

Machine Learning Project

# AIR GUITAR

Team Ronaldo

Lee Hyunsong, Song Minui, Hyewon Cho



# INDEX

- 01      **Introduction - Project Overview**
- 02      **Dataset**
- 03      **Preprocessing**
- 04      **Training (with MLP)**
- 05      **Result & Demo**

## I 01. Introduction

# Problem Definition

Many people dream of enjoying the experience of playing the guitar, but the high cost and low portability of the instrument often make it difficult to start.

1. An expensive piece of equipment
  2. Large and inconvenient to carry
- Spontaneous playing and consistent practice challenging

## I 01. Introduction

# Air-Guitar

lower these barriers and allow *anyone* to enjoy guitar performance easily by creating an **air guitar system** that works with a standard Webcam

## < Project GOAL >

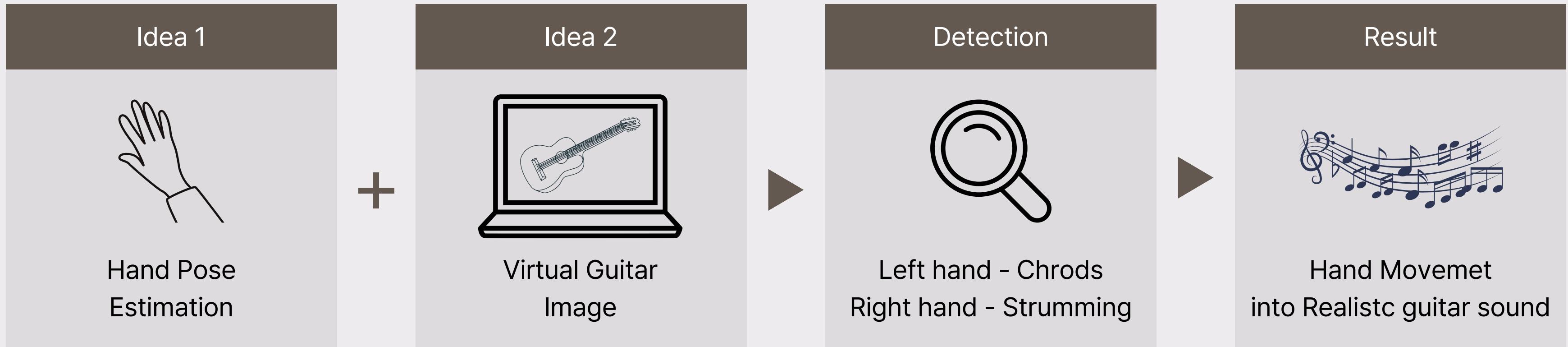
- *Deliver a fun!*
- *Highly accessible!*
- *Portable new style of musical interaction!*



*Image generated by ChatGPT*

# I 01. Introduction

## Air-Guitar



## I 02. Technical Point

### (1) Hand Pose Estimation

OpenCV function - MediaPipe Hands

To track 21 hand landmarks in real time, extracting 2D/3D joint coordinates for chord classification and strumming detection.



*Photo of the actual implementation*

### (2) Guitar Image Synthesis

Overlay a virtual guitar on the webcam feed, updating its position and orientation in real time based on both hands' locations.



## I 02. Technical Point

### (3) Chord Classification - *Left Hand*

#### Coordinate-based classification

Use landmark coordinates to compute distances and angles, then train a **lightweight ML classifier**.

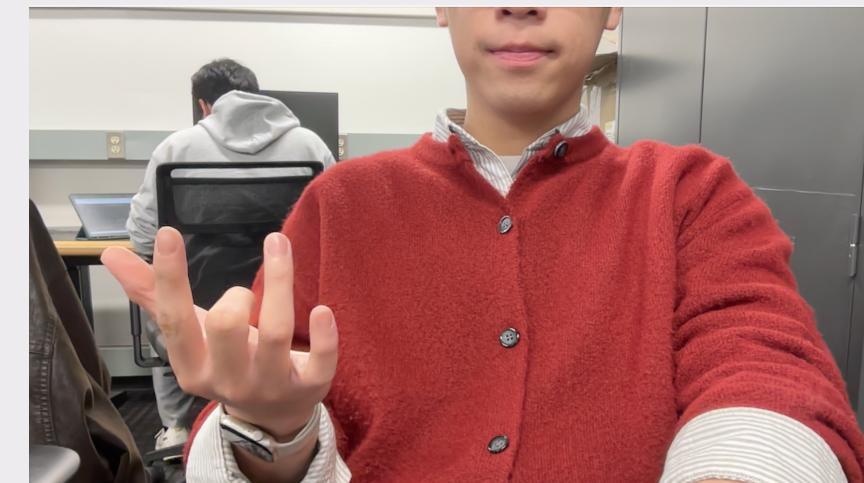
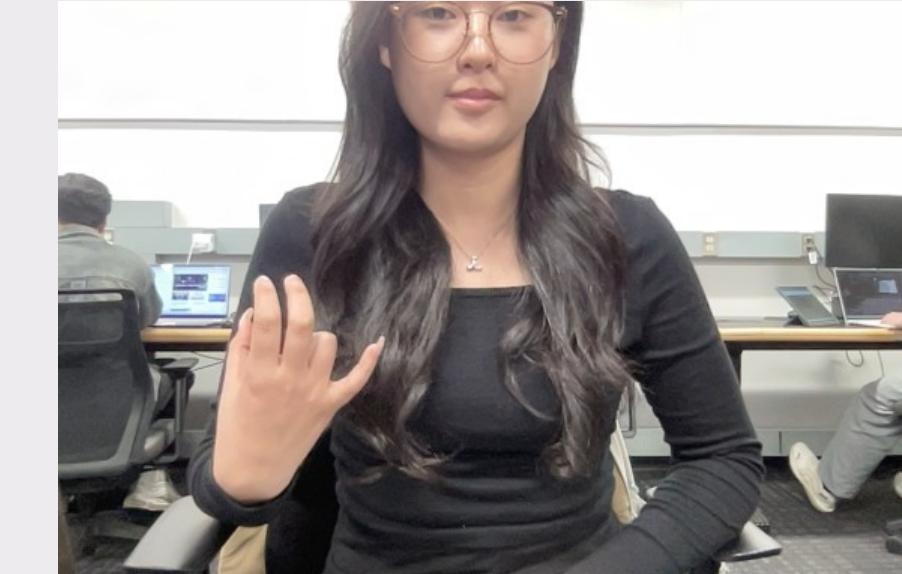
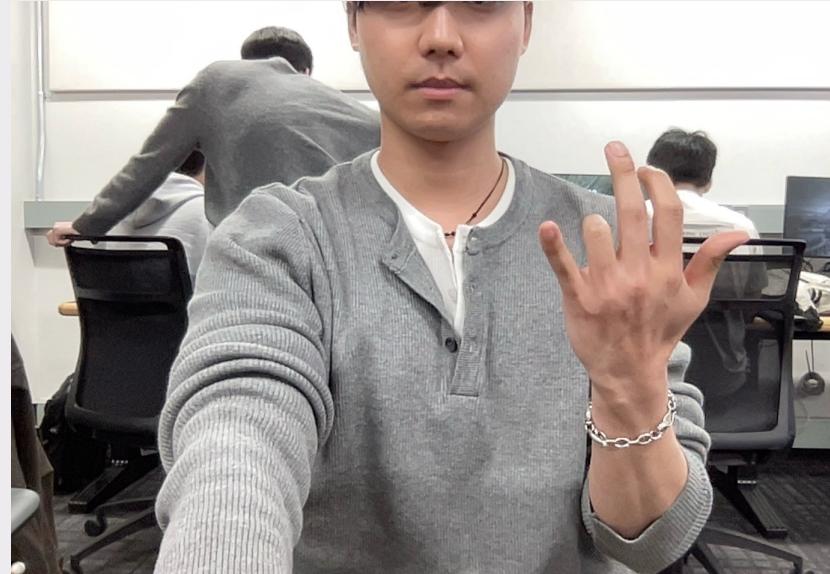
### (4) Strumming Detection - *Right Hand*

We detect strumming by tracking the right hand's vertical movement across frames. Fast, consistent motion is recognized as a strumming event, and optical flow is used to distinguish intentional strums from other hand movements.

## I 02. Dataset

### 1. Collect 60 code samples from each team member.

3 chord types × 60 samples × 3 members = 540 total samples.



## I 03. Preprocessing

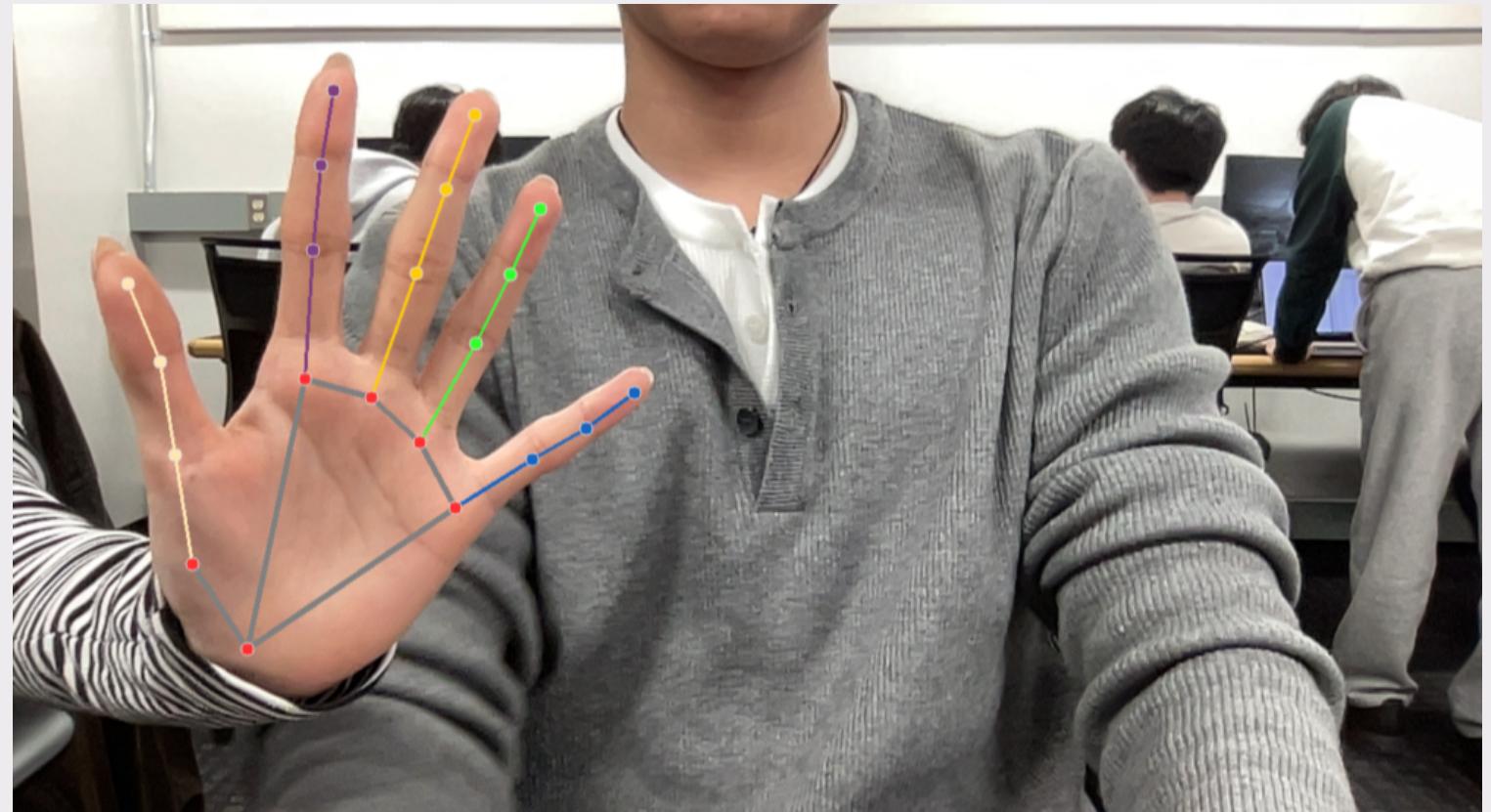
### 2. Extract hand landmark coordinates using MediaPipe & Save them to a CSV file

#### 1. Origin Alignment:

Sets the wrist (Landmark 0) as the origin (0, 0). This ensures the model focuses on the hand's shape rather than its absolute position in the image.

#### 2. Scaling:

Adjusts for the hand's distance from the camera. By dividing by the bounding box size, all landmark coordinates are normalized to a 0 to 1 range, making the data size-invariant.



# I 03. Preprocessing Result

## 2. Extract hand landmark coordinates using MediaPipe and save them to a CSV file

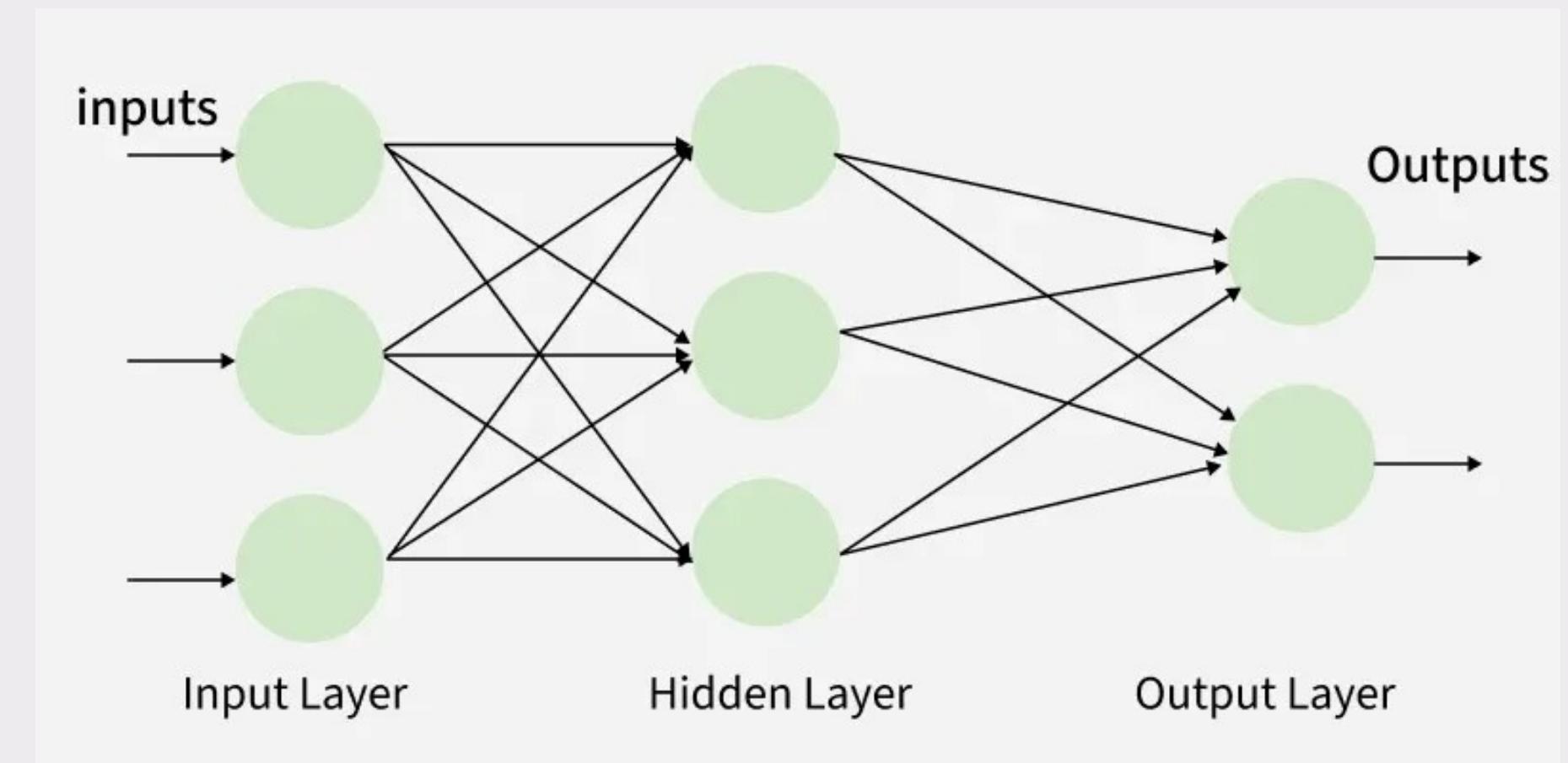
3 coordinates for each finger joint and 6 for the palm  
→ a total of 21 coordinates per sample.

# I 04. Model Training

## (1) MLP

*Multi-Layer Perceptron*

An MLP is a type of artificial neural network composed of multiple layers of neurons. It uses nonlinear activation functions to learn complex patterns and relationships in data.



## I 04. Model Training

### 3. Train an **MLP** model using the coordinates stored in the CSV file

- Hidden layers: (128, 64), (256, 128, 64), (512, 256, 128, 64)
- Dataset split
  - Train : Validation : Test = 8 : 1 : 1

```
pipe = Pipeline([  
    ("scaler", StandardScaler()),  
    ("mlp", MLPClassifier(  
        hidden_layer_sizes=tuple(args.hidden),  
        activation="relu",  
        solver="adam",  
        alpha=1e-4,  
        batch_size=64,  
        learning_rate_init=1e-3,  
        max_iter=500,  
        early_stopping=True,  
        validation_fraction=1/9, # train:val:test=8:1:1  
        n_iter_no_change=20, # 20 epoch 연속으로 향상되지 않으면, 학습이 500  
        에포크에 도달하지 않았더라도 훈련을 중단  
        random_state=args.seed,  
        verbose=True  
    ))  
])
```

# | 05. Result

Hidden Layer  
(128, 64)

[TEST] acc = 0.8367				
	precision	recall	f1-score	support
0	0.8667	0.8667	0.8667	15
1	0.8571	0.7500	0.8000	16
2	0.8000	0.8889	0.8421	18
accuracy			0.8367	49
macro avg	0.8413	0.8352	0.8363	49
weighted avg	0.8391	0.8367	0.8359	49
Confusion matrix:				
[[13 1 1] [ 1 12 3] [ 1 1 16]]				

83%

Hidden Layer  
(256, 128, 64)

[TEST] acc = 0.8571				
	precision	recall	f1-score	support
0	0.8462	0.7333	0.7857	15
1	0.8235	0.8750	0.8485	16
2	0.8947	0.9444	0.9189	18
accuracy			0.8571	49
macro avg	0.8548	0.8509	0.8510	49
weighted avg	0.8566	0.8571	0.8551	49
Confusion matrix:				
[[11 3 1] [ 1 14 1] [ 1 0 17]]				

86%

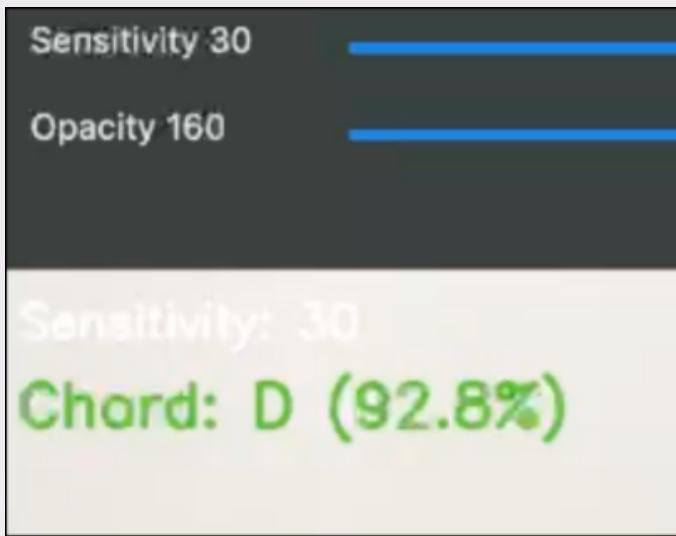
Hidden Layer  
(512, 256, 128, 64)

[TEST] acc = 0.9184				
	precision	recall	f1-score	support
0	1.0000	0.7333	0.8462	15
1	0.8421	1.0000	0.9143	16
2	0.9474	1.0000	0.9730	18
accuracy			0.9184	49
macro avg	0.9298	0.9111	0.9111	49
weighted avg	0.9291	0.9184	0.9150	49
Confusion matrix:				
[[11 3 1] [ 0 16 0] [ 0 0 18]]				

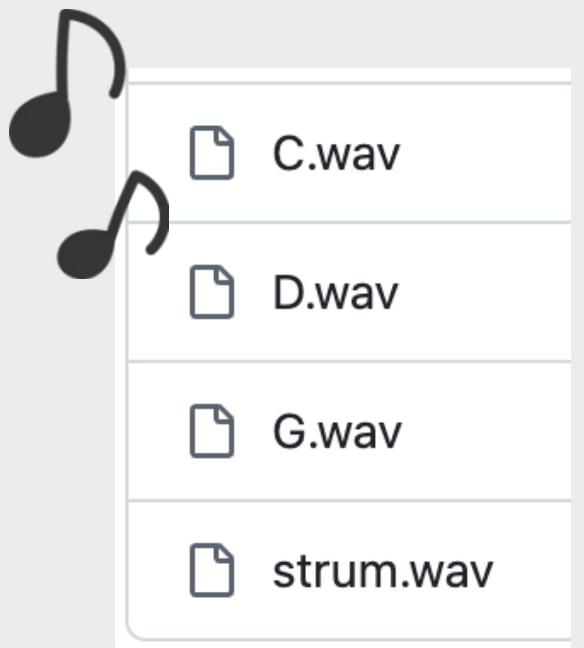
91%

# Program execution Screen

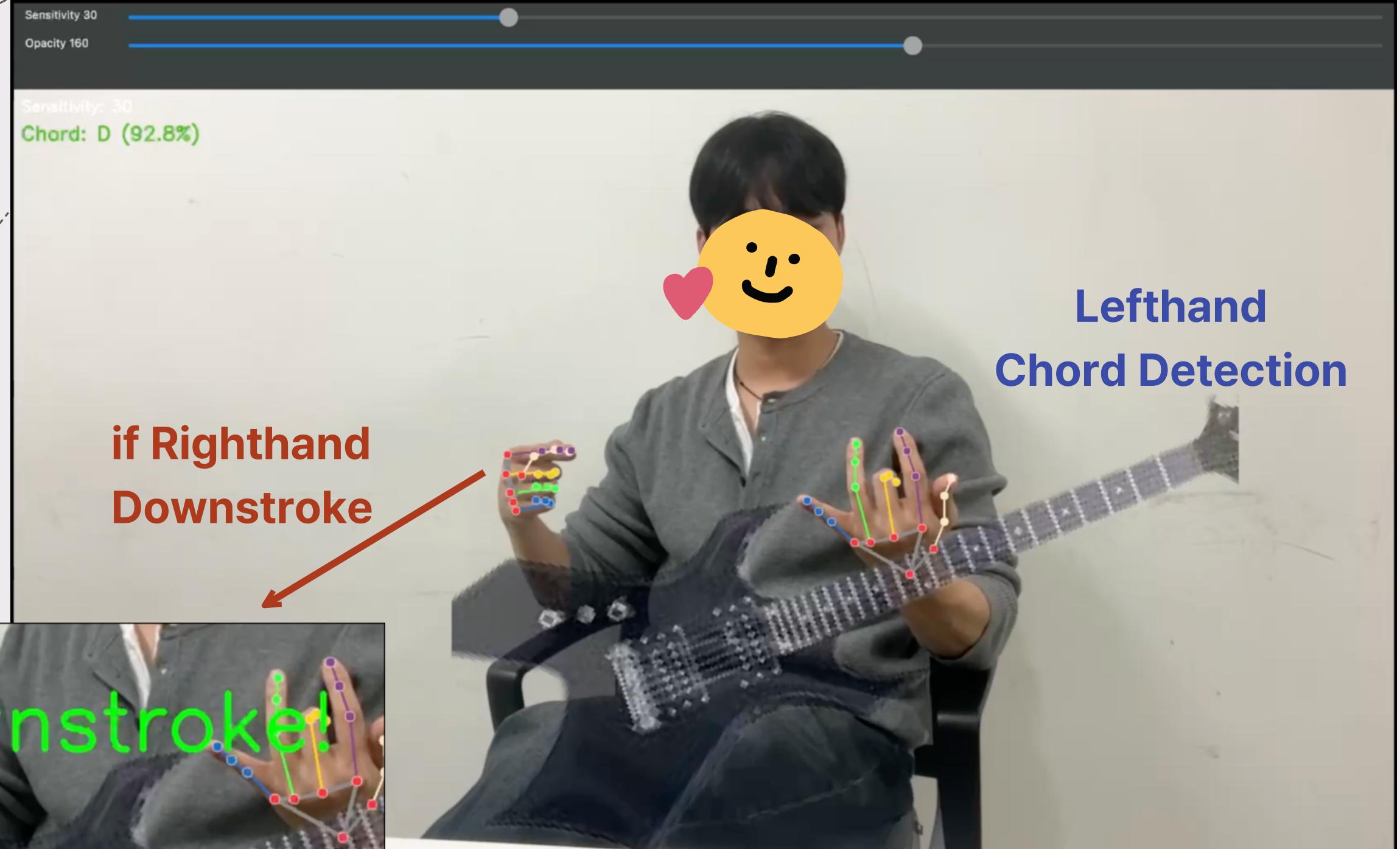
## Sensitivity & Opacity



## Chord : C / D / G Accuracy



if Righthand  
Downstroke



# I Challenges & Future Work

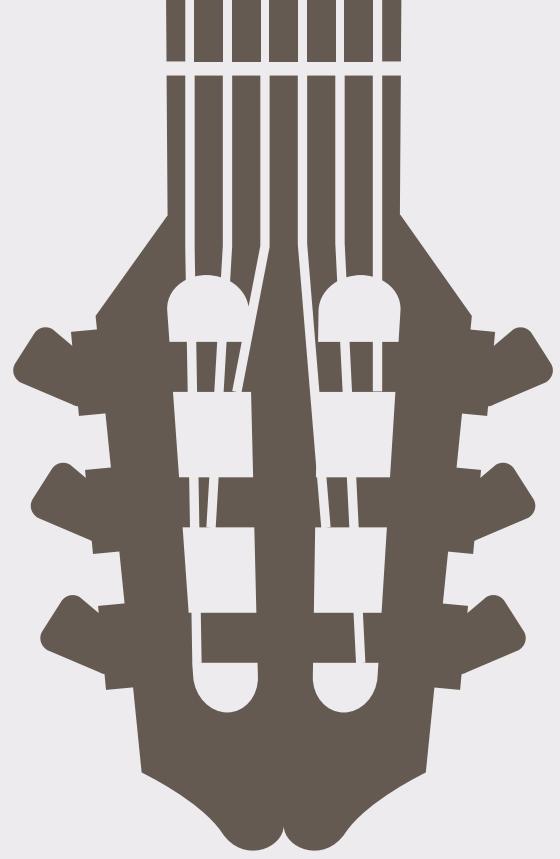
## < Challenge >

Since *no dataset existed* in the first place, we had to create it ourselves.

It's a bit disappointing because with a wider variety of chords and a larger amount of data, we could have built a much richer and more robust program.

## < Future Work >

- We would like to enable upstroke detection as well.
- We also want to implement support for recognizing multiple users simultaneously so that more than one person can play the guitar together.
- I would like to extend this to various instruments and create a real virtual band.



# Thank you

