

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
15	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			
	perception		for system	

Apple

Gesture, Face, Hand.

Speech

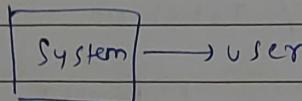
Posture/motion recg.

Modality

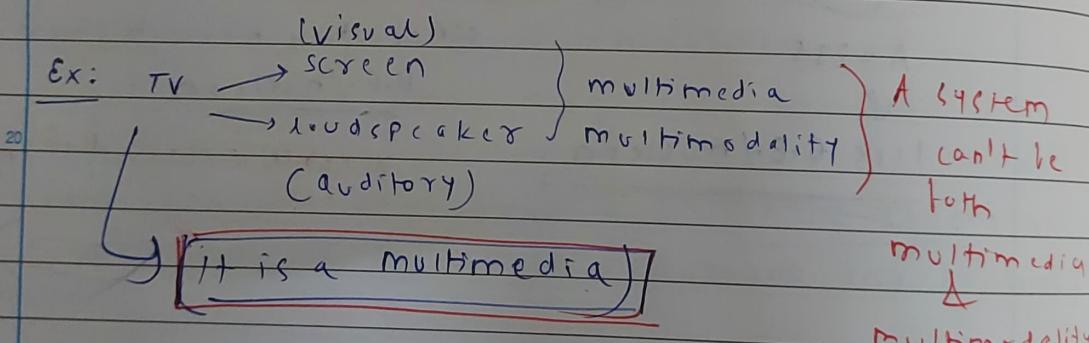
= capability of system to mimic ↪
such human sensory channel.

Multimedia vs. Multimodality

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

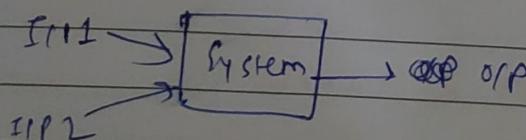


- * Multimodality = system that can stimulate more
than 1 human sensory channel.



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

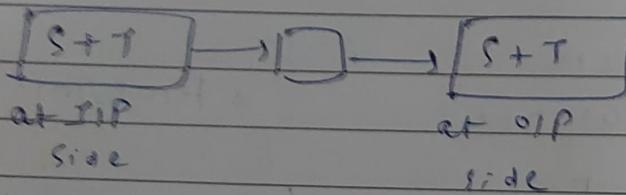
focused on o/p side

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



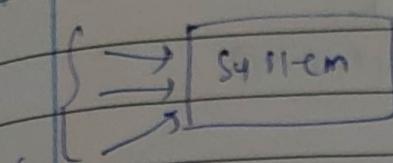
Sharon Orratt
~~(2012)~~ (2012)

- * 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

37th ICF (Benoit)

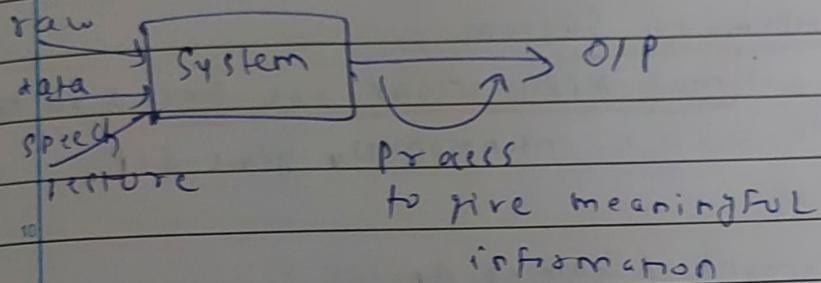
- * A multimodal system represent & ^{manipulate} ~~manipulate~~ information from diff human communication channels. These system automatically extracts meaning from multimodal raw data & gives desired information



human sensory channel } at IIP then it is
multimodal

IIP = human sensory channel

OIP = desired OIP



* Domains of multimodal system :

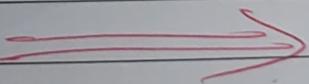
i) Active IIP mode

ii) Passive IIP mode

iii) virtual environment

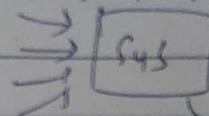
* 20 Put that there = 1979 MIT

↳ Speech + Gesture
(Voice)



Raw data

1002



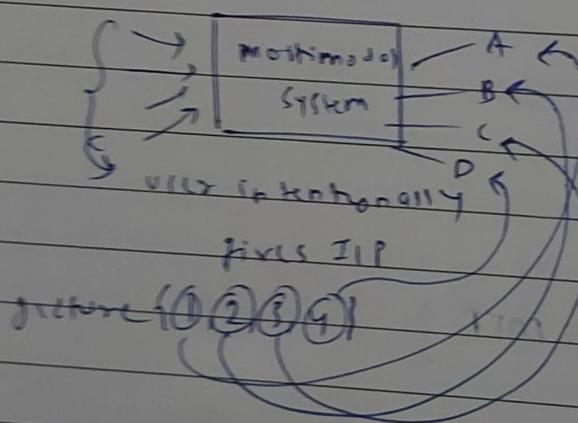
convert to
meaningful O/P

* Multimodal systems are subset of multimedia systems.

* Domains of multimodal System

- ① Active I/P mode → max. Applications are based on this
- ② Passive I/P mode
- ③ Virtual environment (Gaming Appln).

① User can intentionally give information to the system.
For ex: fixed set of commands, gesture



"We know which type of I/P to provide to give a particular O/P"

② A dialog system with diff I/P modes. Chatbot can be of 1 example."

Chatbot

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

* fixed I/I to get derived O/P

P O R
6 7 26

Page:
Date:

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

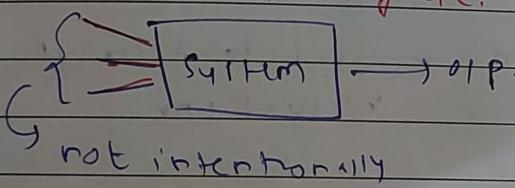
$$\begin{array}{r}
 8 \quad 102 \\
 1100 \\
 \hline
 110 \\
 \hline
 1010 \\
 \hline
 82
 \end{array}$$

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.)

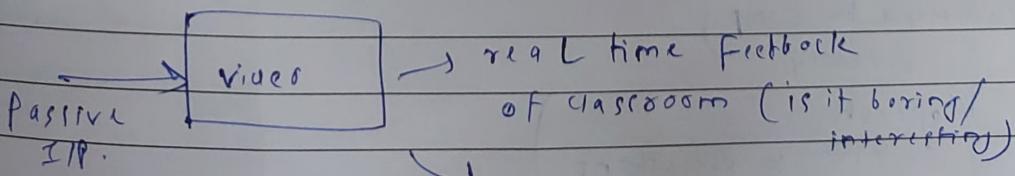


Unlock a phone
with happy
face - intentional

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)

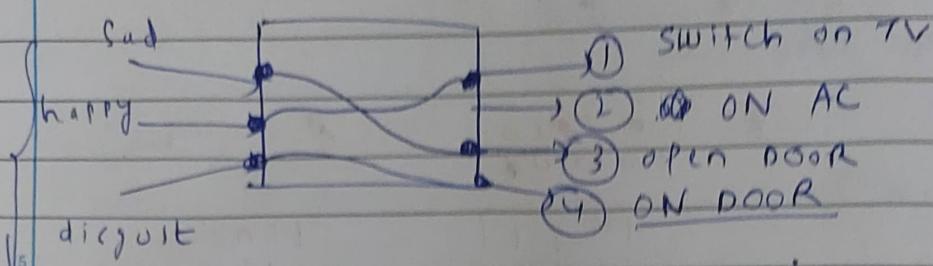
how can classroom

be made interesting / how

whether to continue class

or not / improve skill of
professor

"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.

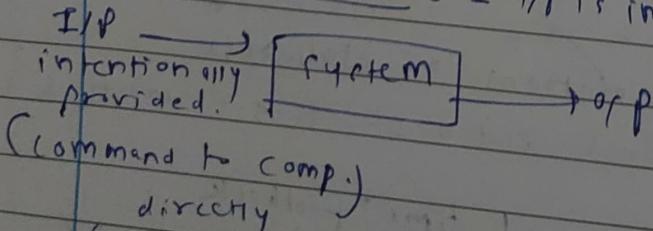
OR

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

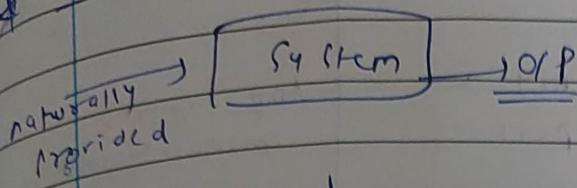
* 25 Active IIP mode = IIP is intentionally provided.



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.

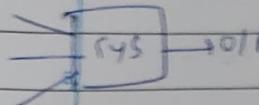
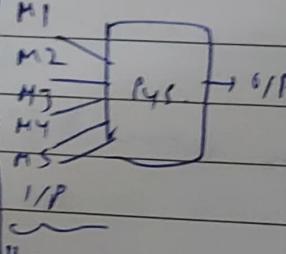
* Process recognition is common for any system.

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 can 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		 Ex: MIT pointing hand + voice	"Sequential" OR "parallel"
M1 + M2 at same time used	"One modality will be selected at a time to give O/P!"	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)

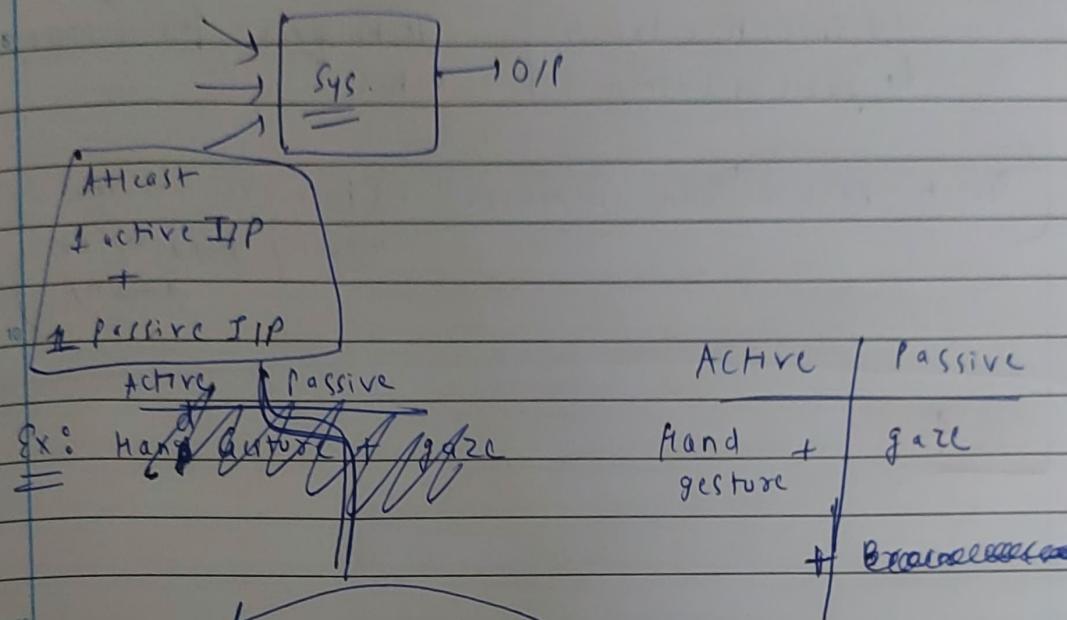
"A system can be of multiple types"

Types of multimodal systems:

Fusion

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

Major application of optical sensors



LIPNET \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

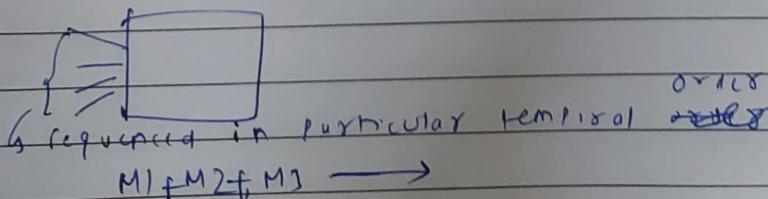
Autonomous
(used in vehicles)

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