

PetBot : A Facial Recognition, Pose Detection, and Ball Tracking

Professor Euntai KIM
T.A Seokwon CHOI, Hongchan JO

Team 15
Dokyu IM, Hokyun IM, Kyunghoon Jung

1

Introduction

Introduction(Pet Robot Market)

The market will be **ACCELERATING**
growing at a **CAGR** of almost

21%

INCREMENTAL
GROWTH
\$ 1.14 bn
2018

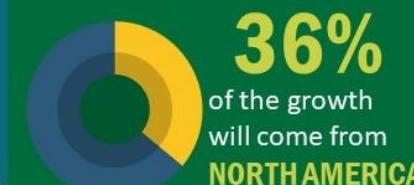
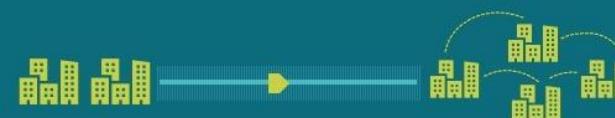
2023

The year-over-year growth rate
for **2019** is estimated at

15.80%

→ **CAGR 21%**
5years Growth
\$1.14 bn (1.6조 원)

The market is **CONCENTRATED** with few
players occupying the market share



One of the **KEY TRENDS** for this
market will be the **INTRODUCTION**
OF ECO-FRIENDLY ROBOTIC PET
DOGS



GLOBAL ROBOTIC PET DOGS MARKET 2022-2026

Market growth will **DECCELERATE**
at a **CAGR** of

11.28%



Incremental growth (\$M)

914.33



The market is **CONCENTRATED**
with few players occupying the
market

Growth Contributed by
NORTH AMERICA



35%

Growth for **2022**



16.34%

→ **CAGR 11.28%**

5years Growth
\$0.914 bn (1.3조 원)

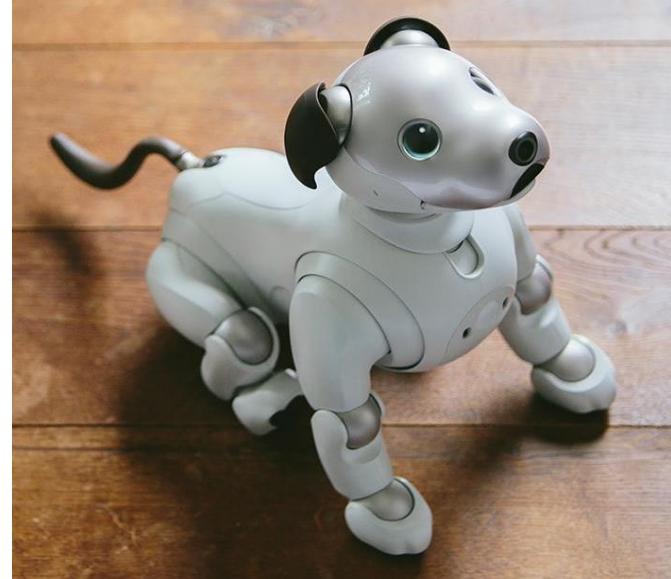
2

Related Works

Related Work



Hasbro's Joy('15.11, \$99)



Sony's AIBO('19.02, \$2,899)



Loona Pet Robot('23.02, \$299)



돌봄(CES 2019)



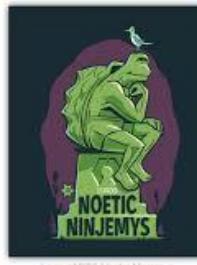
Guru IOT's PEDDY('19.06, \$299)

3

Robot Configuration

Software

Step 1



ROS Environment Setting

Step 2



kobuki

Only released in EOL d



Kobuki Development Setting

Step 3



Migrate to Ros1 Noetic

Step 4



JetPack



[pytorch/vision](#)

Datasets, Transforms and Models specific to Computer Vision



Install DL Libraries

Hardware



RGB / Depth
Camera



WIFI Module



Edge GPU
Device

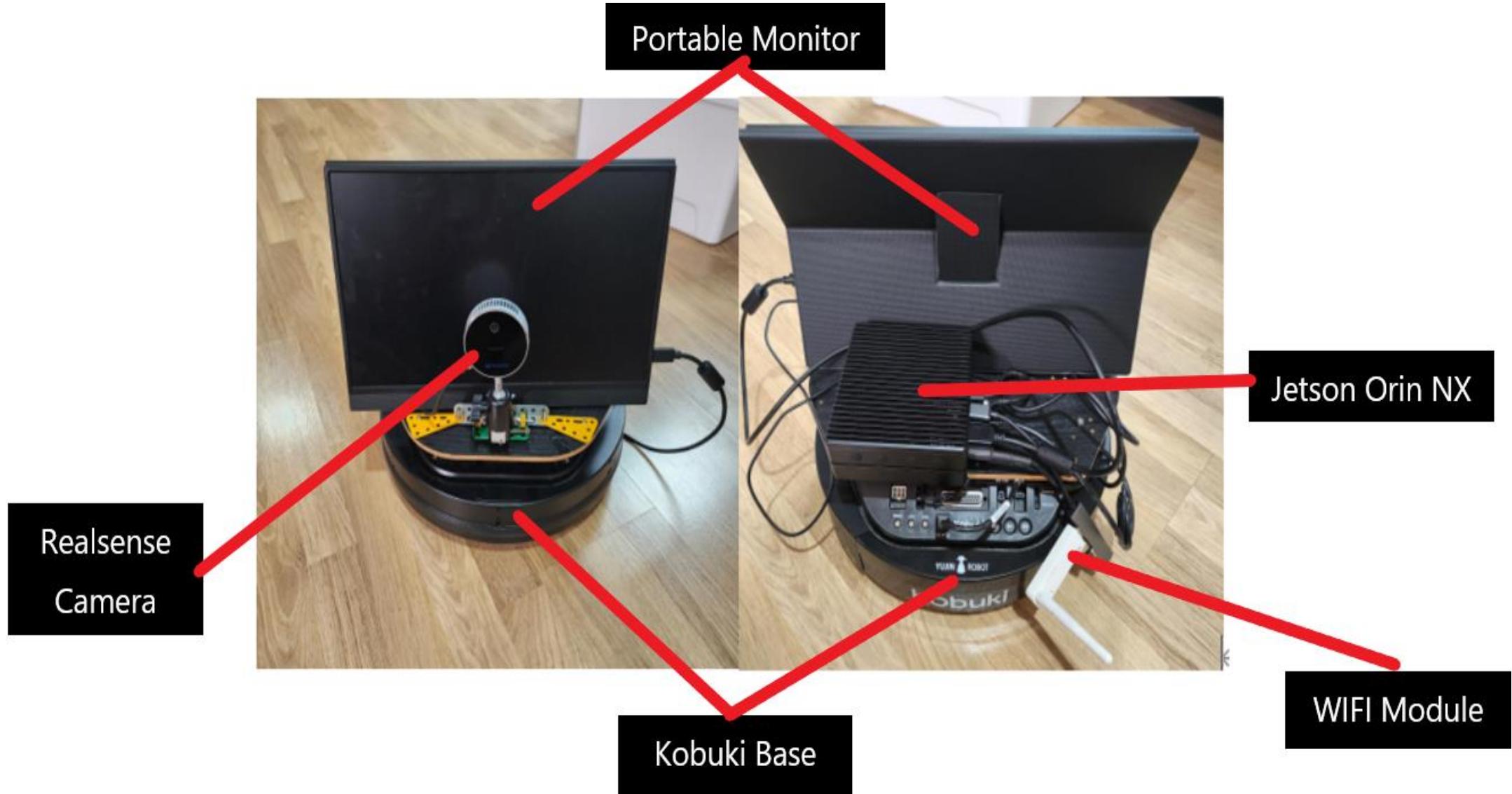


Robot Base



Portable Monitor

Hardware



4

Function Implementation

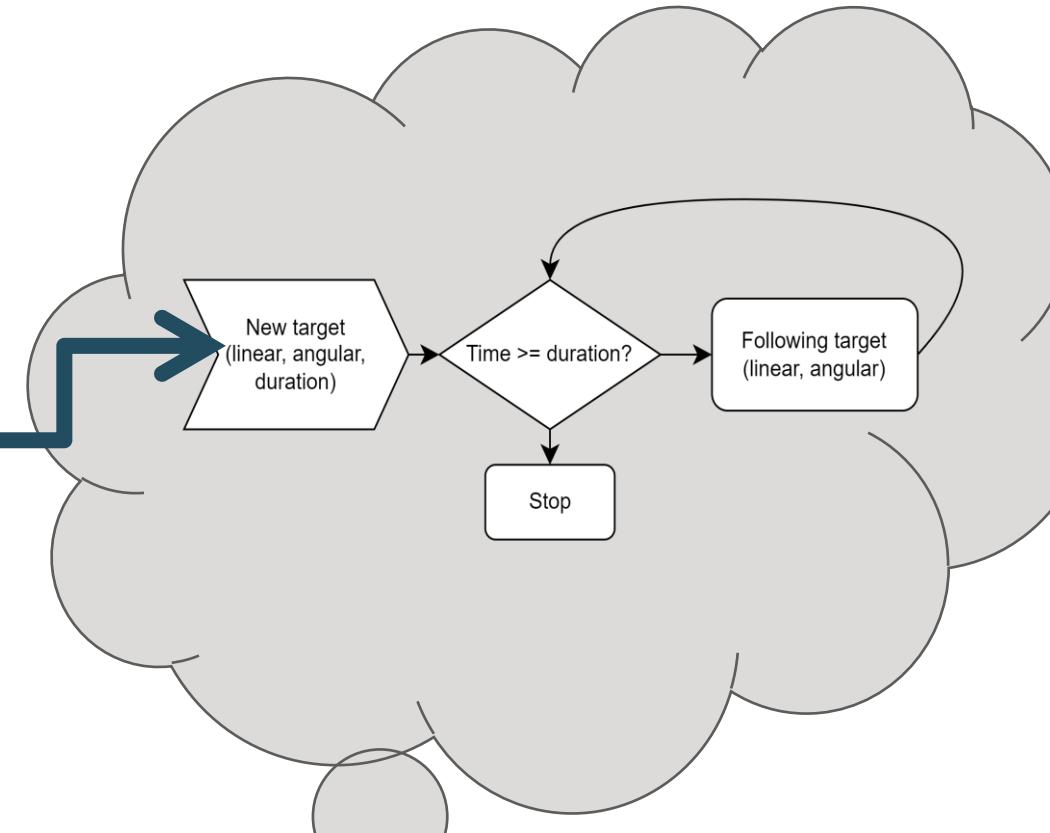
Basic Movements



High-Level Decision

Movement
Command

< 1.0m/s, 0.01rad/s, 2.0s >
<Linear, Angular, Duration>



Low-Level
Controller



Move!

Kobuki API



For the safe development and operation, robot needs fail safe.



If a hazard is detected, this must be operated before any code.

Wall Detection



Wall Detection: Move-Collision-GoBack-Turn←

Cliff Detection



Cliff Detection: Move-Detect-GoBack-Turn←

Emergency Stop
(PickUp Detection)



PickUp Detection: Move-Wrong Move-PickUp-Terminated←



Real Time Pose Recognition/Classification

Pose Estimation Framework

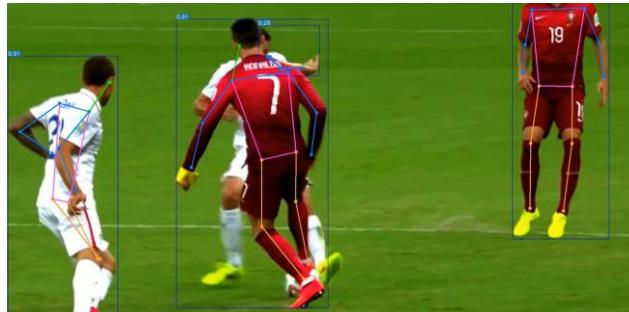
1. Yolo(You Only Look Once) Series
2. MediaPipe



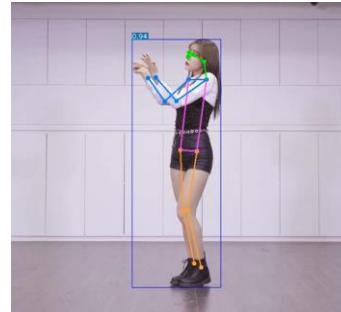
Mediapipe



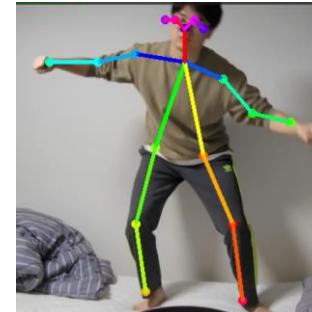
Yolo(You Only Look Once) Series



Test(Yolov7) : Football video



Test(Mediapipe) : Dance video



Test(Open-pose) : Realtime webcam

Measurement Environments

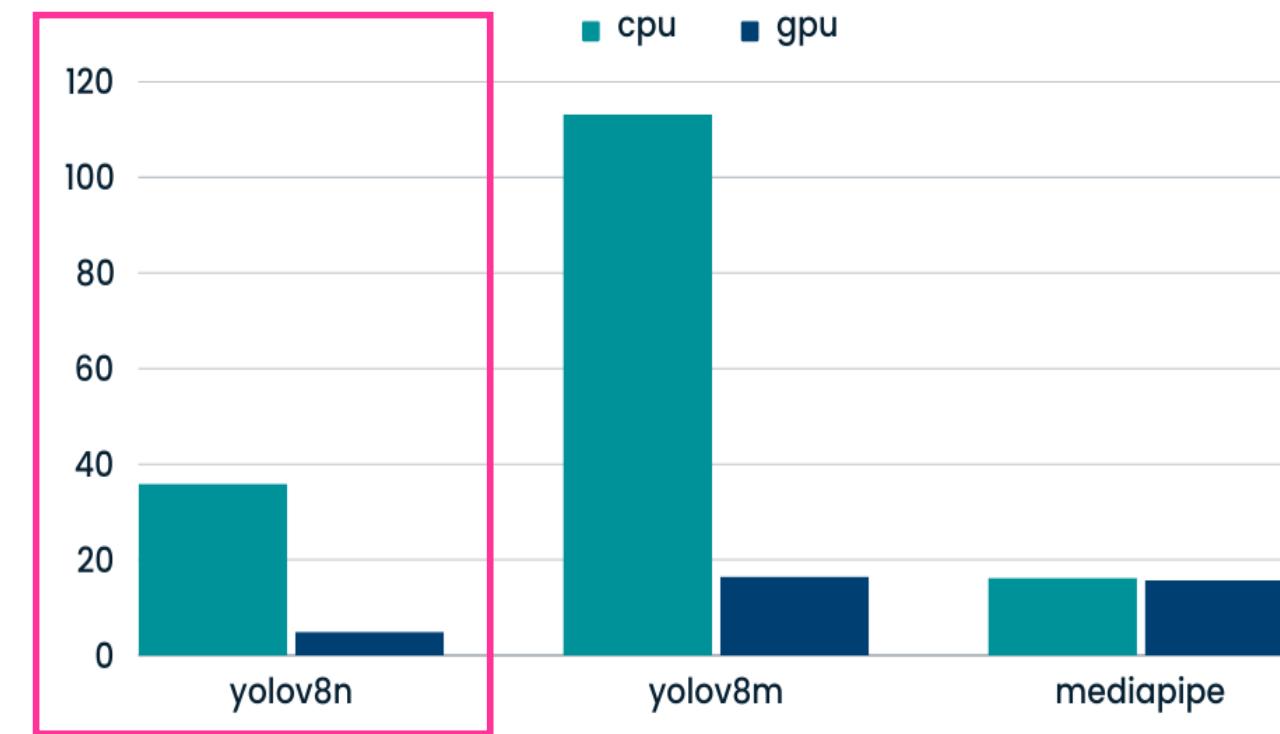
Input Size : 640 x 480

CPU/GPU WarmUp : Using Dummy Example(About 100)

CPU : 13th Gen Intel(R) Core™ i7-13700 2.10GHz
GPU : RTX 3060ti

| | cpu | gpu |
|-----------|--------|-------|
| yolov8n | 35.94 | 4.94 |
| yolov8m | 113.18 | 16.49 |
| mediapipe | 16.25 | 15.7 |

Result Table



Result Graph

Data Collecting 1st

*3 Poses x 3 Person x 10 images(multi-view) = **93 images(Real World)**

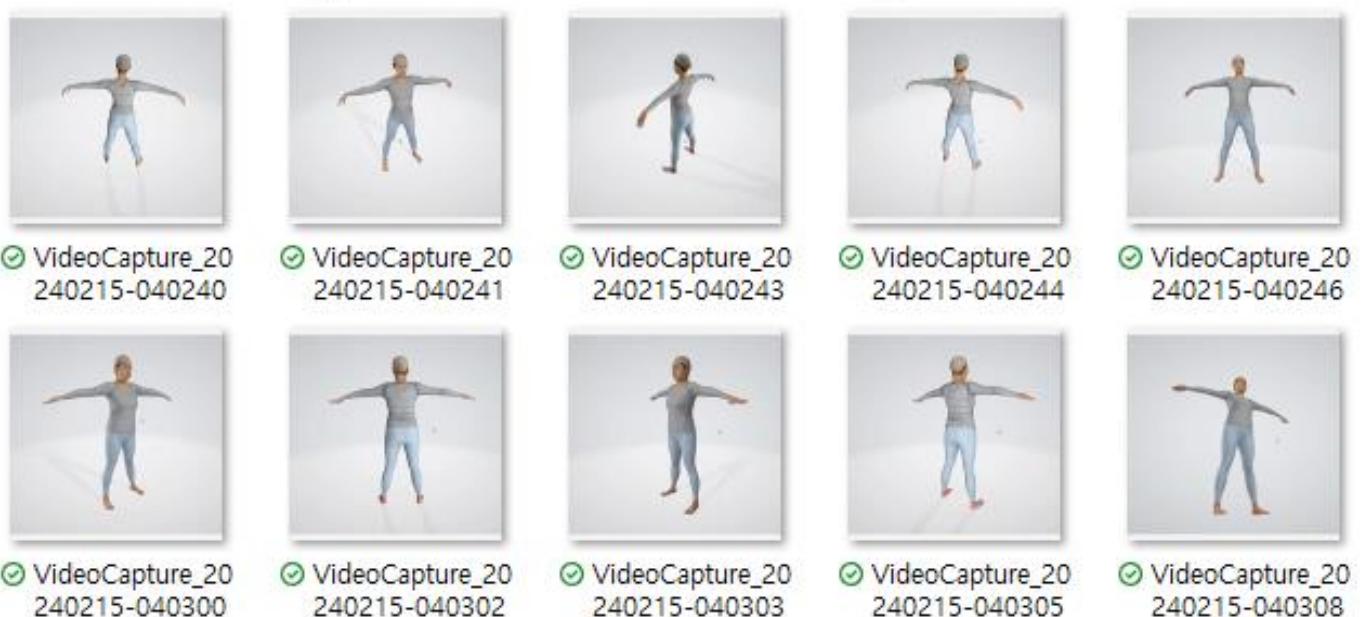


*Attention Position(Stop)/Arm Stretch(Start)/Sit(Spin)

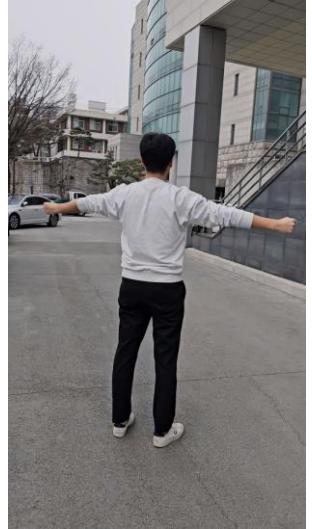
Korean 3D Scan Big Data :
3 Poses x 10 Person x 9 images = **274 images**



Total 367 images



Data Collecting 2nd : Data Augmentation



Data Augmentation



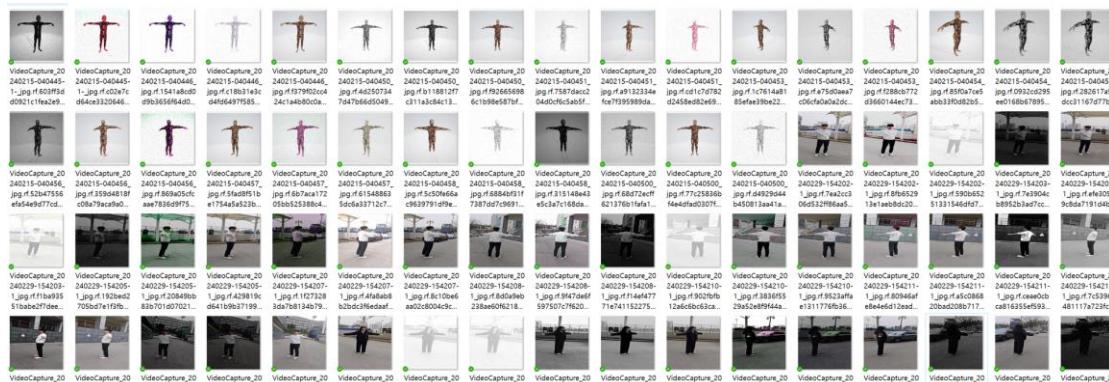
Horizontal Stretching



Adding Noise & Color Filtering



Adding Noise & Removing Color



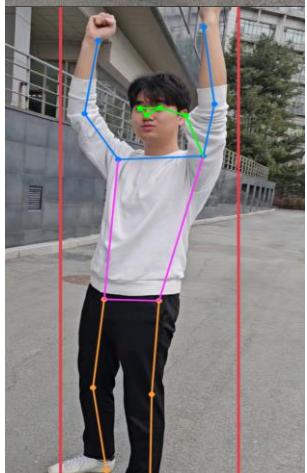
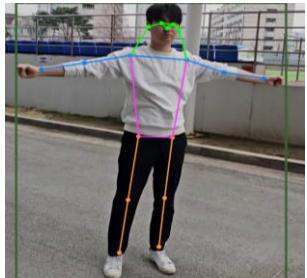
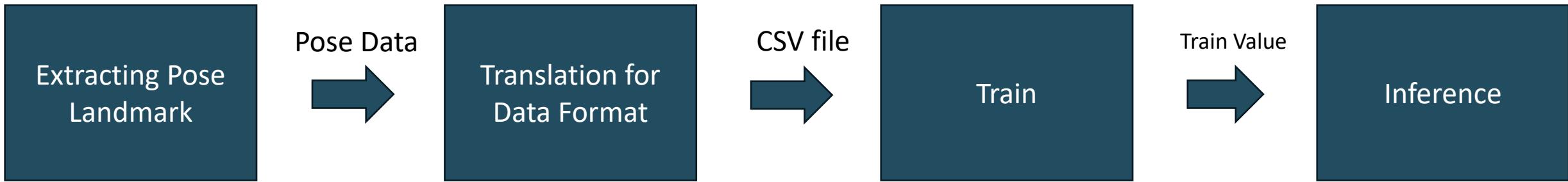
TestSet

(367 – 60) x 3.14 times = Total 966 images

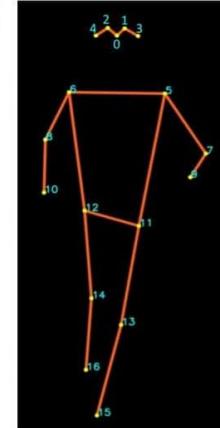
Data 1st

Augmentation

Pipeline



| Index | Key point |
|-------|----------------|
| 0 | Nose |
| 1 | Left-eye |
| 2 | Right-eye |
| 3 | Left-ear |
| 4 | Right-ear |
| 5 | Left-shoulder |
| 6 | Right-shoulder |
| 7 | Left-elbow |
| 8 | Right-elbow |
| 9 | Left-wrist |
| 10 | Right-wrist |
| 11 | Left-hip |
| 12 | Right-hip |
| 13 | Left-knee |
| 14 | Right-knee |
| 15 | Left-ankle |
| 16 | Right-ankle |



| | BU | BV | BW | BX | BY | BZ | CA | | | |
|----|---------------|---------------|---------------|--------------|----------|----------|-----------|---------|----------|---|
| 1 | RIGHT_PINKY_X | RIGHT_PINKY_Y | RIGHT_PINKY_Z | -0.827996503 | -1.75142 | 0.986837 | 0.107401 | -0.8056 | -0.12362 | 0 |
| 2 | -0.0060938201 | -0.830895322 | -1.36522 | 0.986289 | 0.201143 | -0.80967 | -0.19008 | 0.6 | | |
| 3 | 0.000661393 | -0.830895322 | -1.36522 | 0.986289 | 0.201143 | -0.80967 | -0.19008 | 0.6 | | |
| 4 | 0.011974001 | -0.829517011 | -1.33139 | 0.985127 | 0.214762 | -0.81581 | -0.02978 | 0.6 | | |
| 5 | 0.061895834 | -0.830893198 | -1.21958 | 0.984827 | 0.265268 | -0.80584 | -0.105 | 0.6 | | |
| 6 | 0.083531456 | -0.846780017 | -1.2752 | 0.984739 | 0.298476 | -0.80092 | -0.27082 | 0.6 | | |
| 7 | 0.089020216 | -0.845600609 | -1.60286 | 0.984044 | 0.361031 | -0.81022 | -0.43974 | 0 | | |
| 8 | 0.09586136 | -0.873492843 | -1.53807 | 0.983557 | 0.397564 | -0.82125 | -0.64119 | 0 | | |
| 9 | 0.098163472 | -0.900626347 | -1.62103 | 0.983494 | 0.423056 | -0.83504 | -0.676917 | 0.6 | | |
| 10 | 0.096119163 | -0.903610154 | -1.37828 | 0.982563 | 0.432304 | -0.84286 | -0.94591 | 0.6 | | |
| 11 | -0.113924796 | -0.909183951 | -1.34089 | 0.979972 | 0.40636 | -0.8712 | -0.19005 | 0.6 | | |
| 12 | -0.180425932 | -0.890726069 | -1.36494 | 0.976898 | 0.338464 | -0.87062 | -1.41904 | 0.6 | | |
| 13 | -0.207748481 | -0.874283721 | -1.44103 | 0.97428 | 0.301401 | -0.8622 | -1.5258 | 0.6 | | |
| 14 | -0.238965673 | -0.875736291 | -1.4484 | 0.973952 | 0.252116 | -0.86426 | -1.62166 | 0.6 | | |
| 15 | -0.275716671 | -0.873891067 | -1.56368 | 0.973068 | 0.201943 | -0.86884 | -1.70093 | 0.6 | | |
| 16 | -0.287847896 | -0.877702753 | -1.5709 | 0.973586 | 0.184656 | -0.87212 | -1.73114 | 0.6 | | |
| 17 | -0.300571186 | -0.875959511 | -1.63006 | 0.973792 | 0.16047 | -0.88129 | -1.76289 | 0.5 | | |
| 18 | -0.3159537456 | -0.87229568 | -1.57698 | 0.973864 | 0.13822 | -0.87971 | -1.69932 | 0 | | |
| 19 | -0.333023801 | -0.858728774 | -1.57139 | 0.973383 | 0.125046 | -0.87951 | -1.65213 | 0 | | |

CSV file



Train

Train Value



Inference

Optimizer : adam

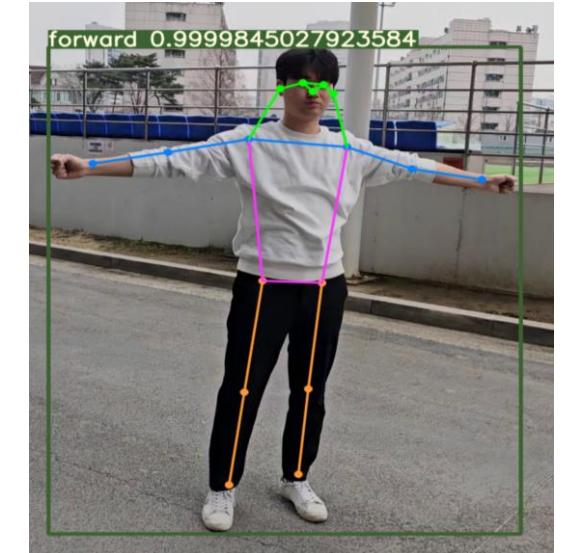
Loss : categorical crossentropy

Model : Sequential

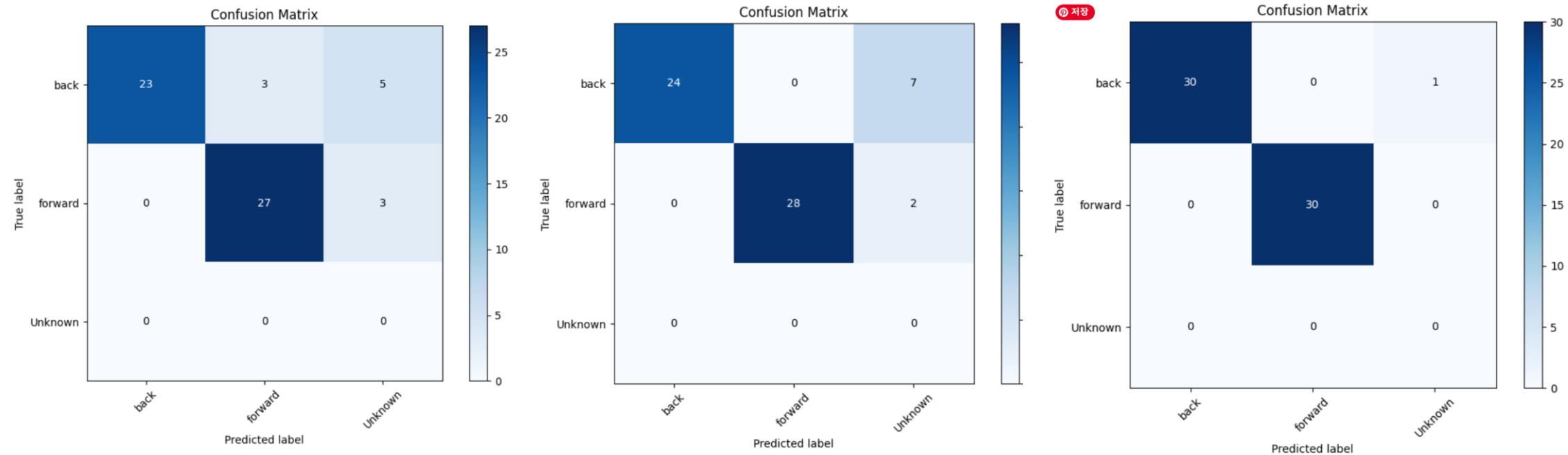
512, relu

256, relu

class_num, softmax



Result



Only 3D Scan Dataset

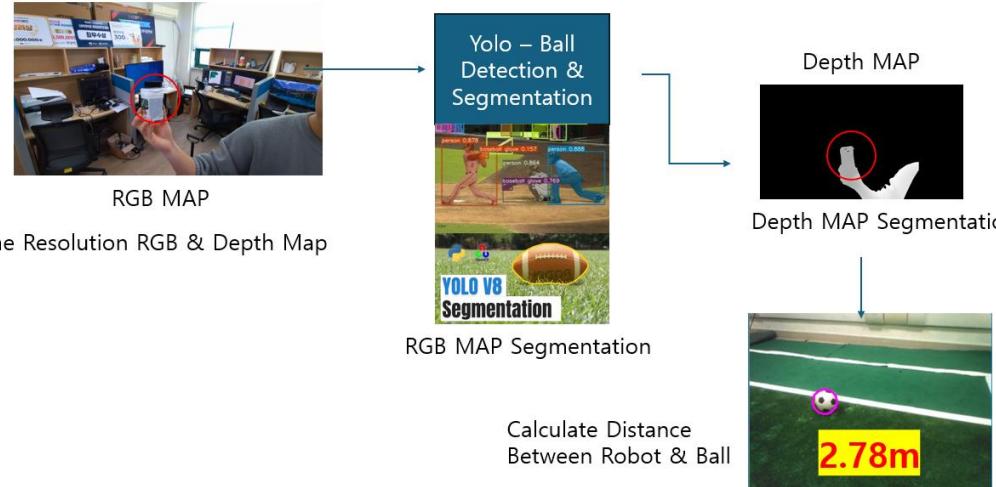
3D Scan + Real World Dataset

After Data Augmentation

| | Before Augmentation | After Augmentation |
|------------|------------------------------|--------------------|
| Accuracy | 81.97% | 100% |
| Error Rate | 18.03% | 0% |
| Precision | 100% (back), 90% (forward) | 100% |
| Recall | 74.19% (back), 90% (forward) | 100% |
| F1 Score | 85.19% (back), 90% (forward) | 100% |

Trial And Error1(RGB-D Map Processing)

Ball Tracking PipeLine



Same Resolution RGB & Depth Map

Trial And Error2(Marker Tracking)

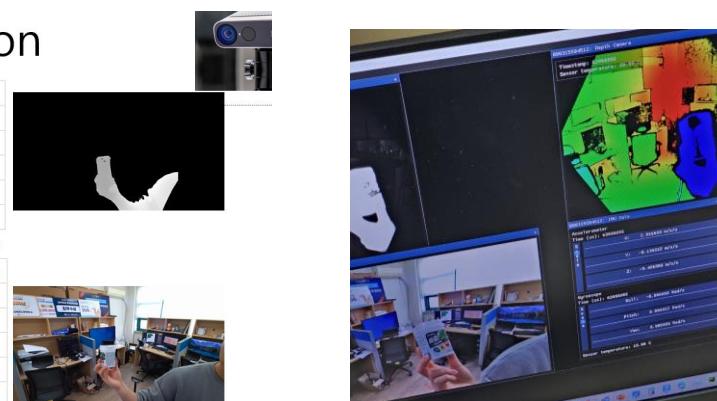


Aruco Marker Size is constant
So, 3D localization possible!

RGB-D MAP Super-resolution

| Mode | Resolution | FoV | FPS | Operating range* | Exposure time |
|----------------------|------------|-------------|--------------|------------------|---------------|
| NFOV unbinne | 640x576 | 75°x65° | 0.5, 15, 30 | 0.5 - 3.86 m | 12.8 ms |
| 320x288 | 75°x65° | 0.5, 15, 30 | 0.5 - 5.46 m | | 12.8 ms |
| WFOV 2x2 binned (SW) | 512x512 | 90°x120° | 0.5, 15, 30 | 0.25 - 2.88 m | 12.8 ms |
| WFOV unbinne | 1024x1024 | 20°x120° | 0.5, 15 | 0.25 - 2.21 m | 20.3 ms |
| Passive IR | 1024x1024 | N/A | 0.5, 15, 30 | N/A | 1.6 ms |

| RGB Camera Resolution (HxV) | Aspect Ratio | Format Options | Frame Rates (FPS) | Nominal FOV (HxV) (post-processed) |
|-----------------------------|--------------|-----------------|-------------------|------------------------------------|
| 3840x2160 | 16:9 | MJPEG | 0.5, 15, 30 | 90°x59° |
| 2560x1440 | 16:9 | MJPEG | 0.5, 15, 30 | 90°x59° |
| 1920x1080 | 16:9 | MJPEG | 0.5, 15, 30 | 90°x59° |
| 1280x720 | 16:9 | MJPEG/YUY2/NV12 | 0.5, 15, 30 | 90°x59° |
| 4096x2072 | 4:3 | MJPEG | 0.5, 15 | 90°x74.3° |
| 2048x1536 | 4:3 | MJPEG | 0.5, 15, 30 | 90°x74.3° |



Can we detect ball in 3d space?
With Constant Ball Size?

Result

$$\text{Distance} = \frac{R_{\text{real}} \times f \times W_{\text{image}}}{R_{\text{img}} \times W_{\text{sensor}}}$$

W_{image} : Width of the image

R_{real} : Actual radius of the Ball(0.15m)

R_{img} : Radius of the ball detected in the image

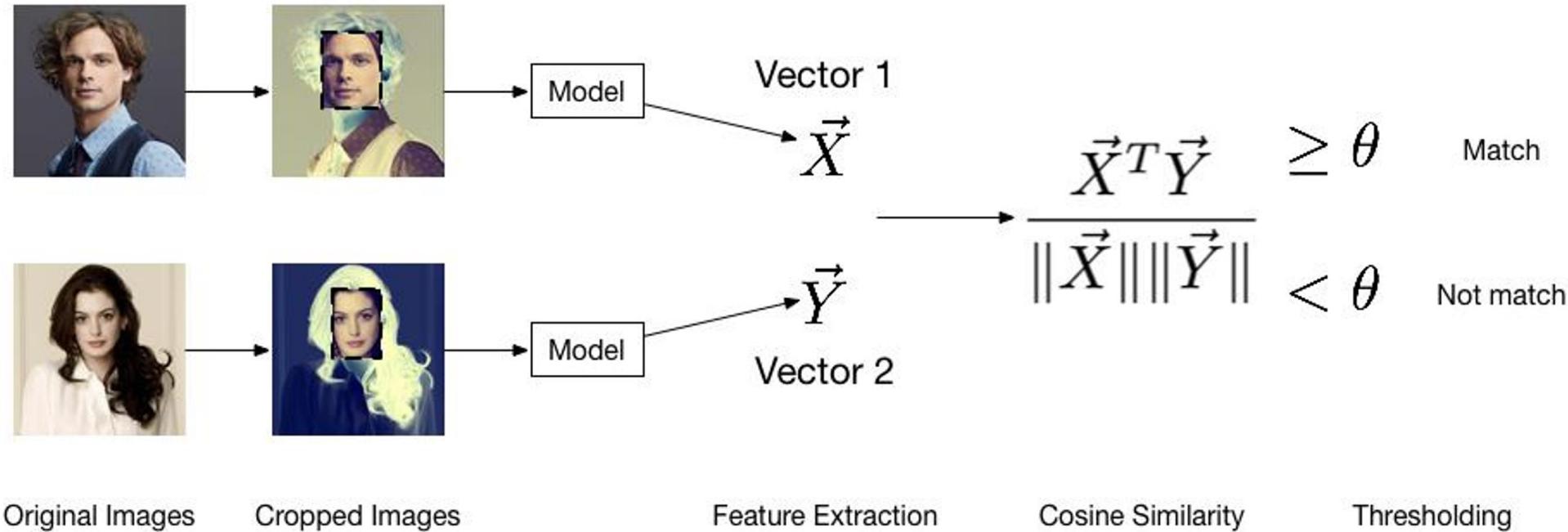
f : Focal length of the camera

W_{sensor} : Width of the camera sensor (0.036m, typical for a 35mm camera sensor)

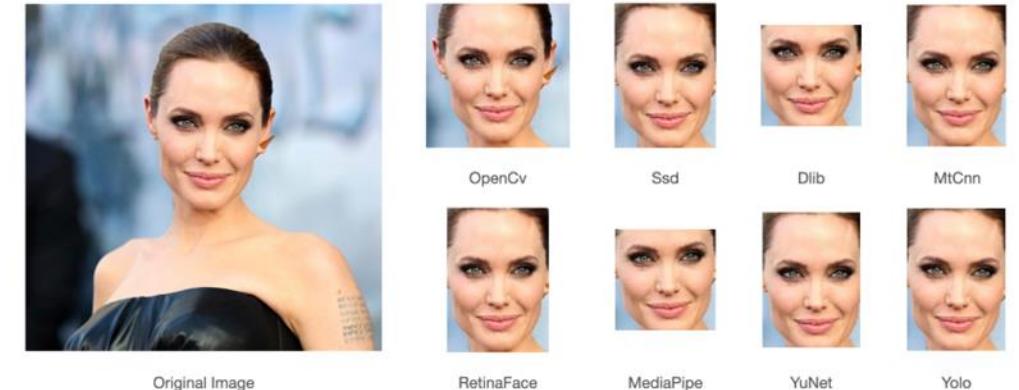
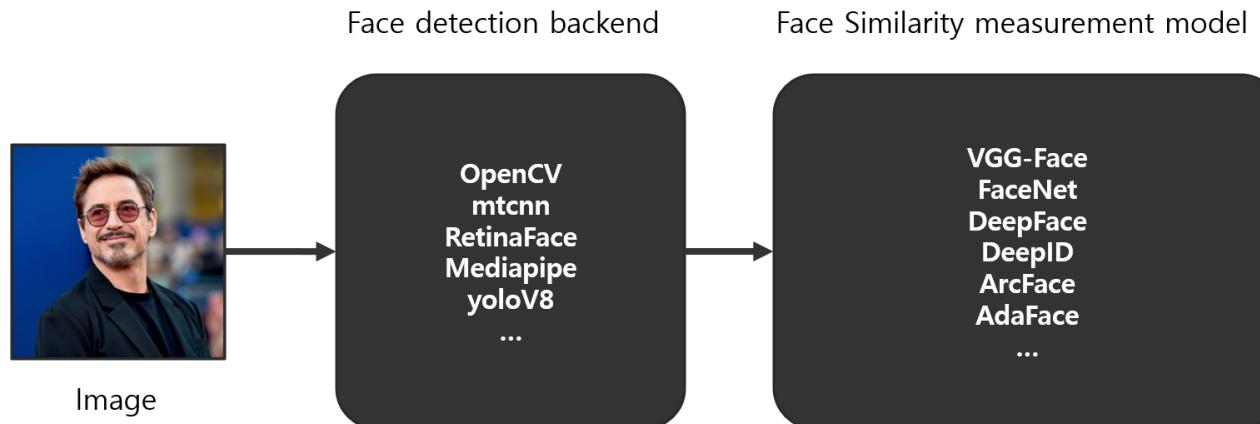
D : Distance between the camera and the ball (in meters)



Face Recognition Pipeline

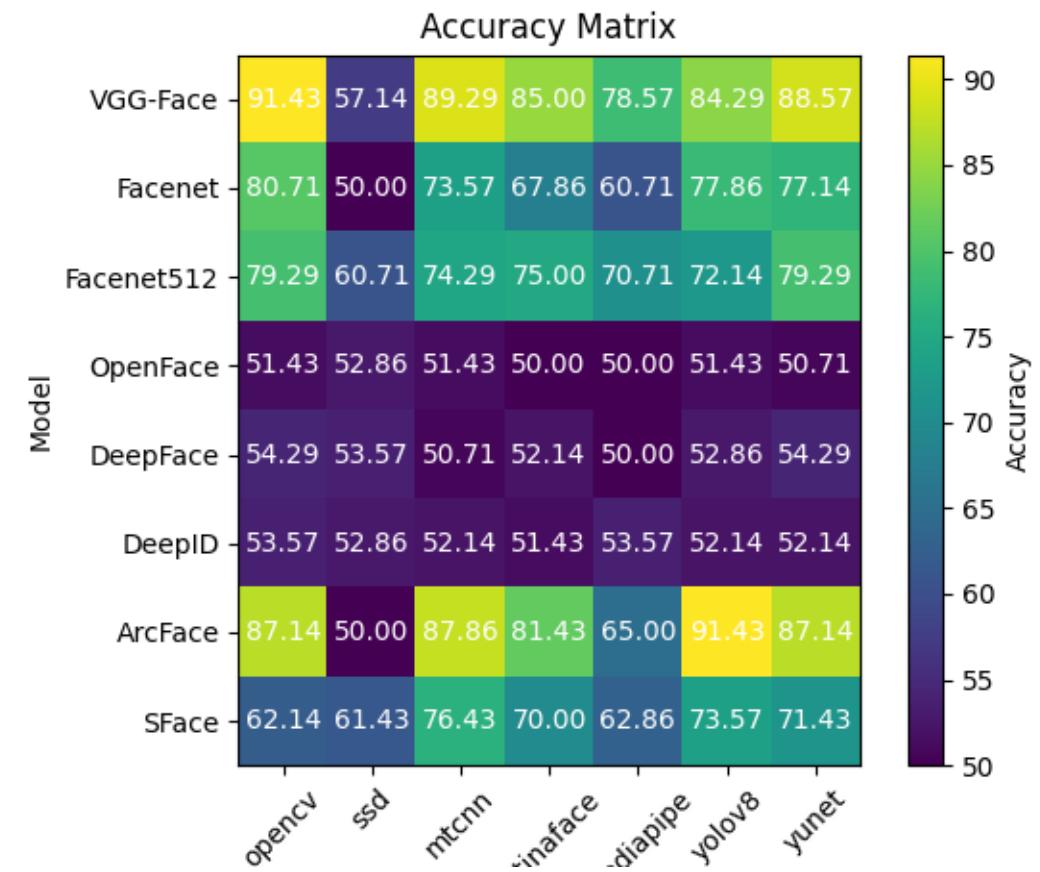


Face Recognition Pipeline

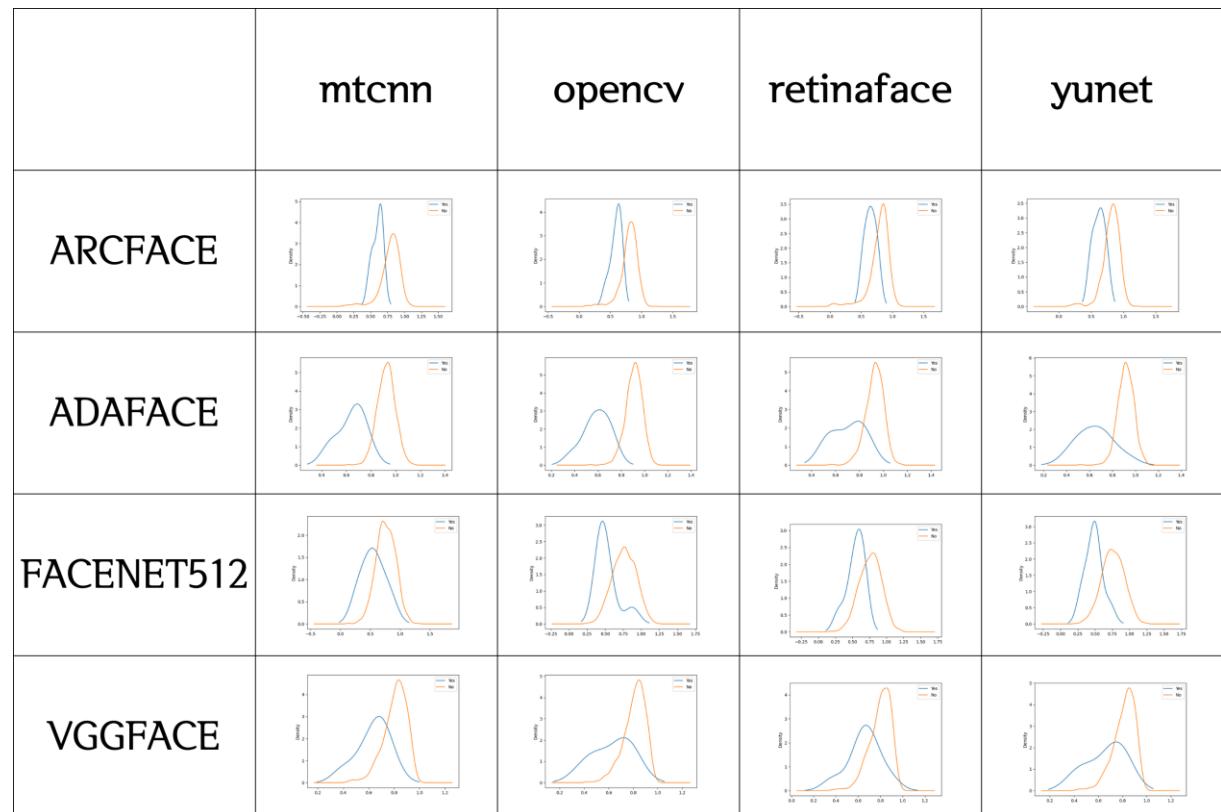


| Model | LFW Score | YTF Score |
|---------------------|-----------|-----------|
| Facenet512 | 99.65% | - |
| SFace | 99.60% | - |
| ArcFace | 99.41% | - |
| Dlib | 99.38 % | - |
| Facenet | 99.20% | - |
| VGG-Face | 98.78% | 97.40% |
| <i>Human-beings</i> | 97.53% | - |
| OpenFace | 93.80% | - |
| DeepID | - | 97.05% |

Quantitative Evaluation of Face Recognition Detection Model Combination

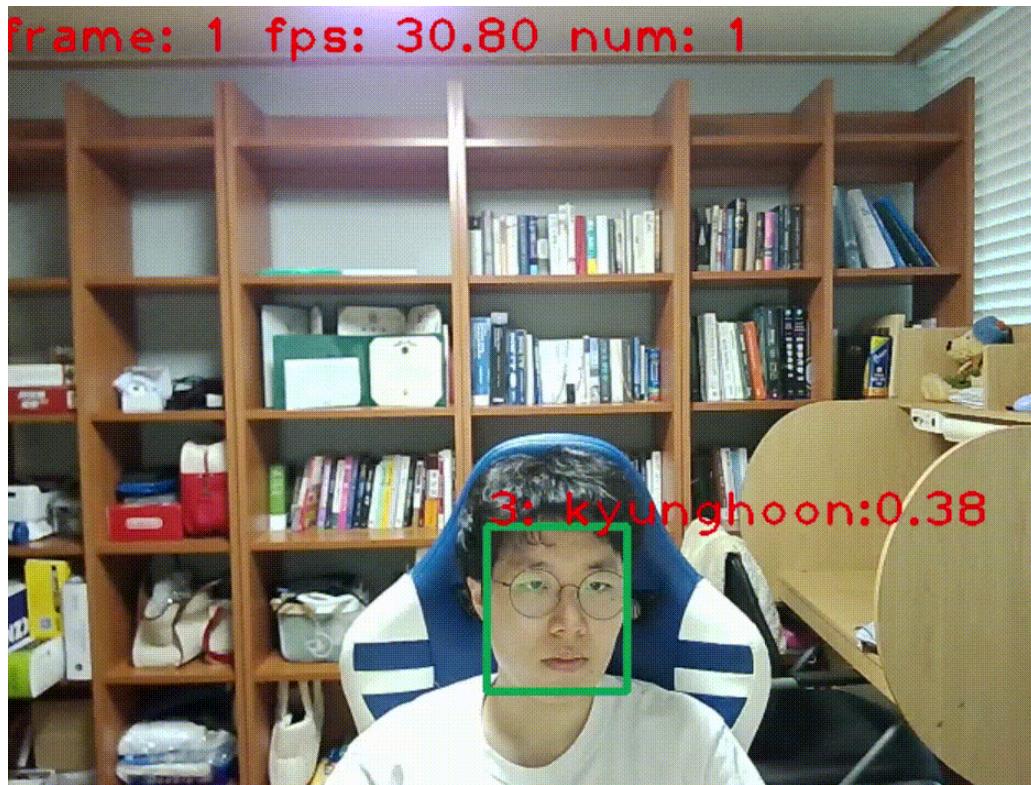


Qualitative Evaluation of Face Recognition Detection Model Combination

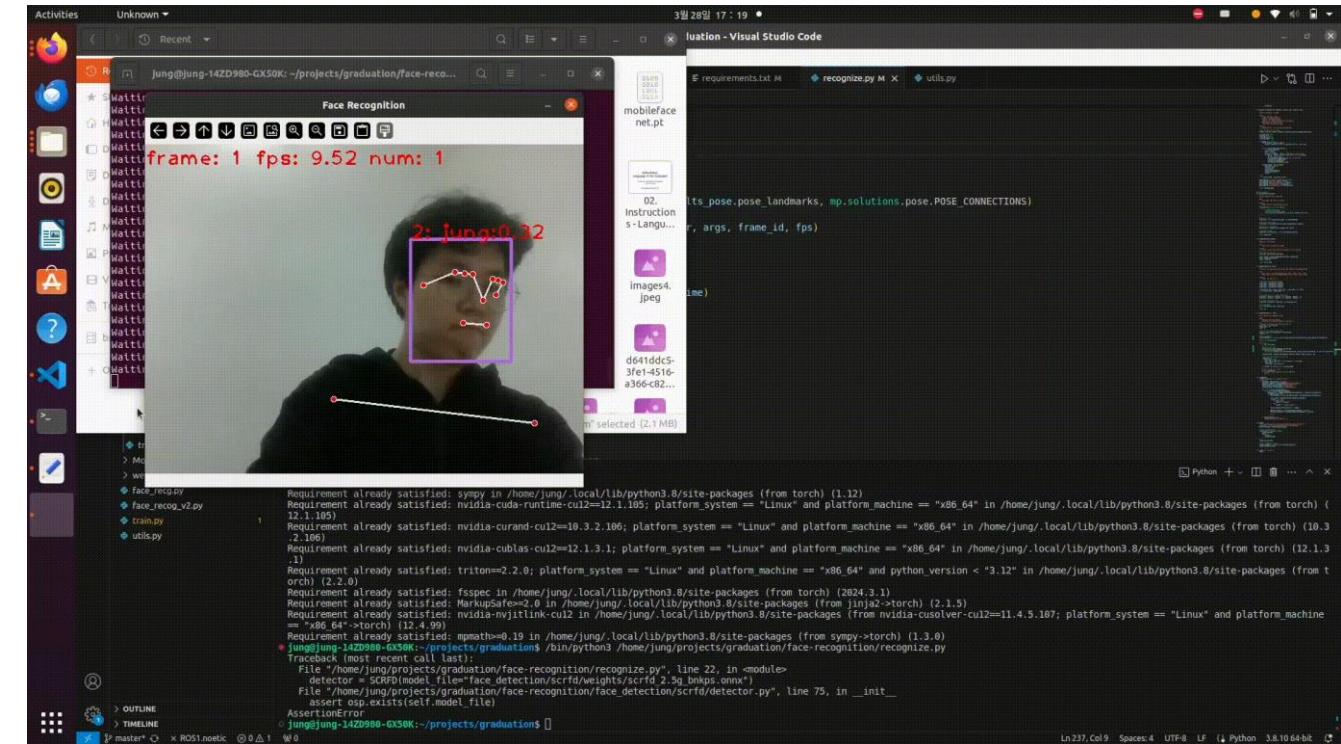


Real Time Face Recognition

- RTX 4090 (fps 30 or above)

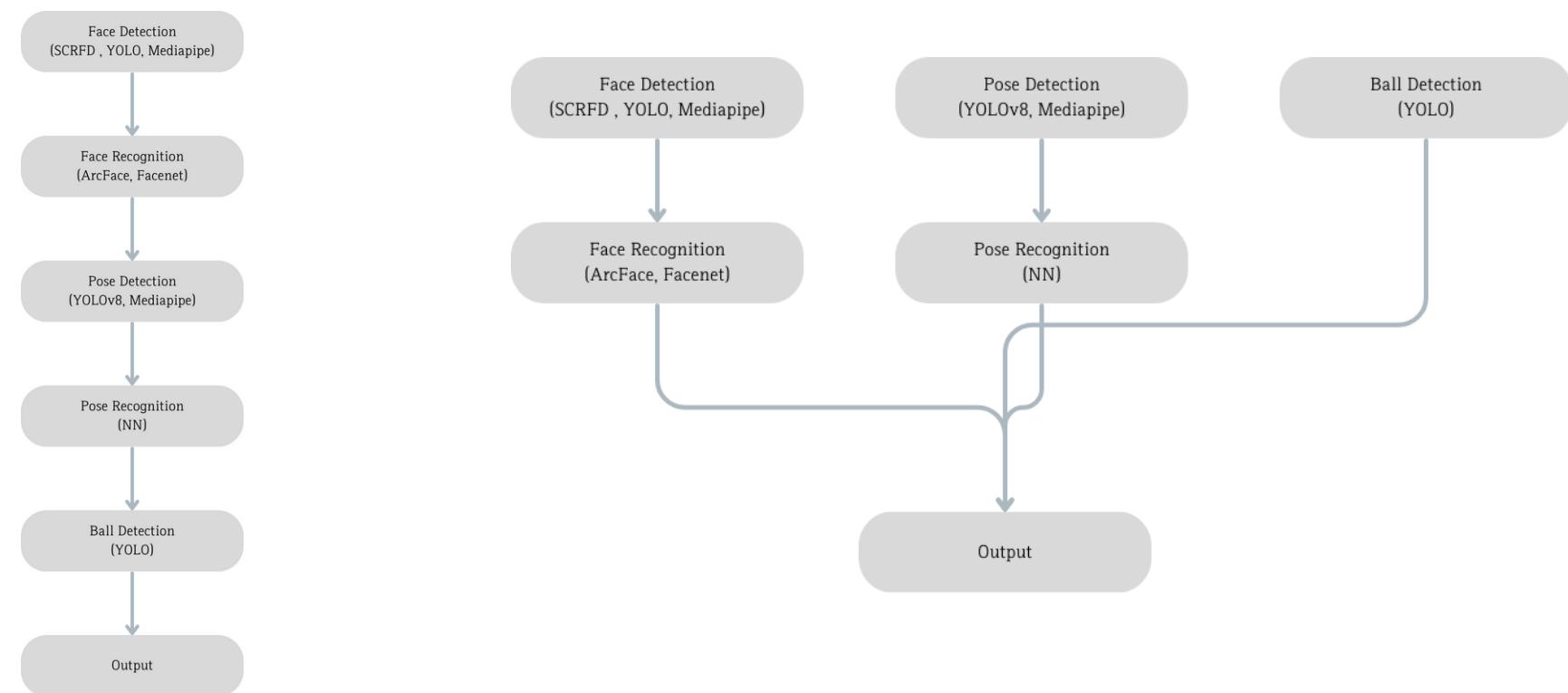


- No GPU (fps 10 or below)



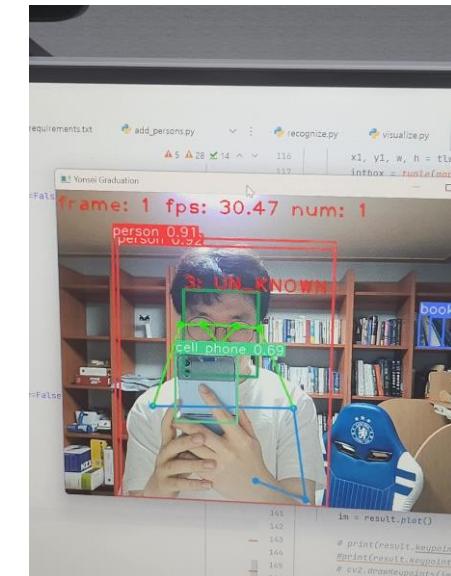
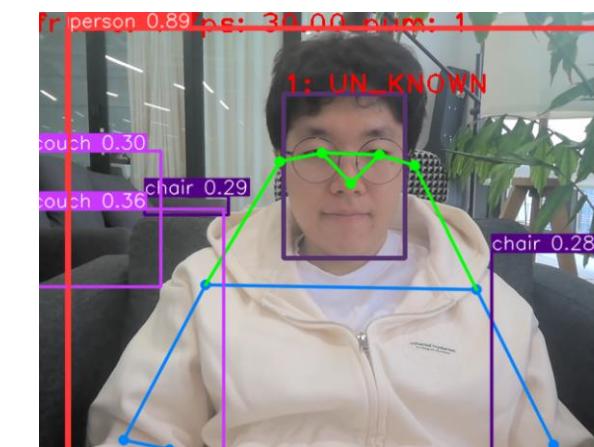
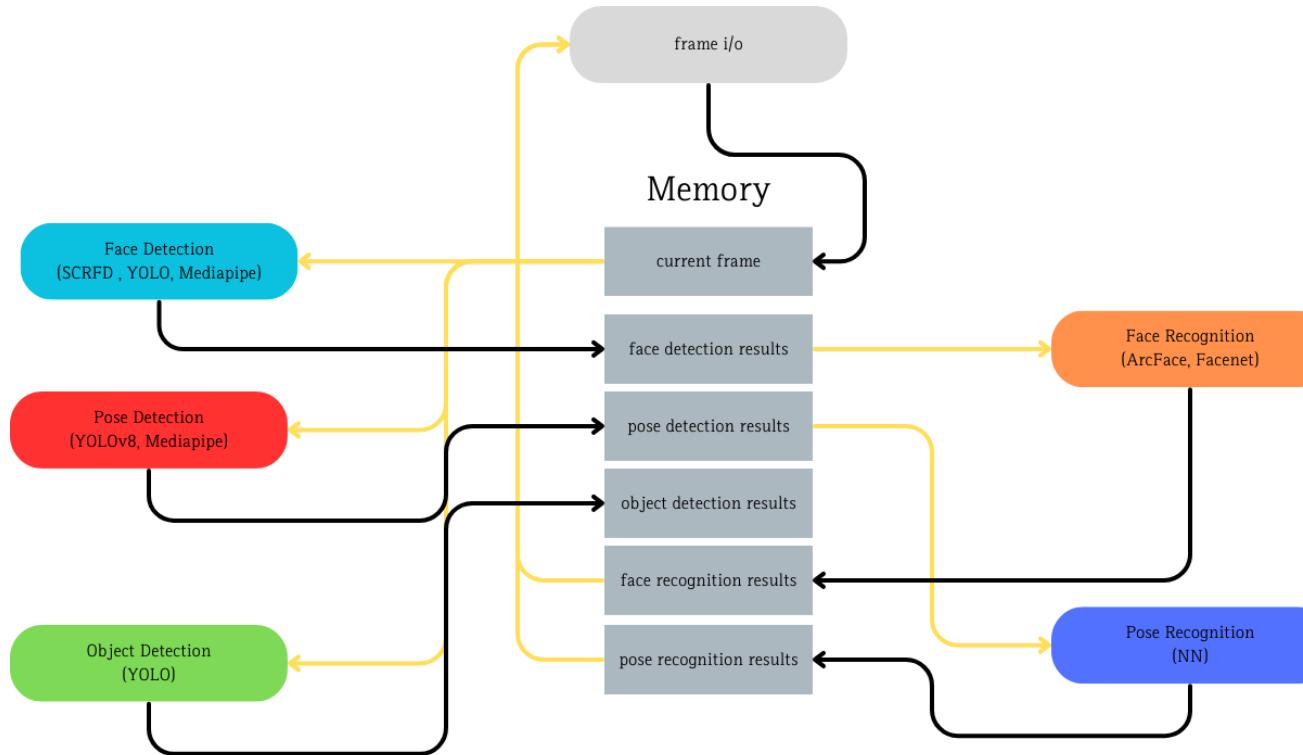
Real Time Face Recognition

- the drop in FPS in the face accuracy measurement model can be addressed by running a lighter model
- However, it is hard to solve the FPS drop when multiple models such as pose detection object detection run together



Multithreading Multiple Models

- Each model operates in parallel by reading from or writing to its allocated memory
- Achieve 30 FPS in Laptop CPU, whereas without multithreading it would be under 1 FPS

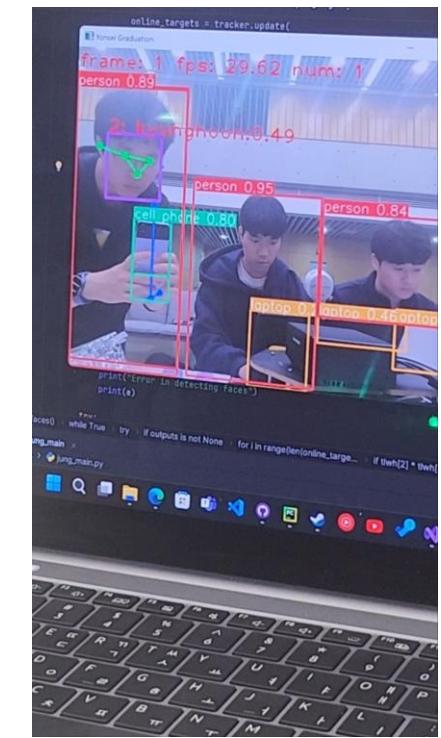
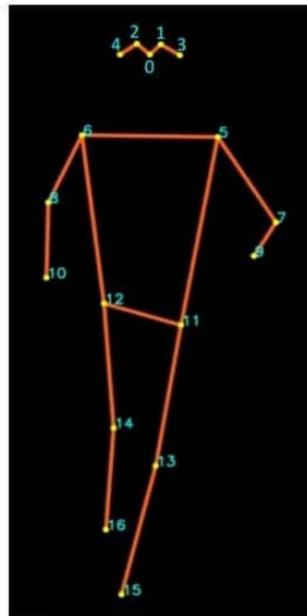


Face-Pose Matching

- Since we are using independent face detection and pose detection models, it is necessary to align the same face and pose
- "Based on the fact that the YOLOv8 pose detection model extracts keypoints for the nose and eyes, if the nose and eye keypoints are within the bounding box of the face detection model, they are considered to belong to the same person"



| Index | Key point |
|-------|----------------|
| 0 | Nose |
| 1 | Left-eye |
| 2 | Right-eye |
| 3 | Left-ear |
| 4 | Right-ear |
| 5 | Left-shoulder |
| 6 | Right-shoulder |
| 7 | Left-elbow |
| 8 | Right-elbow |
| 9 | Left-wrist |
| 10 | Right-wrist |
| 11 | Left-hip |
| 12 | Right-hip |
| 13 | Left-knee |
| 14 | Right-knee |
| 15 | Left-ankle |
| 16 | Right-ankle |



Human/Ball Tracking and Following

How to follow detected human / ball?

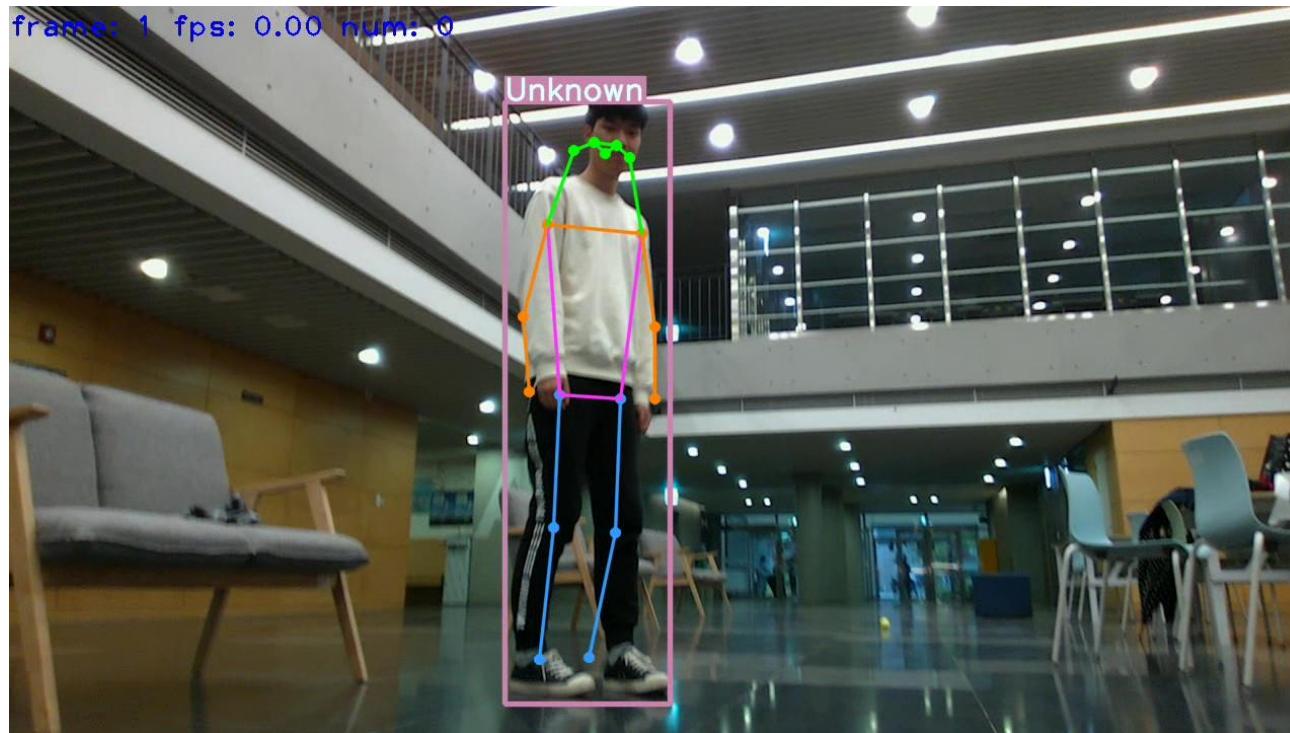


Image Space Detection

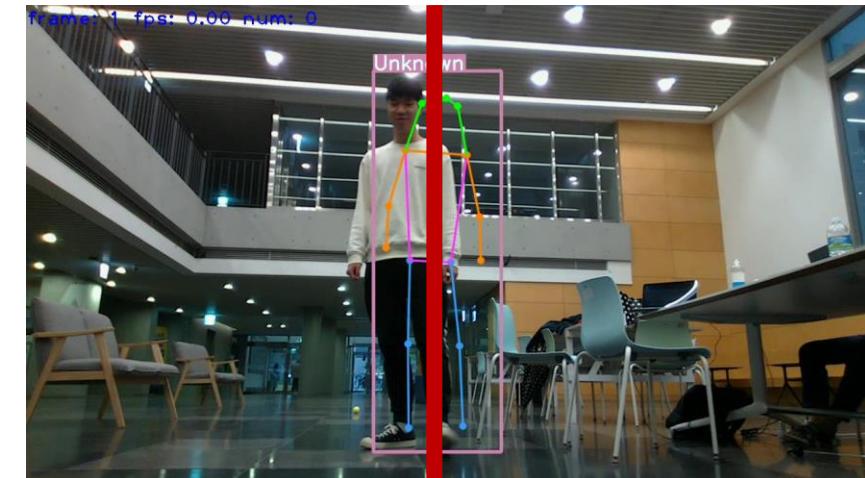
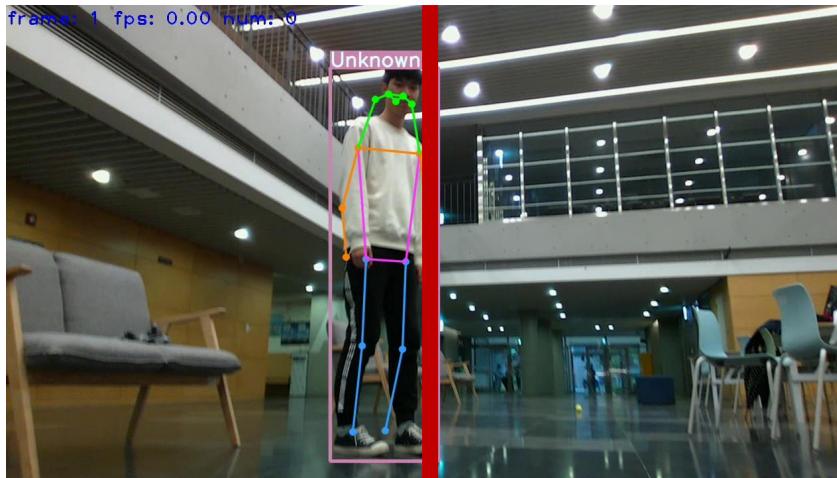


< 1.0m/s, 0.01rad/s, 2.0s >
<Linear, Angular, Duration>

Action

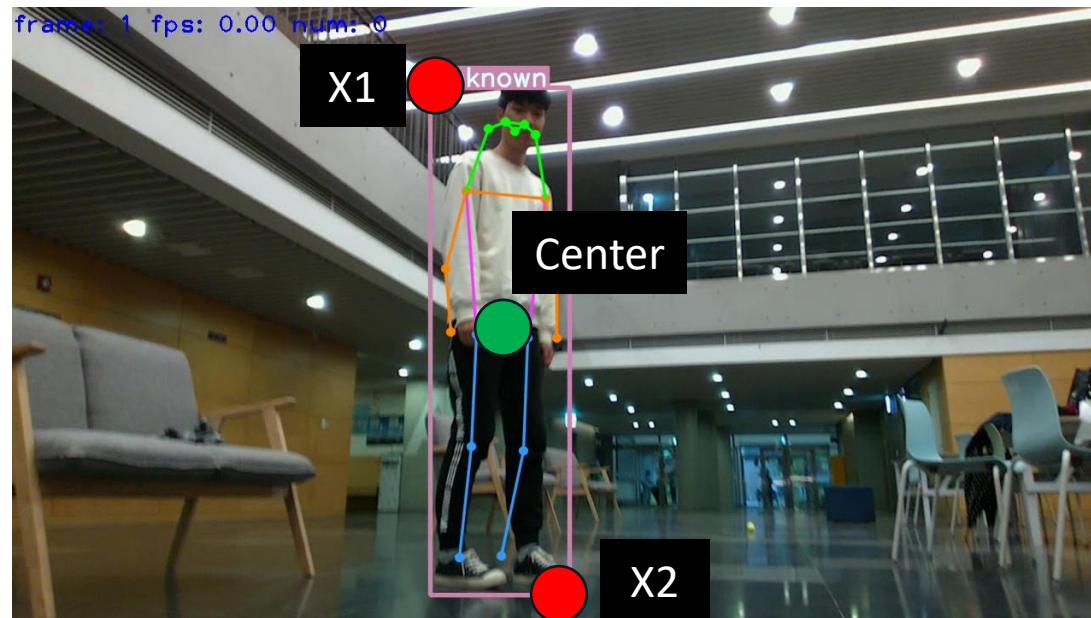
Human/Ball Tracking and Following

Goal : Keep the object in the middle of the screen

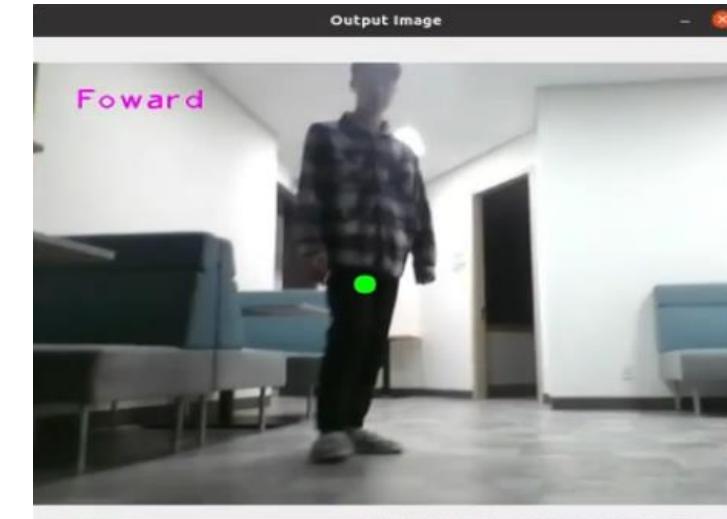


Human/Ball Tracking and Following

How far away is the object from the center of image?

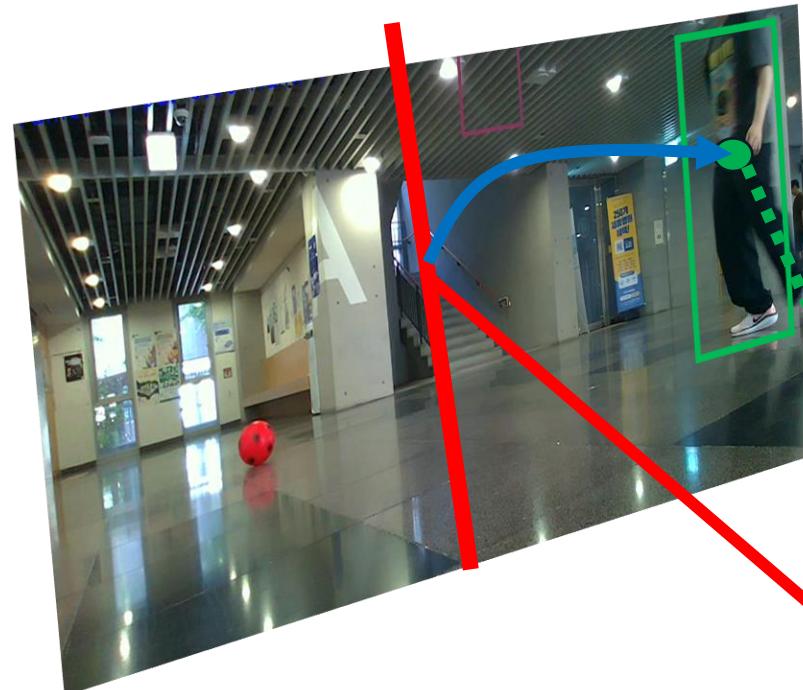


$$\text{Center}_{obj} = (\text{box}['x1'] + \text{box}['x2']) / 2$$



Human/Ball Tracking and Following

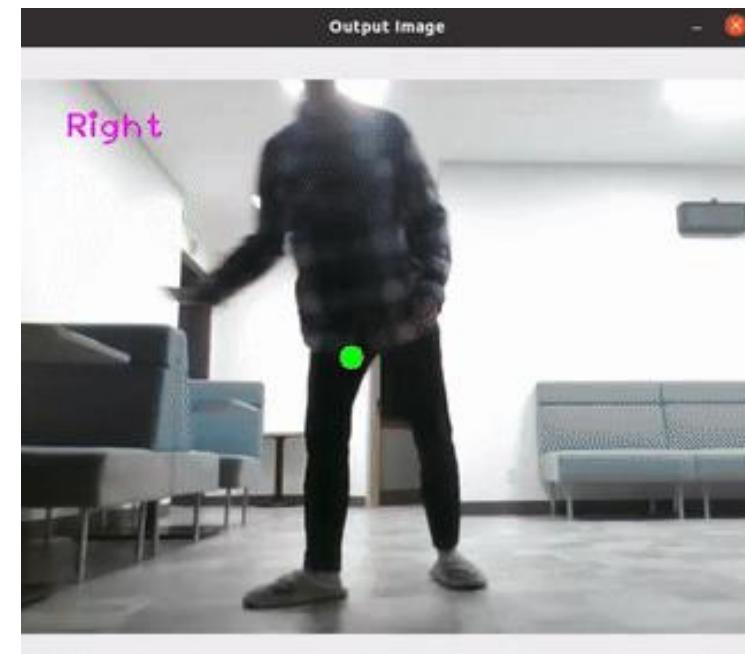
How far away is the object from the center of image?



$$\text{linear} = 0.01 \text{m/s}$$
$$\Delta\text{angular} = \alpha * \frac{\left(\text{center}_{obj} - \frac{\text{width}_{img}}{2} \right)}{\text{width}_{img}}$$
$$\Delta t = 0.01 \text{s}$$

$\langle \text{linear}, \Delta\text{angular}, \Delta t \rangle$

Send this to low-level controller
(Part 4-i)

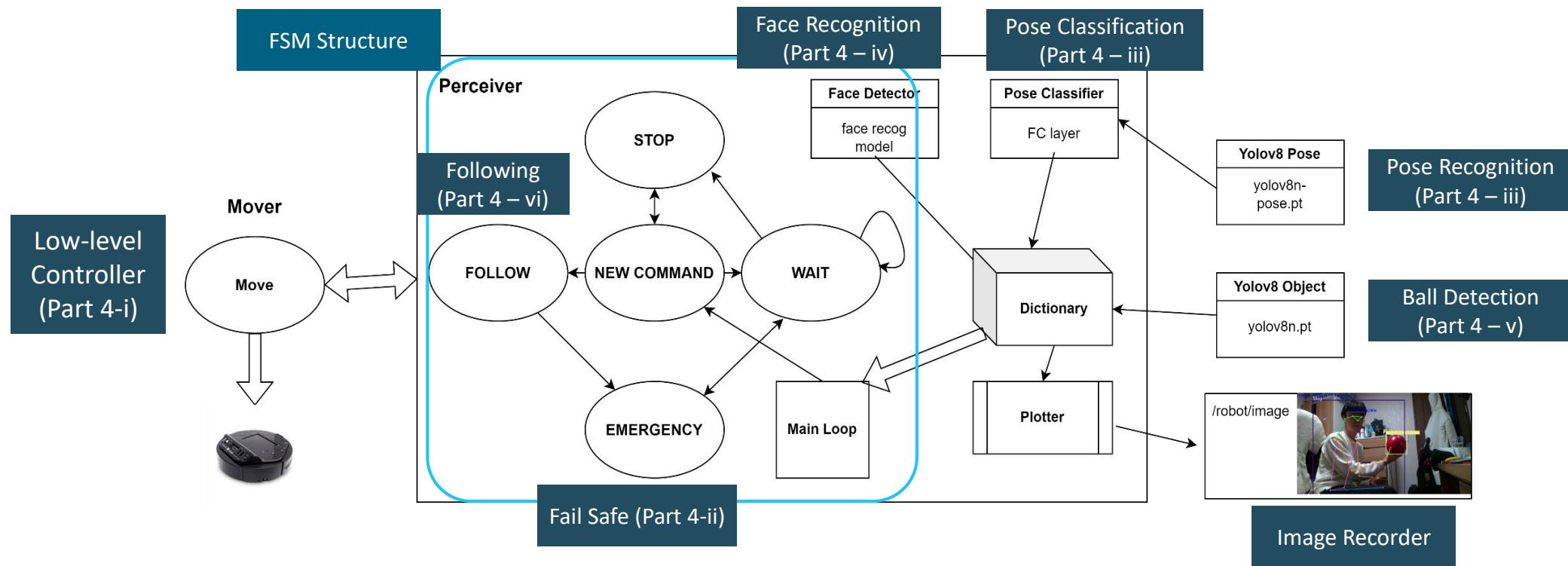


5

Integrating Functions

Node Design

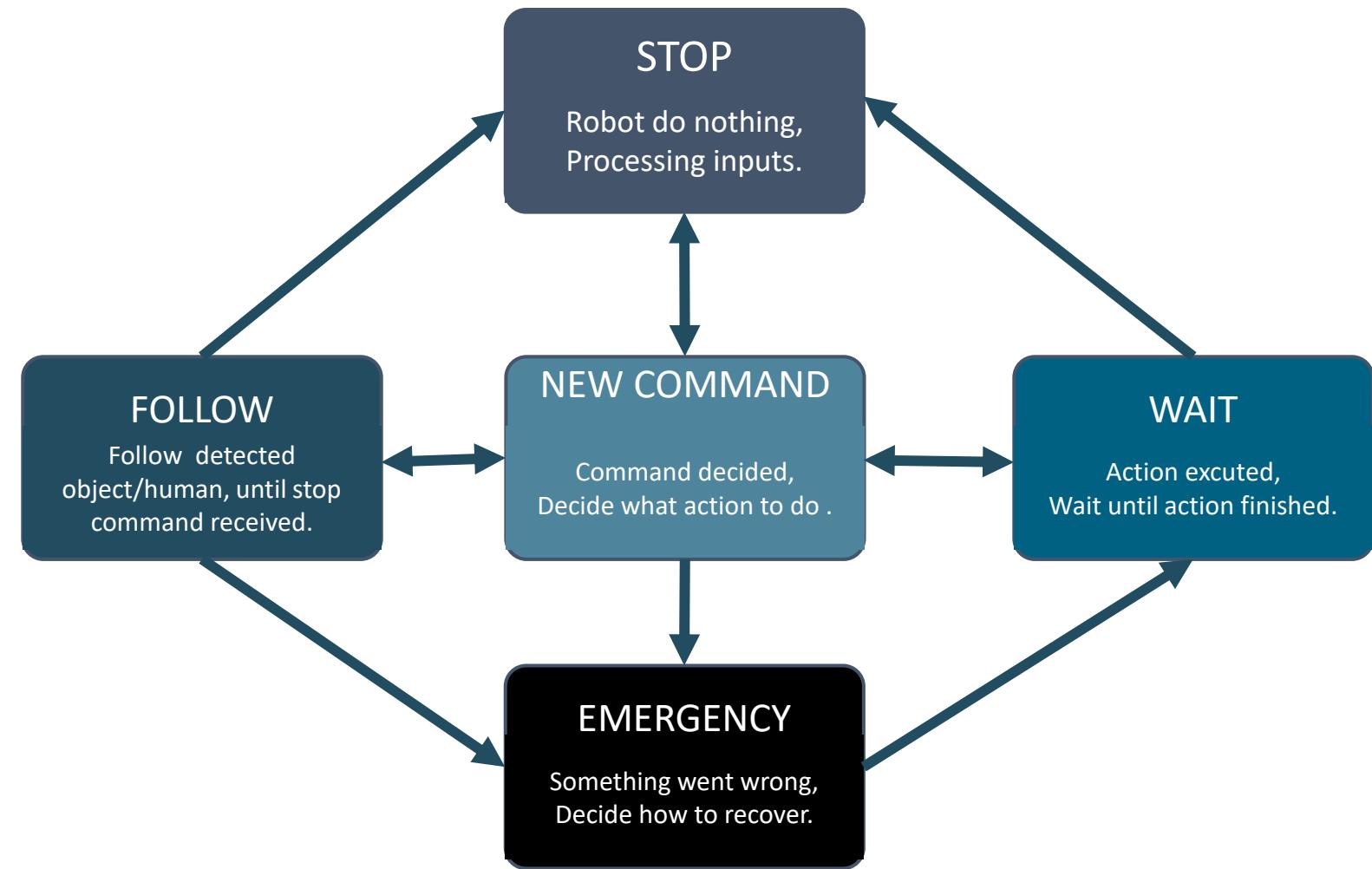
1. Construct ROS nodes for running each functions together.
2. Finite State Machine (FSM) design for internally making decisions.



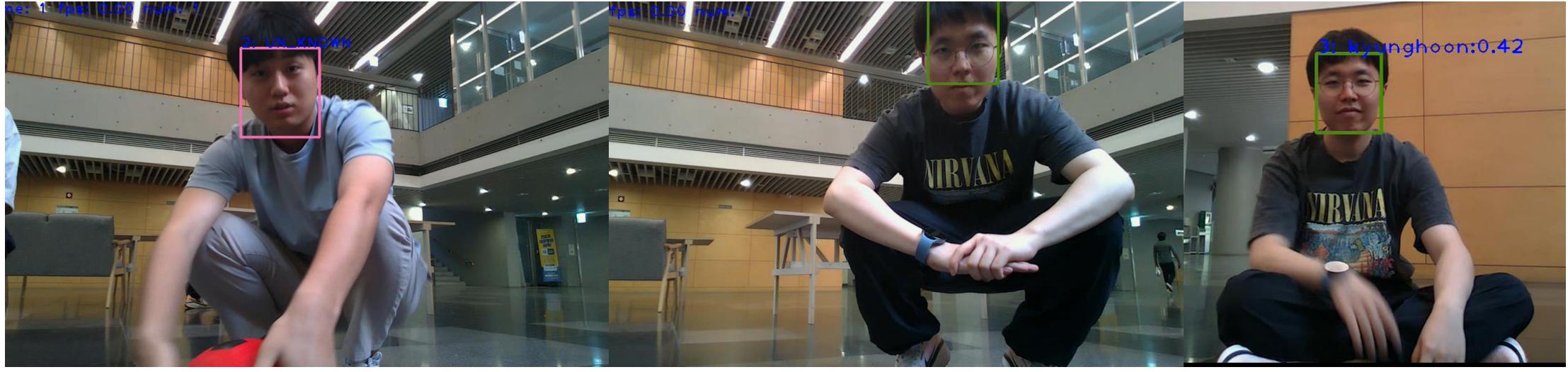
```

def fsm(self):
    # implement FSM
    if self.state == State.STOP or self.state == State.COMMAND:
        if self.state == State.STOP: #STOP
            self.set_target(0, self.current.angular.z)
        if self.command == 0:
            self.state = State.STOP
        elif self.command == 1:
            self.state = State.FOLLOW
            self.set_target(x=0.1)
        elif self.command == 2:
            self.state = State.FOLLOW
            self.set_target(x=0.8)
        elif self.command >= 2:
            if self.command == 3:
                self.action_spin()
            ##### keyop #####
            elif self.command == 10:
                self.action_forward()
            elif self.command == 11:
                self.action_turn_left()
            elif self.command == 12:
                self.action_backward()
            elif self.command == 13:
                self.action_turn_right()
            ##### keyop end #####
            self.state = State.WAIT
        elif self.state == State.EMERGENCY: #EMERGENCY
            if self.check_timer():
                self.set_target(x=0)
                self.set_timer(2)
                self.state = State.WAIT
        elif self.state == State.WAIT: #WAIT
            if self.check_timer():
                self.state = State.COMMAND
        elif self.state == State.FOLLOW:
            if self.command != 1:
                self.state = State.COMMAND
    
```

For reliable priority control and control,
Pet's brain has an FSM structure!



Pose Recognition/Face Recognition/3D Ball Detection



UN_KNOWN

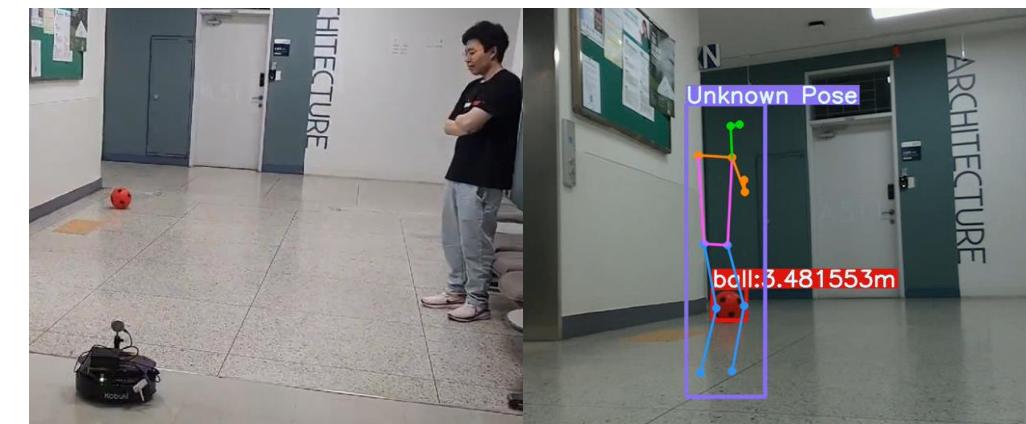
Owner(Kyunghoon)



Start : Arm Stretch

Stop : Attention Pose

Spin : Sit



3D Ball Detection

6

Results/Conclusions

Scene1 : Face Recognition -> Pose Recognition(Start/Stop)



Cam1 (Robot ver)



Cam2 (Global ver)

Scene2 : Face Recognition -> Pose Recognition(Start/Stop/ *Spin) -> Bump and Change Direction -> Ball Tracking

*Add This Week

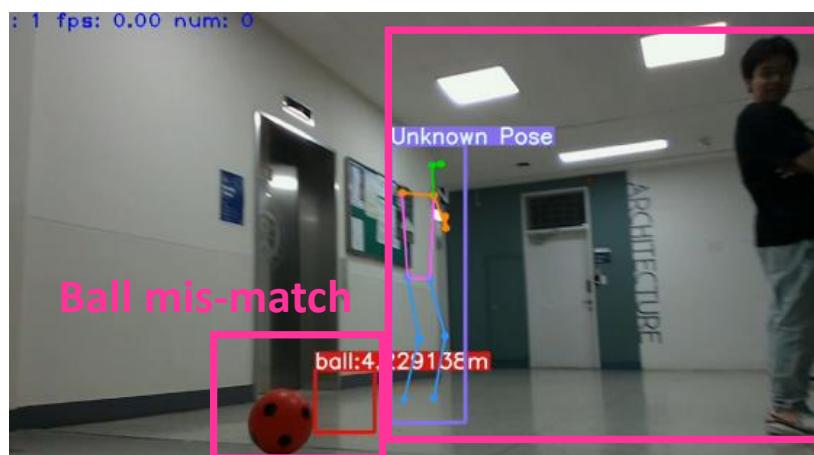
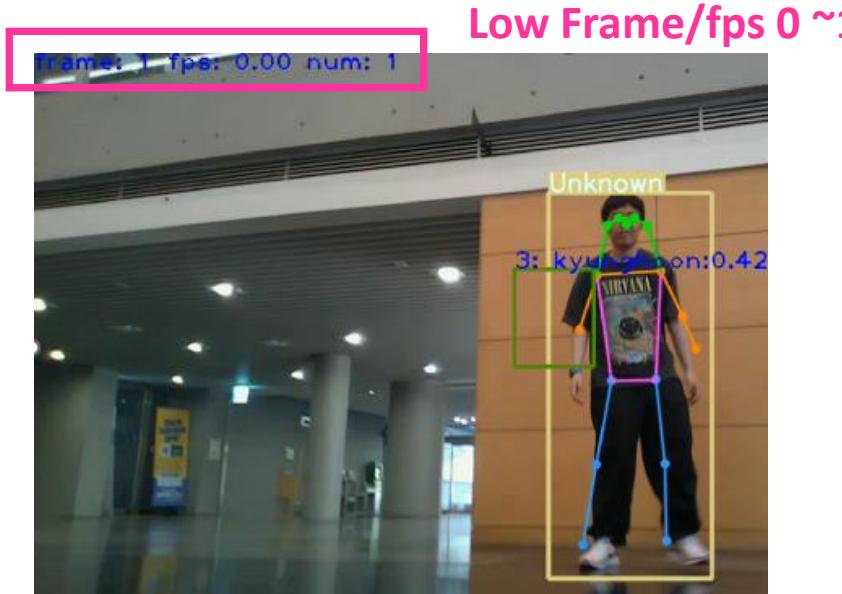


Face Recognition -> Pose Recognition -> Bump and Change Direction

-> Ball Tracking

Conclusions1: Why Robot Seems to be Not Realistic?

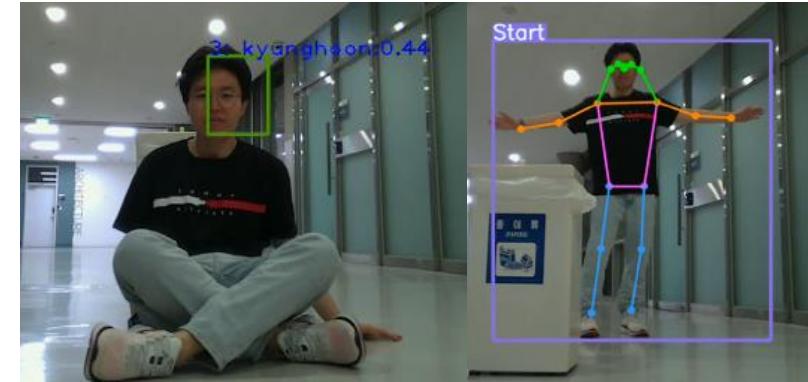
Problem1 : Robot Processor(Jetson)



Human mis-match

**Can not perceive
Abrupt Change**

Problem2 : Camera's Angle is limited

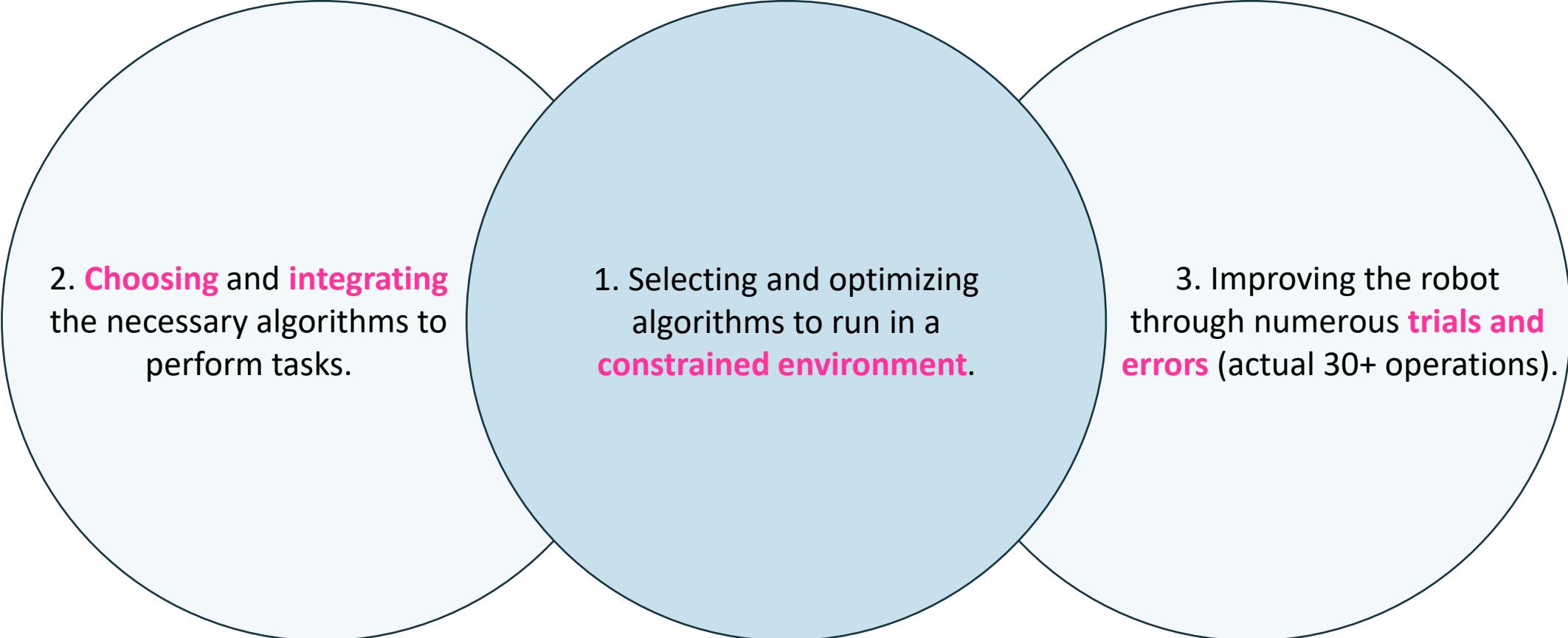


**For Face Recognition(Real-time), Camera towards to up
For Ball Detection, Camera towards to down**

→ **Could only unlock robot from right in front of it**

**Problem3 : Robots should be generalists
as well as specialists?**

Conclusions2: What we have learned



2. **Choosing** and **integrating** the necessary algorithms to perform tasks.

1. Selecting and optimizing algorithms to run in a **constrained environment**.

3. Improving the robot through numerous **trials and errors** (actual 30+ operations).

Thank You
감사합니다