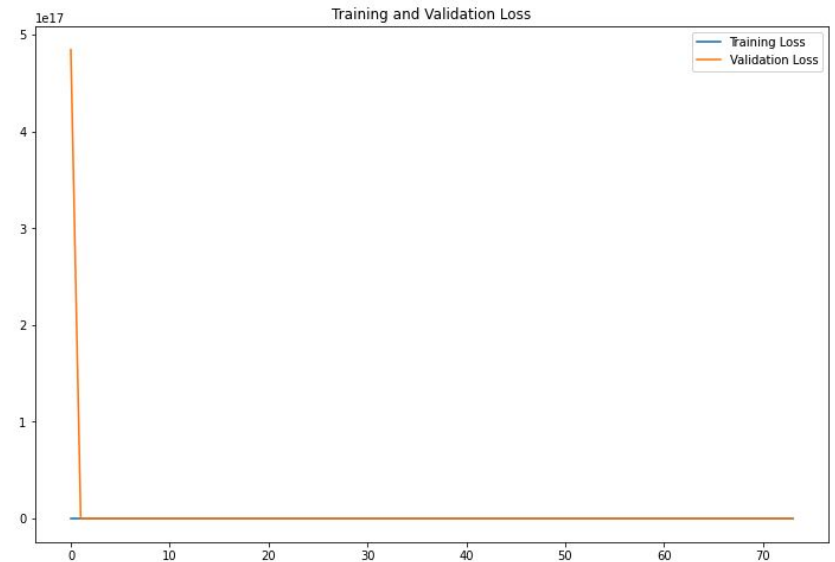
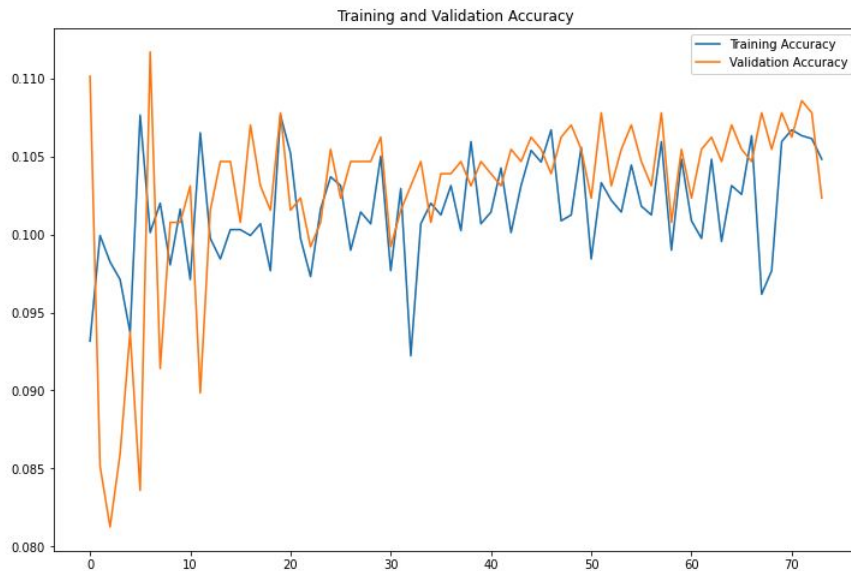
ResNet Architecture

Layer name	Output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112x112	7x7, 64, stride 2				
conv2_x	56x56	3x3, maxpool, stride 2				
		$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28x28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14x14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7x7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1x1	Average pool, 1000-d fc, softmax				
FLOPs		1.8x10 ⁹	3.6x10 ⁹	3.8x10 ⁹	7.6x10 ⁹	11.3x10 ⁹

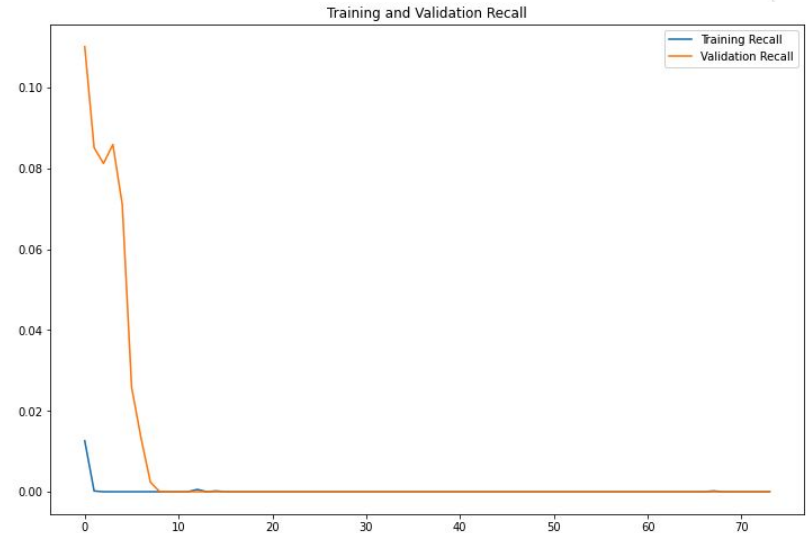
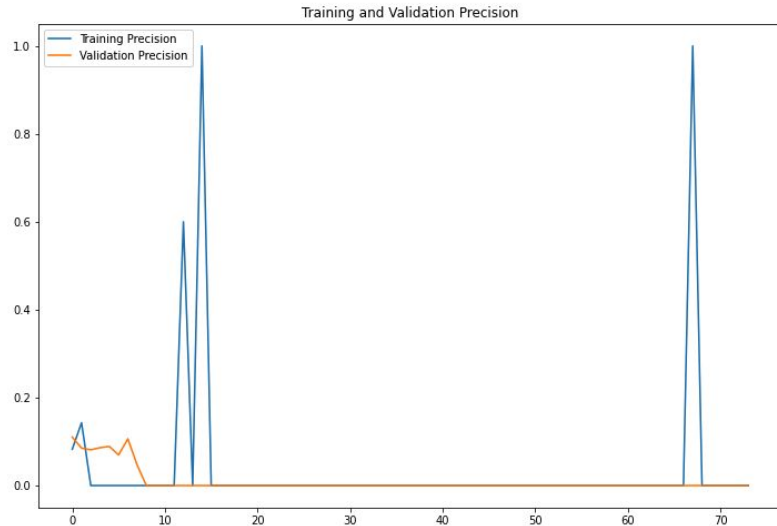
Figure 5.24 Architecture of several ResNet variations from the original paper

GoogleNet Performance on Indoor Scenes

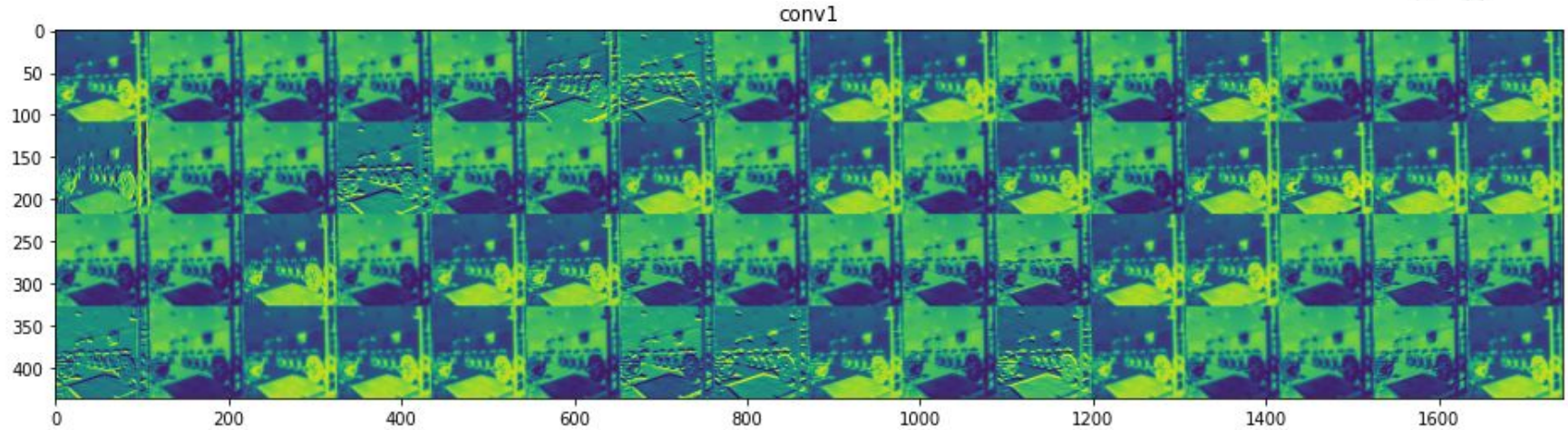
🎯 Top accuracy of around 11%



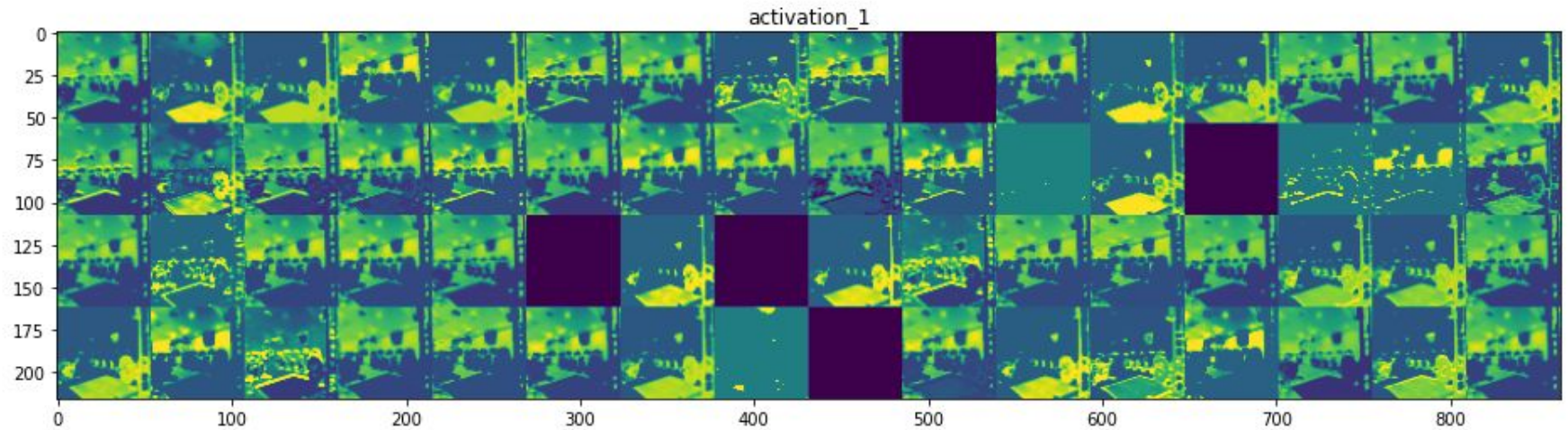
GoogleNet Performance on Indoor Scenes



GoogLeNet Performance on Indoor Scenes

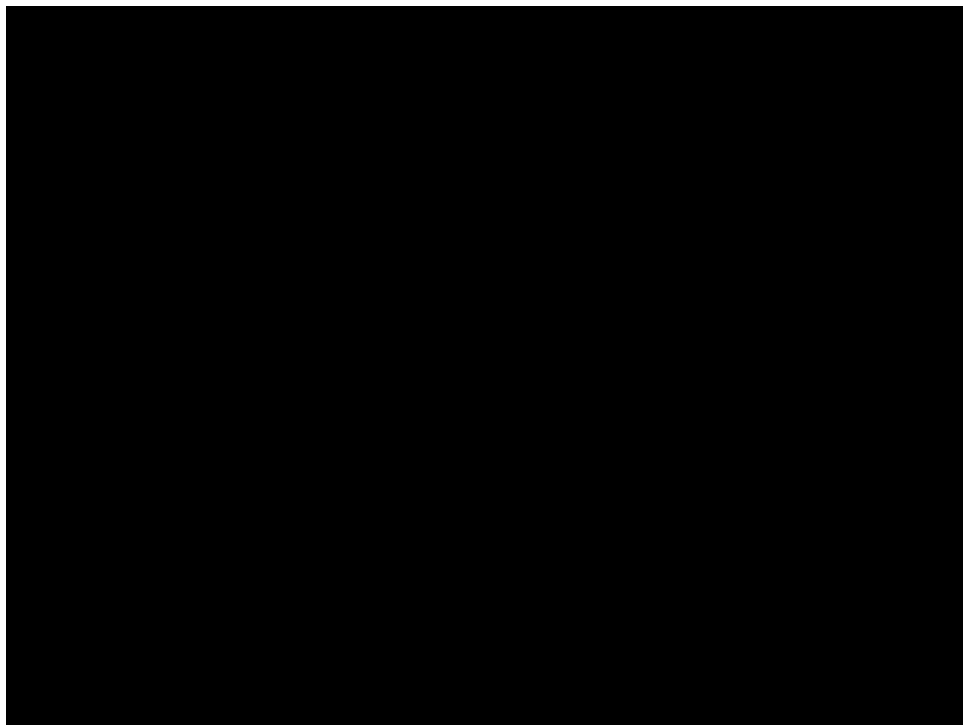


GoogleNet Performance on Indoor Scenes



A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. The nodes are represented by small circles, some of which are larger and have concentric circles, suggesting different levels of connectivity or importance. The lines are thin and gray, creating a mesh-like structure.

5. Demo



Conclusion and Further Steps

- ◎ All 3 architectures performed poorly on the Indoor Scene dataset
- ◎ Each one performed the same on the test data
 - 20 out of 200 correctly classified
 - Correctly classified all 20 bathrooms
- ◎ These are great architectures to pull inspiration from
- ◎ Each dataset needs its own custom architecture
- ◎ Add More Data for training



Thanks!

Any questions?