

# lec 1

- \* medium = technical means to store, transfer & convey information

Physical channel = air / lights.

- \* modality = human sensory channel.

visual channel = see

auditory channel = hear

olfactory = smell

haptic = touch

Gustatory = taste

Vestibular = Balance

temp sensor

*	Senses	organs	Modality	sensor
	Vision	eyes	visual	camera
	Hearing	ears	auditory	microphone
20	Touch	skin	Haptic	Touch screen, gloves
	olfact.	nose	olfactory	
	Taste	tongue	gustatory	
	Balance	equilibrium organ	vestibular	

25  

human for System

Apple

## Gesture, Face, Hand.

## Speech

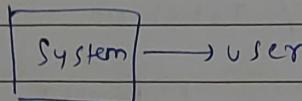
## Posture/motion recg.

## Modality

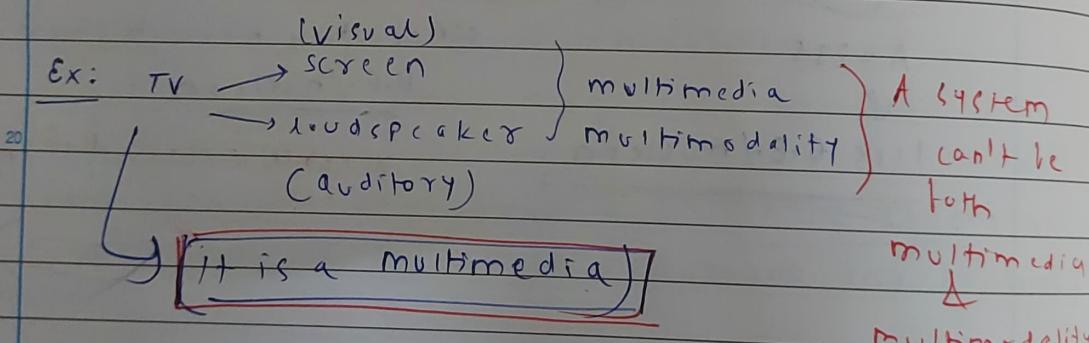
= capability of system to mimic ↪  
such human sensory channel.

## Multimedia vs. Multimodality

- \* Multimedia = multiple media, system with  
more than one way to transfer info  
from system to user.

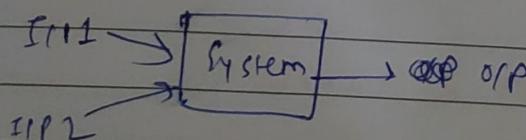


- \* Multimodality = system that can stimulate more  
than 1 human sensory channel.



## Alternate def

Multimodal System = A system that provides more  
than 1 way IIP to interact with it.



TV = no visual  
IIP = no IIP

Smartphone = multimodal. (touchscreen, microphone)

IIP

Camlin

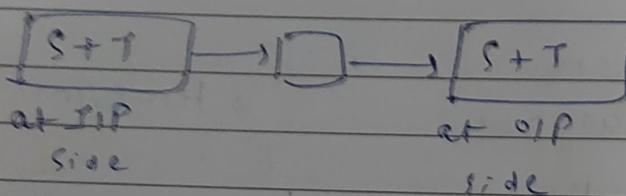
focused on o/p side

and mainly off

- \* multimedia system = interactive system that provides information via several I/O channels (sound, graphics etc.)

focused on I/O side

- \* multimodal system = processes 2 or more combined user I/O modes (such as speech, touch, gestures, hand, head etc. . .) Δ O/P gives that many no. of modes



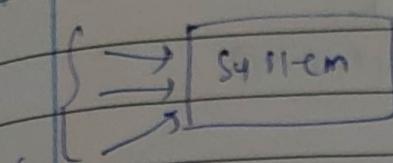
Sharon Orratt  
~~(2012)~~ (2012)

- \* For a truly multimodal system, info we provide at the input must be related & coordinated to give O/P.

- \* multimedia = o/p of the system

37<sup>th</sup> ICF (Benoit)

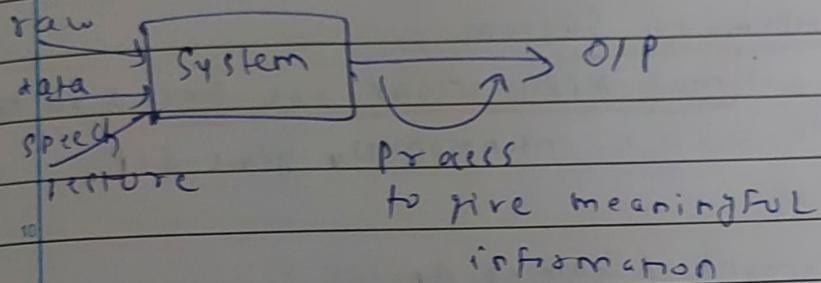
- \* A multimodal system represent & <sup>manipulate</sup> ~~manipulate~~ information from diff human communication channels. These system automatically extracts meaning from multimodal raw data & gives desired information



human sensory channel } at IIP then it is  
multimodal

IIP = human sensory channel

OIP = desired OIP

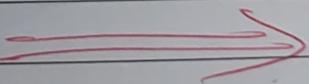


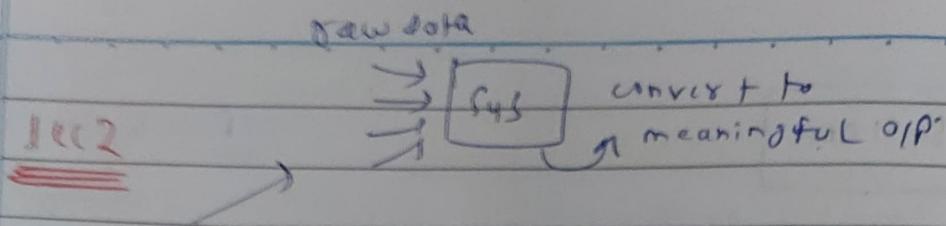
\* Domains of multimodal system :

- i) Active IIP mode
- ii) Passive IIP mode
- iii) virtual environment

\* 20 Put that there = 1979 MIT

↳ Speech + Gesture  
(Voice)





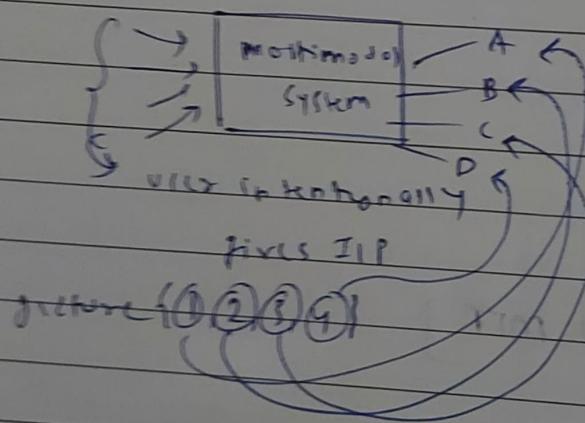
\* multimodal systems are subset of multimedia systems.

## \* Domains of multimodal System

- max. Applications are based on this

  - ① Active I/P mode = user intentionally give I/P
  - ② Passive I/P mode
  - ③ virtual environment (Gaming Appln).

**i** User can intentionally give information to the system  
forex: fixed set of commands, gesture



"We know which type of TIP to provide to give a particular tip."

Q A dialog system with diff IP modes. Chatbot can be of 1 example."

~~Chatbot~~

we input some text & we get fixed set of answers.

To fixed III to get derived off

P O R  
6 7 26

Page:  
Date:

$$\begin{array}{r}
 8 \quad 1010 \\
 10001 \\
 \hline
 10111 \\
 \hline
 6+17+10
 \end{array}$$

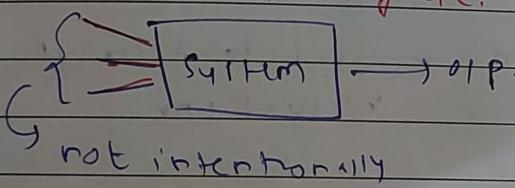
$$\begin{array}{r}
 8 \quad 102 \\
 1100 \\
 \hline
 110 \\
 \hline
 1010 \\
 \hline
 82
 \end{array}$$

8 4 1216

$$\begin{array}{r}
 1000 \\
 0100 \\
 \hline
 100 \\
 000
 \end{array}$$

Q 15 Passive IIP mode = All kinds of info that user do not produce intentionally.

Ex: (motions, gazing at something  
= sad, happy etc.)

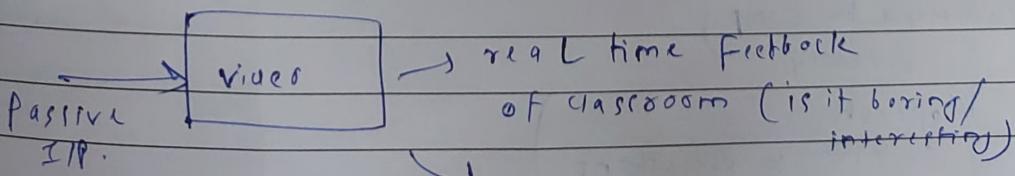


Unlock a phone  
with happy  
face - intentional  
not intentionally

provided by the user.  
(no fixed pattern at  
IIP end)

facial exp  
if students  
classroom

Ex: Depending upon emotions of students  
(how → he captures video of class & clicking)



real time feedback

of classroom (is it boring/  
interesting)

how can classroom

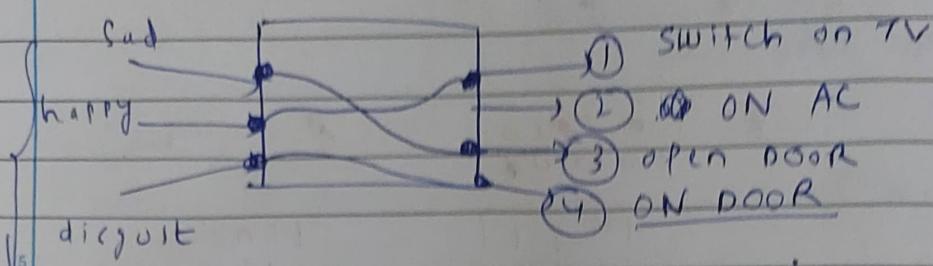
be made interesting / how

whether to continue class

or not / improve skill of  
professor

Carolin

"fixed pattern"



intentionally to get OIP (so this is not passive IIP but active IIP).

③ 10 virtual env = User can behave naturally to perform a task.

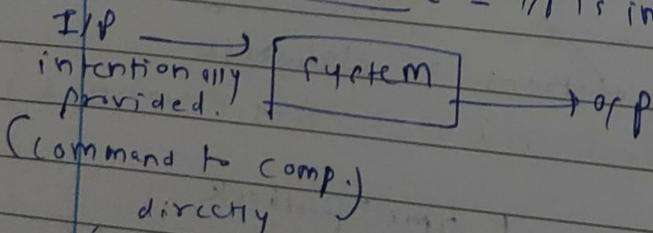
OR

Users can use all kind of information, behave naturally as he/she is in natural environment & the system would be able to analyse all this information & act accordingly.

\* 1 multimodal appln. related to Active / Passive / virtual env.

11C3

\* 25 Active IIP mode = IIP is intentionally provided.

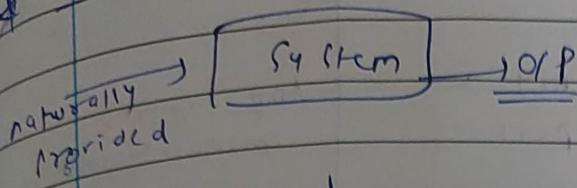


(Command to comp.)  
directly

30 speech  
Pointing to something  
gesture  
writing

i/p

Passive mode



(not give command  
to computer directly)

(perform this task X)

gaze, brain wave patterns, lip movement, facial exp

virtual env (can take up both active & passive i/p)

i/p →   
may be anything  
(Active/Passive)

virtual exp  
at o/p end. (real life experience  
kind of thing we get)

Recognition process will be done for both active or passive input.

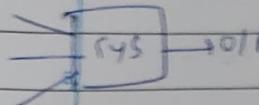
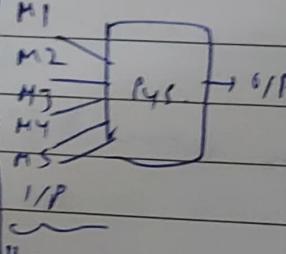
\* Process recognition is common for any system.

### modality relations

CARE = Complementarity, Assignment, Redundancy, Equivalence

# Properties of multimodal systems

Page:  
Date:

Complementarity	Assignment	Redundancy	equivalence
I/P 	ONLY 1 modality will be selected for a certain piece of information	ONLY a part of information will be used.	Any available modality <del>can</del> be used with combination of anything
multiple modalities are used together at a time to reach off			
Ex: Speech + gesture both will work at same time	M1 M2 M3 M4 → SYS → O/P		"Sequential": OR "parallel"
M1 + M2 at same time used	"one modality will be selected at a time to give O/P!"	Ex: MIT pointing hand + voice PUT THAT red circle in triangle if only 1 triangle & circles, then gesture (pointing hand) is redundant, voice is enough	"using Assignment property"
Ex: 1979 MIT "Put that there" (pointing hand + voice) at same time	"when we operate 1 modality, other modalities will be switched off"		
All I/P work together	1 of modality work at a time		Any kind of system can work with any combination of inputs (whether Sequential, parallel or Assignment)

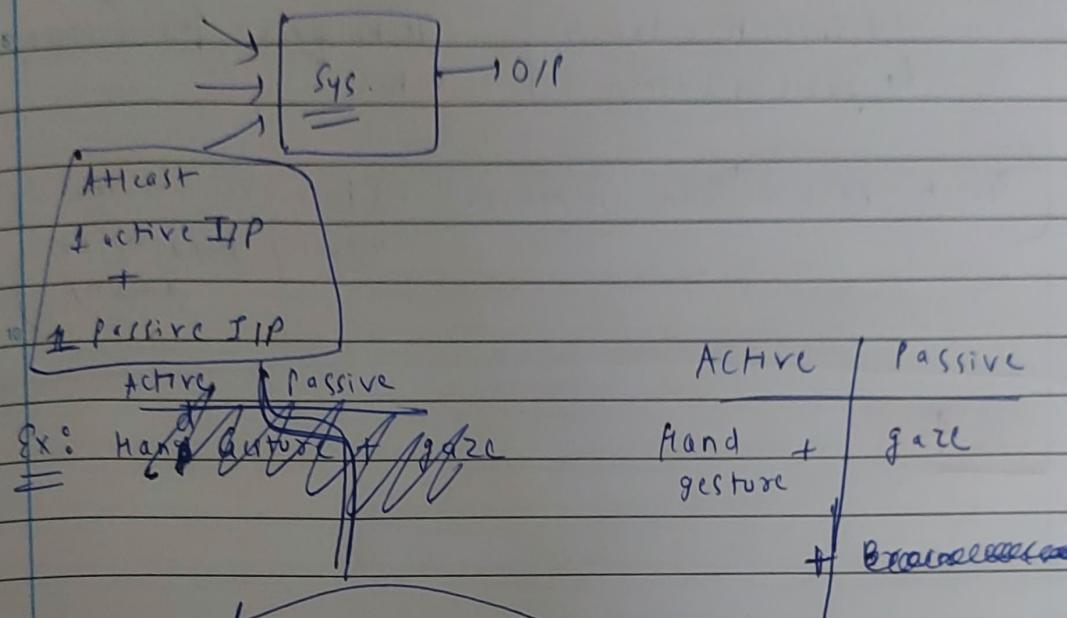
"A system can be of multiple types"

Types of multimodal systems:

Fusion

① ~~Feature based~~ mm / Blended mm system/interface:

Major application of optical sensors



Lip movement  $\Rightarrow$  works on CNN

$\hookrightarrow$  Lip movement = Passive

$\hookrightarrow$  speech = Active.

(recognize lip movements even if noise in background)

(MIT media lab)

Freeze release see decreased

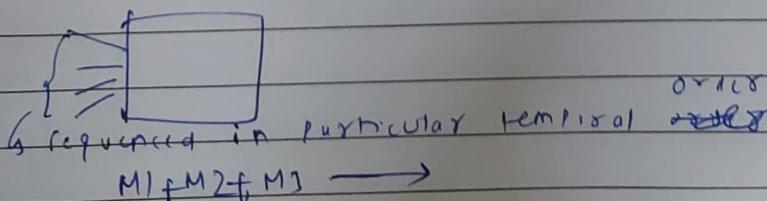
$\hookrightarrow$  detect lip movement + speech

②

Temporally cascaded MM system

much less explored

Process 2 or more modalities sequenced in particular temporal order.



Ex: first gaze then gesture, speech recognition.

Ex: A cascaded multimodal natural user interface to reduce driver distraction (IEEE paper)

- ↳ Speech + button
  - ↳ Speech + touch
  - ↳ gaze + button
- } different cascading techniques for autonomous vehicles.

Ex: (Speech + gaze = "Turn on music" + (increase volume by moving hand up or down))

lec 4

used to capture diff information (faster, more accurate)  $\Rightarrow$  for recognition purposes

Ex: A cascaded multimodal natural user interface to reduce driver distraction (IEEE paper)

- ↳ speech + button }
- ↳ speech + touch } different cascading techniques for autonomous vehicles.
- ↳ gaze + button }

Ex: (Speech + gaze = "Turn on music" + (increase volume by moving hand up or down))

## Face Recognition

If ~~bad~~ lighting conditions = problem may be there

facial recog  $\Rightarrow$  3D model based techniques, variation in angle of hand is used & detection is used.

Hand Gesture  $\Rightarrow$  video is taken & then frames of video taken to capture hand

\* 1 gesture can be provided in several ways.

## Application of facial recog

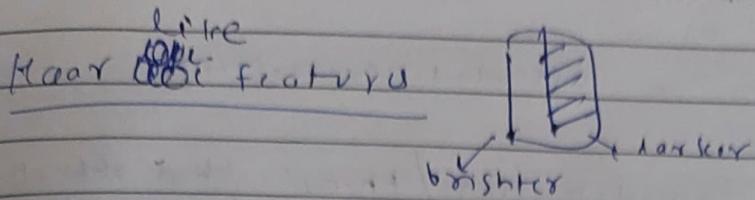
↳ Viola-Jones algorithm

30 Viola-Jones = to extract facial features using haar-like features.

can be used to train diff. objects/body parts.

Challenge of face detection  $\Rightarrow$  Olden classifier for  
George. recognition.

Page :  
Date :

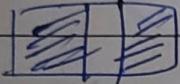


8-10 haar like features

(combination of black & white patches)

these patches made to detect  
features from scaled face

\* to extract information  
about nose



\* different haar like features for eyes, nose, lips, mouth

\* to extract features from facial regions

\* By default, it select 4 Facial Features:

① eyes

② nose

③ mouth

④ mouth + nose

(II) Creating an integral image = facial features are  
combined to extract information.

(III) Adaboost training

(IV) Cascading

can also detect other part of body

→ if 8/10 features matched, then it is face of

that particular person A

## lec 6

occluded = overlapping

\* Challenges of face detection / things that affect face detection also

- ① Pose
- ② Presence or absence of structural components
- ③ facial expression
- ④ occlusion (many faces are there)
- ⑤ Image orientation
- ⑥ Imaging (lighting) conditions

To deal with challenges

"not much spread"

knowledge based method = rule based method,

mathematical based method

(diff b/w eyes, nose, ~~cheek~~ <sup>nose</sup> lip)

Problems: if we take mathematical calculations, then

these calculations may change with orientation  
change of face (if orientation variation,  
change of camera)

"if face is static (face is not moving much)"

25 feature-invariant approaches

invariant to expression  
face ~~expression~~

→ Some features generated from face which will be  
invariant to lightning conditions, pose angle & so on

→ diff methods to find those feature points.

→ Used because of its invariant prop.

→ widely used

→ edge detection, texture of face, . . . . .

→ All features combined to detect the face.

## ~~Template matching~~

not much used

- Match template with current image
- correlation based method

Cons = template & image must comprise of almost similar image otherwise detection may not work correctly (if not much diff. b/w input image and template)

- correlate with template in forward.

→ if <sup>whole</sup> image is only rotated, then template will work but if face is rotated by 0, then will not work.

## ~~Apearance based method~~ = widely used

- training based method.

- folder  $\Rightarrow$  lot of illumination variation, diff faces,

$y_{train}$

angle variation.

② ~~detect~~ for detection.

③ SVM, NN

- train dataset & use information to detect face.

## \* DIFF types of features :

① SIFT (Scale invariant Feature transform)

② HOG (Histogram of gradients)  $\Rightarrow$  colour info

③ LBP (Local binary pattern)  $\Rightarrow$  edges & how allocated.

④ SURF (Speeded up robust Features)

⑤ DAISY features.

⑥ ICF (integral channel feature)

\* HOG  $\Rightarrow$  used for ~~colours~~ info + info about  
edge / line.

## lect

→ Knowledge based methods

(ex: Attendance marker (Person has to come at camera  
& mark attendance)

→ Feature invariant

→ Template matching = ~~not~~ used for recognition purpose  
(whether particular student present  
in a group of students  
picture)

feature comb = nowadays many used.