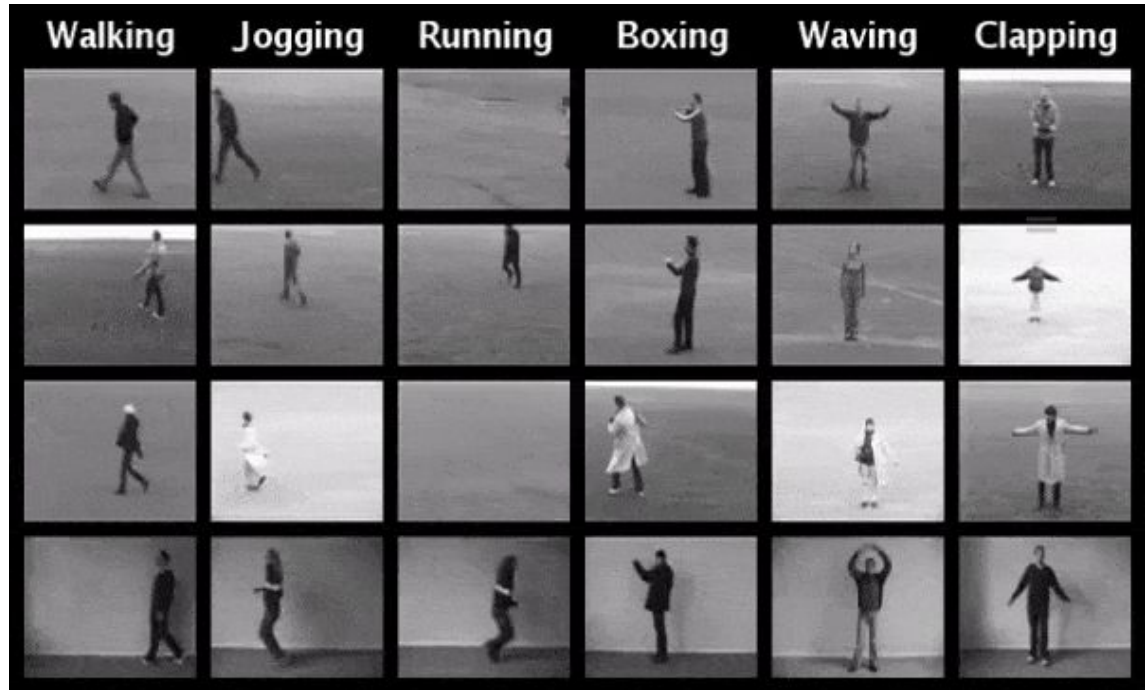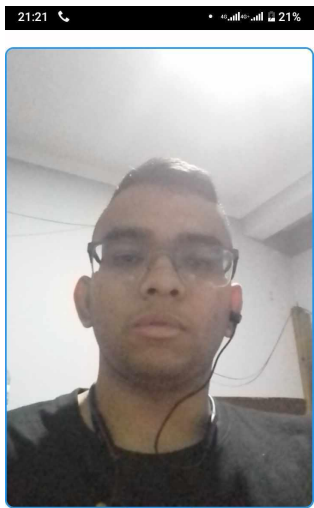# Human Action Recognition

Ashok Prasad Neupane

Anil Shrestha

Jeevan Neupane
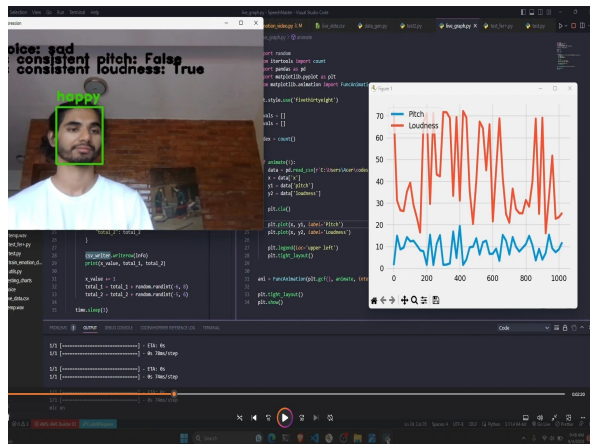
Contacts: neupane.ashok.9696@gmail.com

9818467416

Face Recognition using Siamese Network



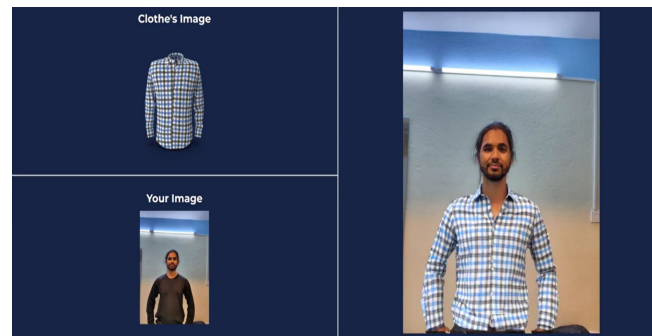Emotion classification on audio and video



Smoking detection using fine tuned YOLO v8



Studied about Full Self Driving Car(Tesla)



Viton(Virtual Try On)

# Research question

How can deep learning models be optimized to accurately and efficiently recognize human actions in real-time from RGB video sequences, considering challenges like occlusion, poor lighting, and motion variability?

Browse State-of-the-Art  Datasets  Methods  More ⌄

Sign In

▦ Videos

# Kinetics-700

☑ Edit

Introduced by Carreira et al. in A Short Note on the Kinetics-700 Human Action Dataset

Kinetics-700 is a video dataset of 650,000 clips that covers 700 human action classes. The videos include human-object interactions such as playing instruments, as well as human-human interactions such as shaking hands and hugging. Each action class has at least 700 video clips. Each clip is annotated with an action class and lasts approximately 10 seconds.
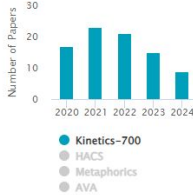
**Homepage**

## Benchmarks

☑ Edit

| Trend | Task | Dataset Variant | Best Model | Paper | Code |
|-------|------|-----------------|------------|-------|------|
| | Action Classification | Kinetics-700 | InternVideo2-6B | 📄 | ◯ |
| | Semantic Object Interaction Classification | Kinetics-700 | 3D ResNet-50 | 📄 | ◯ |
| | Image Clustering | Kinetics-700 | TURTLE | 📄 | ◯ |

## Papers

Search for a paper or author

## Usage △



● Kinetics-700
● HACS
● Metaphorics
● AVA

## License ⓘ

☑ Edit

🔗 Commons Attribution 4.0 International License

Feature extraction techniques like optical flow plus SVM classifier

3D Convolution

Multi-stream net

(a) Learned connectivity between blocks

Long term recurrent convolution net

Input    Visual        Sequence      Output
         Features      Learning

# Deep Learning Based Approach

Figure 2: The overall framework of InternVideo.

Intern video

V JEPA

Video MAE

# Video Foundation models

# Opportunities for us

# Block diagram

# Applications of HAR:



**Fraud Detection**

**Human-Computer Interaction**

**Theft Detections in Supermarkets**

**Health Care Surveillance**

**Crime Surveillance**

**Sports Analysis**

**Accident Detection**

**Video Processing**

# Timeline

## Human Action Recognition

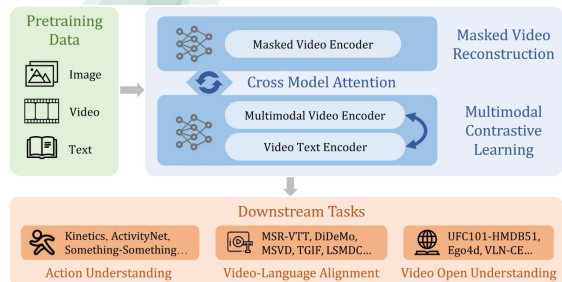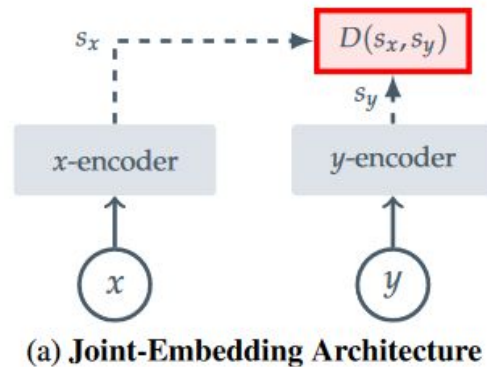| | Ashoj | Kartik | Mangsir | Push | Magh | Falgun |
|---|---|---|---|---|---|---|
| Literature review | ██ | ██ | | | | |
| Methodology selection | ██ | ██ | | | | |
| Data preprocessing | | | ██ | | | |
| Implementation | | | | | ██ | ██ |
| Training | | | | | ██ | ██ |
| Evaluation | | | | | ██ | ██ |
| Documentation | | | ██ | ██ | ██ | ██ |

# References

- **Human Action Recognition and Prediction: A Survey**: https://arxiv.org/abs/1806.11230
- **Revisiting 3D ResNets for Video Recognition**: https://arxiv.org/pdf/2109.01696v1
- **AssembleNet++: Assembling Modality Representations via Attention Connections**: https://arxiv.org/pdf/2008.08072v1
- **InternVideo**: http://arxiv.org/abs/2212.03191
- **VideoMAE**: http://arxiv.org/abs/2203.12602
- **Long Term Recurrent Convolutional Neural Network**: http://arxiv.org/abs/1411.4389
- **V-JEPA**: http://arxiv.org/abs/2301.08243
- **Dataset**: https://paperswithcode.com/dataset/kinetics
- **Shop Lifting dataset**: https://www.kaggle.com/datasets/mateohervas/dcsass-dataset
- **Sapiens**: https://about.meta.com/realitylabs/codecavatars/sapiens/