



Language Technologies Institute



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



Jonathan Francis jmf1@andrew.cmu.edu



Paul Liang pliang@andrew.cmu.edu



Sai Krishna Rallabandi srallaba@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



Carnegie Mellon University

What is Multimodal?





Carnegie Mellon University

Multimodal Communicative Behaviors





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





Carnegie Mellon University

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."



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)



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





"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: <u>http://sspnet.eu/</u>

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



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

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



First Two Core Challenges





Core Challenge 1 – Translation





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



Carnegie Mellon University

Core Challenge 1: Translation

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





Core Challenge 1: Translation - Example

a person jogs a few steps

A person steps forward then turns around and steps forwards again.





Ahuja, C., & Morency, L. P. (2019). Language2Pose: Natural Language Grounded Pose Forecasting. *arXiv preprint arXiv:1907.01108*.



Core Challenge 2: Fusion







Core Challenge 2: Fusion

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



1) Early Fusion



2) Late Fusion





Core Challenge 2: 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



Two More Core Challenges





Carnegie Mellon University

Core Challenge 3: Representation





Core Challenge 3: Early Examples





Language Technologies Institute

Core Challenge 3: Early Examples

Multimodal Vector Space Arithmetic





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



34

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









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





Core Challenge 4: 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

Implicit Alignment

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





Core Challenge 4: Explicit Alignment



Applications:

- Re-aligning asynchronous data

- Finding similar data across modalities (we can estimate the aligned cost)

- Event reconstruction from multiple sources





Core Challenge 4: Explicit Alignment





Carnegie Mellon University

Core Challenge 4: Implicit Alignment



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





One Last Core Challenge





Language Technologies Institute



One Last Core Challenge





Language Technologies Institute

42

Core Challenge 5: Co-Learning

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





Core Challenge 5: Co-Learning



Pham et al., Found in Translation: Learning Robust Joint Representations by Cyclic Translations Between Modalities, https://arxiv.org/abs/1812.07809



Five Multimodal Core Challenges



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



Taxonomy of Multimodal Research

Representation

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

Translation

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

- Encoder-decoder
- Online prediction

Alignment

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

Fusion

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

- Model-based
 - o Kernel-based
 - o Graphical models
 - Neural networks

Co-learning

- Parallel data
 - Co-training
 - o 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





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

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















in in black shirt is playing guitar."

construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

boy is doing backflip on wakeboard.







(a) answer-phone

(a) get-out-car

(b) push-up

(b) cartwheel





(a) fight-person







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

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





Course Syllabus

Three Course Learning Paradigms



Reading assignments and course participation (30% of your grade)
$$\begin{split} i_t &= \sigma \left(W_{xi} x_t + W_{hi} h_{t-1} + W_{ci} c_{t-1} + b_i \right) \\ f_t &= \sigma \left(W_{xf} x_t + W_{hf} h_{t-1} + W_{cf} c_{t-1} + b_f \right) \\ c_t &= f_t c_{t-1} + i_t \tanh \left(W_{xc} x_t + W_{hc} h_{t-1} + b_c \right) \\ o_t &= \sigma \left(W_{xo} x_t + W_{ho} h_{t-1} + W_{co} c_t + b_o \right) \\ h_t &= o_t \tanh(c_t) \end{split}$$

Course project assignments (70% of your grade)



Course lectures (including guest lectures)



Language Technologies Institute

Course Recommendations and Requirements

Ready to read about 20 papers this semester !

- Research papers as part of the weekly reading assignments
- Asked to answer research questions about these papers
- 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
- Motivated to produce a high-quality course project
 - Three course project assignments
 - Designed to enhance state-of-the-art algorithms



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

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 2020, ACL 2020, IJCAI 2020, ICML 2020
- 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
 - <u>https://arxiv.org/abs/1706.07230</u>
- Efficient Low-rank Multimodal Fusion with Modality-Specific Factors
 - <u>https://arxiv.org/abs/1806.00064</u>





Examples of Previous Course Projects

- Multimodal Sentiment Analysis with Word-Level Fusion and Reinforcement Learning
 - https://arxiv.org/abs/1802.00924
- Found in Translation: Learning Robust Joint Representations by Cyclic Translations Between Modalities
 - https://arxiv.org/abs/1812.07809
- Self-Supervised Visual Representations for Cross-Modal Retrieval
 - <u>https://arxiv.org/abs/1902.00378</u>
- Attend and Attack: Attention Guided Adversarial Attacks on Visual Question Answering Models
 - https://nips2018vigil.github.io/static/papers/accepted/33.pdf



Process for Selecting your Course Project

- Thursday 8/29: Lecture describing available multimodal datasets and research topics
- Monday 9/2: Submit a short paragraph listing your top 3 choices
- Tuesday 9/3: During the later part of the lecture, we will have a discussion period to help with team formation
- Wednesday 9/11: Pre-proposals are due. You should have selected your teammates, dataset and task



Carnegie Mellon University

Course Grades



```
\begin{split} i_{t} &= \sigma \left( W_{xi} x_{t} + W_{hi} h_{t-1} + W_{ci} c_{t-1} + b_{i} \right) \\ f_{t} &= \sigma \left( W_{xf} x_{t} + W_{hf} h_{t-1} + W_{cf} c_{t-1} + b_{f} \right) \\ c_{t} &= f_{t} c_{t-1} + i_{t} \tanh \left( W_{xc} x_{t} + W_{hc} h_{t-1} + b_{c} \right) \\ o_{t} &= \sigma \left( W_{xo} x_{t} + W_{ho} h_{t-1} + W_{co} c_{t} + b_{o} \right) \\ h_{t} &= o_{t} \tanh(c_{t}) \end{split}
```

- Reading assignments 20%
- Lecture and course participation 10%
- Project preferences and proposal 5%
- First project assignment
 - Report and presentation 15%
- Mid-term project assignment
 - Report and presentation 20%
- 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	Tuesday Lectures	Thursday Lectures
Week 1 8/27 & 8/29	 Course introduction Research and technical challenges Course syllabus and requirements 	 Multimodal applications and datasets Research tasks and datasets Team projects
Week 2 9/3 & 9/5 ***	 Basic concepts: neural networks Language, visual and acoustic Loss functions and neural networks 	 Basic concep Gradient Project preferences due on Monday night
Week 3 9/10 & 9/12 *Pre-proposal*	 Convolutional neural networks Convolutional kernels and CNNs Residual networks 	 Recurrent ne Gated ne Backpror Pre-proposals due on Sunday 9/11
Week 4 9/17 & 9/19	 Multimodal representation learning Multimodal auto-encoders Multimodal joint representations 	 Coordinated representations Deep canonical correlation analysis Non-negative matrix factorization
Week 5 9/24 & 9/26	 Modular and factorized representations Module networks Factorized representations 	 Unsupervised representations Variational auto-encoder Generative adversarial approaches
Week 6 10/1 & 10/3	First project assignment - Presentations	First assignment due on Sunday 10/6



Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures
Week 7	Multimodal alignment	Alignment and representation
10/8 & 10/10	Explicit - dynamic time warping	Multi-head attention
	Implicit - attention models	Multimodal transformers
Week 8	Reinforcement learning	Multimodal RL
10/15 & 10/17	Markov decision process	Deep Q learning
***	 Q learning and policy gradients 	 Multimodal applications
Week 9	Probabilistic graphical models	Discriminative graphical models
10/22 & 10/24	Dynamic Bayesian networks	Boltzmann distribution and CRFs
***	Coupled and factor HMMs	Continuous and fully-connected CRFs
Week 10	Multimodal fusion and co-learning	New directions in Multimodal ML
10/29 & 10/31	 Multi-kernel learning and fusion 	 Overview of recent approaches in
	 Multimodal transfer learning 	multimodal machine learning
Week 11 11/5 & 11/7	Mid-term project assignment - Presentati	Midterm due on 11/10.



Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures	
Week 12 11/12 & 11/14	 Multi-lingual representations Neural machine translation Guest lecture: Graham Neubig 	KnowledMultGues	ge representation imodal knowledge discovery st lecture
Week 13 11/19 & 11/21	Thanksgiving week (no classes)		
Week 14 11/26 & 11/28	 Language, vision and action Neural machine translation Guest lecture 	 Multimodal affective computing Emotion and sentiment analysis Guest lecture 	
Week 15 12/3 & 12/5	Final project assignment - Presentations		Final due on 12/8.



Piazza https://piazza.com/cmu/fall2019/11777/home

C Secure https://piazza.com/cmu/fall2	018/11777/home		야 ☆ 🖾
ZQ 11-777 •	Q & A Resources	Statistics Manage Class	Louis-Phillippe Morency
Carnegie Mellon University - Fall 2018 11-777: Multimoda Syllabus 1	l Machin	e Learning	Ê
Course Information Staff Resources			
Description	-	Announcements	+ Add
Multimodal machine learning (MMML) is a vibrar field which addresses some of the original goals integrating and modeling multiple communicative acoustic and visual messages. With the initial re recognition and more recently with language & v and video captioning, this research field brings so multimodal researchers given the heterogeneity often found between modalities. This course will mathematical concepts related to MMML includin fusion, heterogeneous representation learning a probabilistic models and computational algorithm current and upcoming challenges. The course will present the fundamental concept deep neural networks relevant to the five main ci- machine learning: (1) multimodal representation mapping, (3) modailty alignment, (4) multimodal These include, but not limited to, multimodal affect captioning and cross-modal multimedia retrieval. General Information Tuesdays and Thursday, 4:30pm-5:50pm Lecture location (Tuesdays and some Thur DH A302	It multi-disciplinary res or artificial intelligence or artificial intelligence or artificial intelligence or artificial intelligence or artificial intelligence or artificial intelligence or artificial intelligence of the data and the cor teach fundamental g multimodal alignme and multi-stream tempo cribing state-of-the-art g multimodal alignme and multi-stream tempo cribing state-of-the-art is for MMML and discu- ts of machine learning (2) translatio fusion and (5) co-lear hallenges in multimodal learning, (2) translatio fusion and (5) co-lear - oncoder, deep canor on models and multim iscuss many of the rec the cognition, image art 	Add an Announcement Click the Add button to add speech click the Add button to add click the Add button to add cli	 Announcements Question/Answers Lecture slides Reading assignment Project resources





Gradescope – Entry Code: 9NGKXX

THINGS TO DO

C Secure https://www.gradescope.com/courses/22361

III gradescope <≡

11777

Advanced Multimodal Machine Learning

Dashboard

Assignments

Roster

Course Settings

INSTRUCTOR

Louis-Philippe Morency

11777 Fall 2018 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 crossmodal multimedia retrieval.

Add students or staff to your course from the Roster page

☆ 🔼

Entry Code: M576ZG

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

Account



Project Preferences – Due Monday 9/3

- Post your project preferences:
 - List of your ranked preferred projects
 - Use alphanumeric code of each dataset
 - Detailed dataset list in the "Lecture1.2-datasets" slides
 - Previous unimodal/multimodal experience
 - Available CPU / GPU resources
- For topics or datasets not in the list:
 - Include a description with links (for other students)

https://piazza.com/cmu/fall2019/11777/home



Carnegie Mellon University