# Multimodal Machine Learning

## Louis-Philippe (LP) Morency

**CMU Multimodal Communication and Machine Learning Laboratory** [MultiComp Lab]

# CMU Course 11-777: Multimodal Machine Learning

# Lecture Objectives

- ## What is Multimodal?

- ## Multimodal: Core technical challenges
    - ### Representation learning, translation, alignment, fusion and co-learning

- ## Multimodal representation learning
    - ### Joint and coordinated representations
    - ### Multimodal autoencoder and tensor representation
    - ### Deep canonical correlation analysis

- ## Fusion and temporal modeling
    - ### Multi-view LSTM and memory-based fusion
    - ### Fusion with multiple attentions

# What is Multimodal?

# What is Multimodal?



**Multimodal distribution**

➢ Multiple modes, i.e., distinct "peaks" (local maxima) in the probability density function

# What is Multimodal?



**Sensory Modalities**

# Multimodal Communicative Behaviors

## **V**erbal

**Lexicon**
> Words

**Syntax**
> Part-of-speech
> Dependencies

**Pragmatics**
> Discourse acts

## **V**ocal

**Prosody**
> Intonation
> Voice quality

**Vocal expressions**
> Laughter, moans

## **V**isual

**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

**Carnegie Mellon University**

# 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.

Language Technologies Institute

Carnegie Mellon University

# Multiple Communities and Modalities



Psychology



Medical



Speech



Vision



Language



Multimedia



Robotics



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

# 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 tutorial

1970        1980        1990        2000        2010

Language Technologies Institute          Carnegie Mellon University

# The McGurk Effect (1976)



[Hearing lips and seeing voices – Nature](#)

1970    1980    1990    2000    2010

Language Technologies Institute

Carnegie Mellon University

# The McGurk Effect (1976)



[Hearing lips and seeing voices – Nature](#)

# ➢ The "Computational" Era(Late 1980s until 2000)

## 1) Audio-Visual Speech Recognition (AVSR)

Language Technologies Institute

Carnegie Mellon University

# 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 away that exploits the complementarity and redundancy.

**Ⓐ Joint representations:**

# Joint Multimodal Representation



"Wow!"

"I like it!"

Joyful tone

Joint Representation
(Multimodal Space)

Tensed voice

Language Technologies Institute

Carnegie Mellon University

# 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 Representation**

Depth$_{Multimodal}$

Depth$_{Video}$

Depth$_{Verbal}$

**Visual**

**Verbal**

# Multimodal Vector Space Arithmetic



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

Language Technologies Institute

Carnegie Mellon University

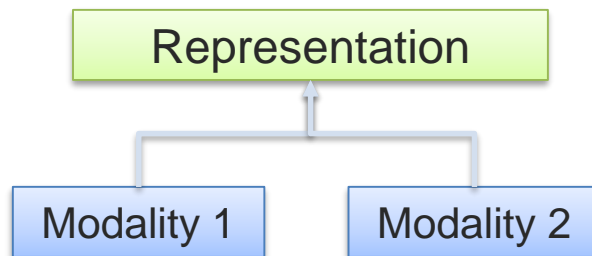# Core Challenge 1: Representation
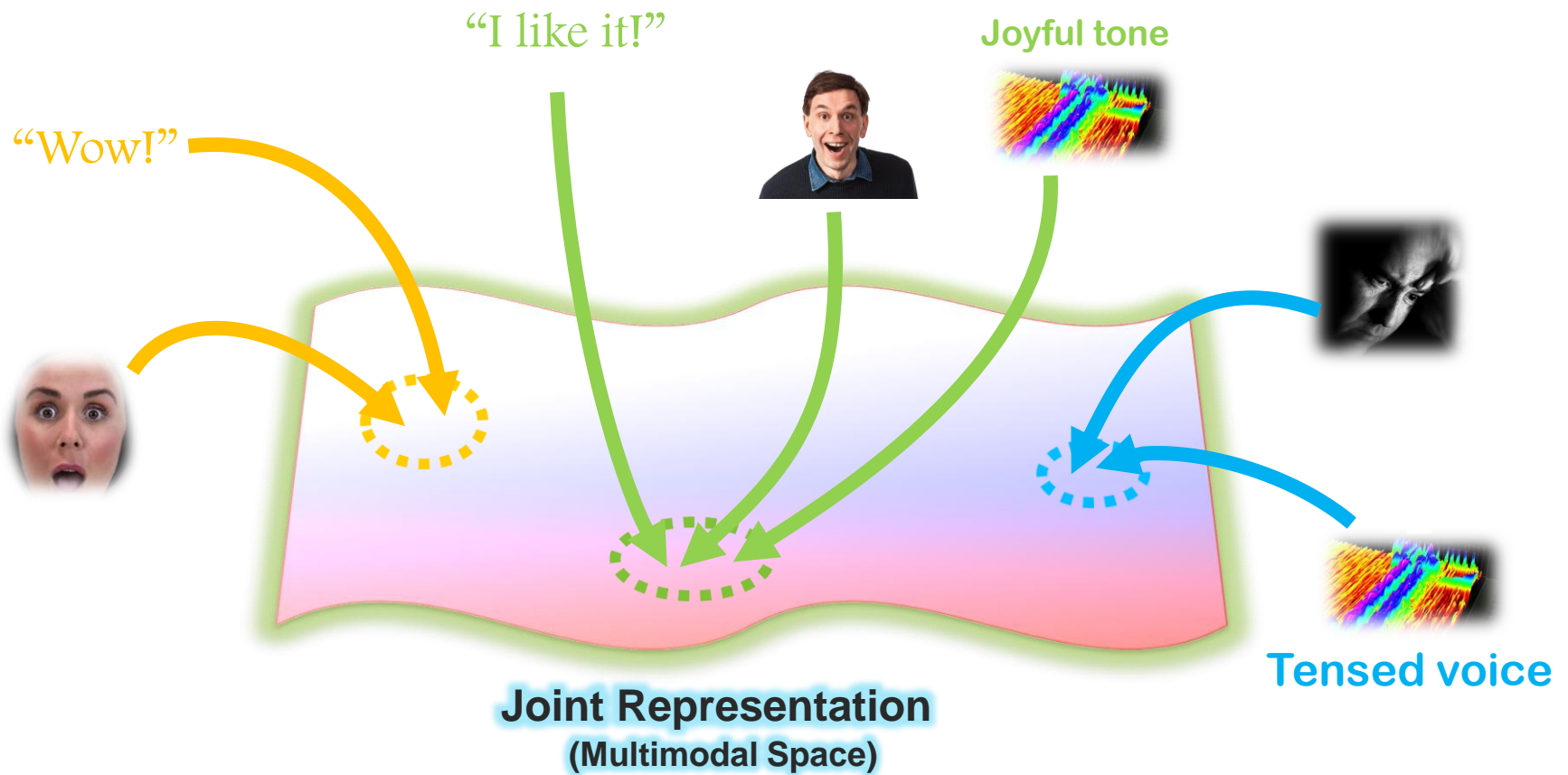
**Definition:** Learning how to represent and summarize multimodal data in away that exploits the complementarity and redundancy.

Ⓐ **Joint representations:**

| Representation |
|:---:|

| Modality 1 | Modality 2 |

Ⓑ **Coordinated representations:**

| Repres. 1 | ⬌ | Repres 2 |

| Modality 1 | | Modality 2 |

Language Technologies Institute

Carnegie Mellon University

# Coordinated Representation: Deep CCA

Learn linear projections that are maximally correlated:

$$(\boldsymbol{u}^*, \boldsymbol{v}^*) = \underset{\boldsymbol{u}, \boldsymbol{v}}{\text{argmax}} \; corr(\boldsymbol{u}^T \boldsymbol{X}, \boldsymbol{v}^T \boldsymbol{Y})$$



Andrew et al., ICML 2013

Language Technologies Institute

Carnegie Mellon University

# Core Challenge 2: Alignment

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

Modality 1          Modality 2



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

Language Technologies Institute

Carnegie Mellon University

# 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

Language Technologies Institute

Carnegie Mellon University

# Core Challenge 3: Fusion

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

**Ⓐ Model-Agnostic Approaches**

**1) Early Fusion**



Modality 1

Modality 2

...

Classifier

**2) Late Fusion**



Modality 1 → Classifier →

Modality 2 → Classifier →

# 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.

Language Technologies Institute

Carnegie Mellon University

# 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: T**ransfer knowledge between modalities, including their representations and predictive models.

Language Technologies Institute

Carnegie Mellon University

# Core Challenge 5: Co-Learning

# 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

Language Technologies Institute

Carnegie Mellon University

# Multimodal Applications

| 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

# Multimodal Representations

# Core Challenge: Representation

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

Ⓐ **Joint representations:**

Ⓑ **Coordinated representations:**

# Deep Multimodal autoencoders

- A deep representation learning approach
- A bimodal auto-encoder
    - Used for Audio-visual speech recognition



[Ngiam et al., Multimodal Deep Learning, 2011]

# Deep Multimodal autoencoders - training

- Individual modalities can be pretrained
  - RBMs
  - Denoising Autoencoders
- To train the model to reconstruct the other modality
  - Use both
  - Remove audio



[Ngiam et al., Multimodal Deep Learning, 2011]

# Deep Multimodal autoencoders - training

- Individual modalities can be pretrained
  - RBMs
  - Denoising Autoencoders
- To train the model to reconstruct the other modality
  - Use both
  - Remove audio
  - Remove video



[Ngiam et al., Multimodal Deep Learning, 2011]

# Multimodal Encoder-Decoder

- Visual modality often encoded using CNN

- Language modality will be decoded using LSTM

  - A simple multilayer perceptron will be used to translate from visual (CNN) to language (LSTM)



Text
$X$

Image
$Y$

# **Multimodal Joint Representation**

- For supervised learning tasks
- Joining the unimodal representations:
  - Simple concatenation
  - Element-wise multiplication or summation
  - Multilayer perceptron
- How to explicitly model both unimodal and bimodal interactions?

e.g. Sentiment

# Multimodal Sentiment Analysis

**MOSI dataset (Zadeh et al, 2016)**



- 2199 subjective video segments
- Sentiment intensity annotations
- 3 modalities: text, video, audio

**Multimodal joint representation:**

$$h_m = f(W \cdot [h_x, h_y, h_z])$$

Sentiment Intensity [-3,+3]

softmax

$h_m$

$h_x$   $h_y$   $h_z$

Text
*X*

Image
*Y*

Audio
*Z*

# Unimodal, Bimodal and Trimodal Interactions



**Speaker's behaviors**

**Sentiment Intensity**

**Unimodal**
- "This movie is sick" → ? → *Ambiguous !*
- "This movie is fair" → ✚
- Smile → ✚ → *Unimodal cues*
- Loud voice → ? → *Ambiguous !*

**Bimodal**
- "This movie is sick" | Smile → ✚✚ → *Resolves ambiguity (bimodal interaction)*
- "This movie is sick" | Frown → ━ ━
- "This movie is sick" | Loud voice → ? → *Still Ambiguous !*

**Trimodal**
- "This movie is sick" | Smile | Loud voice → ✚✚✚ → *Different trimodal interactions !*
- "This movie is fair" | Smile | Loud voice → ✚

# Multimodal Tensor Fusion Network (TFN)

Models both unimodal and bimodal interactions:

$$h_m = \begin{bmatrix} h_x \\ 1 \end{bmatrix} \otimes \begin{bmatrix} h_y \\ 1 \end{bmatrix} = \begin{bmatrix} h_x & h_x \otimes h_y \\ 1 & h_y \end{bmatrix}$$

*Important !*

[Zadeh, Jones and Morency, **EMNLP 2017**]

e.g. Sentiment

softmax

**Bimodal**

**Unimodal**

$h_m$

$h_x$

$h_y$

Text
$X$

Image
$Y$

Language Technologies Institute

Carnegie Mellon University

# Multimodal Tensor Fusion Network (TFN)

Can be extended to three modalities:

$$h_m = \begin{bmatrix} h_x \\ 1 \end{bmatrix} \otimes \begin{bmatrix} h_y \\ 1 \end{bmatrix} \otimes \begin{bmatrix} h_z \\ 1 \end{bmatrix}$$

**Explicitly models unimodal, bimodal and trimodal interactions !**

[Zadeh, Jones and Morency, **EMNLP 2017**]

Language Technologies Institute

**Carnegie Mellon University**

# Experimental Results – MOSI Dataset

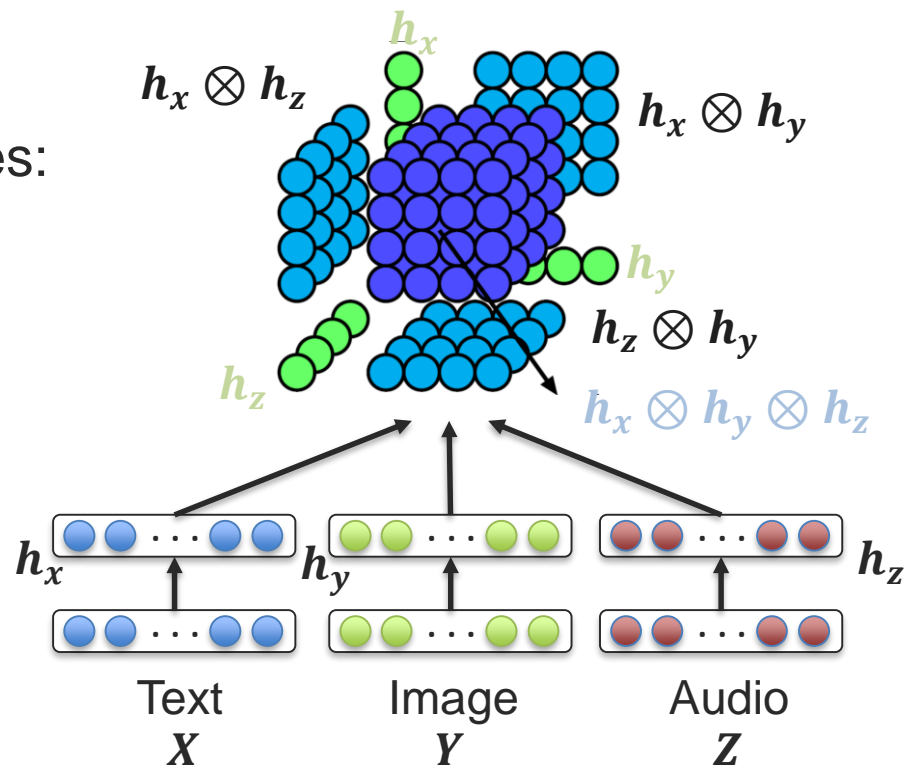| Multimodal Baseline | Binary | | 5-class | Regression | |
|---|---|---|---|---|---|
| | Acc(%) | F1 | Acc(%) | MAE | $r$ |
| Random | 50.2 | 48.7 | 23.9 | 1.88 | - |
| C-MKL | 73.1 | 75.2 | 35.3 | - | - |
| SAL-CNN | 73.0 | - | - | - | - |
| SVM-MD | 71.6 | 72.3 | 32.0 | 1.10 | 0.53 |
| RF | 71.4 | 72.1 | 31.9 | 1.11 | 0.51 |
| TFN | **77.1** | **77.9** | **42.0** | **0.87** | **0.70** |
| Human | 85.7 | 87.5 | 53.9 | 0.71 | 0.82 |
| $\Delta^{SOTA}$ | ↑ 4.0 | ↑ 2.7 | ↑ 6.7 | ↓ 0.23 | ↑ 0.17 |

**Improvement over State-Of-The-Art**

| Baseline | Binary | | 5-class | Regression | |
|---|---|---|---|---|---|
| | Acc(%) | F1 | Acc(%) | MAE | $r$ |
| $TFN_{language}$ | 74.8 | 75.6 | 38.5 | 0.99 | 0.61 |
| $TFN_{visual}$ | 66.8 | 70.4 | 30.4 | 1.13 | 0.48 |
| $TFN_{acoustic}$ | 65.1 | 67.3 | 27.5 | 1.23 | 0.36 |
| $TFN_{bimodal}$ | 75.2 | 76.0 | 39.6 | 0.92 | 0.65 |
| $TFN_{trimodal}$ | 74.5 | 75.0 | 38.9 | 0.93 | 0.65 |
| $TFN_{notrimodal}$ | 75.3 | 76.2 | 39.7 | 0.919 | 0.66 |
| TFN | **77.1** | **77.9** | **42.0** | **0.87** | **0.70** |
| $TFN_{early}$ | 75.2 | 76.2 | 39.0 | 0.96 | 0.63 |

Language Technologies Institute

Carnegie Mellon University

# Multimodal VAE (MVAE)

- Introduce a multimodal variational autoencoder (MVAE) with a new training paradigm that learns a joint distribution and is robust to missing data



(a)       (b)       (c)

[Wu, Mike, and Noah Goodman. "Multimodal Generative Models for Scalable Weakly-Supervised Learning.", NIPS 2018]

Language Technologies Institute

Carnegie Mellon University

# Multimodal VAE (MVAE)

- Transform unimodal datasets into "multi-modal" problems by treating labels as a second modality

$$z \sim p(z) \qquad z \sim p(z|x_2 = 5) \qquad z \sim p(z) \qquad z \sim p(z|x_2 = ankle\ boot)$$



[Wu, Mike, and Noah Goodman. "Multimodal Generative Models for Scalable Weakly-Supervised Learning.", NIPS 2018]

Language Technologies Institute

Carnegie Mellon University

# Coordinated Multimodal Representations

# Coordinated Multimodal Representations

Learn (unsupervised) two or more coordinated representations from multiple modalities. A loss function is defined to bring closer these multiple representations.



Similarity metric (e.g., cosine distance)

Text $X$

Image $Y$

Language Technologies Institute

Carnegie Mellon University

# Coordinated Multimodal Embeddings



Distance(s,t)

Image features s

Text: *a parrot rides a tricycle*

[Huang et al., Learning Deep Structured Semantic Models for Web Search using Clickthrough Data, 2013]

Carnegie Mellon University

# Canonical Correlation Analysis

*"canonical": reduced to the simplest or clearest schema possible*

(1) Learn two linear projections, one for each view, that are maximally correlated:

$$(\boldsymbol{u}^*, \boldsymbol{v}^*) = \underset{\boldsymbol{u}, \boldsymbol{v}}{\operatorname{argmax}}\, corr\left(\boldsymbol{H_x}, \boldsymbol{H_y}\right)$$

$$= \underset{\boldsymbol{u}, \boldsymbol{v}}{\operatorname{argmax}}\, corr\left(\boldsymbol{u^T X}, \boldsymbol{v^T Y}\right)$$

projection of Y

projection of X

$\boldsymbol{H_x}$   $\boldsymbol{H_y}$

$U$   $V$

Text   Image
$X$   $Y$

# Correlated Projection

**①** Learn two linear projections, one for each view, that are maximally correlated:

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



Two views $X, Y$ where same instances have the same color

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

We want to learn multiple projection pairs $\left(\boldsymbol{u}_{(i)}\boldsymbol{X}, \boldsymbol{v}_{(i)}\boldsymbol{Y}\right)$:

$$\left(\boldsymbol{u}^*_{(i)}, \boldsymbol{v}^*_{(i)}\right) = \underset{\boldsymbol{u}_{(i)}, \boldsymbol{v}_{(i)}}{\mathrm{argmax}} \; corr\left(\boldsymbol{u}^T_{(i)}\boldsymbol{X}, \boldsymbol{v}^T_{(i)}\boldsymbol{Y}\right) \quad \approx \boldsymbol{u}^T_{(i)}\boldsymbol{\Sigma}_{XY}\boldsymbol{v}_{(i)}$$

② We want these multiple projection pairs to be orthogonal ("canonical") to each other:

$$\boldsymbol{u}^T_{(i)}\boldsymbol{\Sigma}_{XY}\boldsymbol{v}_{(j)} = \boldsymbol{u}^T_{(j)}\boldsymbol{\Sigma}_{XY}\boldsymbol{v}_{(i)} = \boldsymbol{0} \qquad \text{for } i \neq j$$

$$\boldsymbol{U}\boldsymbol{\Sigma}_{XY}\boldsymbol{V} = tr(\boldsymbol{U}\boldsymbol{\Sigma}_{XY}\boldsymbol{V}) \qquad \text{where } \boldsymbol{U} = [\boldsymbol{u}_{(1)}, \boldsymbol{u}_{(2)}, \dots, \boldsymbol{u}_{(k)}]$$

$$\text{and } \boldsymbol{V} = [\boldsymbol{v}_{(1)}, \boldsymbol{v}_{(2)}, \dots, \boldsymbol{v}_{(k)}]$$

Language Technologies Institute

Carnegie Mellon University

# Canonical Correlation Analysis

**3** Since this objective function is invariant to scaling, we can constraint the projections to have unit variance:

$$U^T \Sigma_{XX} U = I \qquad V^T \Sigma_{YY} V = I$$

**Canonical Correlation Analysis:**

maximize: $tr(U^T \Sigma_{XY} V)$

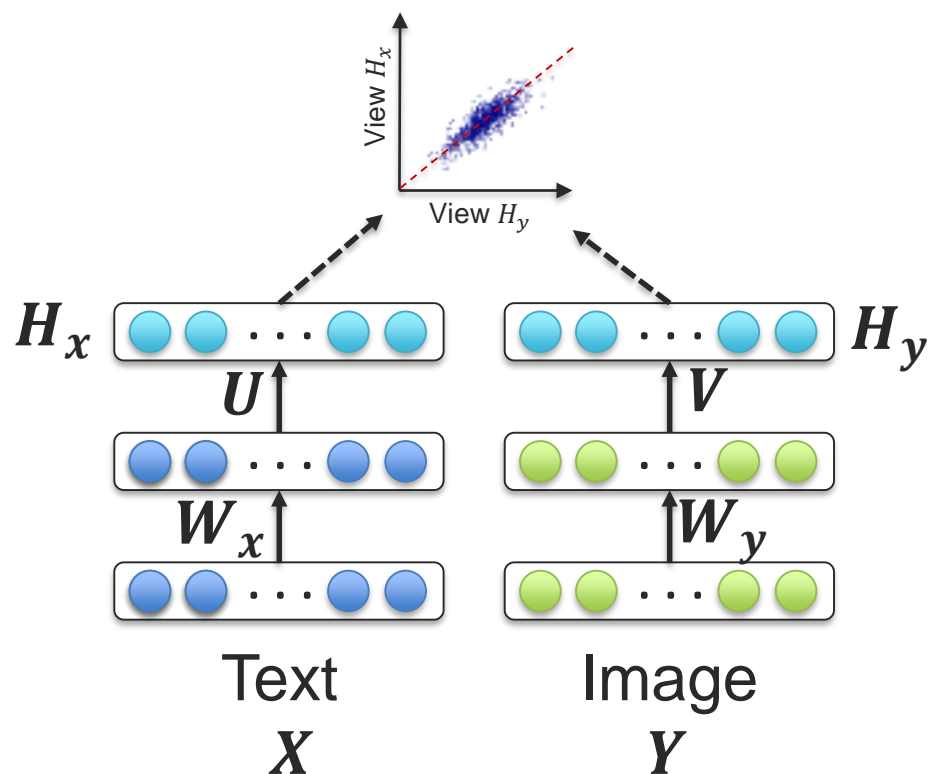subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$

Language Technologies Institute

Carnegie Mellon University

# Deep Canonical Correlation Analysis

Same objective function as CCA:

$$\underset{V,U,W_x,W_y}{\operatorname{argmax}} \ corr(H_x, H_y)$$

① Linear projections maximizing correlation

② Orthogonal projections
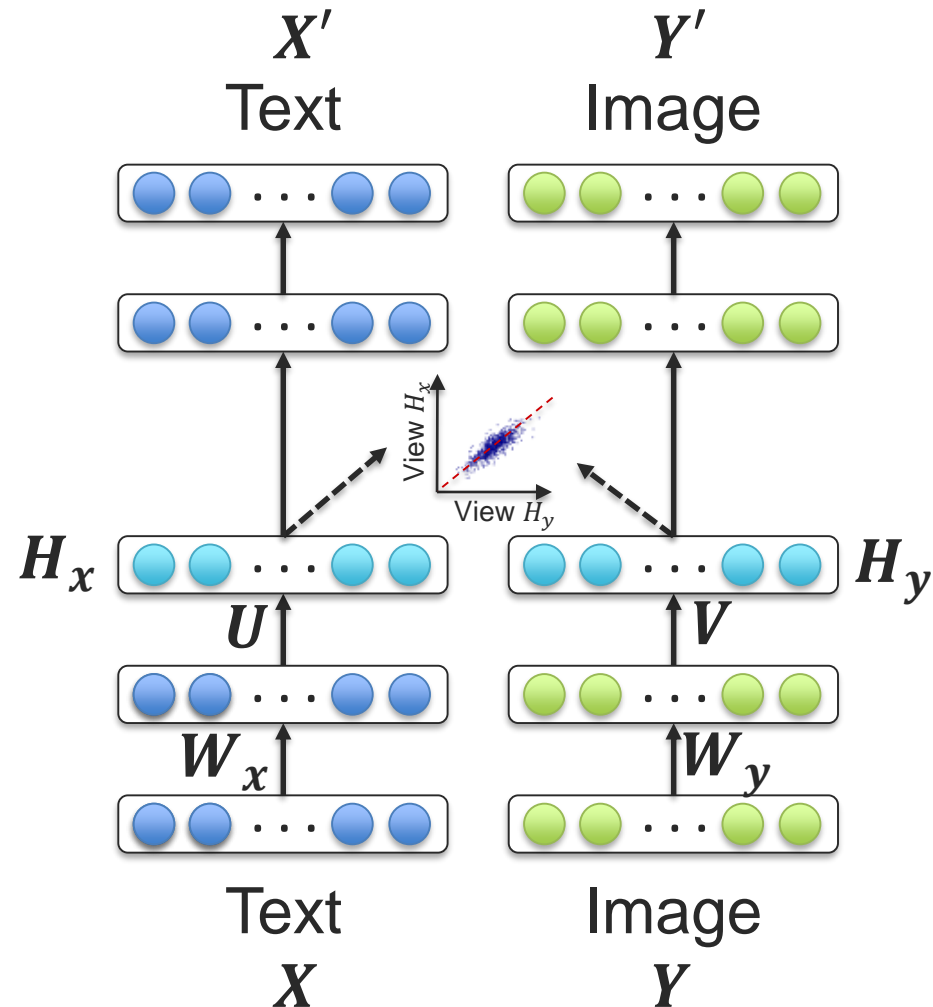
③ Unit variance of the projection vectors

Andrew et al., ICML 2013

Language Technologies Institute

Carnegie Mellon University

# Deep Canonically Correlated Autoencoders (DCCAE)

Jointly optimize for DCCA and autoencoders loss functions

➤ A trade-off between multi-view correlation and reconstruction error from individual views
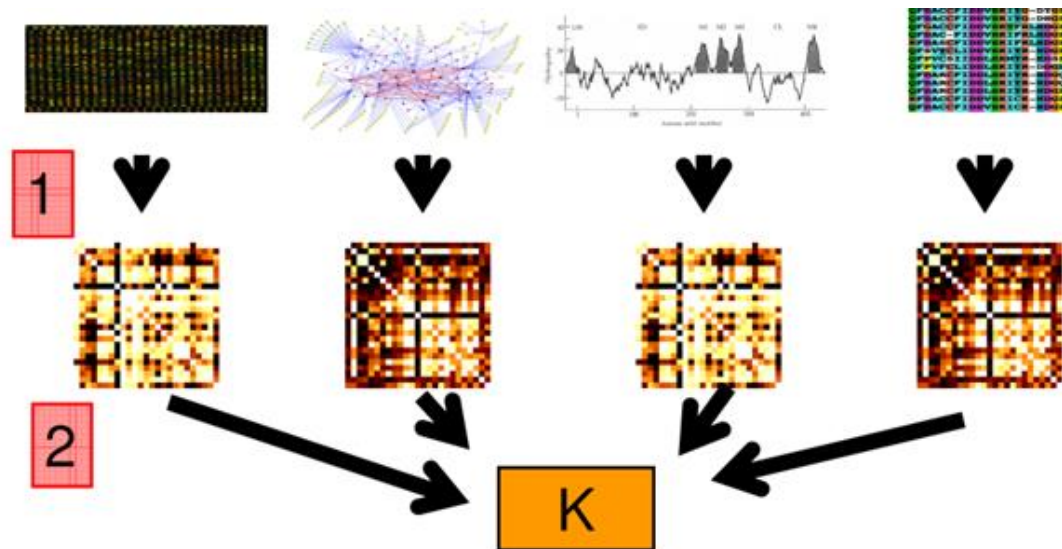


Wang et al., ICML 2015

# Multimodal Fusion

# Multiple Kernel Learning

- Pick a family of kernels for each modality and learn which kernels are important for the classification case
- Generalizes the idea of Support Vector Machines
- Works as well for unimodal and multimodal data, very little adaptation is needed



[Lanckriet 2004]
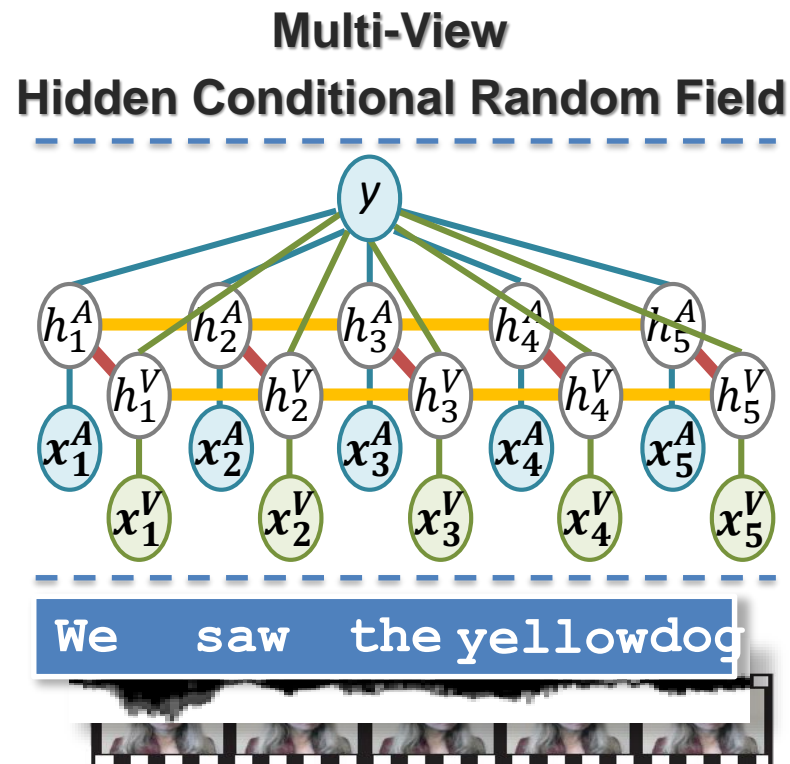
# Multimodal Fusion for Sequential Data

Modality-*private* structure

- Internal grouping of observations

Modality-*shared* structure

- Interaction and synchrony

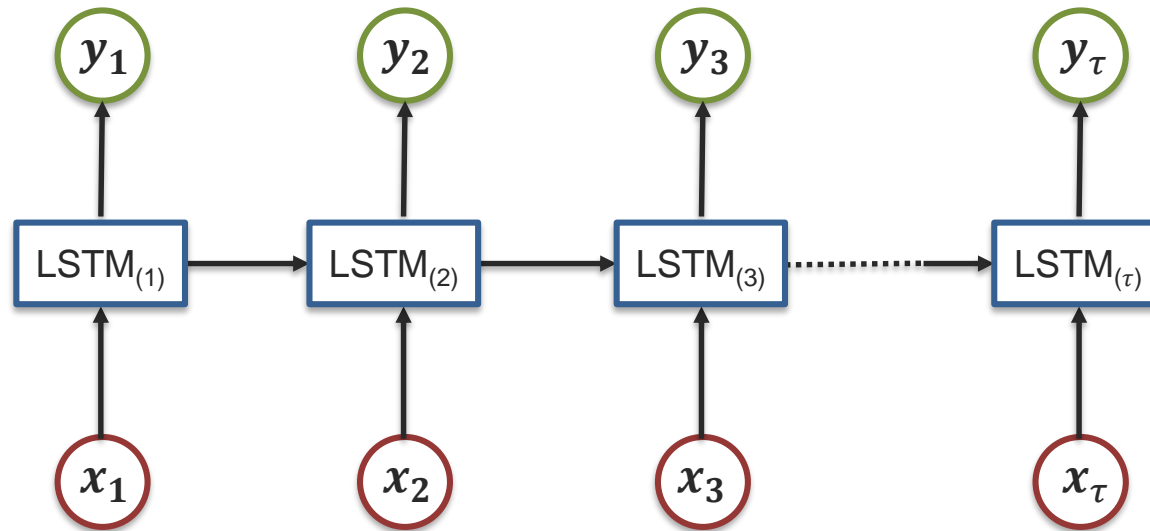$$p(y \mid x^A, x^V; \theta) = \sum_{h^A, h^V} p(y, h^A, h^V \mid x^A, x^V; \theta)$$

**Multi-View
Hidden Conditional Random Field**



We saw the yellowdog

➢ Approximate inference using loopy-belief

[Song, Morency and Davis, CVPR 2012]

Language Technologies Institute

Carnegie Mellon University

# Sequence Modeling with LSTM

Language Technologies Institute

Carnegie Mellon University

# Multimodal Sequence Modeling – Early Fusion

Language Technologies Institute

Carnegie Mellon University

# Multi-View Long Short-Term Memory (MV-LSTM)



[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, **ECCV**, 2016]

Language Technologies Institute

Carnegie Mellon University

# Multi-View Long Short-Term Memory
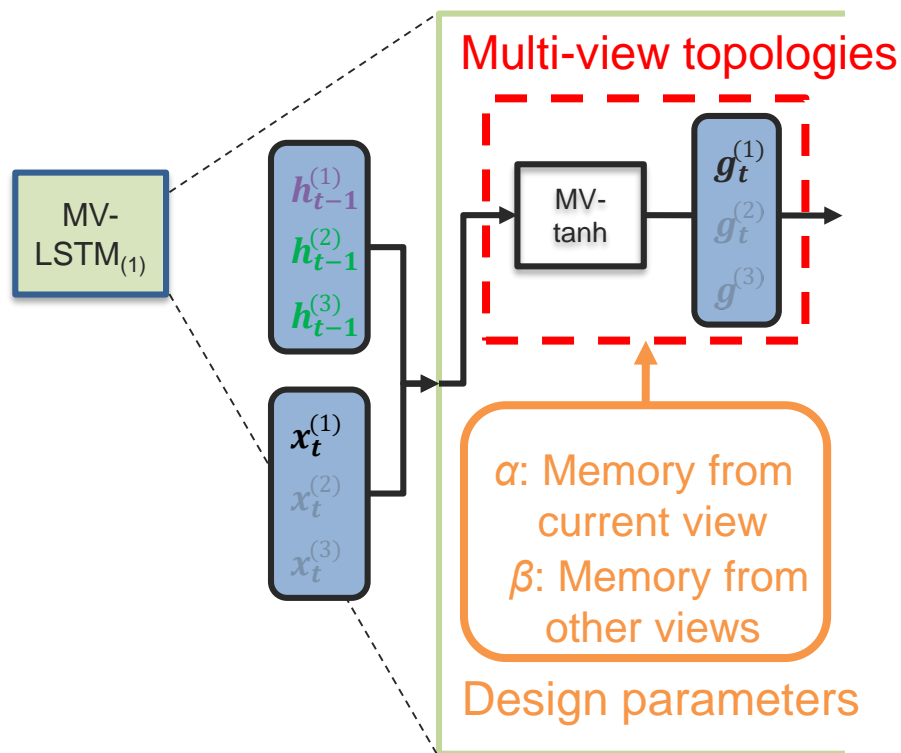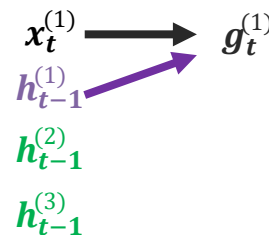


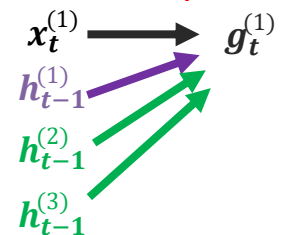[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, **ECCV**, 2016]
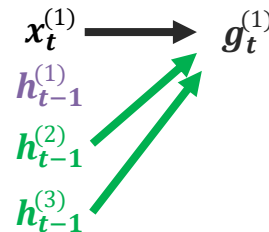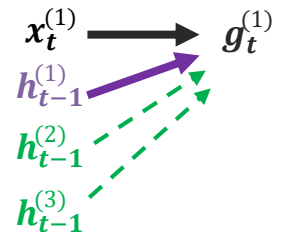
# Topologies for Multi-View LSTM



Multi-view topologies

$h_{t-1}^{(1)}$ $h_{t-1}^{(2)}$ $h_{t-1}^{(3)}$

$x_t^{(1)}$ $x_t^{(2)}$ $x_t^{(3)}$

MV-LSTM$_{(1)}$

MV-tanh

$g_t^{(1)}$ $g_t^{(2)}$ $g^{(3)}$

α: Memory from current view
β: Memory from other views

Design parameters

**View-specific**
*α=1, β=0*

$x_t^{(1)} \rightarrow g_t^{(1)}$
$h_{t-1}^{(1)}$
$h_{t-1}^{(2)}$
$h_{t-1}^{(3)}$

**Fully-connected**
*α=1, β=1*

$x_t^{(1)} \rightarrow g_t^{(1)}$
$h_{t-1}^{(1)}$
$h_{t-1}^{(2)}$
$h_{t-1}^{(3)}$

**Coupled**
*α=0, β=1*

$x_t^{(1)} \rightarrow g_t^{(1)}$
$h_{t-1}^{(1)}$
$h_{t-1}^{(2)}$
$h_{t-1}^{(3)}$

**Hybrid**
*α=2/3, β=1/3*

$x_t^{(1)} \rightarrow g_t^{(1)}$
$h_{t-1}^{(1)}$
$h_{t-1}^{(2)}$
$h_{t-1}^{(3)}$

[Shyam, Morency, et al. Extending Long Short-Term Memory for Multi-View Structured Learning, **ECCV**, 2016]

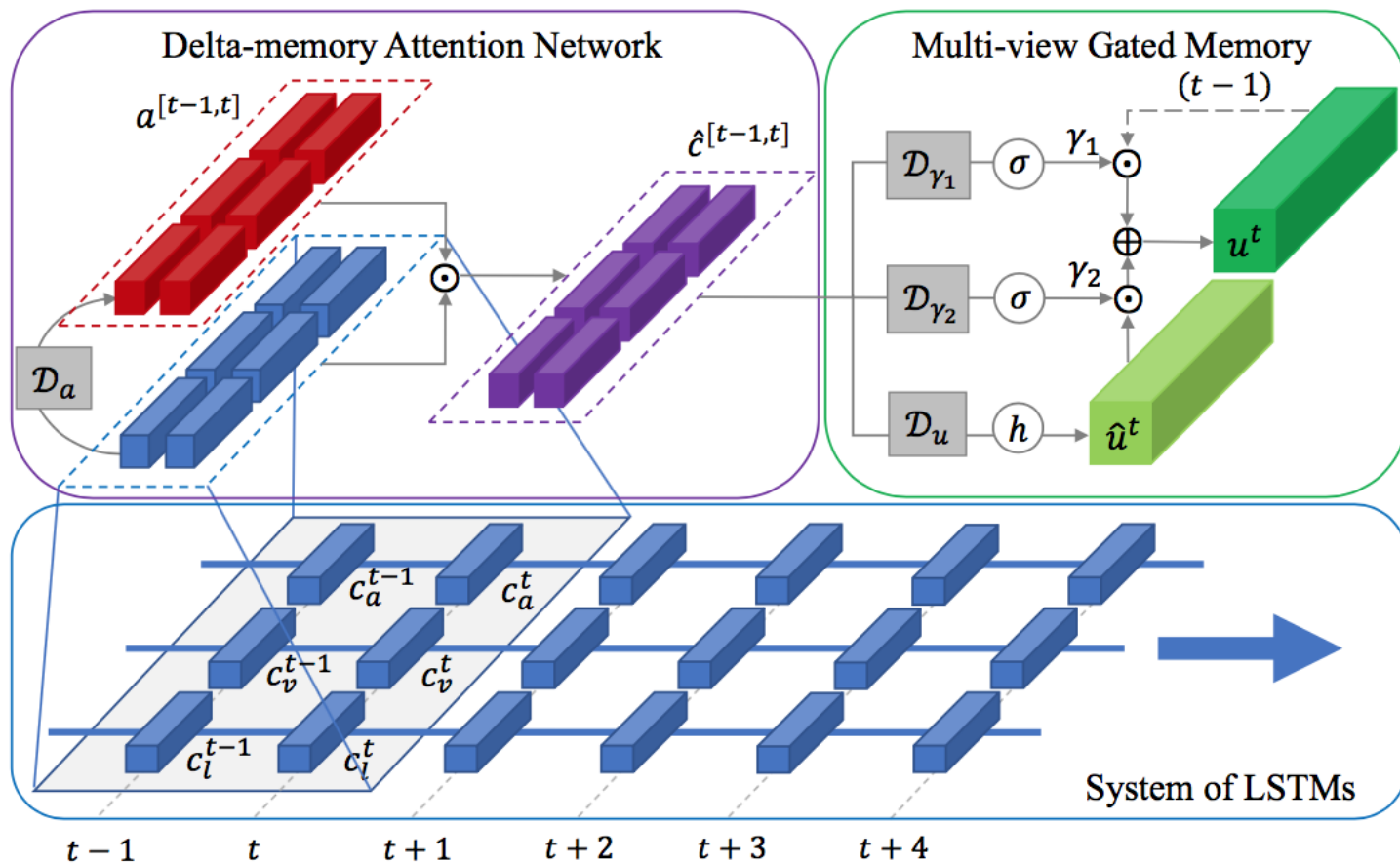Language Technologies Institute

Carnegie Mellon University

# Memory Based

- A memory accumulates multimodal information over time.

- From the representations throughout a source network.

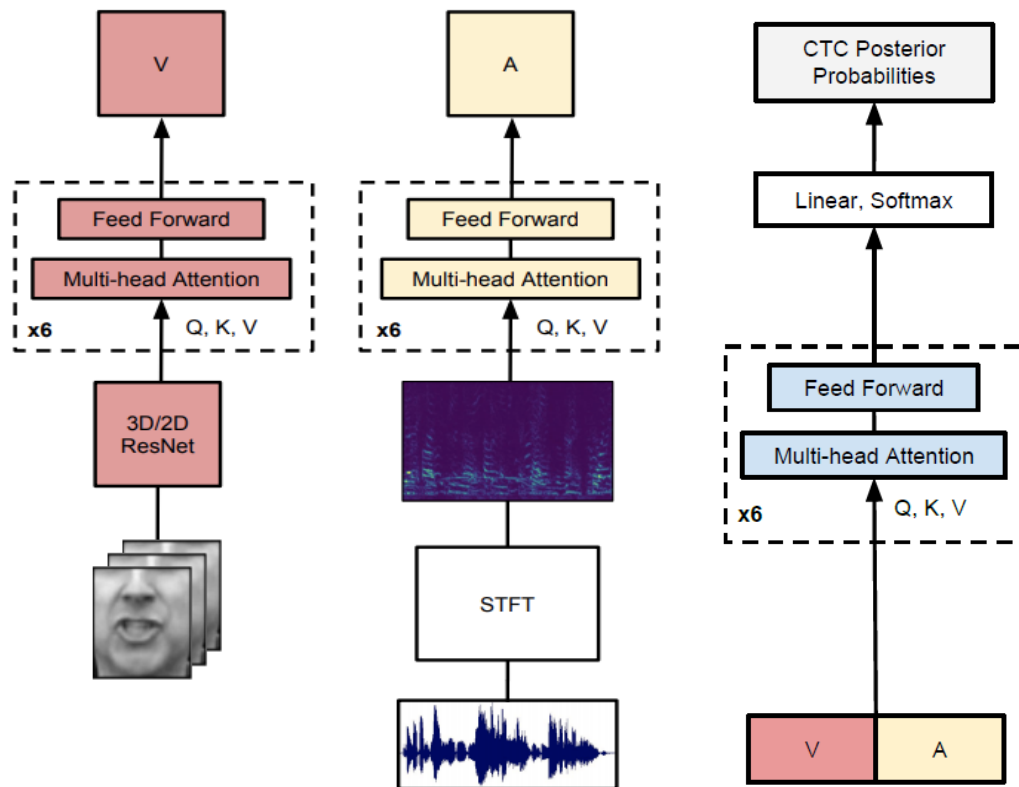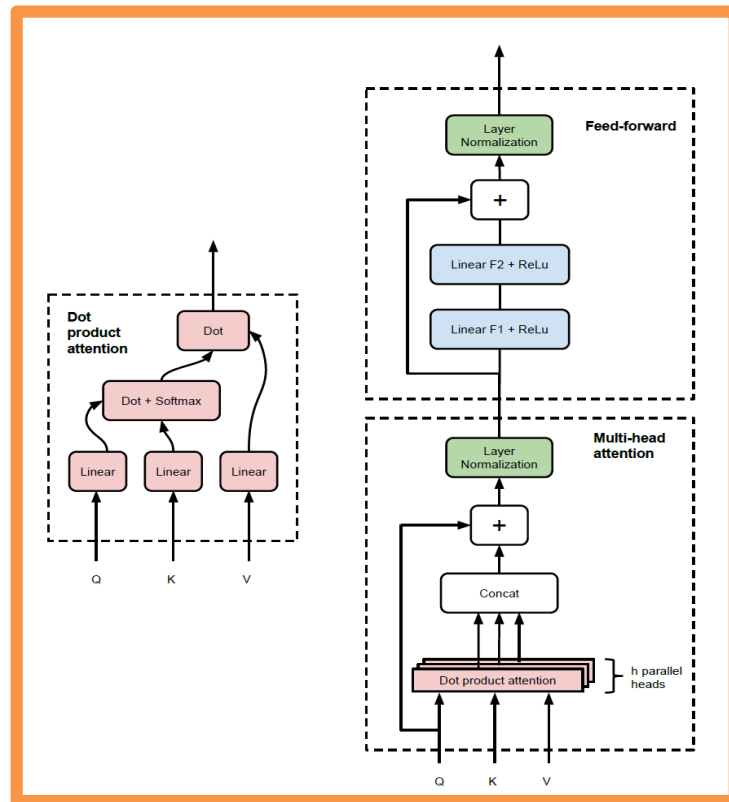- No need to modify the structure of the source network, only attached the memory.

Language Technologies Institute

Carnegie Mellon University

# Memory Based



[Zadeh et al., Memory Fusion Network for Multi-view Sequential Learning, AAAI 2018]
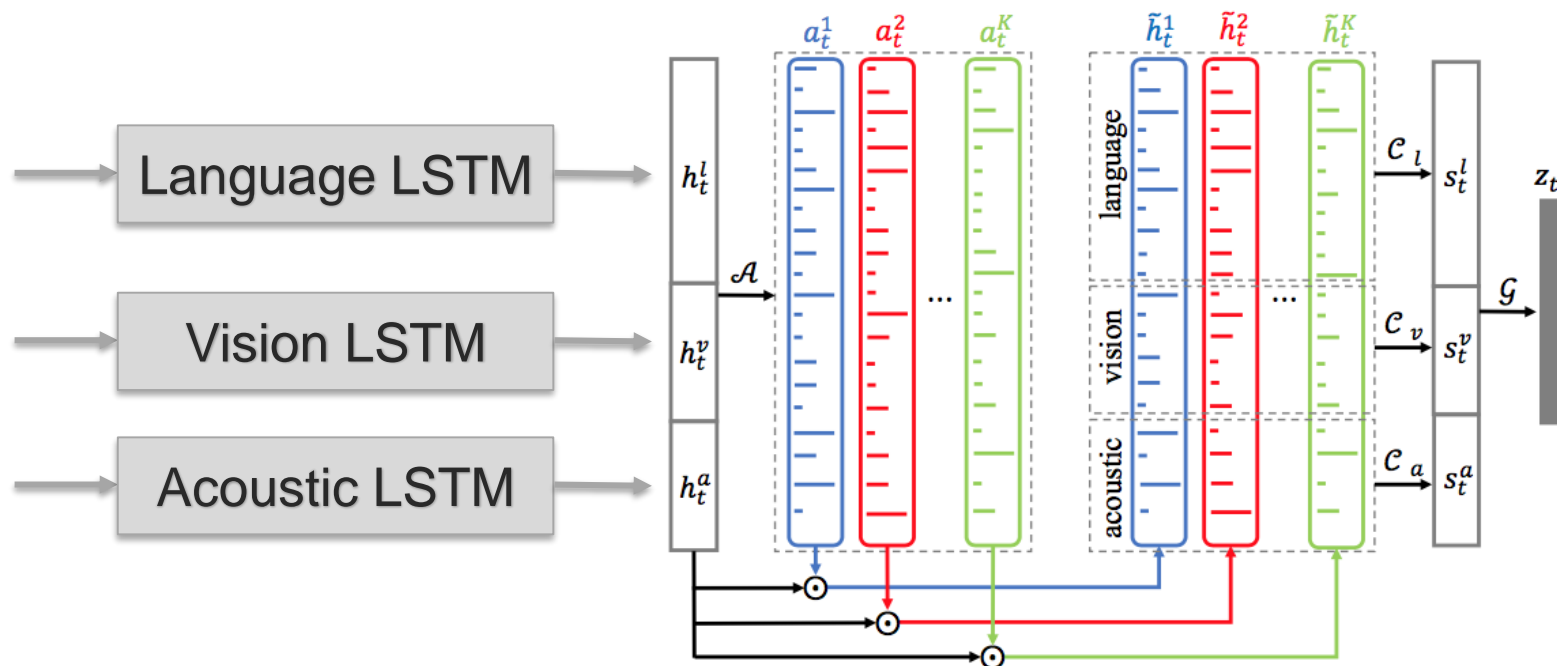
# Multi-Head Attention for AVSR

Afouras, Triantafyllos, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. "Deep audio-visual speech recognition." *arXiv preprint arXiv:1809.02108* (Sept 2018).

# Fusion with Multiple Attentions

- Modeling Human Communication – Sentiment, Emotions, Speaker Traits



[Zadeh et al., Human Communication Decoder Network for Human Communication Comprehension, AAAI 2018]

# Multimodal Machine Learning

**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**

Language Technologies Institute

**Carnegie Mellon University**