# Identify Gate Operations using Action Recognition

DISSERTATION

Submitted in partial fulfillment of the requirements of the

M.Tech Data Science and Engineering Degree programme

By

**VISHALI S**

**ID No.**

*Under the supervision of:*



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

Pilani (Rajasthan), INDIA

August 2022

# Identify Gate Operations using Action Recognition

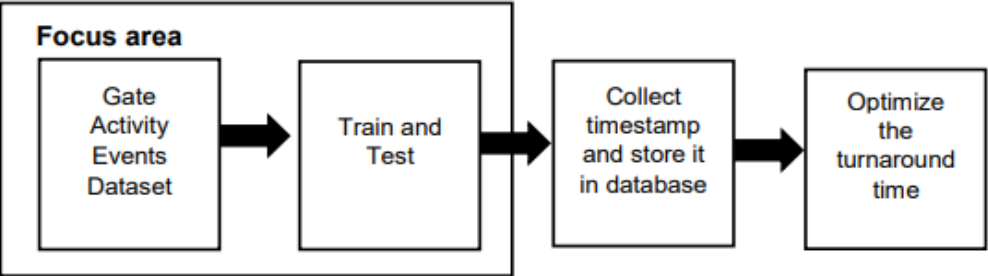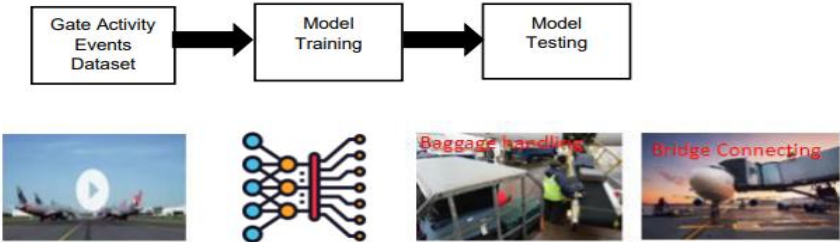| | |
|---|---|
| **Problem Statement:**<br><br>Airport ground operations are one of the main causes of late departures. A common mantra in the aviation industry is "The plane is not making any money while it is on the ground". This statement represents the underlying hurdles faced by the airline industry. There is a need to strive for efficiency, as airlines wish to keep their planes grounded for as short as possible. | **Benefit:**<br><br>While stationary, the airline is not earning any revenue with its plane, while facing many costs at the same time. Not having a quick turnaround time can cost business. In many cases, the amount of time it takes to complete represents cost and can present risk. Fast turnaround time means less money spent on manpower. |
| **Objective of the Project:**<br><br>Build a system which will take video as an input and recognize the ground activities which are bridge connecting, baggage handling and miscellaneous. To achieve this, we must build a model based on visual-analytical approach that can automatically recognize the activities. | **Resource:**<br><br>Operating System: High end Machine with 8VRAM or above<br>Languages/Scripting: Python<br>Libraries: PyTorch, Tensorflow, Matplotlib, OpenCV, Albumentations |
| **Solution Architecture:**<br><br>The spatiotemporal features of the videos are extracted using the image processing techniques. The preprocessed data will be given to the models. In Training phase, deep learning algorithms are used to extract the features and preciously recognize the actions. | **Risk:**<br><br>Data collection was at potential risk. We observed the open-source datasets for action recognition like ActivityNet, UCF101 and HMDB51. It has an average of 100 to 150 videos for each action. Downloaded videos from YouTube & performed trimming operation to focus only on the actions we required. Collected around 130 videos for each category. |



Fig 1: Overall Gate Monitoring System



Fig 2: Proposed Solution Architecture

# Project Plan & Deliverables:

| SL NO | Task | Planned duration (In weeks) | Name of Deliverables | Status |
|-------|------|------------------------------|----------------------|--------|
| 1 | Understanding the project requirements and specifications | 2 weeks | | Done |
| 2 | Data collection | 2 weeks | Dataset | Done |
| 3 | Data Preprocessing | 3 weeks | Processed Dataset | Done |
| 4 | Exploration of algorithms and initial training | 5 weeks | Literature study & Model | Done |
| 5 | Adding miscellaneous data as one more category | 1 week | Dataset | Done |
| 6 | Inference the model | 1 weeks | Model | Done |
| 7 | Documentation for final report | 2 weeks | Code & Report | Done |

Fig 3: Representation of workflow

# Data Acquistion:

## Data Collection:



## Dataset Folder Structure:



## Data Visualization:



## Data Exploration:

| Action Categories | 3 |
|---|---|
| Average Videos per Action Category | 130 |
| Average Number of Frames per Video | 812.85 |
| Average Frames Width per Video | 1154.58 |
| Average Frames Height per Video | 713.00 |
| Average Frames per Seconds per Video | 27.37 |

Statistic of overall dataset



Statistic of overall dataset

## Data Preprocessing:

Input Video Frames

Optical Flow X Coordinates Frames

Optical Flow Y Coordinates Frames

RGB Frames

# Exploration of algorithms and model training:

## Single Frame CNN



| Hyperparameters | Value |
|---|---|
| Epoch | 50 |
| Batch size | 4 |
| Loss | categorical_crossentropy |
| Optimizer | Adam |

This network uses single architecture that fuses information from all frames at the last stage. It works by running an image classification model on every single frame of the video and then average all the individual probabilities to get the final probabilities vector.

## ConvLSTM



| Hyperparameters | Value |
|---|---|
| Epoch | 50 |
| Batch size | 4 |
| Loss | categorical_crossentropy |
| Optimizer | Adam |

This is implemented by using a combination of ConvLSTM cells. A ConvLSTM cell is a variant of an LSTM network. It has got convolution embedded in the architecture, which makes it capable of identifying spatial features of the data while keeping into account the temporal relation.

## LRCN Model



| Hyperparameters | Value |
|---|---|
| Epoch | 50 |
| Batch size | 4 |
| Loss | categorical_crossentropy |
| Optimizer | Adam |

This network combines CNN and LSTM layers in a single model. The CNN model can be used to extract spatial features from the frames in the video and LSTM model can then use the spatial features extracted by CNN at each time-steps for temporal sequence modeling. This way the network learns spatiotemporal features directly in an end-to-end training, resulting in a robust model.

## Two Stream Network



| Hyperparameters | RGB Value | Optical Value |
|---|---|---|
| Epoch | 250 | 350 |
| Batch size | 25 | 25 |
| Learning rate | 0.001 | 0.001 |
| Loss | Cross entropy | Cross entropy |
| Optimizer | SGD | SGD |

This network is composed of two separate convolution neural networks. One to handle spatial features and one to handle temporal or motion features. The input to the two-stream architecture only contains a single frame for the spatial stream and a fixed-size group of optical flow maps for the temporal stream.

# Model Training Results:

## Single Frame CNN

### Total Loss vs Total Validation Loss

### Total Accuracy vs Total Validation Accuracy

**Training Results**

| Loss | 0.09544 |
|---|---|
| Accuracy | 0.9787 |

## ConvLSTM

### Total Loss vs Total Validation Loss

### Total Accuracy vs Total Validation Accuracy

| Loss | 0.4735 |
|---|---|
| Accuracy | 0.80851 |

## LRCN Model

### Total Loss vs Total Validation Loss

### Total Accuracy vs Total Validation Accuracy

| Loss | 0.4916 |
|---|---|
| Accuracy | 0.84523 |

## Two Stream Network

**RGB frames:**



| Loss | 0.6643 |
|---|---|
| Accuracy | 0.63402 |

**Optical frames:**



| Loss | 0.6546 |
|---|---|
| Accuracy | 0.63402 |

# Model Inference:

| Single Frame CNN | ConvLSTM | LRCN Model | Two Stream Network |
|---|---|---|---|



**RGB frames:**



Single Frame CNN (first row):
CLASS NAME: bridge_connecting    AVERAGED PROBABILITY: 9.6e+01
CLASS NAME: misclaneous    AVERAGED PROBABILITY: 0.13
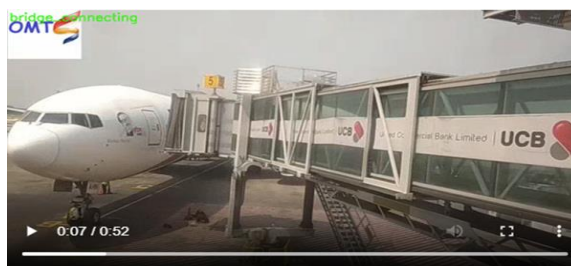CLASS NAME: baggage_handling    AVERAGED PROBABILITY: 0.003

ConvLSTM (first row):
Action Predicted: bridge_connecting
Confidence: 0.8162801265716553

LRCN Model (first row):
Action Predicted: bridge_connecting
Confidence: 0.8911964893341064

Single Frame CNN (second row):
CLASS NAME: baggage_handling    AVERAGED PROBABILITY: 9.2e+01
CLASS NAME: bridge_connecting    AVERAGED PROBABILITY: 6.2
CLASS NAME: misclaneous    AVERAGED PROBABILITY: 1.5

ConvLSTM (second row):
Action Predicted: baggage_handling
Confidence: 0.8300662636756897

LRCN Model (second row):
Action Predicted: baggage_handling
Confidence: 0.8003833293914795

**Optical frames:**

Single Frame CNN (third row):
CLASS NAME: misclaneous    AVERAGED PROBABILITY: 1e+02
CLASS NAME: bridge_connecting    AVERAGED PROBABILITY: 0.13
CLASS NAME: baggage_handling    AVERAGED PROBABILITY: 0.00011
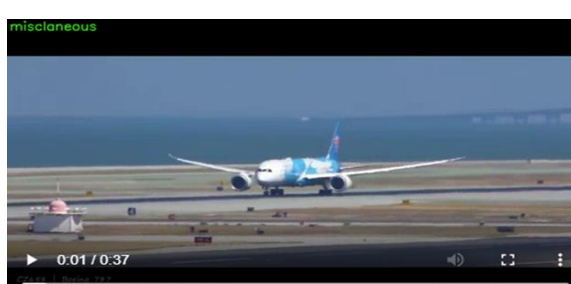
ConvLSTM (third row):
Action Predicted: misclaneous
Confidence: 0.9995587468147278

LRCN Model (third row):
Action Predicted: misclaneous
Confidence: 0.9992884397506714

# Questions?