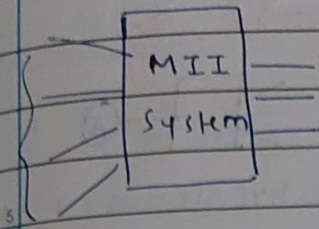


7 9th Feb, 23

Page :

Date :

MII-lec 2 (20th Feb 2023)



↳ not necessary all work at same time

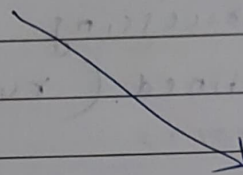
* data level fusion

↳ same modality

S1 S2

G1 G2

FE1 FE2



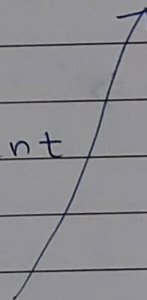
* Feature level fusion

Fusion Processes

↳ similar modality

↳ speech & lip movement

↳ noise affects both



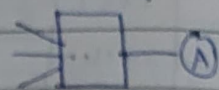
* Decision level fusion

↳ Two or more modalities are different (face & head gesture)

↳ noise affecting one may not affect others

* CARE properties (lec-3)

CARE = Complementarity \Rightarrow



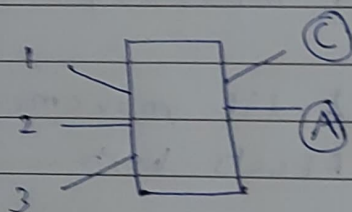
Assignment \Rightarrow 1 modality is used at a time
 redundancy \Rightarrow
 equivalence \Rightarrow

Parallel processing

+ combined (running parallel & combining all modalities)

Ex: Put that there (speaking + pointing together to give O/P)

redundancy \Rightarrow



if I use 1 to run the system,
 2 & 3 are redundant

if I use 3, 1 & 2 = redundant

* redundancy = the modality which is not used

* Equivalence = sequential

(Put that there ~~pointing~~ followed by pointing)

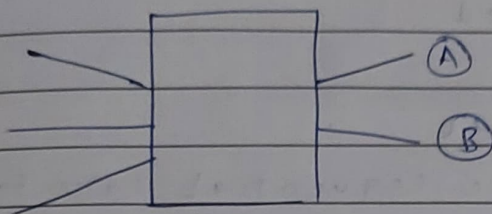
* CASE Properties

C = Concurrent = all commands in 1-go

A = Alternate

S = Synergetic

E = exclusive = commands at a time but they
(1 after another) ↑ five diff O/P



Parallel + Independent

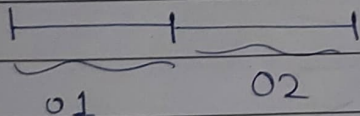
↓
2 modalities
enter at same
time

* Synergetic = Parallel + Combined

↳ benefit ⇒ mismatch % is less
↳ most widely used

* Exclusive = Sequential + Independent

sequential modality



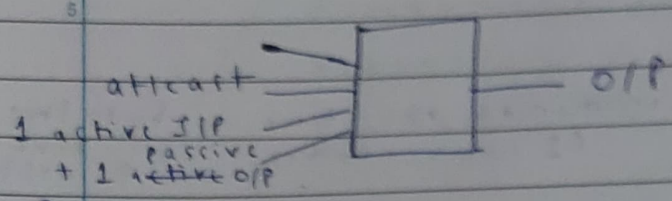
outputs are different

* Types of Multimodal System

① Fusion based

(passive) Unintentional

= gazing + lip movement



② Temporally cascaded

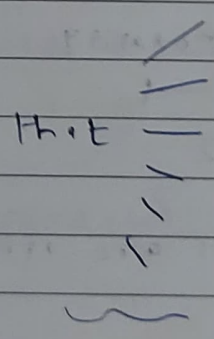
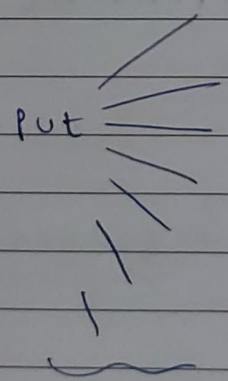
→ sequential
(when we use sequential data to give O/P)

Multimodal

* Fission

Speech IIP = "Put that there"

After receiving, we check different versions of Put



check in dict what comes after Put
(This is Fission Process)

check in dict & accordingly corrective measures taken

→ Detection — Colour
— 3D model = for VR
— motion based = video processing

Page: _____

Date: _____

* Hand Gesture

Data

→ Hand Glove

Adv: — ① Accuracy ↑ (high)
— ② Background can be anything
— ③ correct OIP

Disadv: — ① extra equipment
— ② wired which means we are constrained to 1 place (Majority)
— ③ not a natural way of comm.

→ Coloured marker ⇒ 6th sense (P. Mistry)

on fingers
→ R, G & B markers used to control something

Adv: wireless, multiple Applications

→ Accuracy is fine

Disadv:

→ ① ~~Wear the device~~
① Wear the device at all times

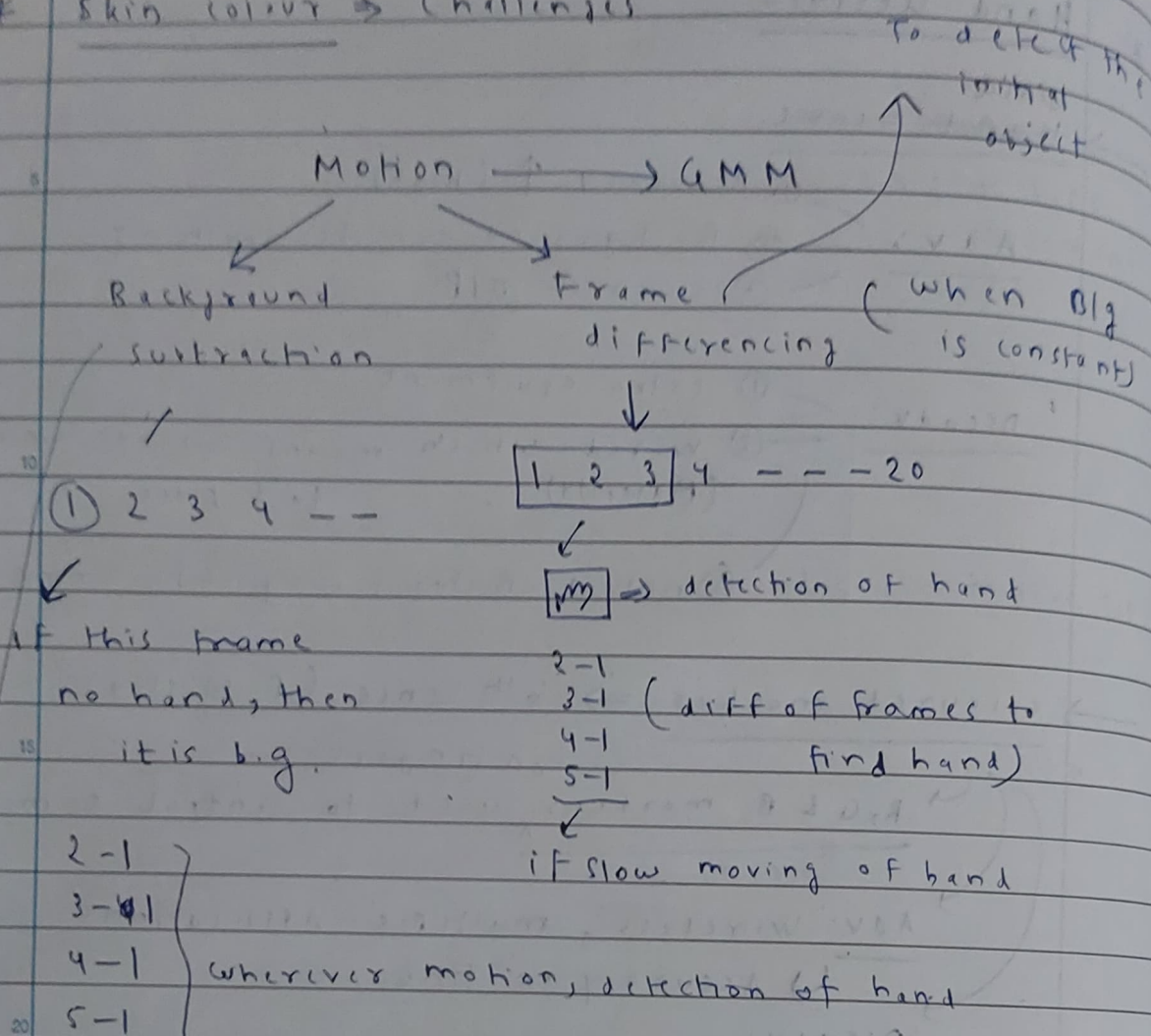
② Similar coloured objects at background (not able to recognize) → not able to detect (overlapping) or track

→ vision based ⇒ use of only ^{hand} ~~video~~ and camera

→ Adv: natural way of comm, wireless,
Accuracy is fine

Disadv: Challenges =

* Skin colour → Challenges



→ "Illumination variation is disadvantage" (lighting conditions)

* Hand Gesture

- ↳ Detection
- ↳ Tracking (To locate the position of hand)
- ↳ Recognition (meaning of the gesture)

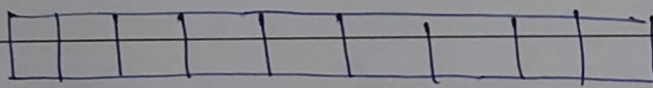
* Limitation of Frame differencing

- ↳ B/g should be stationary

* B/g subtraction

- ↳ The first frame should be b/g
- ↳ challenge = Illumination variation (A slight change in b/g)

* GMM (Gaussian mixture model)

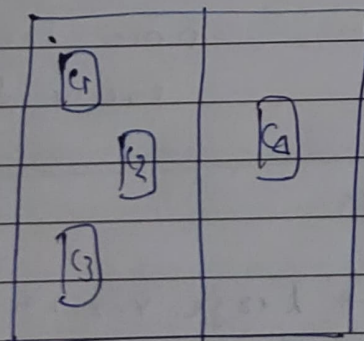


→ identify objects present
(Training Phase)
→ Test Phase

↳ limitation

Training of other things can be done
(like people in case of car-detection)

method 1:



* For first 30-seconds,
training Phase

* later it calculates similarity
matrix to detect objects

This can be set by ourselves

3D model = not reqd.

Page:
Date:

before hand

Method 2: Provide dictionary & then later detect the objects

→ THIS TRACKING

* Colour based = Markov model

* Colour based tracking

(A colour based approach)

1) Coloured marker approach

2) Skin colour object (Vision-based)

RGB → HSV (Hue-saturation-value)

YCbCr

$- \leq H \leq -$

$- \leq S \leq -$

$- \leq V \leq -$

min & max range

decided from

detection of face

skin-colour

since face is easily detected

major disadvantages:

Colour overlapping + large variations

(if multiple persons)

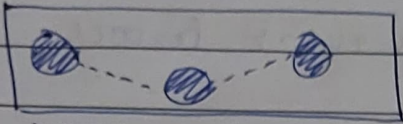
+ illumination variations

to predict the next activity

* Probabilistic approach

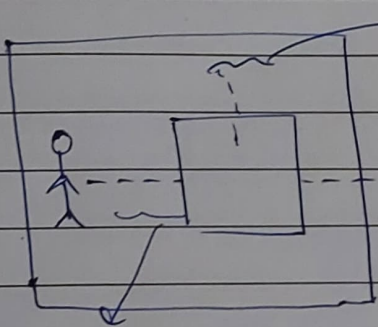
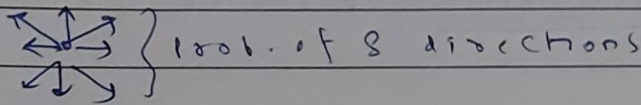
tracking Algo. (study ~~about~~ about this)

→ Kalman Filter (an example of prob. based approach)



first few frames, check where object is moving

gives probability where will object move in next frame (what is next state of object)



but if it changes direction, detection ~~is~~ wrong

} if it comes to this state, detected properly

here also detected properly.

* Appearance based approach

→ most widely used

→ KLT, Camshift
↳ shape + colour

Kenneth, Lucas & Thompson

→ not ^{work} ~~work~~ well for longer videos

* KLT ⇒ feature based technique



→ Features detected from detected object
& use these features in
subsequent frames.

50 features
extracted

edges, lines, corners

→ limitation = if lot of illumination variation