# Multimodality

## Learning from Text, Speech, and Vision

CMU 11-4/611 Natural Language Processing

Lecture 28
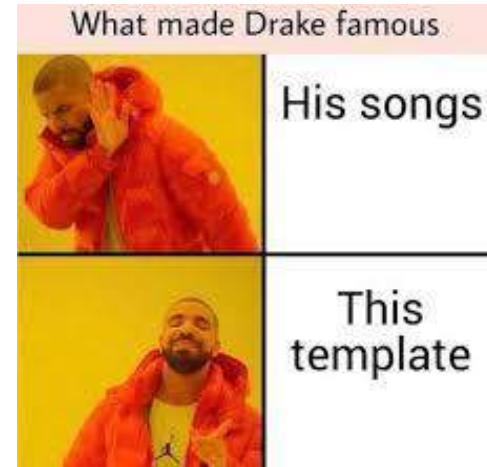April 14, 2020

Shruti Palaskar

# Outline

I.   What is multimodality?

II.  Types of modalities

III. Commonly used Models

IV.  Multimodal Fusion and Representation Learning

V.   Multimodal Tasks: Use Cases

# I. What is Multimodality?

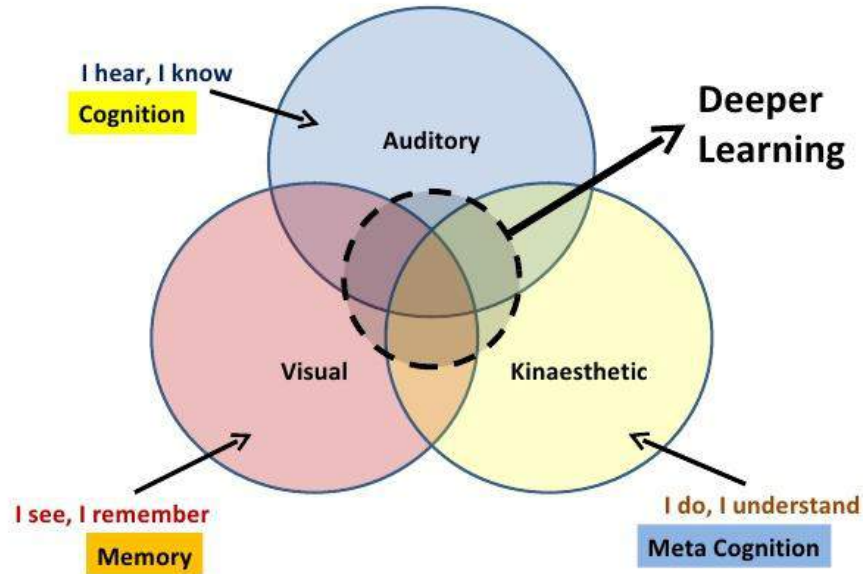# Human Interaction is Inherently Multimodal

# How We Perceive



The curse that afflicts abstract painting

# How We Perceive

## Multi-Modal Learning

I hear, I know
**Cognition**

Auditory

**Deeper Learning**

Visual

Kinaesthetic

I see, I remember
**Memory**

I do, I understand
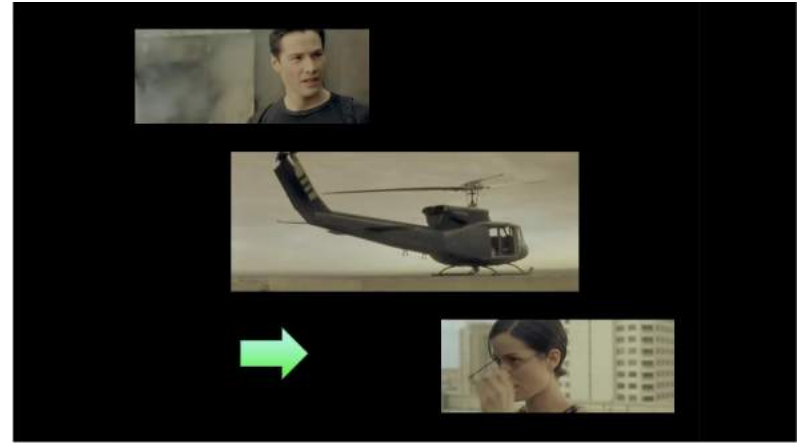**Meta Cognition**

cc Steve Wheeler, University of Plymouth, 2009

# The Dream: Sci-Fi Movies



JARVIS



The Matrix

# Reality?

# Give a caption.

# Give a caption.



Human: A Small Dogs Ears Stick Up As It Runs In The Grass.

Model: A Black And White Dog Is Running On Grass With A Frisbee In Its Mouth

Single sentence image description -> Captioning
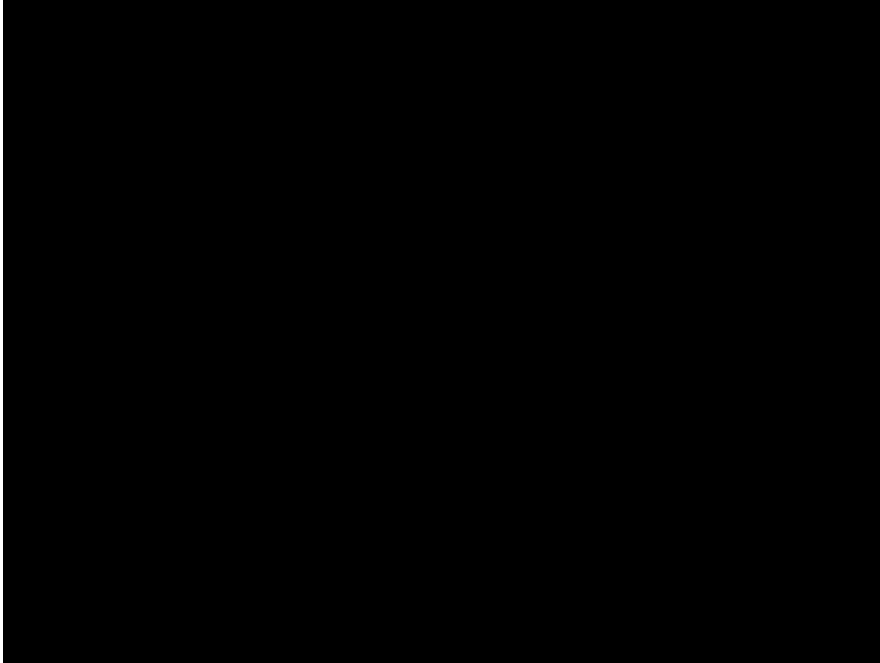
# Give a caption.

# Give a caption.



Human: A Young Girl In A White Dress Standing In Front Of A Fence And Fountain.

Model: Two Men Are Standing In Front Of A Fountain

# Reality?
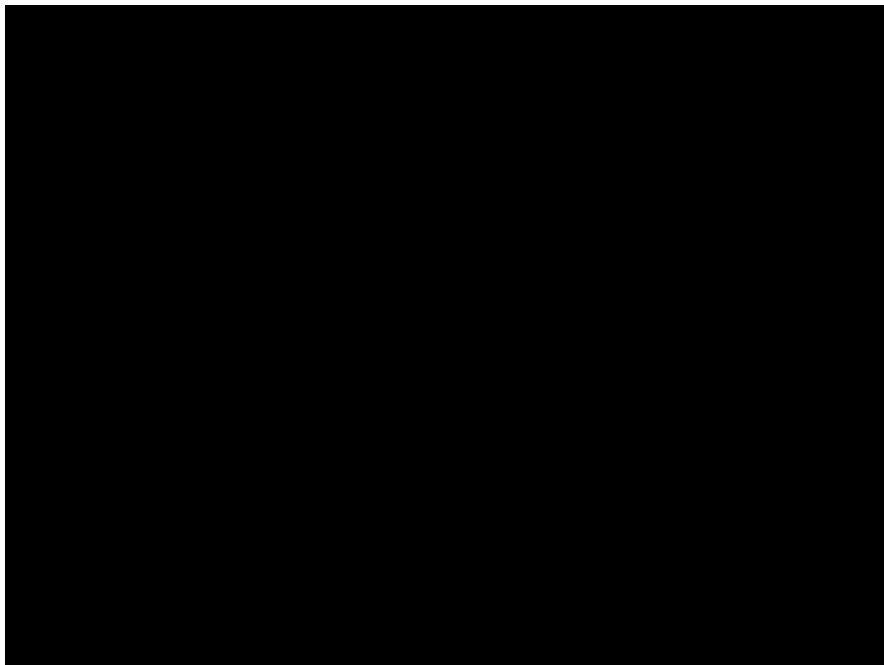
# Watch the video and answer questions.

QUESTIONS

Q. is there only one person ?

Q. does she walk in with a towel around her neck ?

Q. does she interact with the dog ?

Q. does she drop the towel on the floor ?

# Watch the video and answer questions.



QUESTIONS

Q. is there only one person ?
A.     there is only one person and a dog .

Q. does she walk in with a towel around her neck ?
A.     she walks in from outside with the towel around her neck .

Q. does she interact with the dog ?
A.     she does not interact with the dog

Q. does she drop the towel on the floor ?
A.     she dropped the towel on the floor at the end of the video .

Simple questions, simple answers -> Video Question Answering

Reality? Baby Steps. Still a long way to go.

# ...Challenges

Common challenges based on the tasks we just saw

- Training Dataset bias
- Very complicated tasks
- Lack of common sense reasoning within models
- No world knowledge available like humans do
  - Physics, Nature, Memory, Experience

How do we teach machines to perceive?

# Outline

I. What is multimodality?

II. Types of modalities

III. Commonly used Models

IV. Multimodal Fusion and Representation Learning

V. Multimodal Tasks: Use Cases

# II. Types of modalities

# Types of Modalities



IMAGE/VIDEO

The Dad ✔
@thedad

me: it's bedtime now

kid: please let me do just ONE thing

me: ok

kid: *starts watching one movie*

5:59 PM · Oct 25, 2019 · Buffer

TEXT

SPEECH/AUDIO

EMOTION/AFFECT
/SENTIMENT

# Example Dataset: ImageNet

- Object Recognition
- Image Tagging/Categorization
- ~14M images
- Knowledge Ontology
- Hierarchical Tags
  - Mammal -> Placental -> Carnivore -> Canine -> Dog -> Working Dog -> Husky

Deng et al. 2009

# Example Dataset: How2 Dataset



**I'm very close to the green but I didn't get it on the green so now I'm in this grass bunker.**

*Eu estou muito perto do green, mas eu não pus a bola no green, então agora estou neste bunker de grama.*

In golf, get the body low in order to get underneath the golf ball when chipping out of thick grass from a side hill lie.

- Speech
- Video
- English Transcript

- Portuguese Transcript
- Summary

Sanabria et al. 2018

# Example Dataset: Open Pose



- Action Recognition
- Pose Estimation
- Human Dynamic
- Body Dynamics

Wei et al. 2016

# III. Commonly Used Models

# Multilayer Perceptrons



Single Perceptron

# Multilayer Perceptrons



Figure: Hugo Larochelle

Feed Forward

OUTPUT (y)

Feed-forward Weights ($W_i$)

$$y = f(\boldsymbol{W_3} \cdot h_2 + b_3)$$

Hidden States $h_1$ & $h_2$

$$h_2 = f(\boldsymbol{W_2} \cdot h_1 + b_2)$$

$$h_1 = f(\boldsymbol{W_1} \cdot x + b_1)$$

INPUT (x)

28

# Multilayer Perceptrons: Uses in Multimedia

# Multilayer Perceptrons: Limitations

**Limitation #1**

Very large amount of input data samples (xi), which requires a gigantic amount of model parameters.

Figure: Ranzatto

# Convolutional Neural Networks (CNNs)



**Translation invariance:** we can use same parameters to capture a specific "feature" in any area of the image. We can use different sets of parameters to capture different features.

These operations are equivalent to perform **convolutions** with different filters.

# Convolutional Neural Networks (CNNs)

## LeNet-5



LeCun et al. 1998

# Convolutional Neural Networks (CNNs) for Image Encoding



Krizhevsky et al. 2012

# Multilayer Perceptrons: Limitations

**Limitation #1**

Very large amount of input data samples (xi), which requires a gigantic amount of model parameters.



Figure: Ranzatto

**Limitation #2**

Does not naturally handle input data of variable dimension

(eg. audio/video/word sequences)

# Recurrent Neural Networks

Build specific connections capturing the temporal evolution

→ **Shared weights in time**

# Recurrent Neural Networks

Feed-forward
Weights (W)

$$h_t = f(W \cdot x_t + U \cdot h_{t-1} + b)$$

Recurrent
Weights (U)

Updated
state

Previous
state

$$h_t = f(W \cdot x_t + U \cdot h_{t-1} + b)$$

INPUT (x)

# Recurrent Neural Networks for Video Encoding



**Combination method**

CNN   CNN   • • •   CNN

Combination is commonly implemented as a small NN on top of a pooling operation (e.g. max, sum, average).

Recurrent Neural Networks are well suited for processing sequences.

Donahue et al. 2015

37

# Attention Mechanism

Network B focuses on different information from network A at every step.

38

# Loss Function: Softmax

**Cross-entropy loss:**

$$L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}\right)$$

**Softmax function**

**Minimizing the negative log likelihood.**

matrix multiply + bias offset

| 0.01 | -0.05 | 0.1 | 0.05 |
|------|-------|-----|------|
| 0.7 | 0.2 | 0.05 | 0.16 |
| 0.0 | -0.45 | -0.2 | 0.03 |

$W$

| -15 |
|-----|
| 22 |
| -44 |
| 56 |

$x_i$

+

| 0.0 |
|-----|
| 0.2 |
| -0.3 |

$b$

$y_i$ | 2

cross-entropy loss (Softmax)

| -2.85 |
|-------|
| 0.86 |
| 0.28 |

*exp* →

| 0.058 |
|-------|
| 2.36 |
| 1.32 |

*normalize*
(to sum to one)

| 0.016 |
|-------|
| 0.631 |
| 0.353 |

- log(0.353)
=
**0.452**

# IV. Multimodal Fusion & Representation Learning

# Fusion: Model Agnostic



**A** **Model-Agnostic Approaches**

**1) Early Fusion**

Modality 1 →
Modality 2 →
→ Classifier →

**2) Late Fusion**

Modality 1 → Classifier →
Modality 2 → Classifier →
→

# Fusion: Model Based

**Definition:** To join information from two or more modalities to perform a prediction task.

(B) **Model-Based (Intermediate) Approaches**

1) **Deep neural networks**

2) **Kernel-based methods**

3) **Graphical models**



Multiple kernel learning



Multi-View Hidden CRF

# Representation Learning: Encoder-Decoder

# Representation Learning



INPUT    PROJECTION    OUTPUT

w(t-2)

w(t-1)

SUM

w(t+1)

w(t+2)

w(t)

the cat climbed a tree

Given context:

a, cat, the, tree

Estimate prob. of

climbed

**Word2Vec**

Mikolov et al. 2013

# Representation Learning: RNNs



$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

Cho et al. 2014

# Representation Learning: Self-Supervised

Use videos to train a CNN that predicts the audio statistics of a frame.

# Representation Learning: Transfer Learning



Teacher network: Visual Recognition (object & scenes)

# Representation Learning: Joint Learning

# Representation Learning: Joint Learning (Similarity)



Audio-visual correspondence detector network

Vision subnetwork

Audio subnetwork

Fusion layers

**Correspond?**

Yes / No

# V. Common Tasks, Use Cases

# V. Common Tasks

1. Vision and Language
2. Speech, Vision and Language
3. Multimedia
4. Emotion and Affect

- Image/Video Captioning
- Visual Question Answering
- Visual Dialog
- Video Summarization
- Lip Reading
- Audio Visual Speech Recognition
- Visual Speech Synthesis
- …

1. Vision and Language Common Tasks

# Image Captioning



Vinyals et al. 2015

# Image Captioning



"straw" "hat" END

$y_t$

$W_{oh}$

$W_{hh}$ $h_t$

$CNN_{\theta_c}$ $W_{hi}$

$W_{hx}$

$x_t$

START "straw" "hat"

Karpathy et al. 2015

Slides by Marc Bolaños 54

# Image Captioning: Show, Attend and Tell



14x14 Feature Map

A bird flying over a body of water

LSTM

1. Input Image
2. Convolutional Feature Extraction
3. RNN with attention over the image
4. Word by word generation

Xu et al. 2015

# Image Captioning and Detection



a plate of food. food on a plate. a blue cup on a table. a plate of food. a blue bowl with red sauce. a bowl of soup. a cup of coffee. a bowl of chocolate. a glass of water. a plate of food. a silver metal container. a small bowl of sauce. table with food on it. a slice of orange. a table with food on it. a slice of meat. yellow and white cheese.

Johnson et al. 2016

# Video Captioning



Donahue et al. 2015

57

# Video Captioning



LSTM unit (2nd layer)

hidden state at t = T

(a) Stacked LSTM video encoder

(b) Hierarchical Recurrent Neural Encoder

Image

t = 1    Time    t = T

first chunk of data

# Visual Question Answering



What is the mustache made of?

AI System

# Visual Question Answering



What is the mustache made of?

AI System → bananas

# Visual Question Answering



"Yes"

Decode

$[z_1, z_2, \ldots z_N]$ ⊕ $[y_1, y_2, \ldots y_M]$

Encode

ECONOMIC GROWTH
IN PERCENTAGE

Encode

"Is economic growth decreasing ?"

5.74  6.07  6.71  6.23  6.03  6.2  5.7

FY09  FY10  FY11  FY12  FY13  FY14*  FY14**

61

# Visual Question Answering

# Visual Question Answering



63

# Video Summarization

**~1.5 minutes of audio and video**

"Teaser" (33 words on avg)

how to cut peppers to make a spanish omelette ; get expert tips and advice on making cuban breakfast recipes in this free cooking video .

Transcript (290 words on avg)

on behalf of expert village my name is lizbeth muller and today we are going to show you how to make spanish omelet . i 'm going to dice a little bit of peppers here . i 'm not going to use a lot , i 'm going to use very very little . a little bit more then this maybe . you can use red peppers if you like to get a little bit color in your omelet . some people do and some people do n't . but i find that some of the people that are mexicans who are friends of mine that have a mexican she like to put red peppers and green peppers and yellow peppers in hers and with a lot of onions . that is the way they make there spanish omelets that is what she says . i loved it , it actually tasted really good . you are going to take the onion also and dice it really small . you do n't want big chunks of onion in there cause it is just pops out of the omelet . so we are going to dice the up also very very small . so we have small pieces of onions and peppers ready to go .

# Video Summarization: Hierarchical Model



Palaskar et al. 2019

# Action Recognition



ResNeXt



Can Spatiotemporal 3D CNNs Retrace the History of 2D CNNs and ImageNet?

Kensho Hara, Hirokatsu Kataoka, Yutaka Satoh
National Institute of Advanced Industrial Science and Technology (AIST)
Tsukuba, Ibaraki, Japan
{kensho.hara, hirokatsu.kataoka, yu.satou}@aist.go.jp

# 2. Speech, Vision and Language Common Tasks

# Audio Visual Speech Recognition: Lip Reading



t frames | STCNN + Spatial Pooling (x3) | Bi-GRU (x2) | Linear | CTC loss

Assael et al. 2016

# Lip Reading: Watch, Listen, Attend and Spell



Figure 1. *Watch, Listen, Attend and Spell* architecture. At each time step, the decoder outputs a character $y_i$, as well as two attention vectors. The attention vectors are used to select the appropriate period of the input visual and audio sequences.

Chung et al. 2017

# 3. Multimedia Common Tasks

# Multimedia Retrieval

# Multimedia Retrieval

# Multimedia Retrieval: Shared Multimodal Representation

# Multimedia Retrieval

# 4. Emotion and Affect

# Affect Recognition:
# Emotion, Sentiment, Persuasion, Personality

# Outline

I.   What is multimodality?

II.   Types of modalities

III.   Commonly used Models

IV.   Multimodal Fusion and Representation Learning

V.   Multimodal Tasks: Use Cases

# Takeaways

- Lots of multimodal data generated everyday
- Need automatic ways to understand it
  - Privacy
  - Security
  - Regulation
  - Storage
- Different models used for different downstream tasks
  - Highly open-ended research!
- Try it out for fun on Kaggle!

Thank you!

spalaska@cs.cmu.edu