



Language
Technologies
Institute

Carnegie
Mellon
University

Advanced Multimodal Machine Learning

Lecture 1.1: Introduction

Louis-Philippe Morency

* Original version co-developed with Tadas Baltrusaitis

Your Instructor and TAs This Semester (11-777)



Louis-Philippe Morency

morency@cs.cmu.edu

Office: GHC-5411



Carla De Oliveira Viegas

cviegas@andrew.cmu.edu



Harsh Jhamtani

jharsh@andrew.cmu.edu



Varun Lakshminarasimhan

vbl@andrew.cmu.edu



Zhun Liu

zhunl@andrew.cmu.edu



Ying Shen

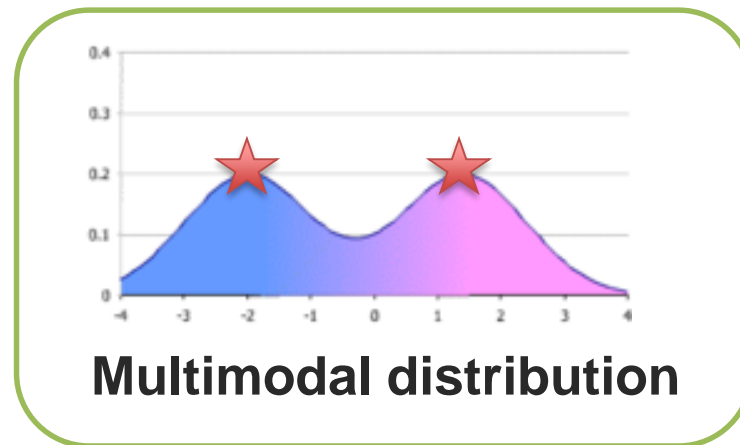
yshen2@andrew.cmu.edu

Lecture Objectives

- Introductions
- What is Multimodal?
 - Multimodal communicative behaviors
- A historical view of multimodal research
- Core technical challenges
 - Representation, translation, alignment, fusion and alignment
- Course syllabus and project assignments
 - Grades and course structure

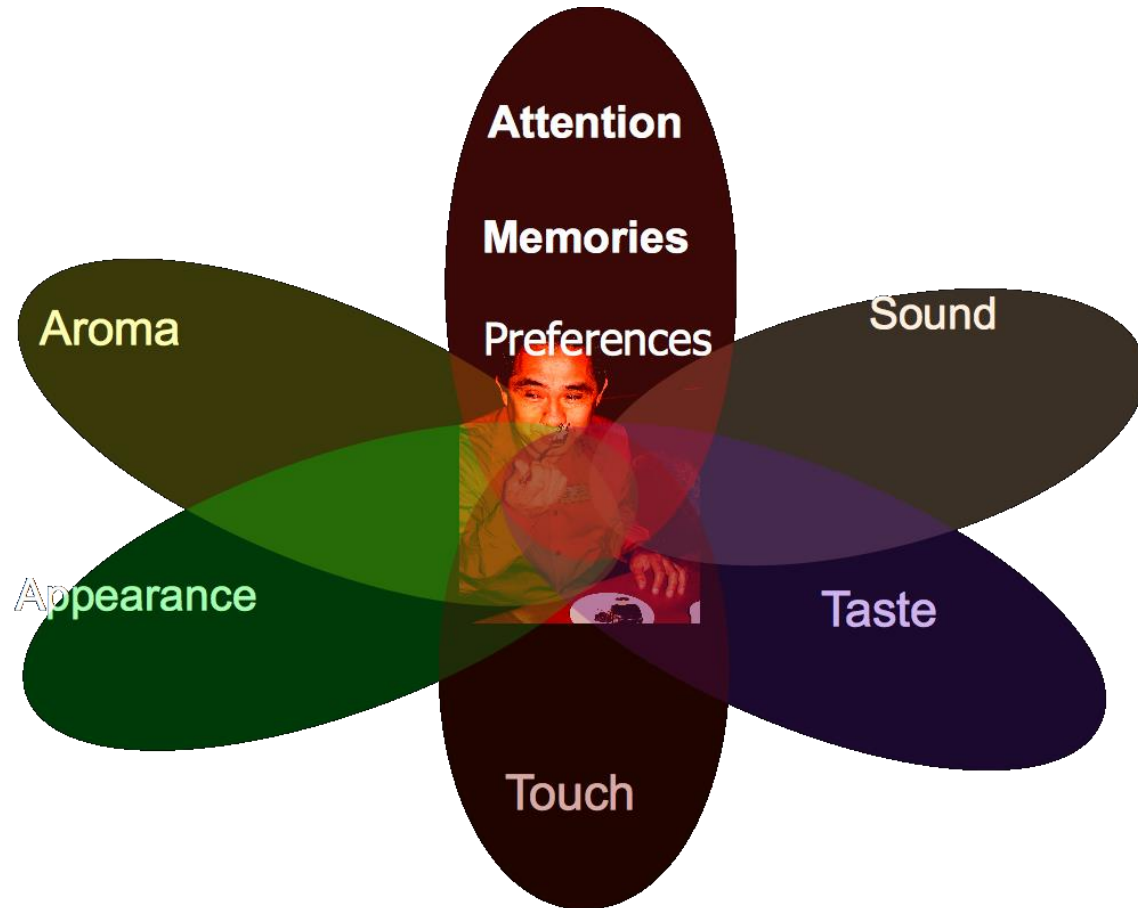
What is Multimodal?

What is Multimodal?



- Multiple modes, i.e., distinct “peaks” (local maxima) in the probability density function

What is Multimodal?



Sensory Modalities

Multimodal Communicative Behaviors

Verbal

Lexicon

Words

Syntax

Part-of-speech

Dependencies

Pragmatics

Discourse acts

Vocal

Prosody

Intonation

Voice quality

Vocal expressions

Laughter, moans

Visual

Gestures

Head gestures

Eye gestures

Arm gestures

Body language

Body posture

Proxemics

Eye contact

Head gaze

Eye gaze

Facial expressions

FACS action units

Smile, frowning



What is Multimodal?

Modality

The way in which something happens or is experienced.

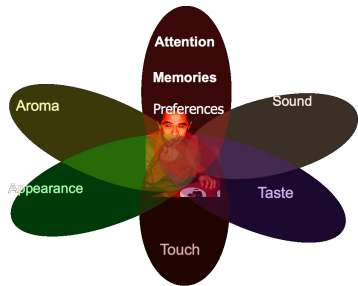
- *Modality* refers to a certain type of information and/or the representation format in which information is stored.
- *Sensory modality*: one of the primary forms of sensation, as vision or touch; channel of communication.

Medium (“middle”)

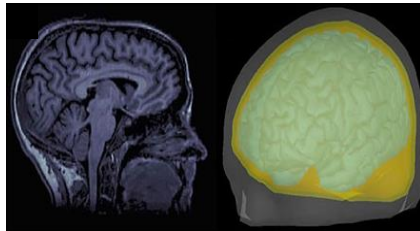
A means or instrumentality for storing or communicating information; system of communication/transmission.

- *Medium* is the means whereby this information is delivered to the senses of the interpreter.

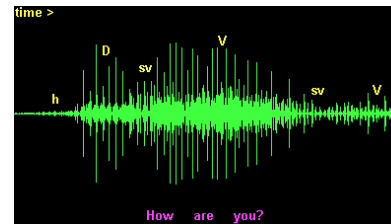
Multiple Communities and Modalities



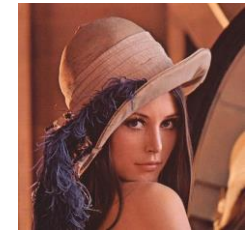
Psychology



Medical



Speech



Vision



Language



Multimedia



Robotics

$$\frac{\partial}{\partial \theta} \log \pi(x, \theta) = \frac{1}{\pi(x, \theta)} \frac{\partial \pi(x, \theta)}{\partial \theta}$$
$$\int \tau(x) \cdot \frac{\partial}{\partial \theta} \log \pi(x, \theta) \cdot \pi(x, \theta) dx = \int \tau(x) \cdot \frac{\partial \pi(x, \theta)}{\partial \theta} dx$$
$$\frac{\partial}{\partial \theta} \log \pi(x, \theta) = \frac{1}{\pi(x, \theta)} \frac{\partial \pi(x, \theta)}{\partial \theta}$$

Learning

Examples of Modalities

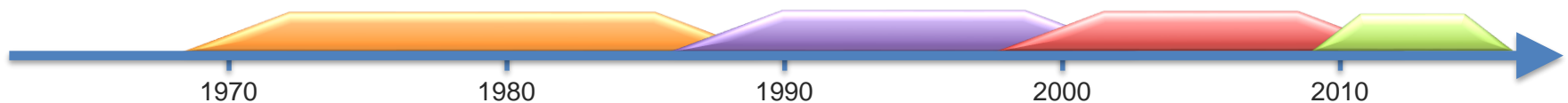
- Natural language (both spoken or written)
- Visual (from images or videos)
- Auditory (including voice, sounds and music)
- Haptics / touch
- Smell, taste and self-motion
- Physiological signals
 - Electrocardiogram (ECG), skin conductance
- Other modalities
 - Infrared images, depth images, fMRI

A Historical View

Prior Research on “Multimodal”

Four eras of multimodal research

- The “behavioral” era (1970s until late 1980s)
- The “computational” era (late 1980s until 2000)
- The “interaction” era (2000 - 2010)
- The “deep learning” era (2010s until ...)
 - ❖ Main focus of this course



Language and Gestures



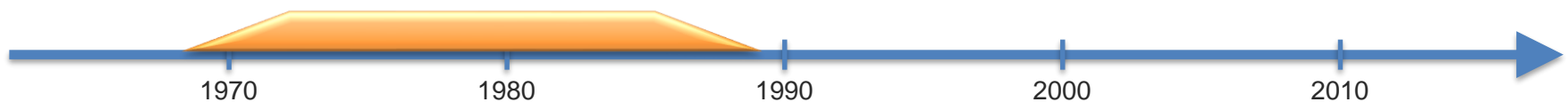
David McNeill

University of Chicago

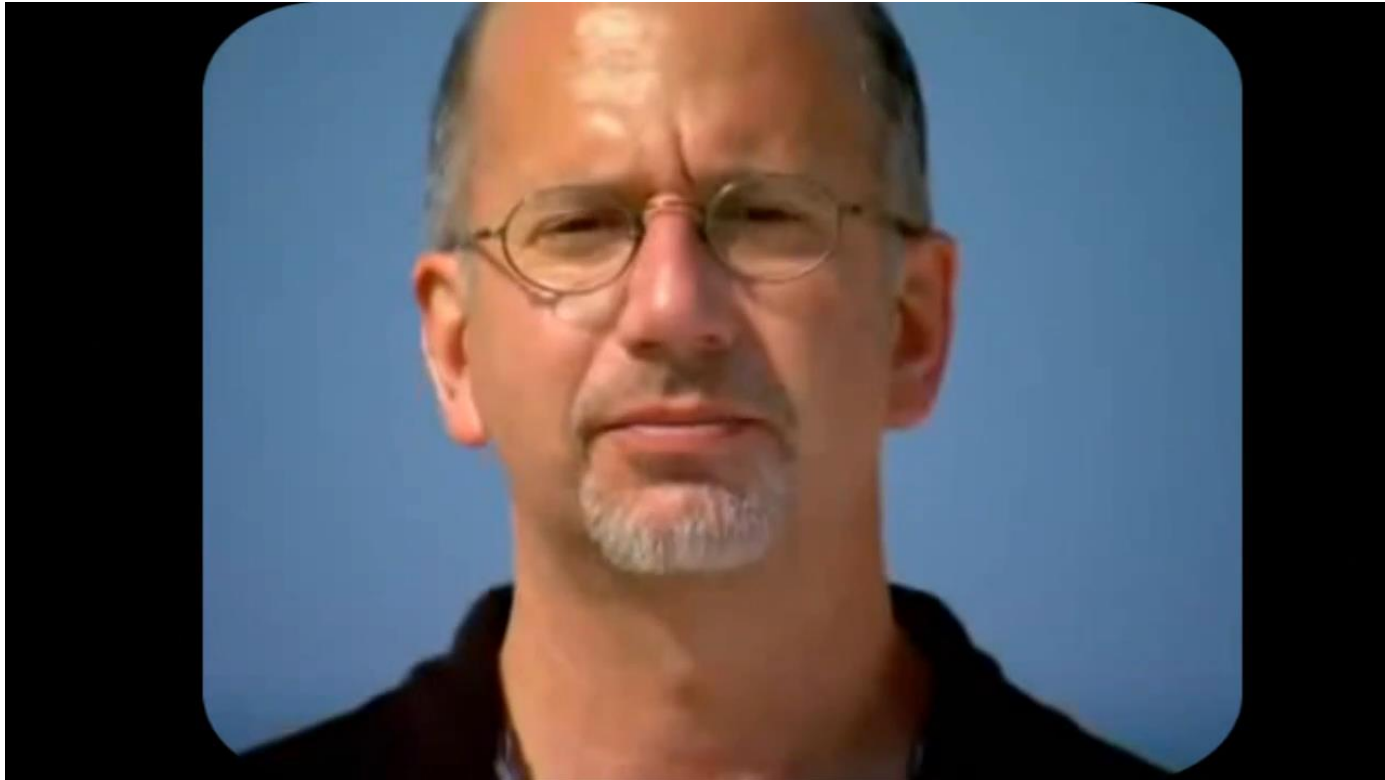
Center for Gesture and Speech Research

“For McNeill, gestures are in effect the speaker’s thought in action, and integral components of speech, not merely accompaniments or additions.”

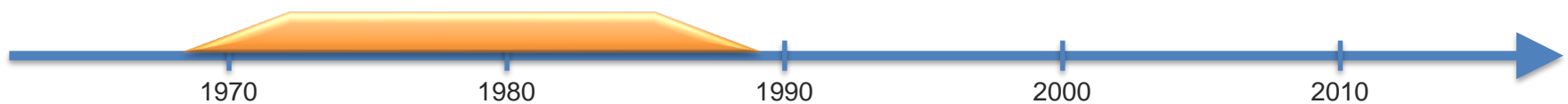
□ TRIVIA: Justine Cassell was a student of David McNeill



The McGurk Effect (1976)



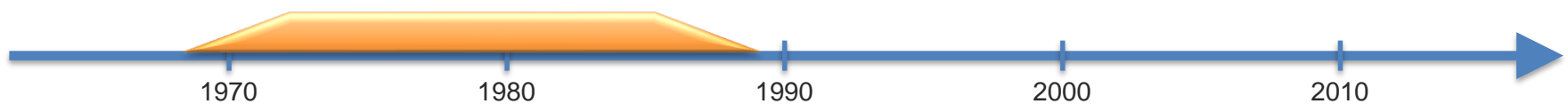
Hearing lips and seeing voices – Nature



The McGurk Effect (1976)

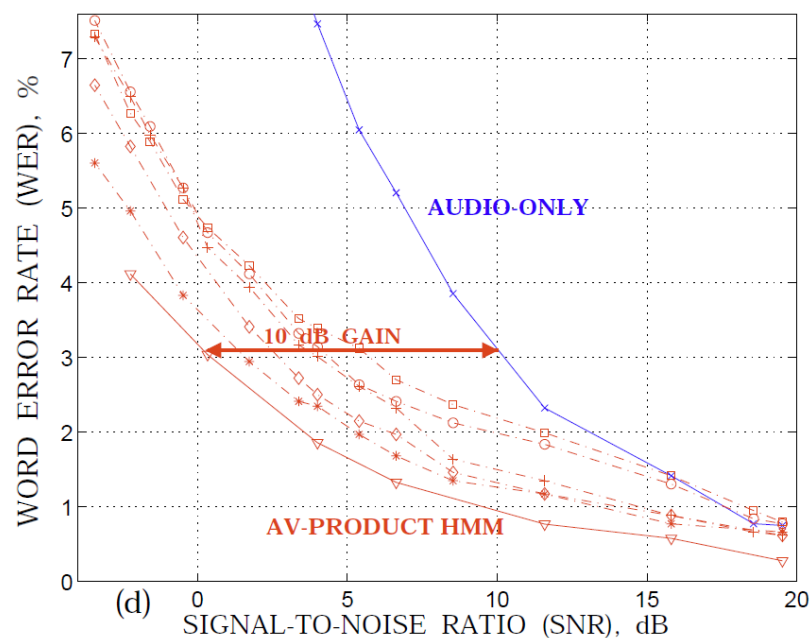


Hearing lips and seeing voices – Nature



➤ The “Computational” Era(Late 1980s until 2000)

1) Audio-Visual Speech Recognition (AVSR)



1970

1980

1990

2000

2010



➤ The “Computational” Era (Late 1980s until 2000)

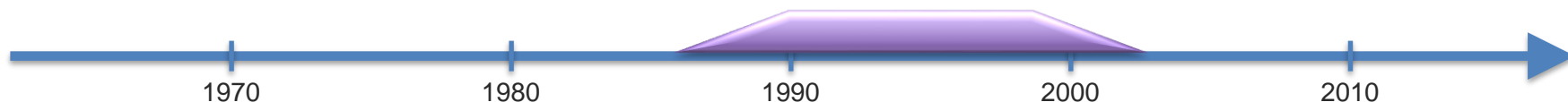
2) Multimodal/multisensory interfaces



Rosalind Picard

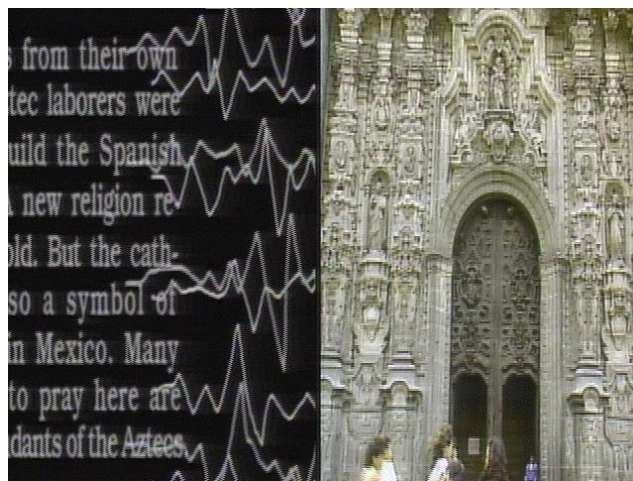
Affective Computing is computing that relates to, arises from, or deliberately influences emotion or other affective phenomena.

- ❑ TRIVIA: Rosalind Picard came from the same group (MIT, Sandy Pentland)



➤ The “Computational” Era (Late 1980s until 2000)

3) Multimedia Computing

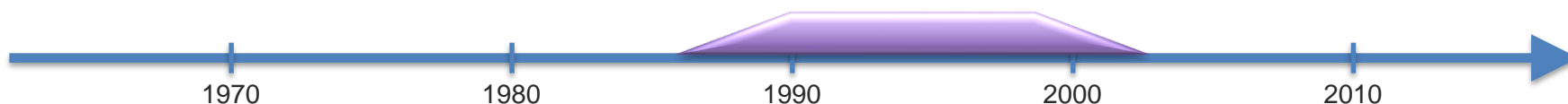


**Carnegie
Mellon
University**



[1994-2010]

“The Informedia Digital Video Library Project automatically combines speech, image and natural language understanding to create a full-content searchable digital video library.”



➤ The “Interaction” Era (2000s)

1) Modeling Human Multimodal Interaction



AMI Project [2001-2006, IDIAP]

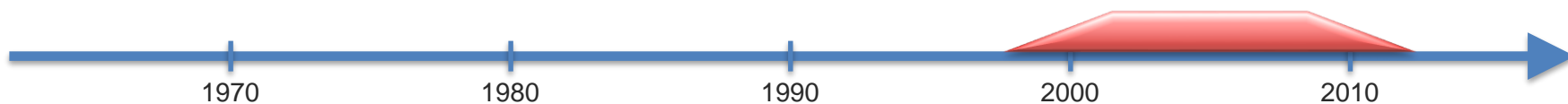
- 100+ hours of meeting recordings
- Fully synchronized audio-video
- Transcribed and annotated



CHIL Project [Alex Waibel]

- Computers in the Human Interaction Loop
- Multi-sensor multimodal processing
- Face-to-face interactions

❑ TRIVIA: Samy Bengio started at IDIAP working on AMI project



➤ The “Interaction” Era (2000s)

1) Modeling Human Multimodal Interaction



CALO Project [2003-2008, SRI]

- Cognitive Assistant that Learns and Organizes
- Personalized Assistant that Learns (PAL)
- Siri was a spinoff from this project



Social Signal Processing Network

SSP Project [2008-2011, IDIAP]

- Social Signal Processing
- First coined by Sandy Pentland in 2007
- Great dataset repository: <http://sspnet.eu/>

❑ TRIVIA: LP's PhD research was partially funded by CALO ☺



1970

1980

1990

2000

2010



➤ The “deep learning” era (2010s until ...)

Representation learning (a.k.a. deep learning)

- Multimodal deep learning [ICML 2011]
- Multimodal Learning with Deep Boltzmann Machines [NIPS 2012]
- Visual attention: Show, Attend and Tell: Neural Image Caption Generation with Visual Attention [ICML 2015]

Key enablers for multimodal research:

- New large-scale multimodal datasets
- Faster computer and GPUS
- High-level visual features
- “Dimensional” linguistic features

Our course focuses on this era!



Core Technical Challenges

Core Challenges in “Deep” Multimodal ML

Representation

Alignment

Fusion

Translation

Co-Learning

Multimodal Machine Learning: A Survey and Taxonomy

By Tadas Baltrusaitis, Chaitanya Ahuja,
and Louis-Philippe Morency

<https://arxiv.org/abs/1705.09406>

- ✓ 5 core challenges
- ✓ 37 taxonomic classes
- ✓ 253 referenced citations

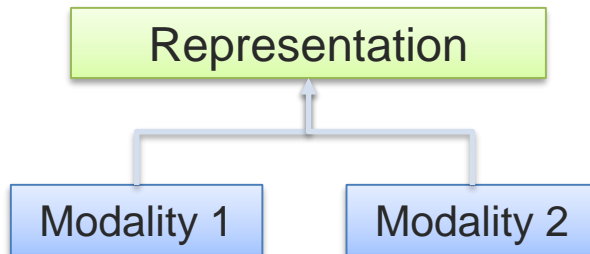
These challenges are non-exclusive.



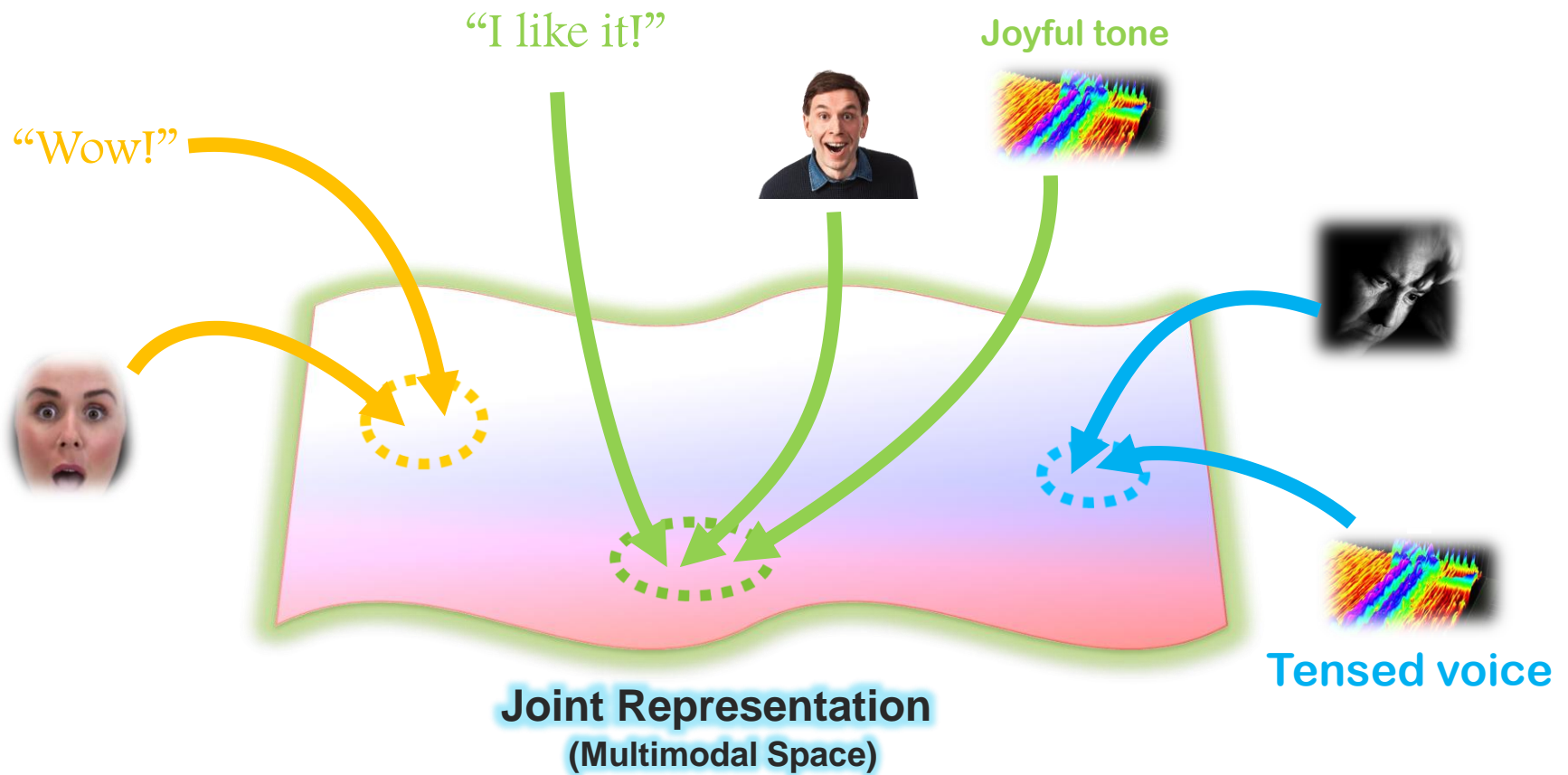
Core Challenge 1: Representation

Definition: Learning how to represent and summarize multimodal data in a way that exploits the complementarity and redundancy.

Ⓐ Joint representations:



Joint Multimodal Representation



Joint Multimodal Representations

Audio-visual speech recognition

[Ngiam et al., ICML 2011]

- Bimodal Deep Belief Network

Image captioning

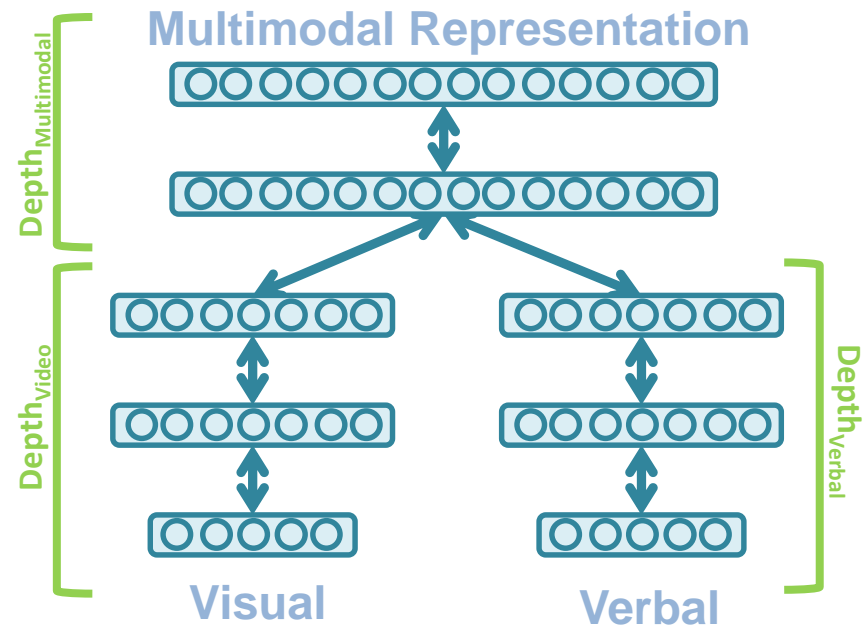
[Srivastava and Salahutdinov, NIPS 2012]

- Multimodal Deep Boltzmann Machine

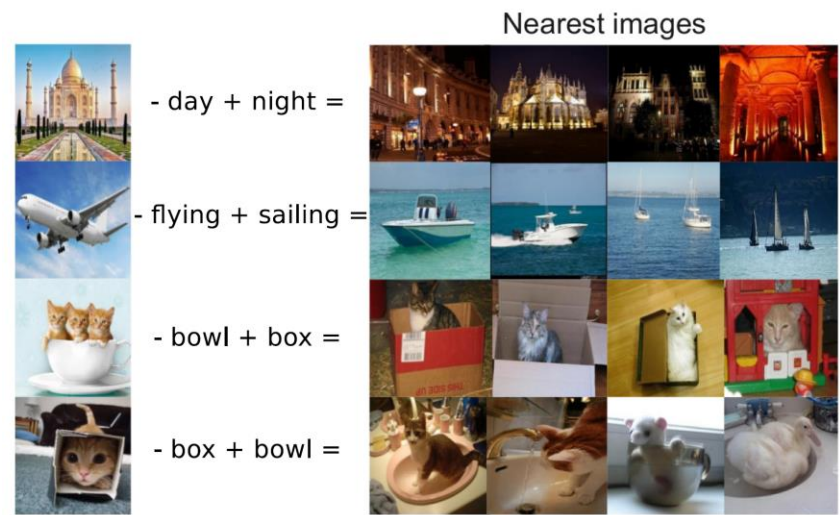
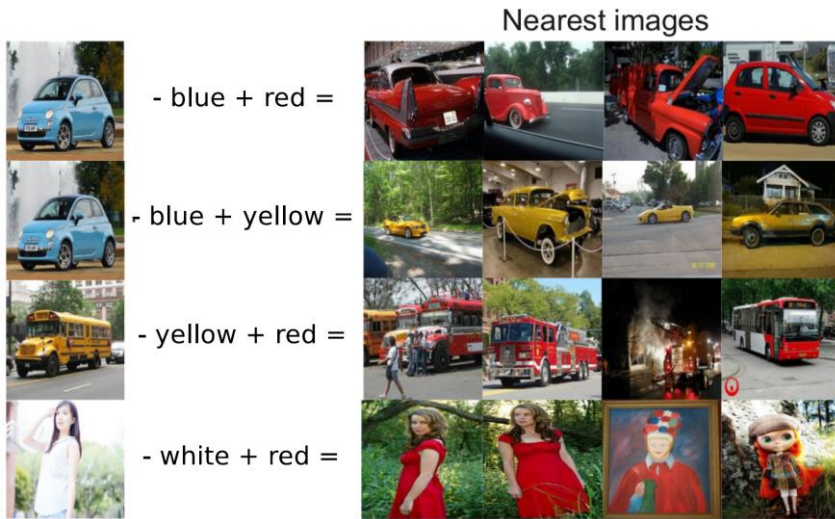
Audio-visual emotion recognition

[Kim et al., ICASSP 2013]

- Deep Boltzmann Machine



Multimodal Vector Space Arithmetic

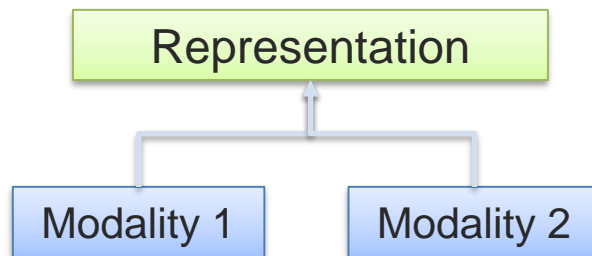


[Kiros et al., Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, 2014]

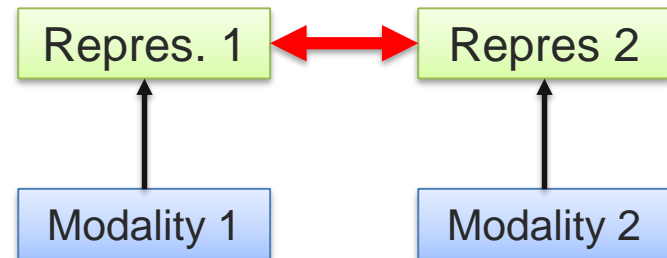
Core Challenge 1: Representation

Definition: Learning how to represent and summarize multimodal data in a way that exploits the complementarity and redundancy.

Ⓐ Joint representations:



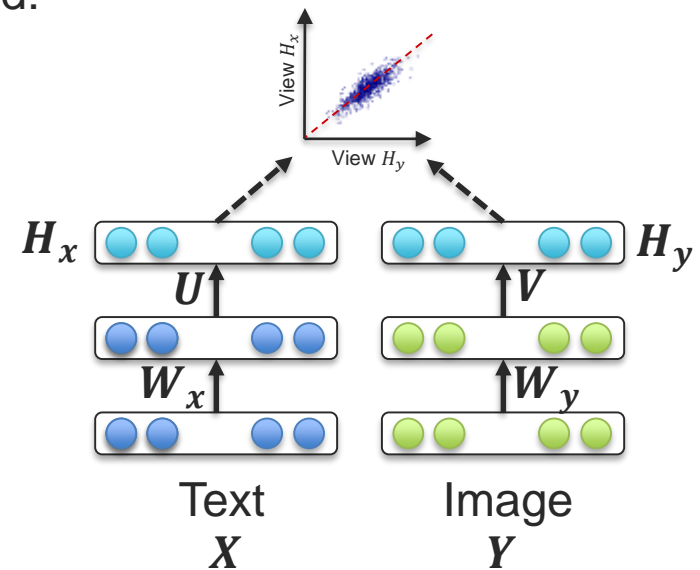
Ⓑ Coordinated representations:



Coordinated Representation: Deep CCA

Learn linear projections that are maximally correlated:

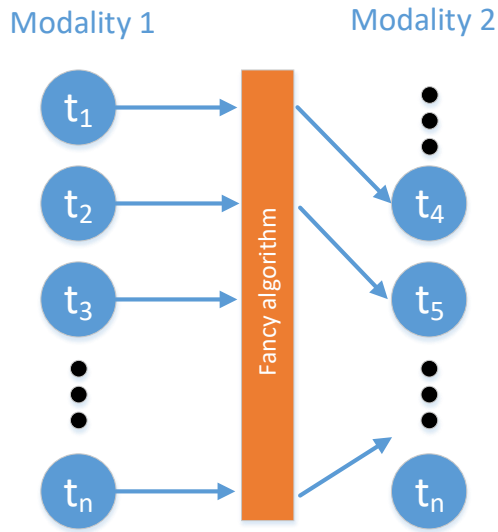
$$(\mathbf{u}^*, \mathbf{v}^*) = \operatorname{argmax}_{\mathbf{u}, \mathbf{v}} \operatorname{corr}(\mathbf{u}^T \mathbf{X}, \mathbf{v}^T \mathbf{Y})$$



Andrew et al., ICML 2013

Core Challenge 2: Alignment

Definition: Identify the direct relations between (sub)elements from two or more different modalities.



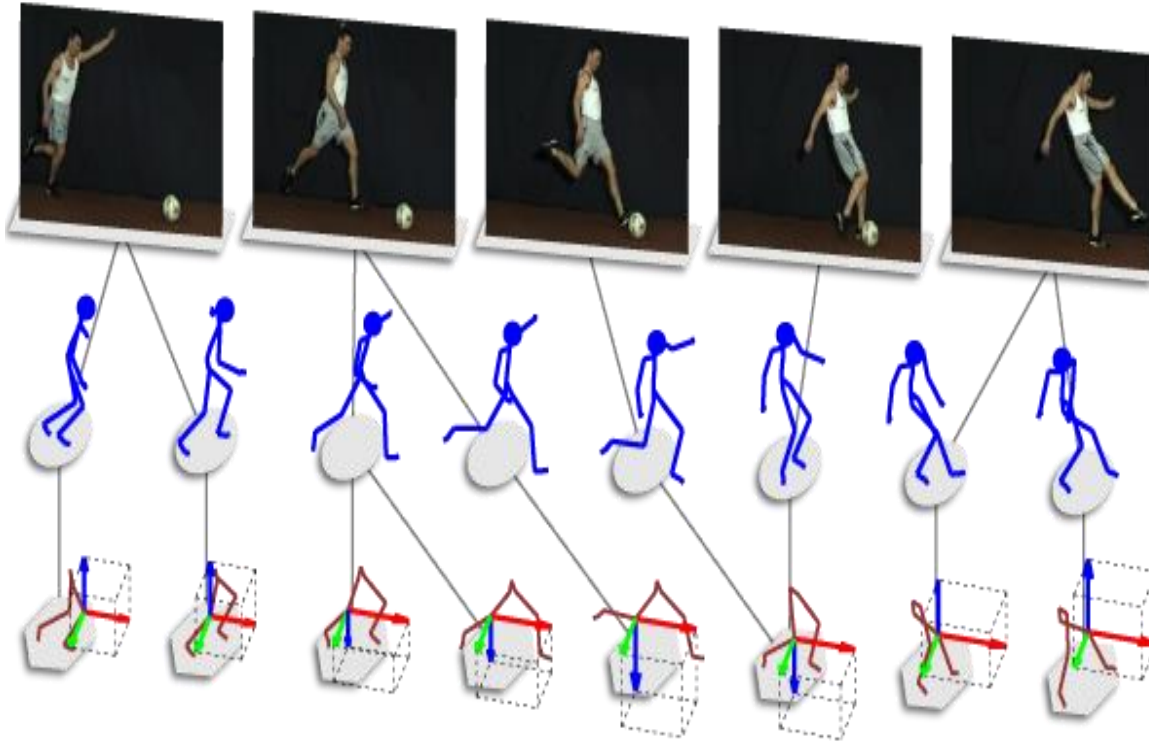
A Explicit Alignment

The goal is to directly find correspondences between elements of different modalities

B Implicit Alignment

Uses internally latent alignment of modalities in order to better solve a different problem

Temporal sequence alignment

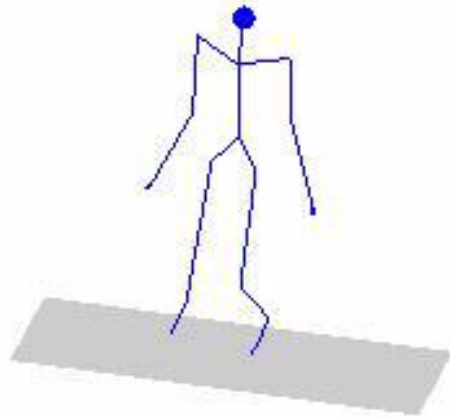


Applications:

- Re-aligning asynchronous data
- Finding similar data across modalities (we can estimate the aligned cost)
- Event reconstruction from multiple sources

Alignment examples (multimodal)

1/273



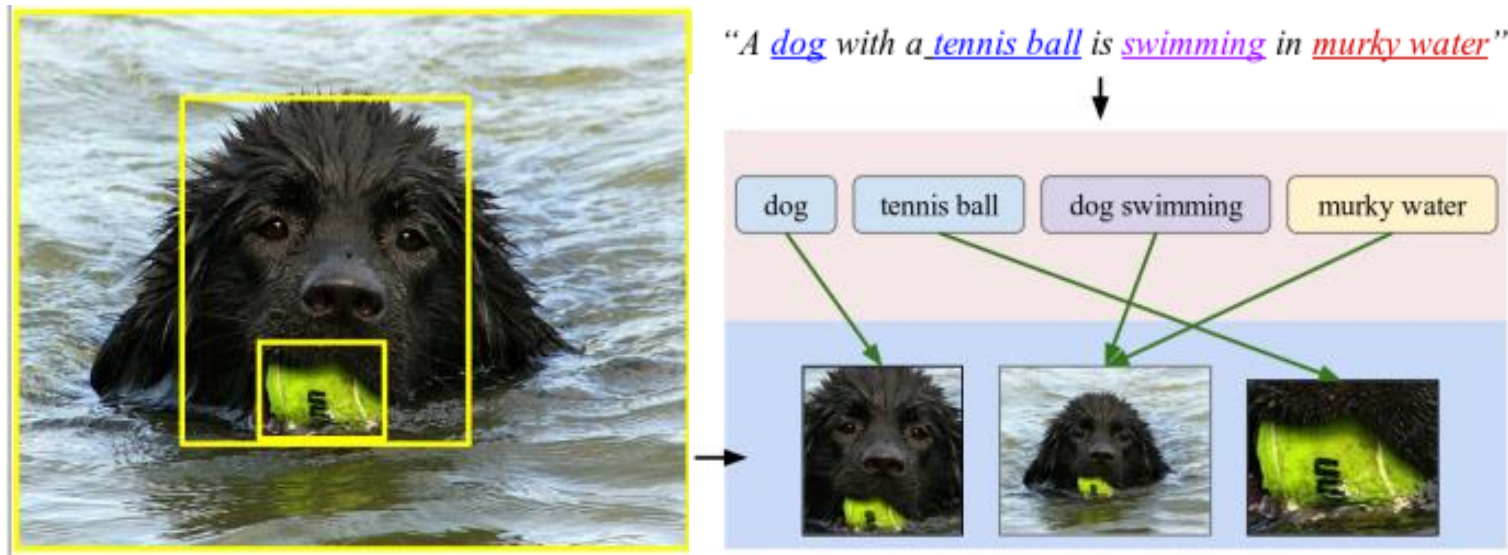
1/51



1/127



Implicit Alignment



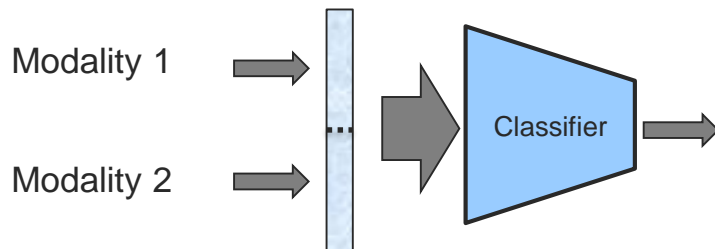
Karpathy et al., Deep Fragment Embeddings for Bidirectional Image Sentence Mapping, <https://arxiv.org/pdf/1406.5679.pdf>

Core Challenge 3: Fusion

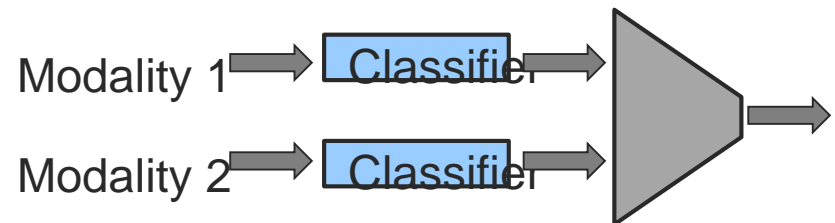
Definition: To join information from two or more modalities to perform a prediction task.

A Model-Agnostic Approaches

1) Early Fusion



2) Late Fusion

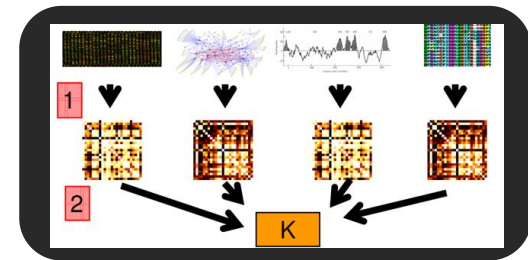


Core Challenge 3: Fusion

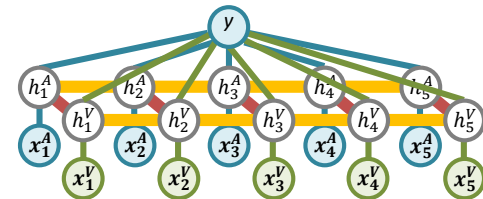
Definition: To join information from two or more modalities to perform a prediction task.

B Model-Based (Intermediate) Approaches

- 1) Deep neural networks
- 2) Kernel-based methods
- 3) Graphical models



Multiple kernel learning

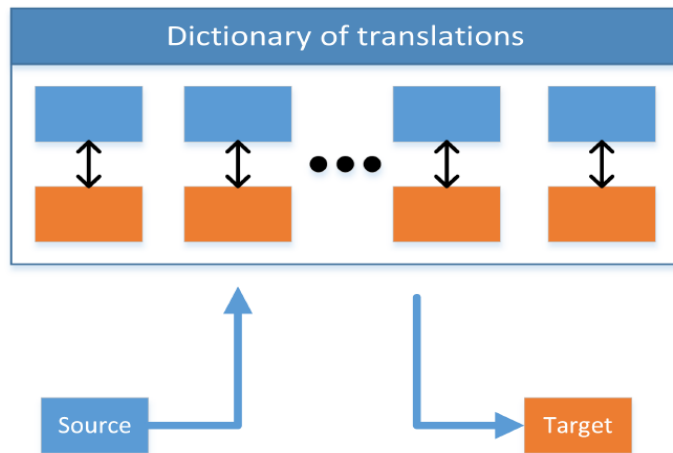


Multi-View Hidden CRF

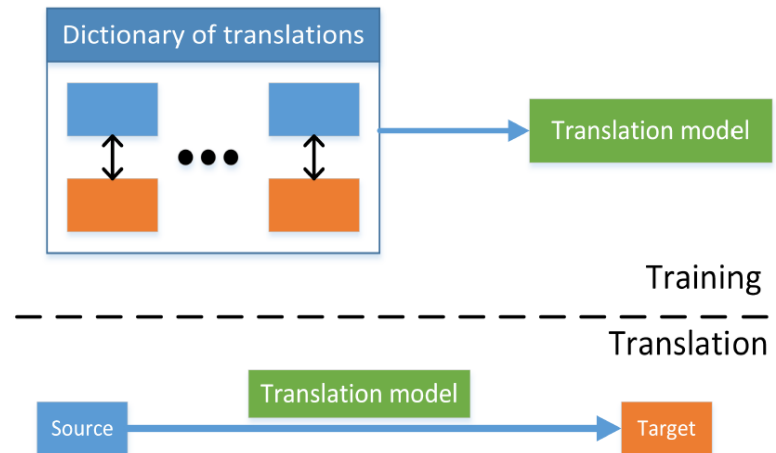
Core Challenge 4: Translation

Definition: Process of changing data from one modality to another, where the translation relationship can often be open-ended or subjective.

A Example-based



B Model-driven



Core Challenge 4 – Translation



Visual gestures
(both speaker and
listener gestures)

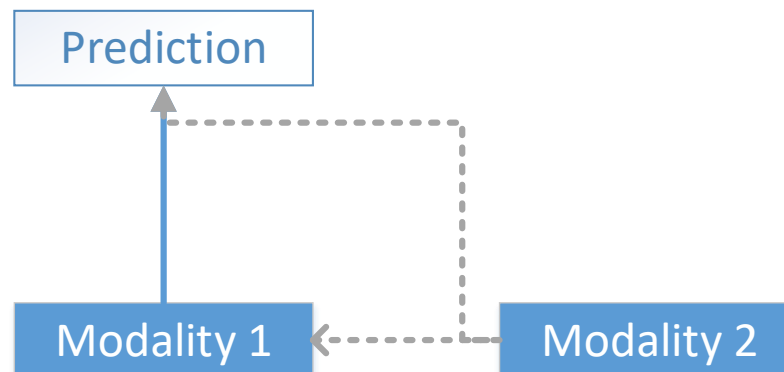


Transcriptions
+
Audio streams

Marsella et al., Virtual character performance from speech, SIGGRAPH/Eurographics Symposium on Computer Animation, 2013

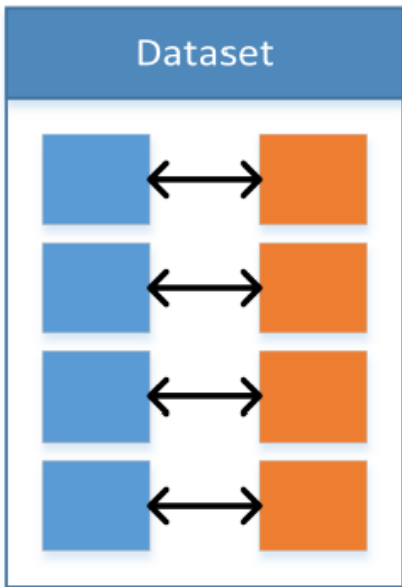
Core Challenge 5: Co-Learning

Definition: Transfer knowledge between modalities, including their representations and predictive models.



Core Challenge 5: Co-Learning

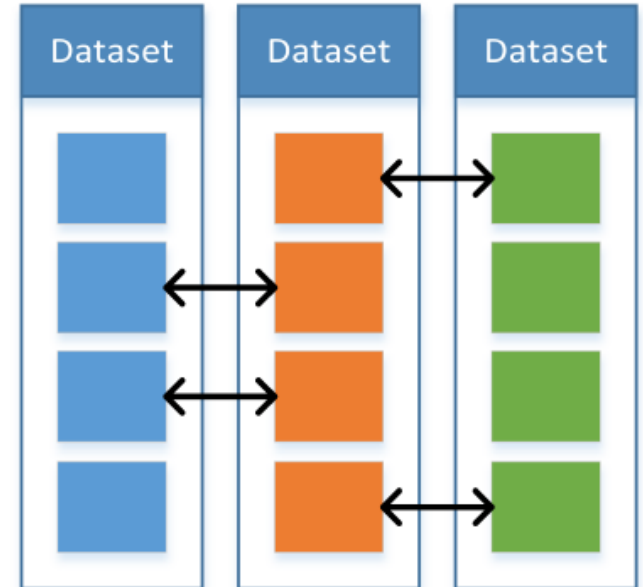
(A) Parallel



(B) Non-Parallel



(C) Hybrid



Taxonomy of Multimodal Research

[<https://arxiv.org/abs/1705.09406>]

Representation

- Joint
 - *Neural networks*
 - *Graphical models*
 - *Sequential*
- Coordinated
 - *Similarity*
 - *Structured*

Translation

- Example-based
 - *Retrieval*
 - *Combination*
- Model-based
 - *Grammar-based*

- *Encoder-decoder*
- *Online prediction*

Alignment

- Explicit
 - *Unsupervised*
 - *Supervised*
- Implicit
 - *Graphical models*
 - *Neural networks*

Fusion

- Model agnostic
 - *Early fusion*
 - *Late fusion*
 - *Hybrid fusion*

Model-based

- *Kernel-based*
- *Graphical models*
- *Neural networks*

Co-learning

- Parallel data
 - *Co-training*
 - *Transfer learning*
- Non-parallel data
 - *Zero-shot learning*
 - *Concept grounding*
 - *Transfer learning*
- *Hybrid data*
 - *Bridging*

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency, Multimodal Machine Learning: A Survey and Taxonomy

Real world tasks tackled by MMML

- Affect recognition
 - Emotion
 - Persuasion
 - Personality traits
- Media description
 - Image captioning
 - Video captioning
 - Visual Question Answering
- Event recognition
 - Action recognition
 - Segmentation
- Multimedia information retrieval
 - Content based/Cross-media



Multimodal Applications

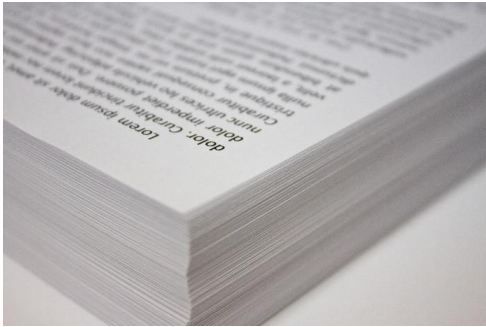
[<https://arxiv.org/abs/1705.09406>]

APPLICATIONS	CHALLENGES				
	REPRESENTATION	TRANSLATION	FUSION	ALIGNMENT	CO-LEARNING
Speech Recognition and Synthesis					
Audio-visual Speech Recognition	✓		✓	✓	✓
(Visual) Speech Synthesis	✓	✓			
Event Detection					
Action Classification	✓		✓		✓
Multimedia Event Detection	✓		✓		✓
Emotion and Affect					
Recognition	✓		✓	✓	✓
Synthesis	✓	✓			
Media Description					
Image Description	✓	✓		✓	✓
Video Description	✓	✓	✓	✓	✓
Visual Question-Answering	✓		✓	✓	✓
Media Summarization	✓	✓	✓		
Multimedia Retrieval					
Cross Modal retrieval	✓	✓		✓	✓
Cross Modal hashing	✓				✓

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency, Multimodal Machine Learning: A Survey and Taxonomy

Course Syllabus

Three Course Learning Paradigms



Research paper reading
and group discussion
(40% of your grade)

$$\begin{aligned}i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\h_t &= o_t \tanh(c_t)\end{aligned}$$

Course project assignments
(60% of your grade)



Course lectures
(including guest lectures)

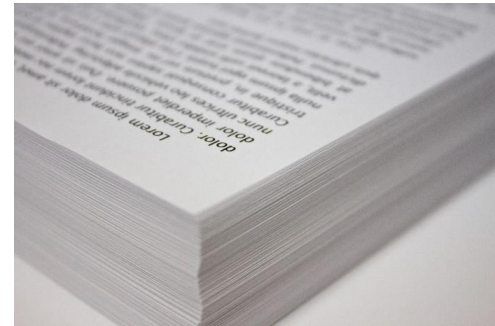
Course Structure

Tuesdays



Course lectures

Thursdays



Group discussion and student presentations

Course Recommendations and Requirements

- 1 Ready to read at least 9 papers this semester !**
 - 9 research papers as part of the weekly reading assignments
 - Asked to answer research questions about each paper
- 2 Already taken a machine learning course**
 - Strongly recommended for students to have taken an introduction machine learning course
 - 10-401, 10-601, 10-701, 11-663, 11-441, 11-641 or 11-741
- 3 Motivated to produce a high-quality course project**
 - Three course project assignments
 - Designed to enhance state-of-the-art algorithms

$$\begin{aligned}i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\h_t &= o_t \tanh(c_t)\end{aligned}$$

Course Project

- Pre-proposal (in 2 weeks)
 - Define your dataset, research task and teammates
- First project assignment (in 5 weeks)
 - Experiment with unimodal representations
 - Explore/discuss simple baseline model(s)
- Midterm project assignment (in 10 weeks)
 - Implement and evaluate state-of-the-art model(s)
 - Discuss new multimodal model(s)
- Final project assignment (in 14 weeks)
 - Implement and evaluate new multimodal model(s)
 - Discuss results and possible future directions

Course Project Guidelines

- Dataset should have at least two modalities:
 - Natural language and visual/images
- Teams of 3 or 4 students
- The project should explore algorithmic novelty
- Possible venues for your final report:
 - NAACL 2019, ACL 2019, IJCAI 2019, ICML 2019
- We will discuss on Thursday about project ideas
- GPU resources available:
 - Amazon AWS and Google Cloud Platform

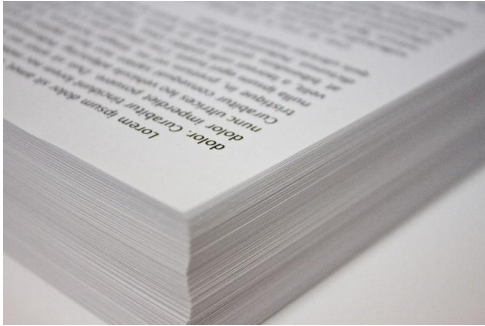
Examples of Previous Course Projects

- Select-Additive Learning: Improving Generalization in Multimodal Sentiment Analysis
 - <https://arxiv.org/abs/1609.05244>
- Preserving Intermediate Objectives: One Simple Trick to Improve Learning for Hierarchical Models
 - <https://arxiv.org/abs/1706.07867>
- Gated-Attention Architectures for Task-Oriented Language Grounding
 - <https://arxiv.org/abs/1706.07230>
- Efficient Low-rank Multimodal Fusion with Modality-Specific Factors
 - <https://arxiv.org/abs/1806.00064>
- Multimodal Sentiment Analysis with Word-Level Fusion and Reinforcement Learning
 - <https://arxiv.org/abs/1802.00924>

Process for Selecting your Course Project

- Thursday 8/30: Lecture describing available multimodal datasets and research topics
- Monday 9/3: Submit a short paragraph listing your top 3 choices
- Tuesday 9/4: During the later part of the lecture, we will have a discussion period to help with team formation
- Sunday 9/12: Pre-proposals are due. You should have selected your teammates, dataset and task

Course Grades



$$\begin{aligned}i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\c_t &= f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \\h_t &= o_t \tanh(c_t)\end{aligned}$$

- Discussion and lecture participation 20%
- Reading assignments 20%

- First project assignment
 - Report and presentation 15%
- Mid-term project assignment
 - Report and presentation 15%
- Final project assignment
 - Report and presentation 30%

Equal Contribution by All Teammates!

- Each team will be required to create a GitHub repository which will be accessible by TAs
- Each report should include a description of the task from each teammate
- Please let us know soon if you have concerns about the participation levels of your teammates

Lecture Schedule

Classes	Lectures	
Week 1 8/28 & 8/30	Course introduction <ul style="list-style-type: none">• Research and technical challenges• Multimodal applications and datasets	Thursday 8/30 in DH A302
Week 2 9/4 & 9/6	Basic mathematical concepts <ul style="list-style-type: none">• Language, image and audio representation• Loss functions and basic neural networks	Project preferences due on Monday night
Week 3 9/11 & 9/13	Convolutional neural networks and optimization <ul style="list-style-type: none">• Neural network optimization• Convolutional neural networks	Pre-proposal due on Sunday 9/12
Week 4 9/18 & 9/20	Recurrent neural networks <ul style="list-style-type: none">• Backpropagation Through Time• Gated networks and LSTM	

Lecture Schedule

Classes	Lectures
Week 5 9/25 & 9/27	Multimodal representation learning <ul style="list-style-type: none">• Multimodal auto-encoders• Multimodal joint representations
Week 6 10/3 & 10/5	<i>First project assignment - Presentation</i> <div style="border: 2px solid red; padding: 5px; display: inline-block; margin-left: 20px;">Thursday in DH A302 . Proposal due: 10/7.</div>
Week 7 10/9 & 10/11	Multivariate statistics and coordinated representations <ul style="list-style-type: none">• Deep canonical correlation analysis• Non-negative matrix factorization
Week 8 10/16 & 10/18	Multimodal alignment and attention models <ul style="list-style-type: none">• Explicit alignment and dynamic time warping• Implicit alignment and attention models

Lecture Schedule

Classes	Lectures
Week 9 10/23 – 10/25	Multimodal optimization <ul style="list-style-type: none">• Practical deep model optimization• Variational approaches
Week 10 10/30 & 11/1	Probabilistic graphical models <ul style="list-style-type: none">• Boltzmann distribution and CRFs• Continuous and fully-connected CRFs
Week 11 11/6 & 11/8	<i>Mid-term project assignment - Pre</i> <div data-bbox="1286 751 1879 932" style="border: 2px solid red; padding: 5px; display: inline-block;">Thursday in DH A302. Midterm due on 11/11.</div>
Week 12 11/13 & 11/15	Multimodal fusion and new directions <ul style="list-style-type: none">• Multi-kernel learning and fusion• New directions in multimodal machine learning

Lecture Schedule

Classes	Lectures
Week 13 11/20 & 11/22	<i>Thanksgiving week (+ Project preparation)</i>
Week 14 11/27 & 11/29	Multi-lingual representations and models <ul style="list-style-type: none">• Neural machine translation• Guest lecture: Graham Neubig
Week 15 TBD <i>* Final *</i>	<i>Final project assignment - Present</i> <div data-bbox="1286 689 1879 872" style="border: 2px solid red; padding: 5px; margin-top: 10px;">Thursday in DH A302. Final project due: 12/9.</div>

Piazza <https://piazza.com/cmu/fall2018/11777/home>

The screenshot shows a web browser window with the URL <https://piazza.com/cmu/fall2018/11777/home>. The page header includes the Piazza logo, navigation links for '11-777', 'Q & A', 'Resources', 'Statistics', and 'Manage Class', and a user profile for 'Louis-Philippe Morency'. The main content area displays 'Carnegie Mellon University - Fall 2018' and the course title '11-777: Multimodal Machine Learning'. Below the title are three icons: a download icon for 'Syllabus', an edit icon, and a trash icon.

Course Information Staff Resources

Description

Edit

Multimodal machine learning (MMML) is a vibrant multi-disciplinary research field which addresses some of the original goals of artificial intelligence by integrating and modeling multiple communicative modalities, including linguistic, acoustic and visual messages. With the initial research on audio-visual speech recognition and more recently with language & vision projects such as image and video captioning, this research field brings some unique challenges for multimodal researchers given the heterogeneity of the data and the contingency often found between modalities. This course will teach fundamental mathematical concepts related to MMML including multimodal alignment and fusion, heterogeneous representation learning and multi-stream temporal modeling. We will also review recent papers describing state-of-the-art probabilistic models and computational algorithms for MMML and discuss the current and upcoming challenges.

The course will present the fundamental concepts of machine learning and deep neural networks relevant to the five main challenges in multimodal machine learning: (1) multimodal representation learning, (2) translation & mapping, (3) modality alignment, (4) multimodal fusion and (5) co-learning. These include, but not limited to, multimodal auto-encoder, deep canonical correlation analysis, multi-kernel learning, attention models and multimodal recurrent neural networks. The course will also discuss many of the recent applications of MMML including multimodal affect recognition, image and video captioning and cross-modal multimedia retrieval.

General Information

Edit

Time

Tuesdays and Thursday, 4:30pm-5:50pm

Lecture location (Tuesdays and some Thursdays)

DH A302

Discussion location (Thursdays, with some exceptions)

WEH 5304 and WEH 5320

Announcements

Add

Add an Announcement
Click the Add button to add an announcement.

- Announcements
- Question/Answers
- Lecture slides
- Reading assignments
- Project resources
- Room assignments



Gradescope – Entry Code: M576ZG

gradescope <≡

11777
Advanced Multimodal Machine Learning

Dashboard
Assignments
Roster
Course Settings

INSTRUCTOR
Louis-Phillippe Morency

Account

11777 | Fall 2018

Entry Code: **M576ZG**

DESCRIPTION

Multimodal machine learning (MMML) is a vibrant multi-disciplinary research field which addresses some of the original goals of artificial intelligence by integrating and modeling multiple communicative modalities, including linguistic, acoustic and visual messages. With the initial research on audio-visual speech recognition and more recently with language & vision projects such as image and video captioning, this research field brings some unique challenges for multimodal researchers given the heterogeneity of the data and the contingency often found between modalities. This course will teach fundamental mathematical concepts related to MMML including multimodal alignment and fusion, heterogeneous representation learning and multi-stream temporal modeling. We will also review recent papers describing state-of-the-art probabilistic models and computational algorithms for MMML and discuss the current and upcoming challenges. The course will present the fundamental concepts of machine learning and deep neural networks relevant to the five main challenges in multimodal machine learning: (1) multimodal representation learning, (2) translation & mapping, (3) modality alignment, (4) multimodal fusion and (5) co-learning. These include, but not limited to, multimodal auto-encoder, deep canonical correlation analysis, multi-kernel learning, attention models and multimodal recurrent neural networks. The course will also discuss many of the recent applications of MMML including multimodal affect recognition, image and video captioning and cross-modal multimedia retrieval.

THINGS TO DO

- 1 Add students or staff to your course from the Roster page.
- 1 Create your first assignment from the Assignments page.

- Submit your weekly answers to reading questions
- Submit your project reports
- View the comments from your graded reports