



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal AI: Understanding Human Behaviors

Louis-Philippe (LP) Morency



MultiComp Lab

PhD students:

Chaitanya Ahuja, Volkan Cirik, Paul Liang,
Hubert Tsai, Alexandria Vail, Torsten Wörtwein
and Amir Zadeh

Postdoctoral researcher:

Jeffrey Girard

Lab coordinator:

Nicole Siverling

Project assistant:

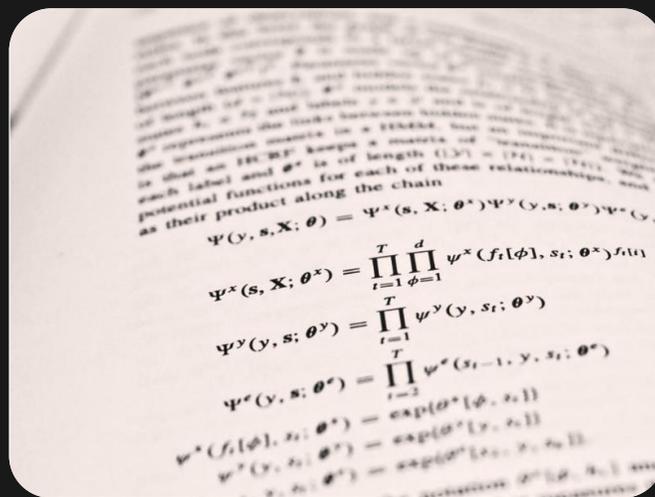
John Friday



MultiComp Lab



Communication Dynamics



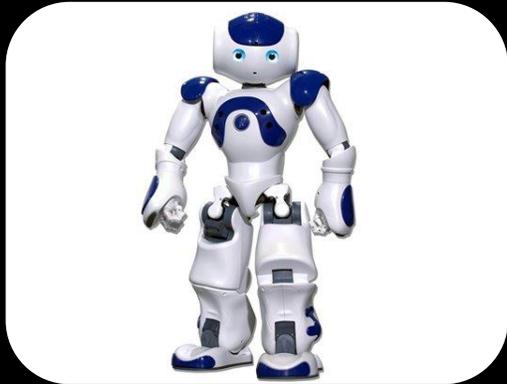
Multimodal AI



Mental Health

Multimodal AI Technologies

Robots



Virtual Humans



Ubiquitous



Mobile



Online



Wearable

Multimodal AI Technologies

Robots



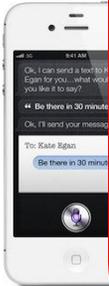
Virtual Humans

Ubiquitous

Video Conferencing



M



Multimodal Communicative Behaviors



Verbal

- **Lexicon**
 - Words
- **Syntax**
 - Part-of-speech
 - Dependencies
- **Pragmatics**
 - Discourse acts

Vocal

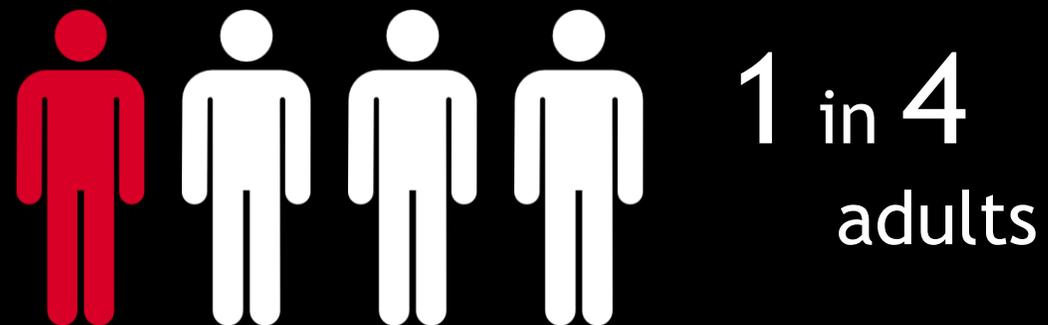
- **Prosody**
 - Intonation
 - Voice quality
- **Vocal expressions**
 - Laughter, moans

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



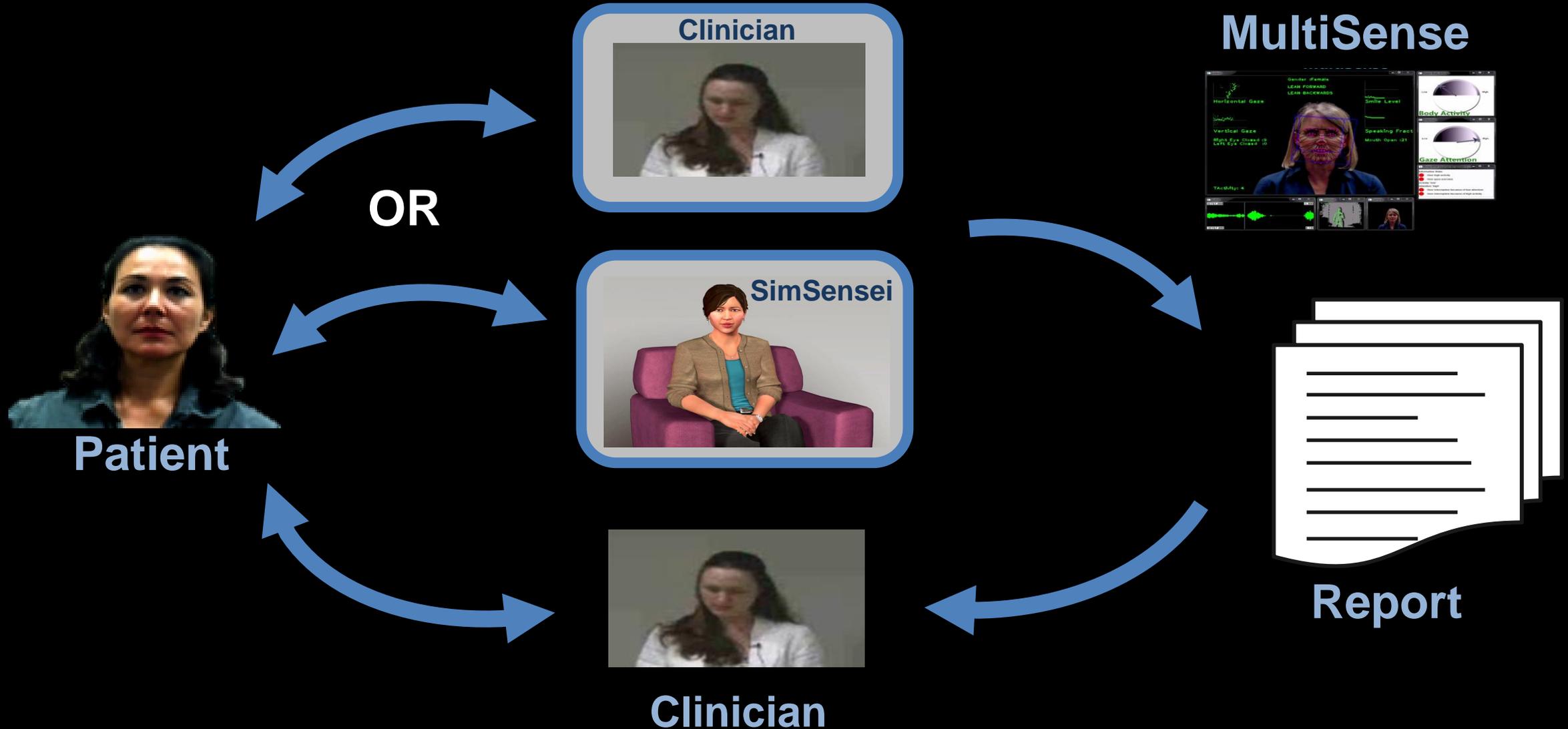
Mental Health by the Numbers



#1 cause
of disability



Multimodal AI for Mental Health Assessment



Behavioral Indicators of Psychological Distress



Distress Assessment
Interview Corpus



Data



Dictionary of Multimodal Behavior Markers



[IVA 2011, FG 2013, ICASSP 2013, 2015, ACII 2013, 2017, Interspeech 2013, 2017, ICMI 2013, 2014, 2018]

Psychosis

- Speech disfluency
- Language structure
- Gaze patterns
- ...

Depression

- Smile dynamics
- Gaze aversion
- Vowel space
- ...

Suicidal Ideation

- Lexical markers
- Voice quality
- Prosodic cues
- ...



UPMC



Scientific Discoveries

① Smile Dynamics - Behavior Indicators



Depressed vs Non-depressed

Number of smiles



Surprising!

Smile duration



Smile intensity



Scientific Discoveries

② Negative Expressions - Behavior Indicators



PTSD vs Non-PTSD

Overall population



Men only



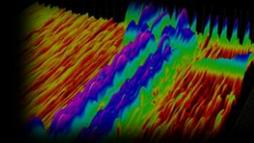
Women only



Opposite!

Scientific Discoveries

③ Speech Patterns - Behavior Indicators



Suicidal vs Non-suicidal

First person pronouns
(e.g., me, my, mine, I)



Repeater vs Non-repeater

Voice tenseness



Important!

Multimodal AI: Automatic Distress Level Prediction



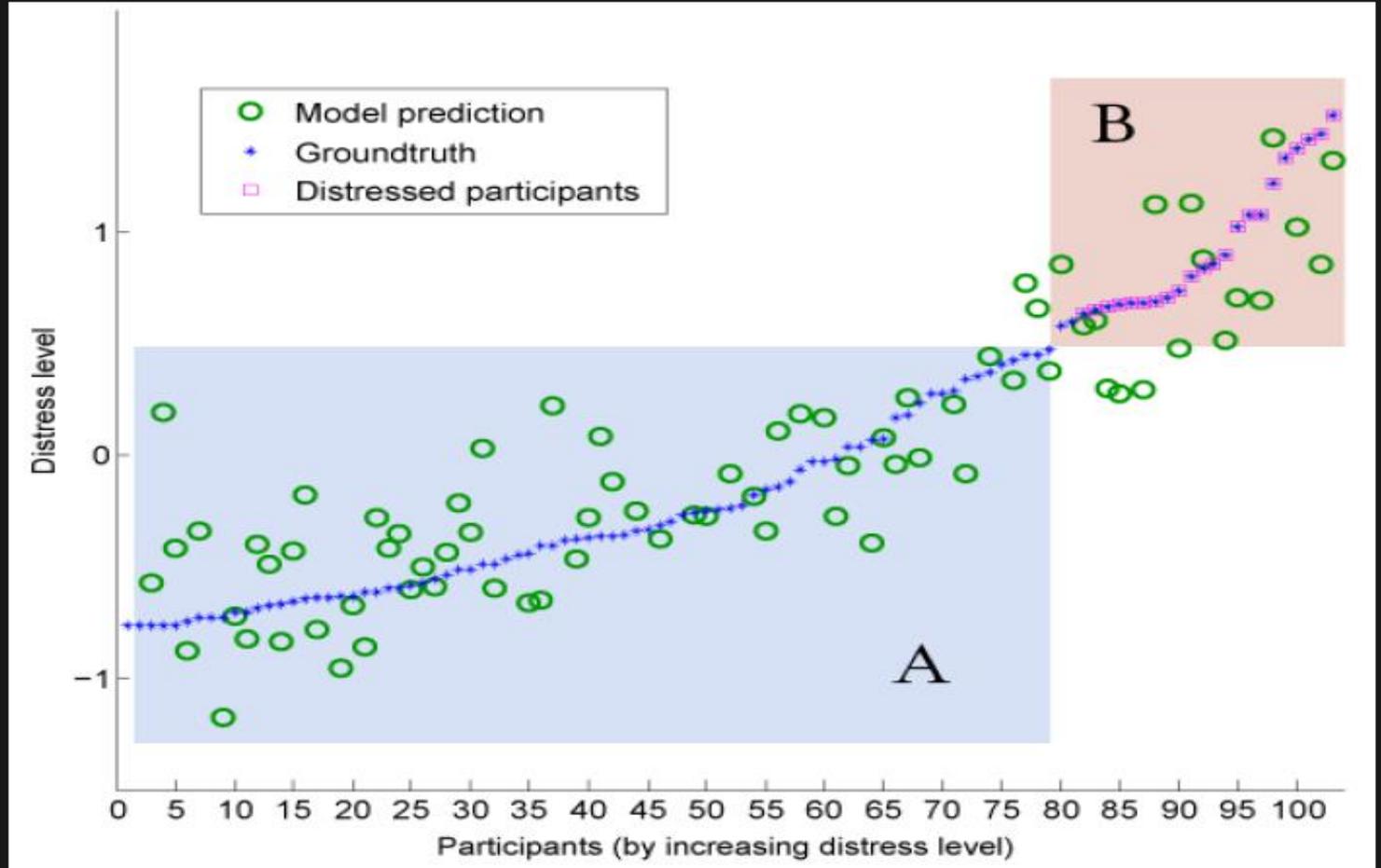
Verbal

We saw the yellow dog

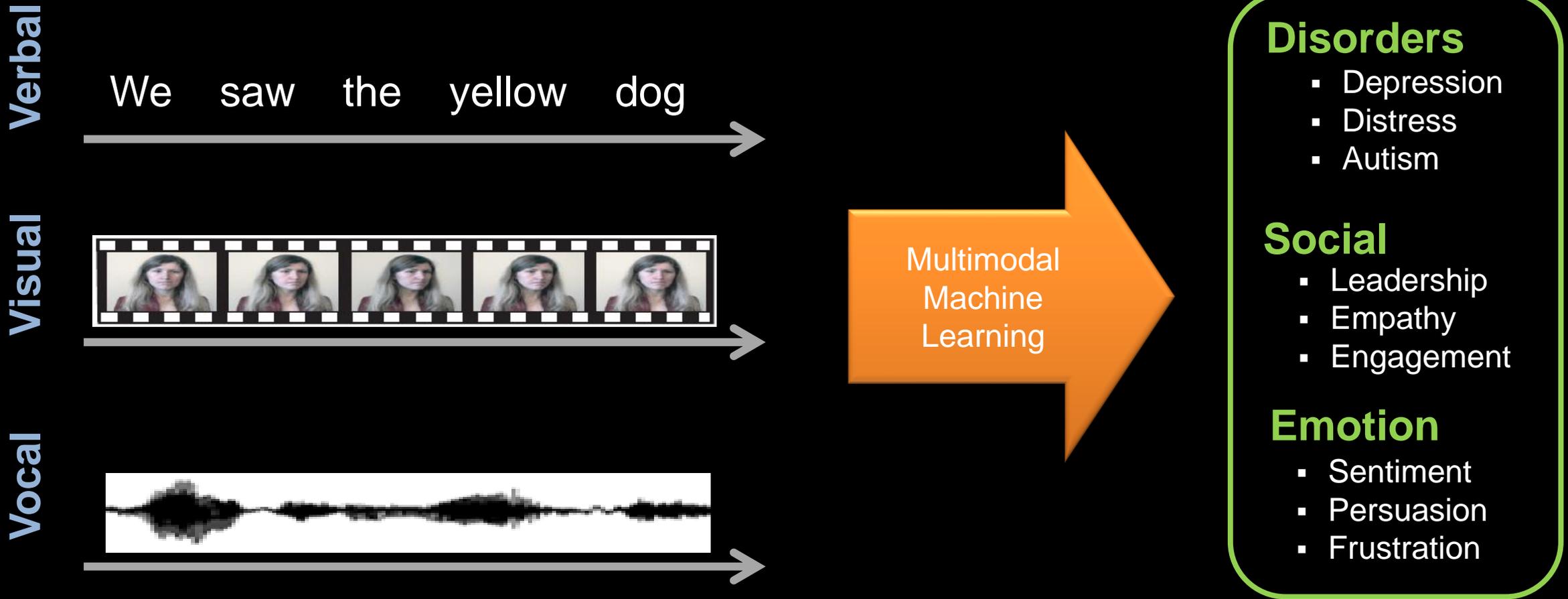
Vocal



Visual



Multimodal Machine Learning



Social-IQ: A QA Benchmark for Artificial Social Intelligence

[CVPR 2019]

Social Intelligence

1250

Videos

7500

Questions

30k

Correct Answers

22.5k

Incorrect Answers



Social Gathering



Intimate Moments



Debates



Discussion



Opinion Sharing



TV Shows

Social Phenomena

- a. Are people getting along?*
- b. How is the atmosphere in the room?*

Mental State and Attitude

- a. Was the man hurt when insulted?*
- b. Was the woman brave?*

Multimodal Behavior

- a. How did the man show his discontent?*
- b. How did the woman respond to the rude person?*

Core Challenges in 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>

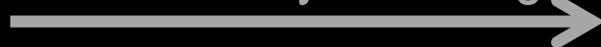
- ✓ 5 core challenges
- ✓ 37 taxonomic classes
- ✓ 253 referenced citations

Core Challenge 1: Multimodal Representation

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

Verbal

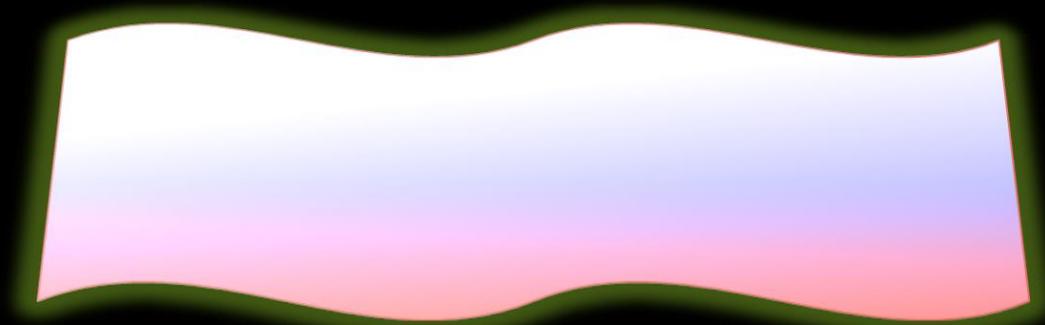
We saw the yellow dog



Visual

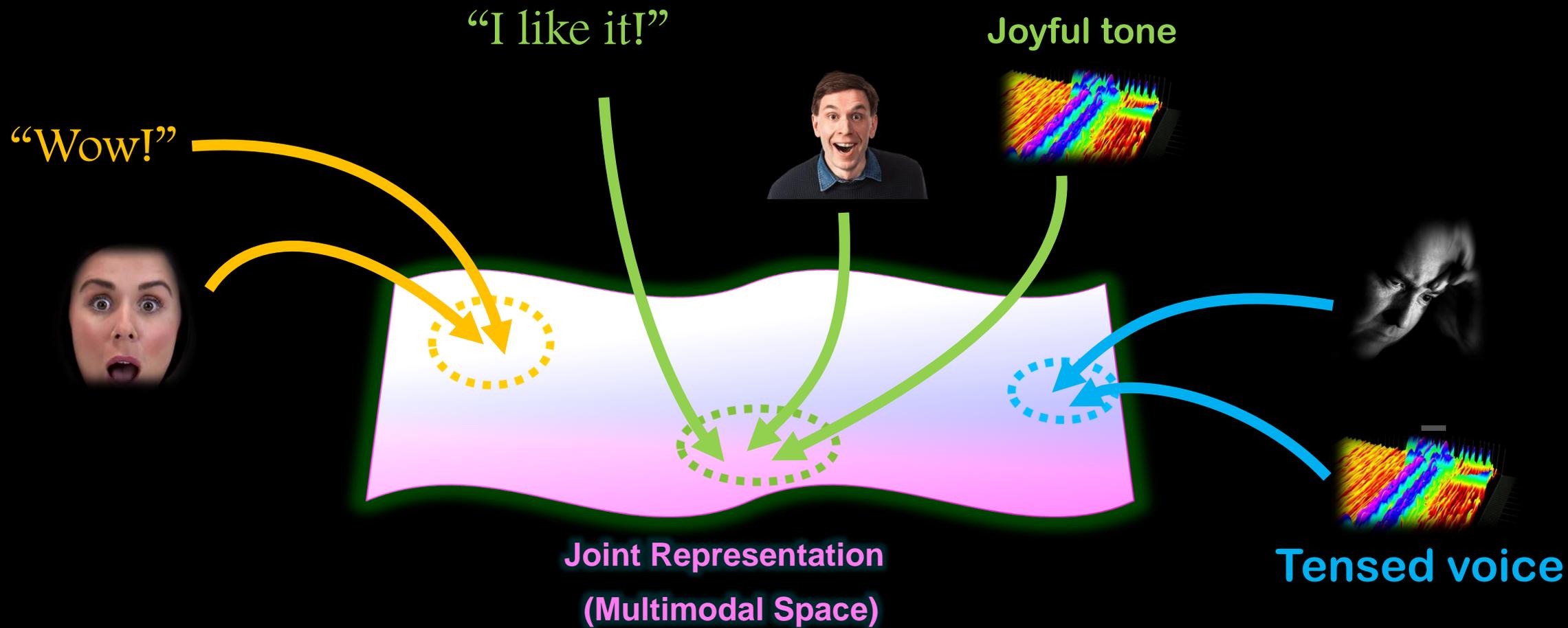


Vocal



Joint Representation
(Multimodal Space)

Joint Multimodal Representation



Multimodal Sentiment Analysis

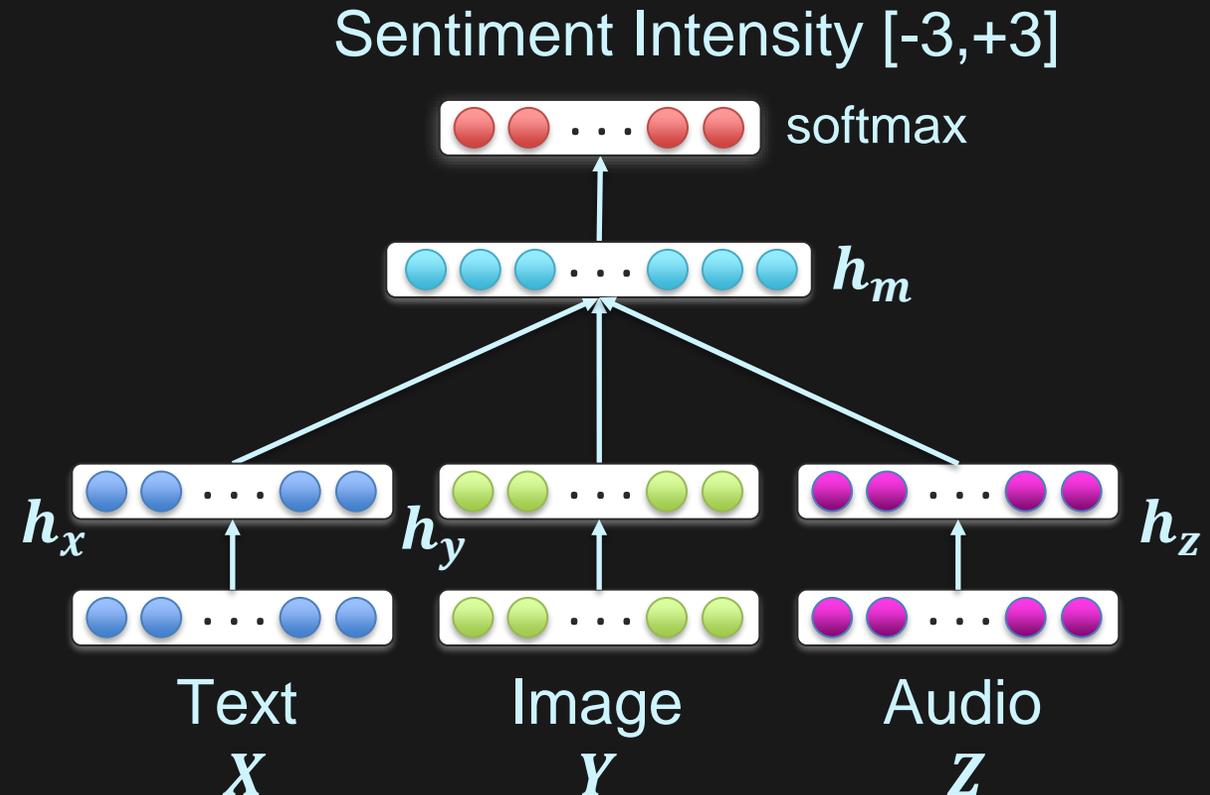
MOSI dataset (Zadeh et al, 2016)



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

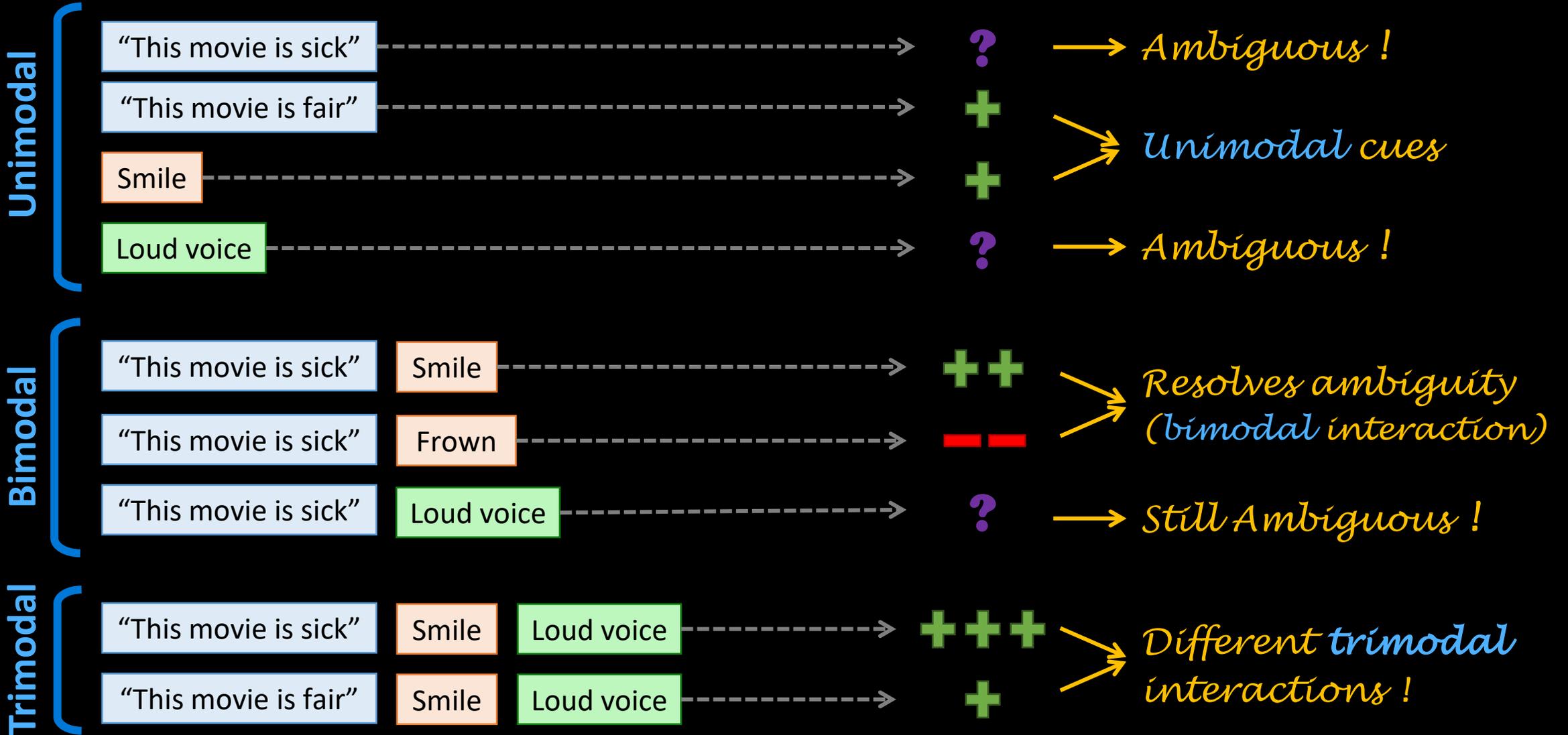
Multimodal joint representation:

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



Speaker's behaviors

Sentiment Intensity

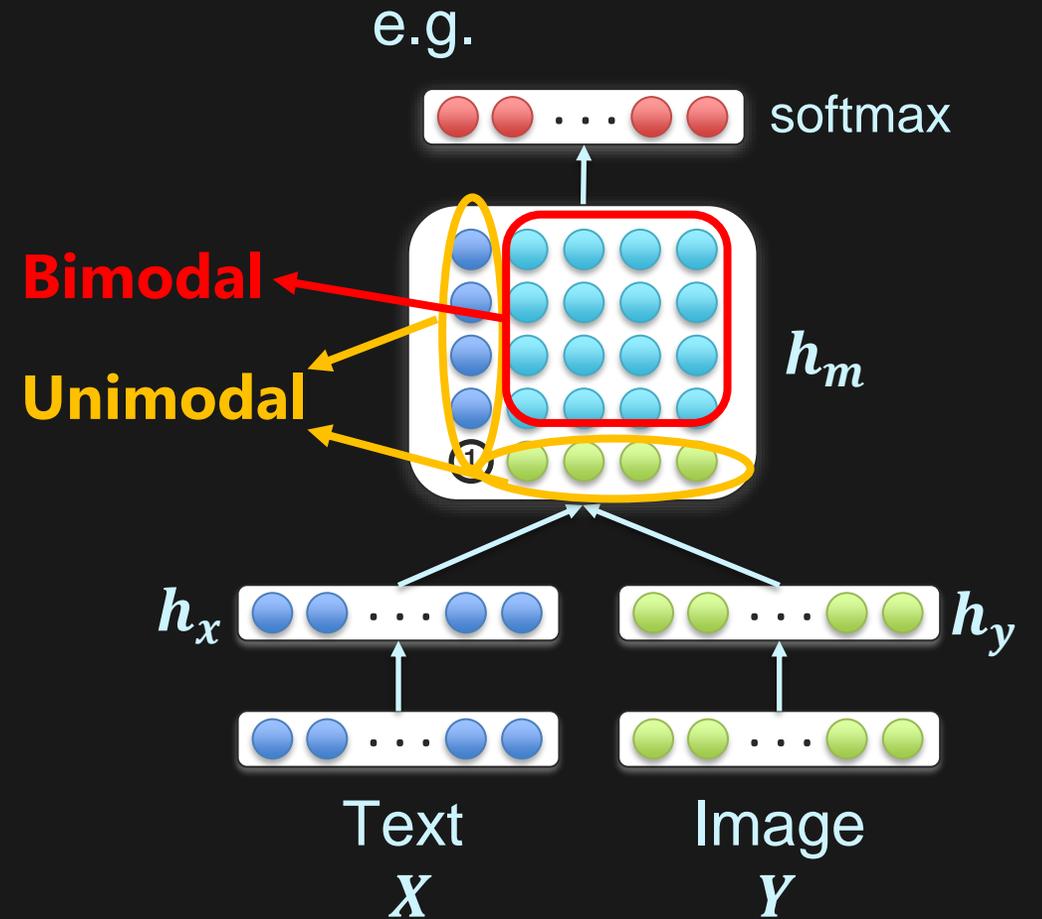


Multimodal Tensor Fusion Network (TFN)

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!



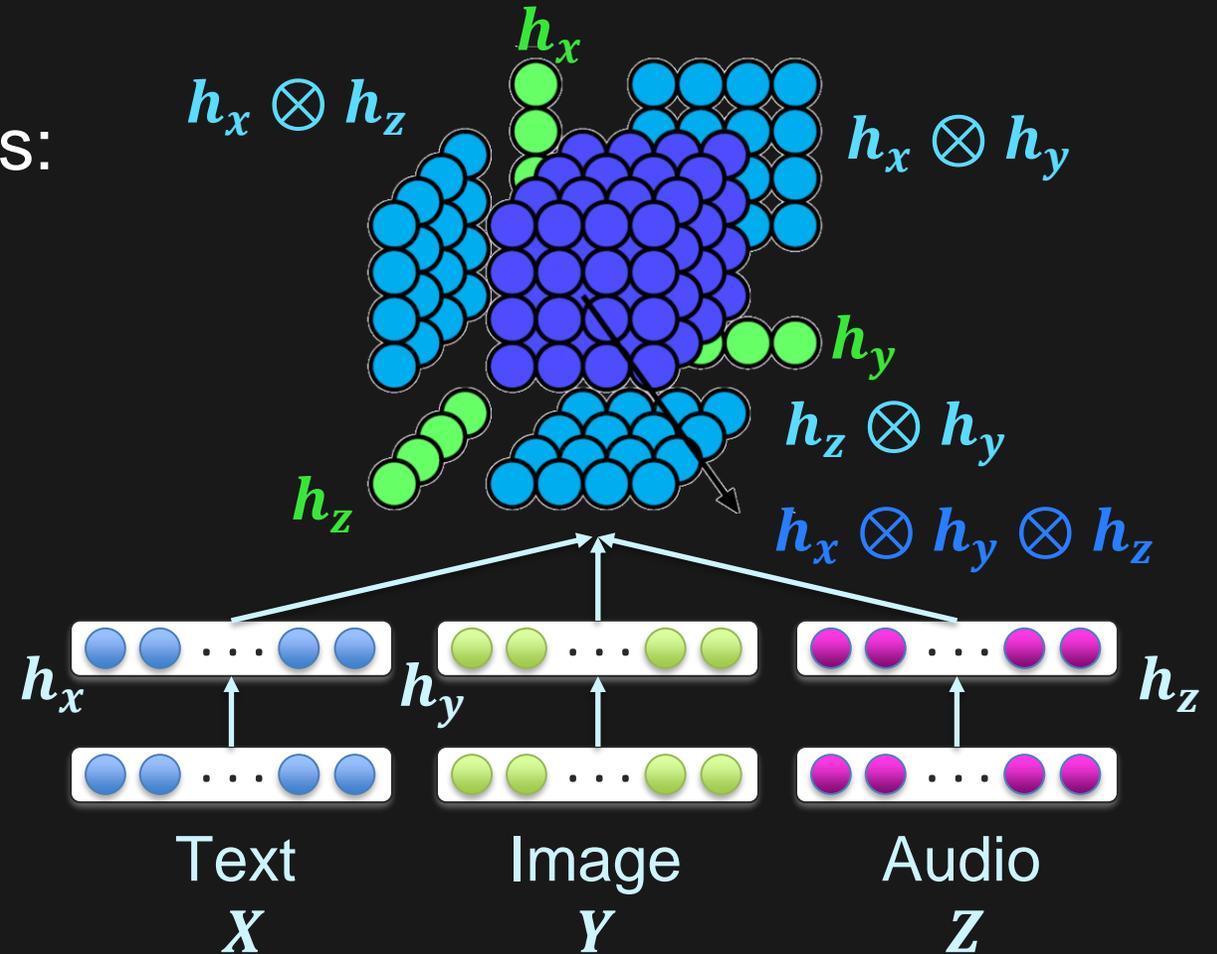
[Zadeh, Jones and Morency, EMNLP 2017]

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 !

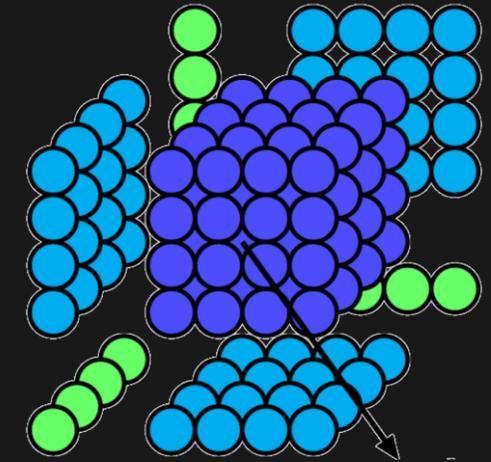


[Zadeh, Jones and Morency, EMNLP 2017]

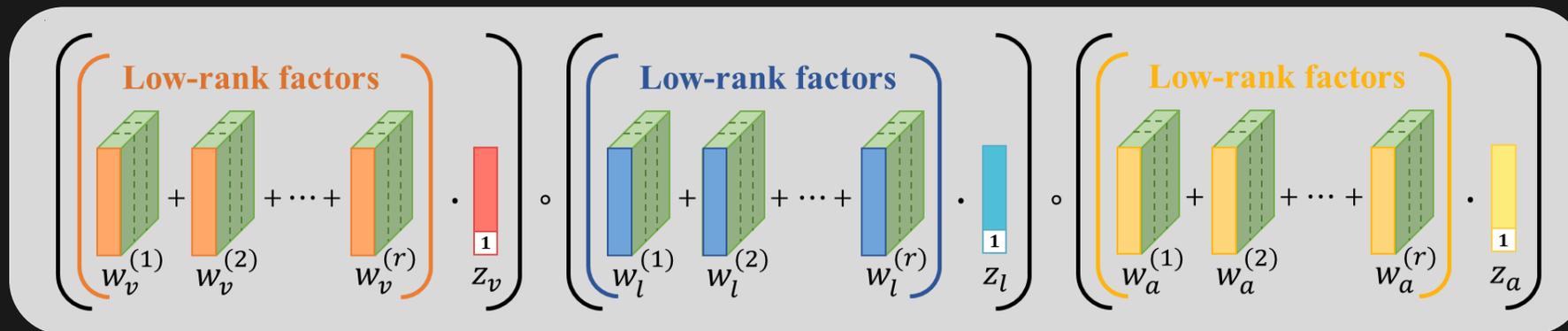
Improving Efficiency of Multimodal Representations [ACL 2018]

Tensor Fusion Network: Explicitly models unimodal, bimodal and trimodal interactions

[Zadeh, Jones and Morency, EMNLP 2017]



Efficient Low-rank Multimodal Fusion



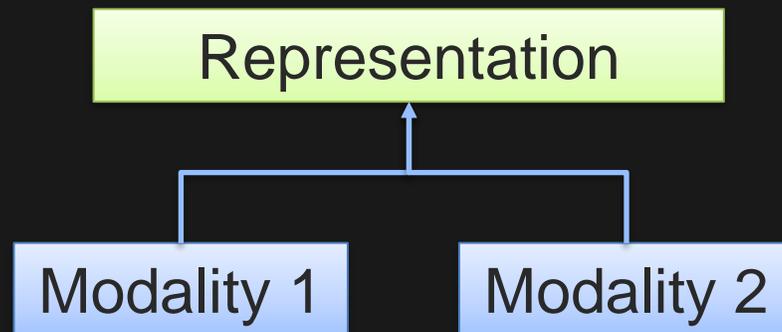
Canonical Polyadic Decomposition

[Liu, Shen, Bharadwaj, Liang, Zadeh and Morency, ACL 2018]

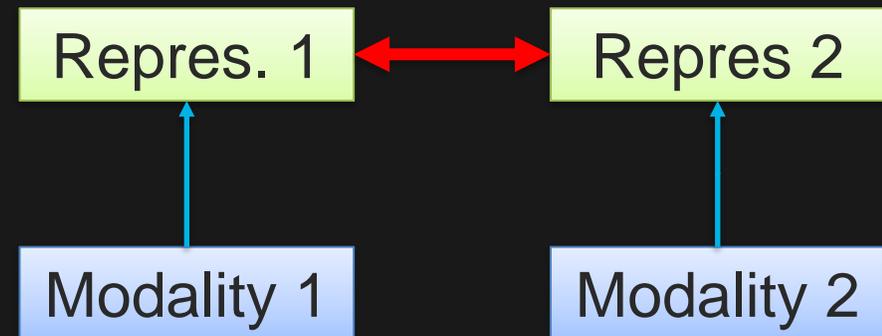
Core Challenge 1: Representation

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

(A) Joint representations:



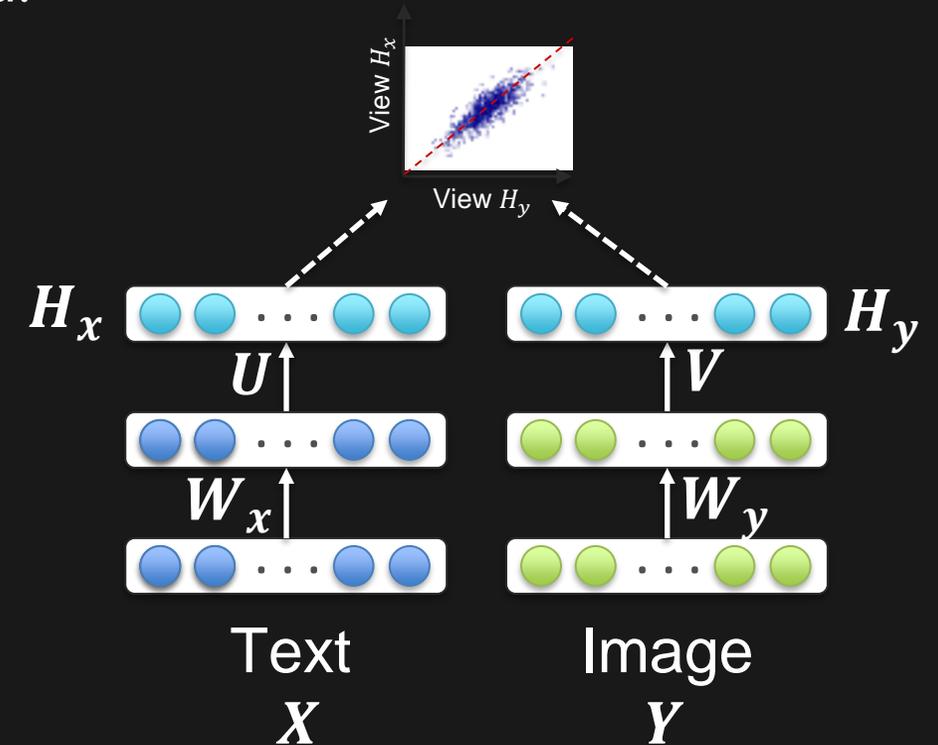
(B) Coordinated representations:



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

Toward Debiasing Sentence Representations [ACL 2020]

“The boy is coding.” OR “The girl is coding.”

“The boys at the playground.”

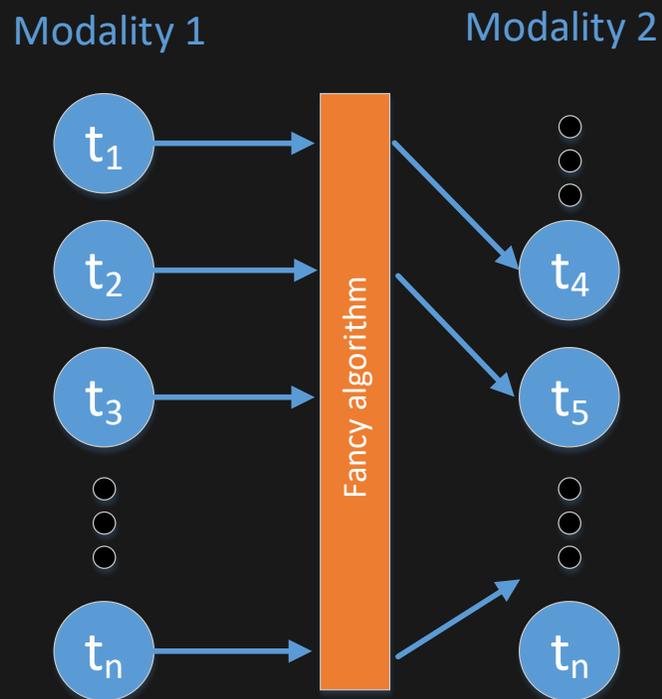
OR

“The girls at the playground.”

RESEARCH QUESTION: How to debias multimodal representations?

Core Challenge 2: Alignment

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



A Explicit Alignment

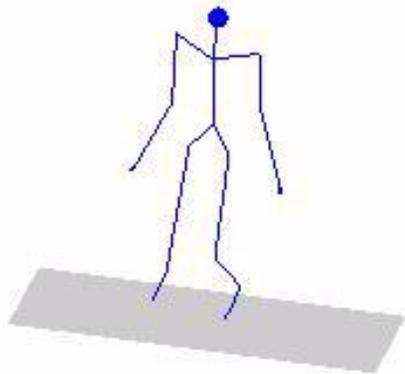
The goal is to directly find correspondences between elements of different modalities

B Implicit Alignment

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

Explicit (Multimodal) Alignment

1/273



1/51



1/127



Alignment and Representation: Multimodal Transformer

[ACL 2019]

Visual



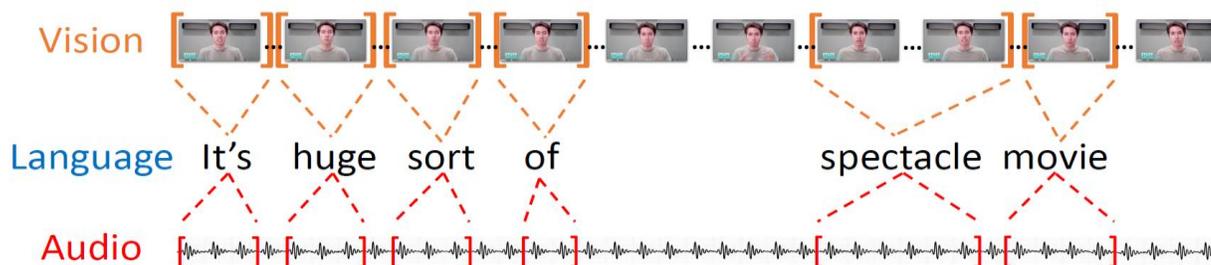
Vocal



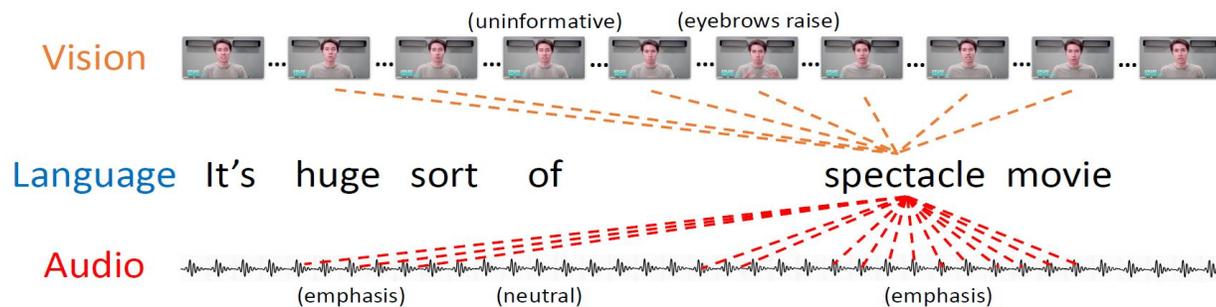
Verbal

"I like..."

Predefined Word-level alignment



Automatic Cross-Modal alignment



Multimodal representation

Representation

Alignment

time

Multimodal Transformer [ACL 2019]

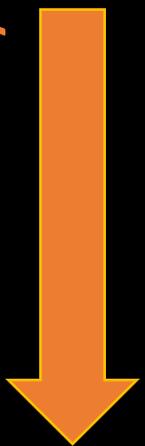
Alignment

Representation

Visual



Visually contextualizing
the verbal modality

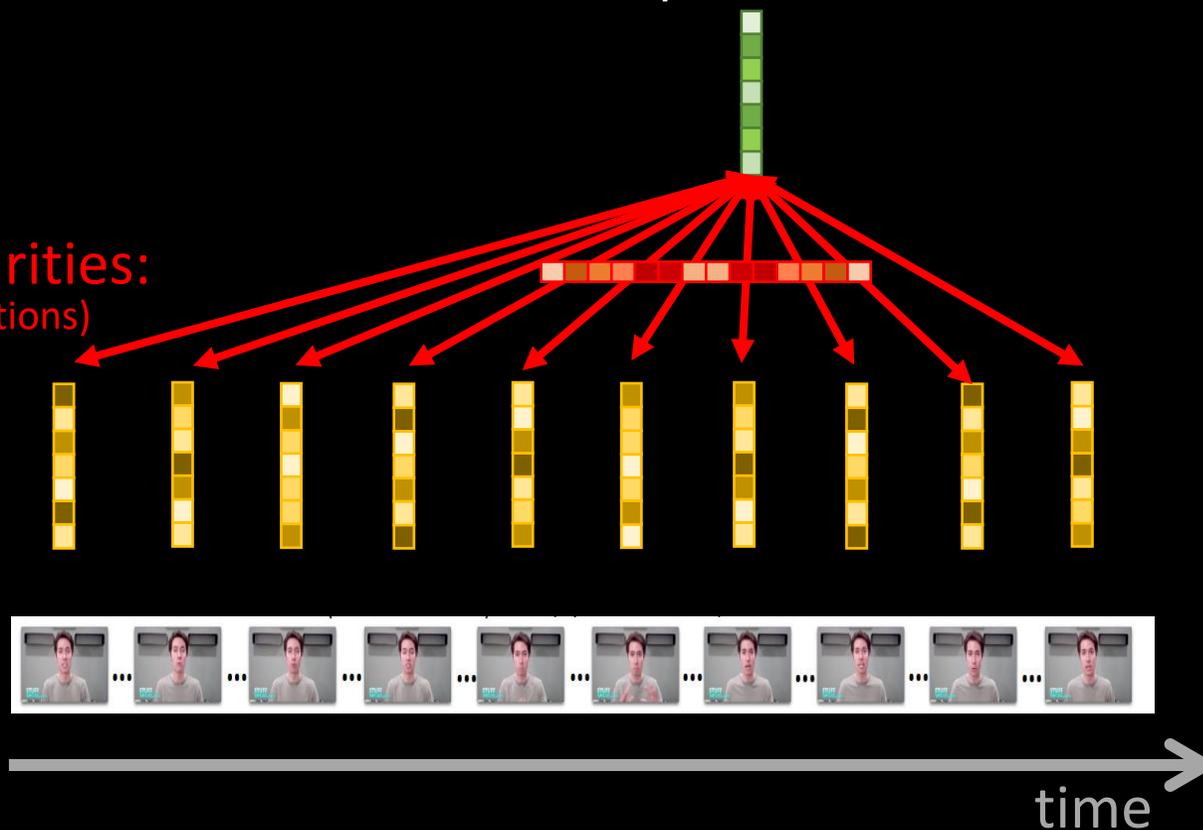


Verbal

“I like...”

“spectacle”

Similarities:
(attentions)



Multimodal Transformer [ACL 2019]

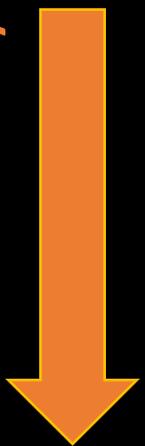
Alignment

Representation

Visual



Visually contextualizing
the verbal modality

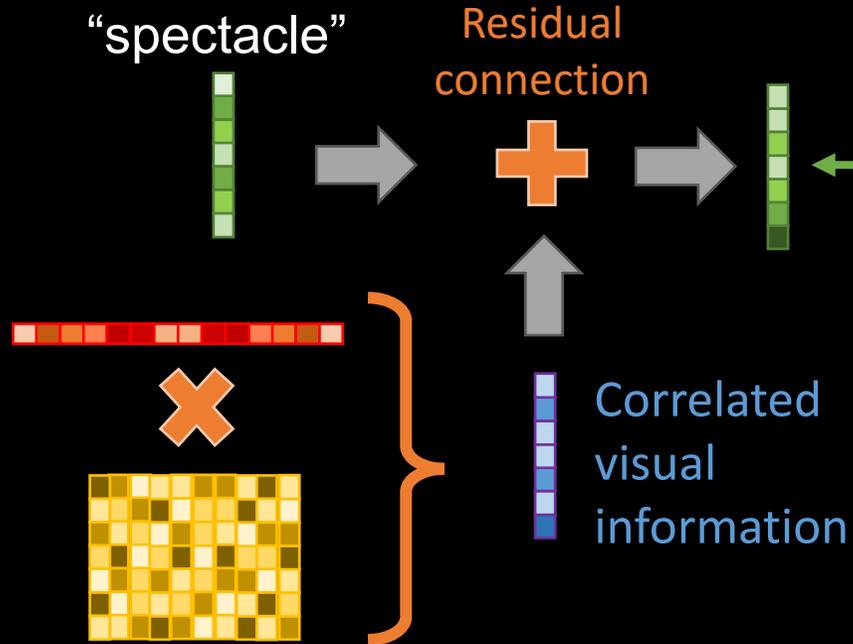


Verbal

"I like..."

Similarities:
(attentions)

Visual embeddings:



time

Multimodal Transformer [ACL 2019]

Alignment

Representation

Visual

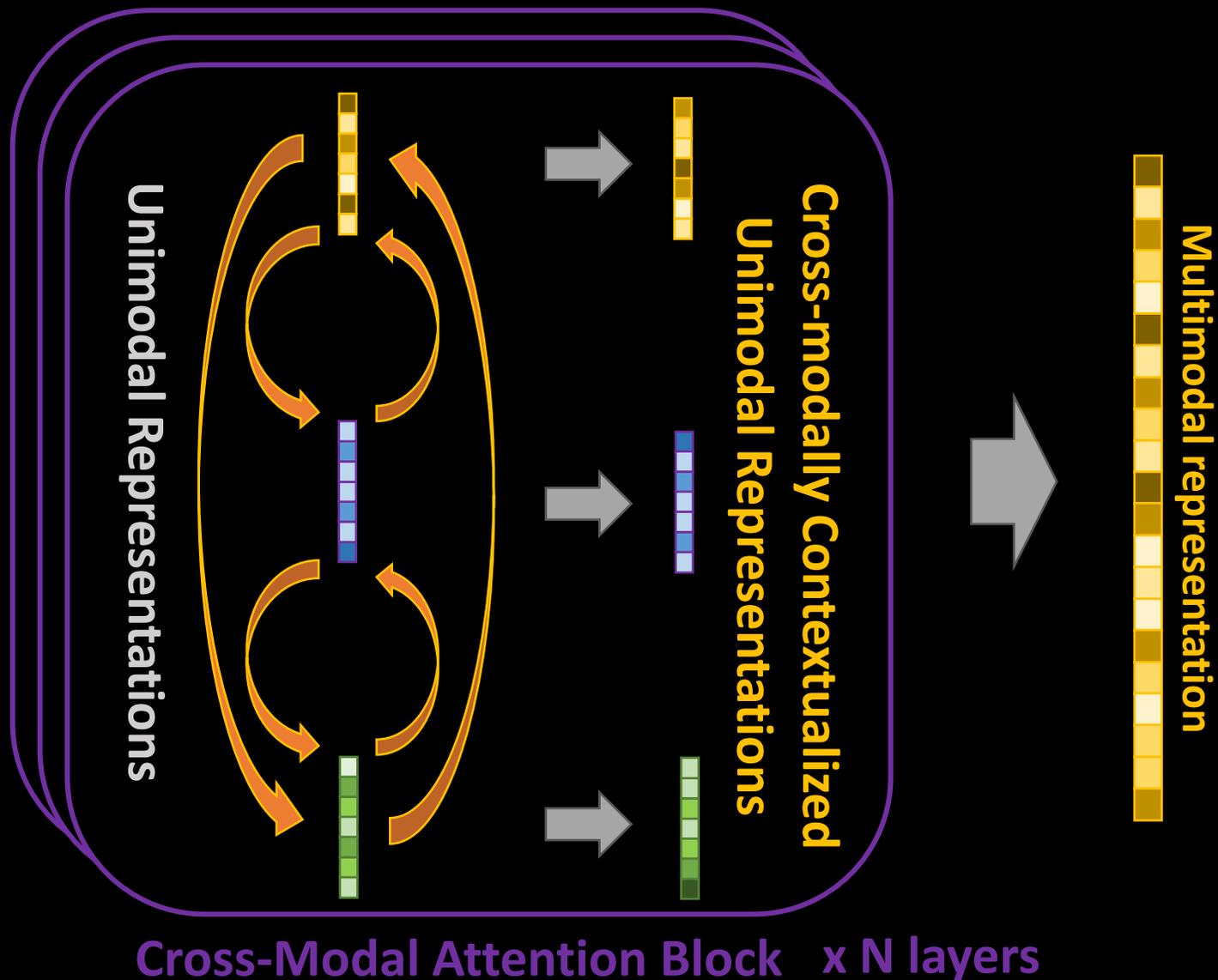


Vocal



Verbal

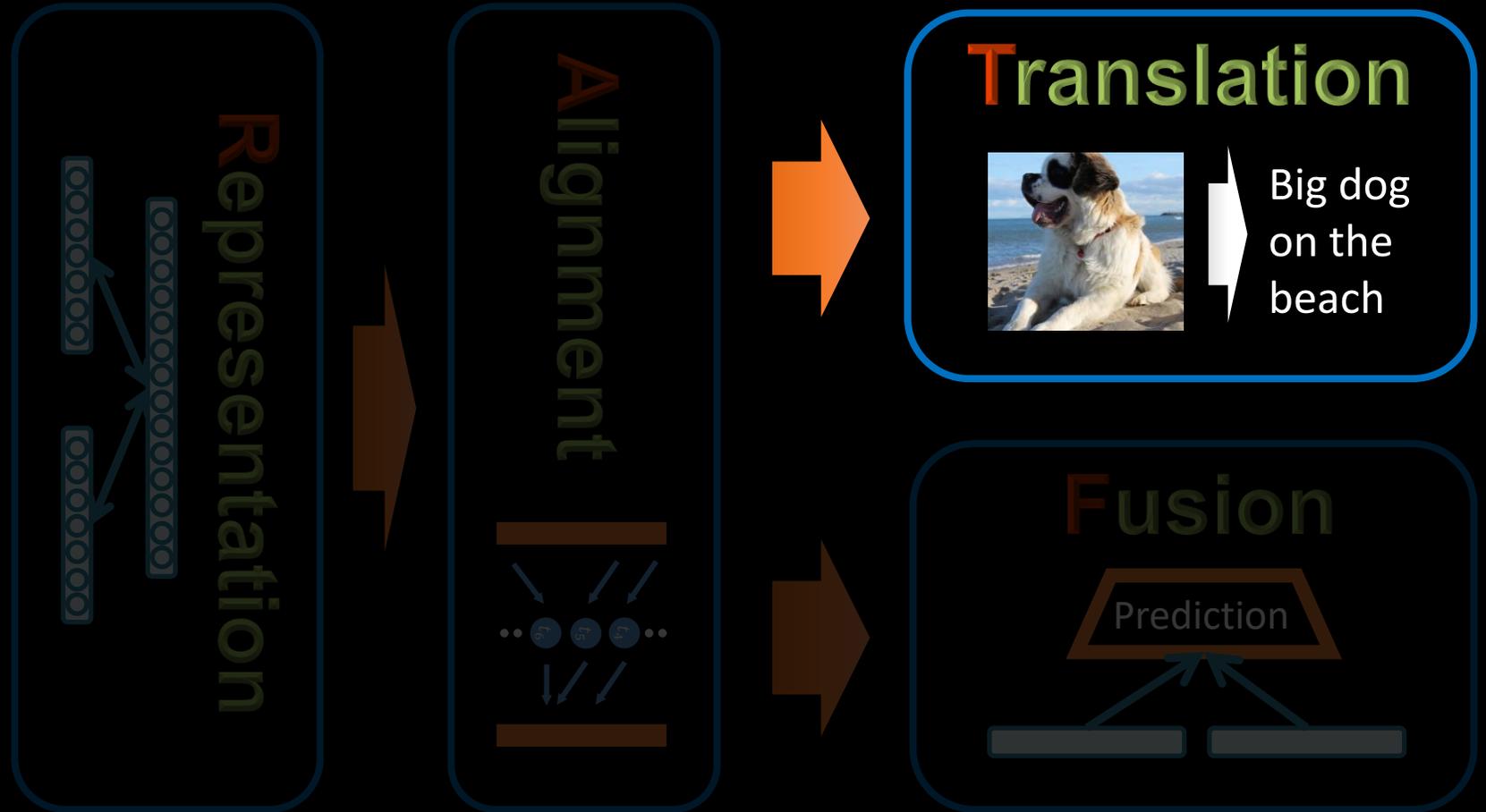
"I like..."



Cross-Modal Attention Block x N layers

Multimodal AI – Core Challenges

[Survey: TPAMI 2019]

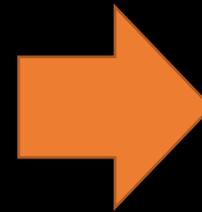


Language-to-Pose [3DV 2019]

Translation

Story Narrative

“Characters walk hands-in-hands slowly while talking and then decide to leave the road and explore the forest...”



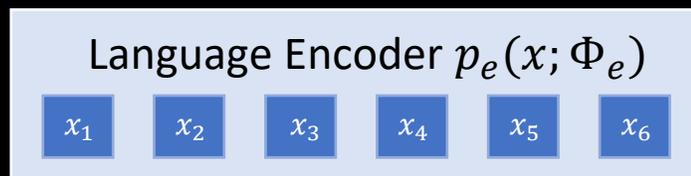
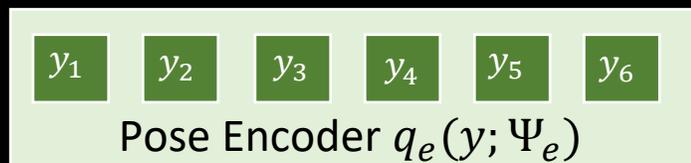
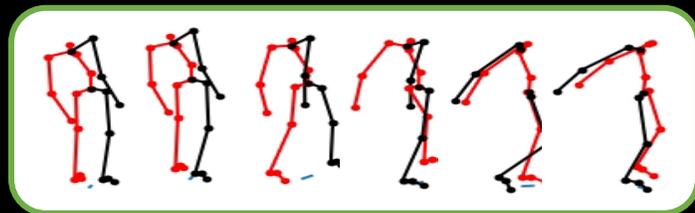
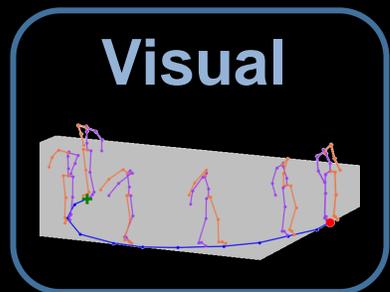
Movie Animation



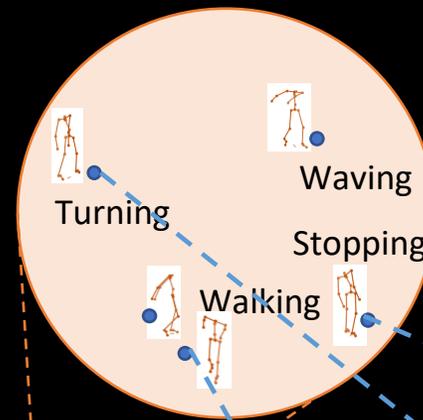
More than a week for professional animators to create 15 seconds !

Language-to-Pose [3DV 2019]

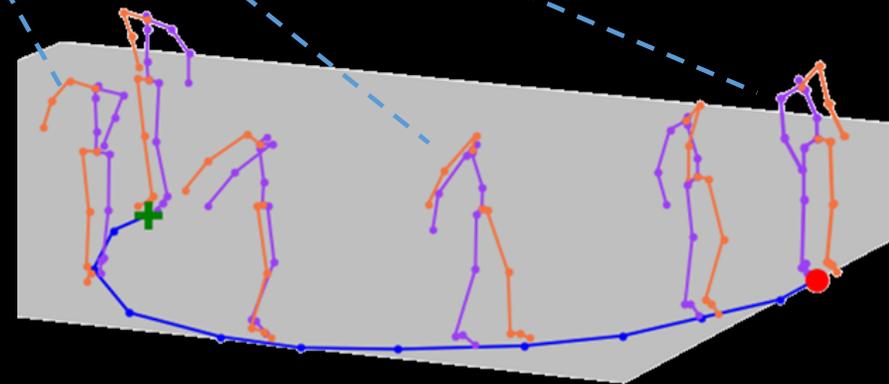
Translation



“A person walks in a circle”



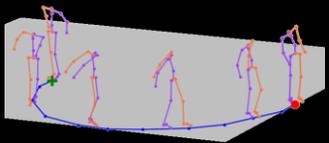
Joint Embedding
Space \mathbb{Z}



Language-to-Pose [3DV 2019]

Translation

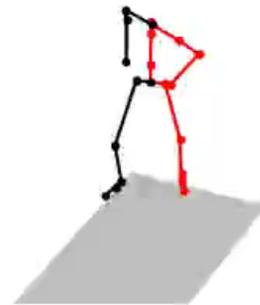
Visual



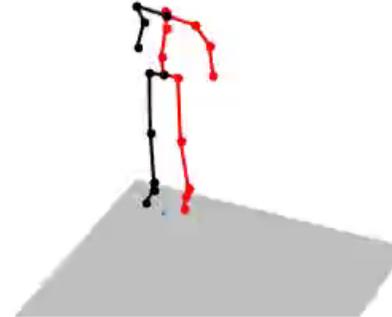
Verbal

“A person
walks in a
circle”

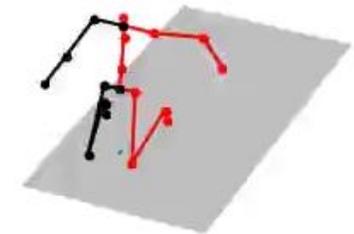
Results from our JL2P model (Joint Language-to-Pose)



a person jogs
a few steps



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



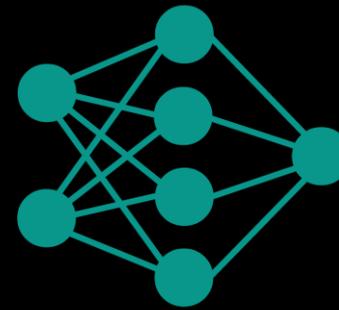
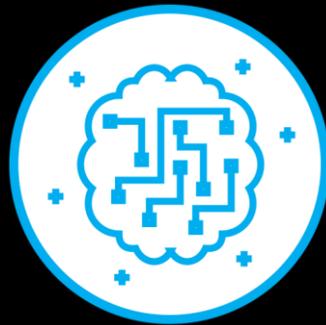
A kneeling person raises
their arms to the sides
and stand up.

Language-to-Network [ACL 2020]



Rhinoceros Auklet: Breeding adults are **cloudy gray**, ... a thick **orange-yellow bill** ...

Acadian Flycatcher: **olive-green** above with a **whitish eye-ring** ... **two distinct white wing-bars**.



- ✓ Rhinoceros Auklet
- ? Acadian Flycatcher
- ✗ Acadian Flycatcher

Language-to-Action [ACL 2020]



Multi-step instruction:

- 1 **Go** to the **entrance** of the **lounge area**.
- 2 **On your right** there will be **a bar**.
- 3 **On top** of the **counter**, you will see **a box**.
Bring me **that**.

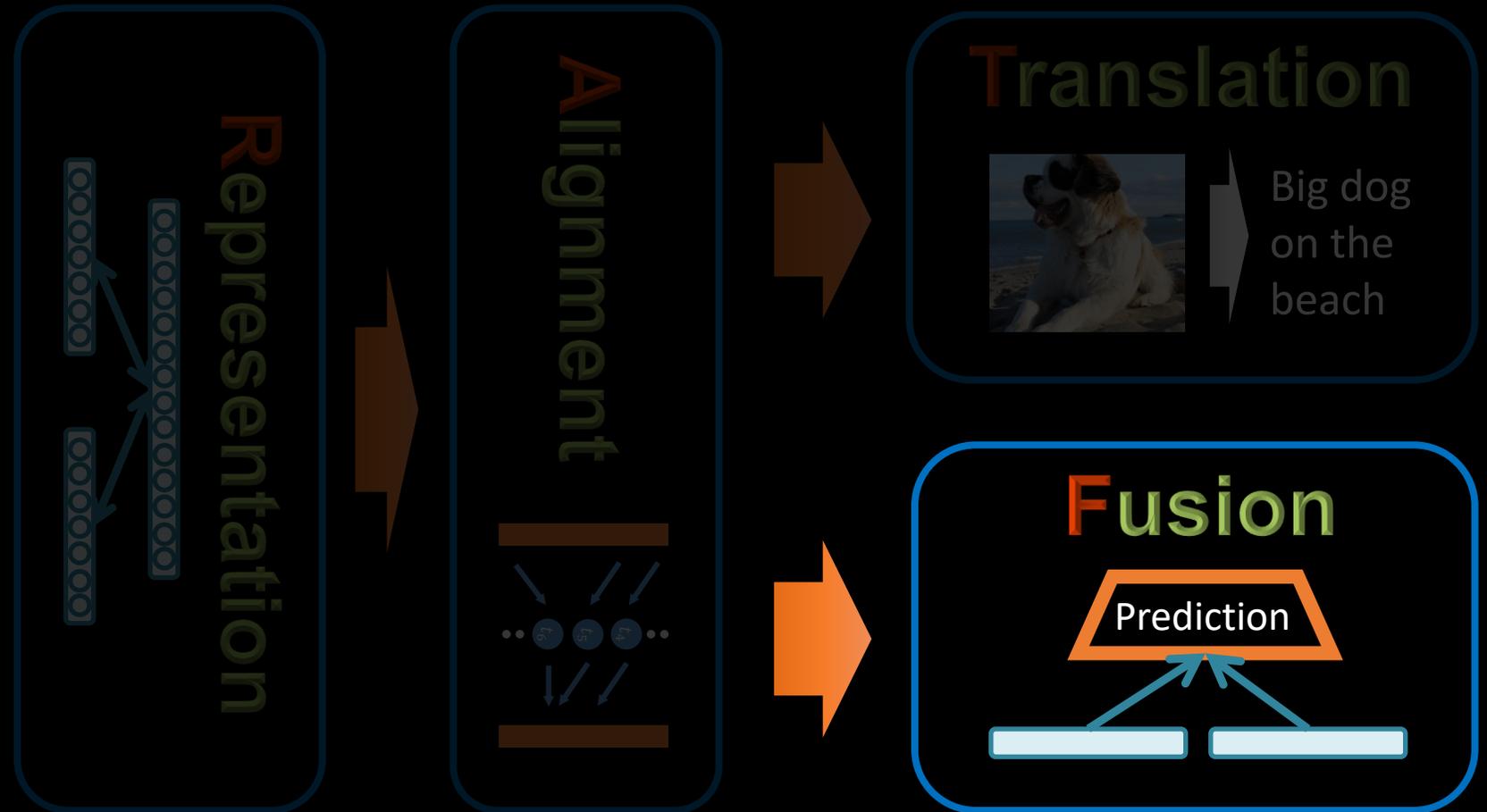
Refer360 dataset

- 17,135 annotated instances
- 2,000 panoramic 360 degrees scenes
- 43.8 average number of words per instructions

<https://github.com/volkancirik/refer360>

Multimodal AI – Core Challenges

[Survey: TPAMI 2019]



Multimodal Sentiment Analysis [ACL 2018]

CMU- MOSEI dataset

Visual



Vocal



Verbal

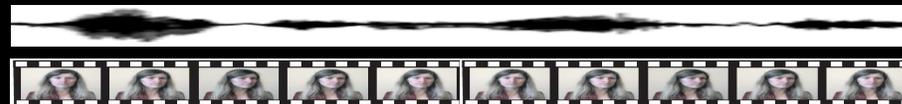
“I like...”



- ➔ 23,000 video segments
- ➔ 1,000 speakers (from vlogs)
- ➔ 3 modalities (verbal, vocal, visual)

<https://github.com/A2Zadeh/CMU-MultimodalSDK>

“He’s average presenter when...”



Multimodal Sentiment Analysis [AAAI 2018]

Visual



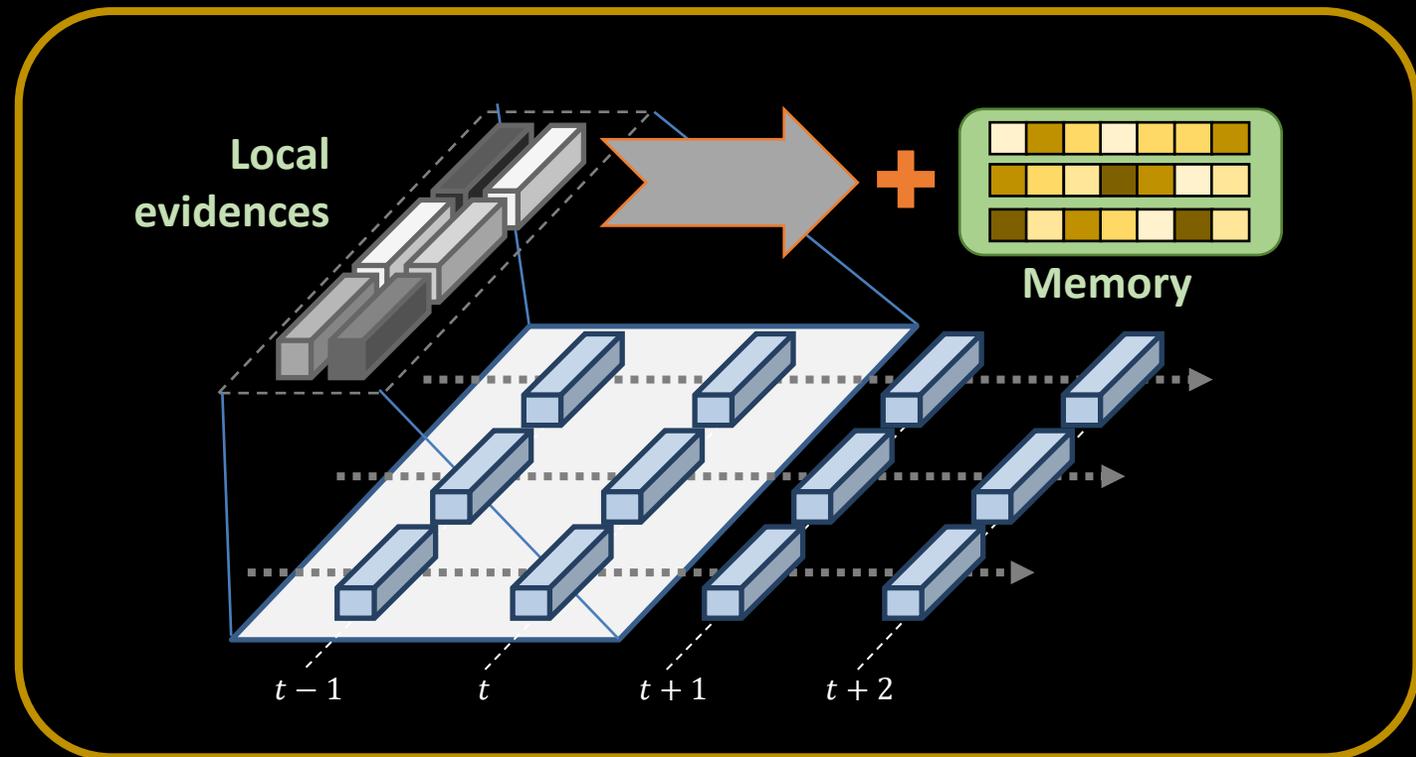
Vocal



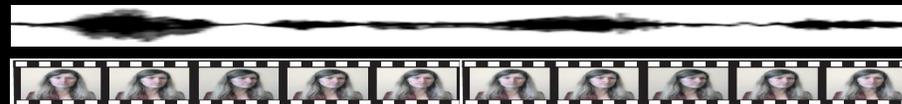
Verbal

“I like...”

Memory Fusion Network

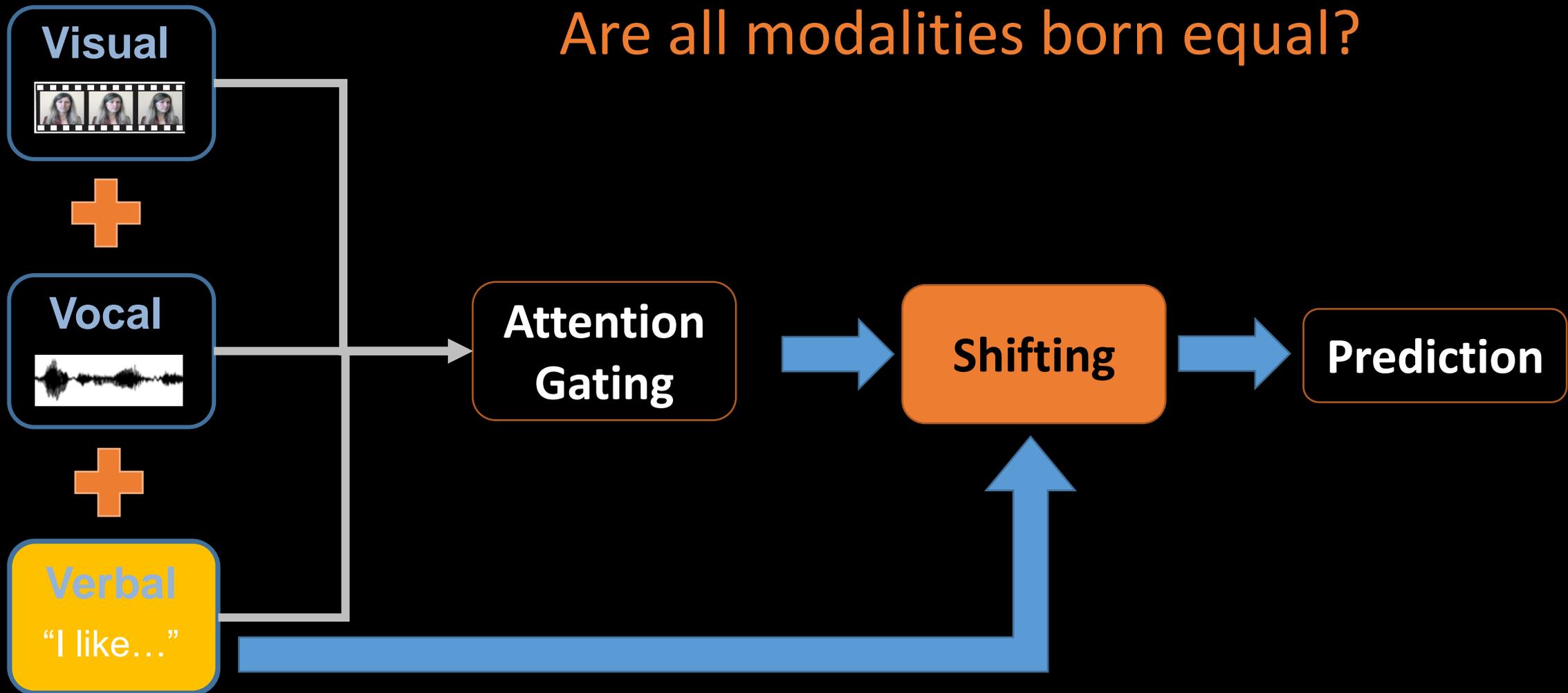


“He’s average presenter when...”



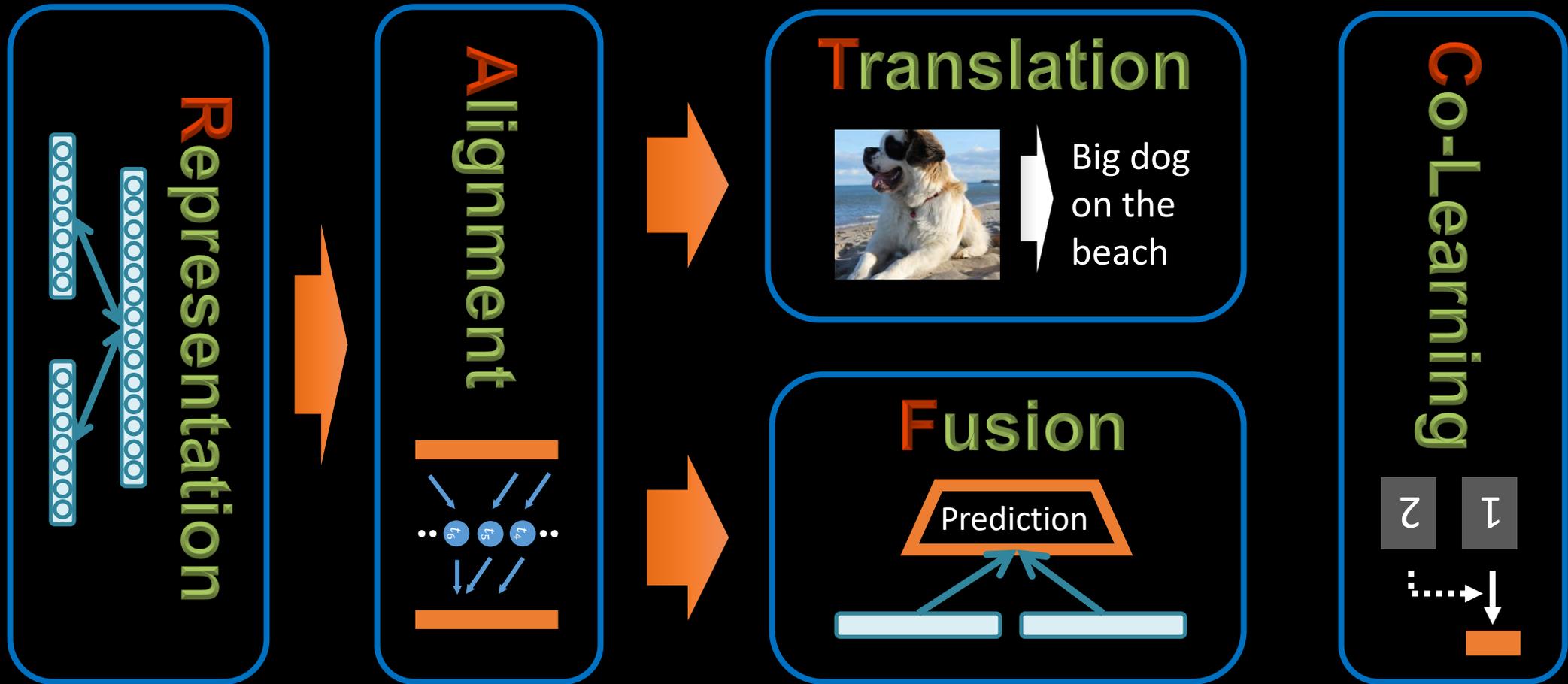
Multimodal Adaptation Gate [ACL 2020]

Are all modalities born equal?



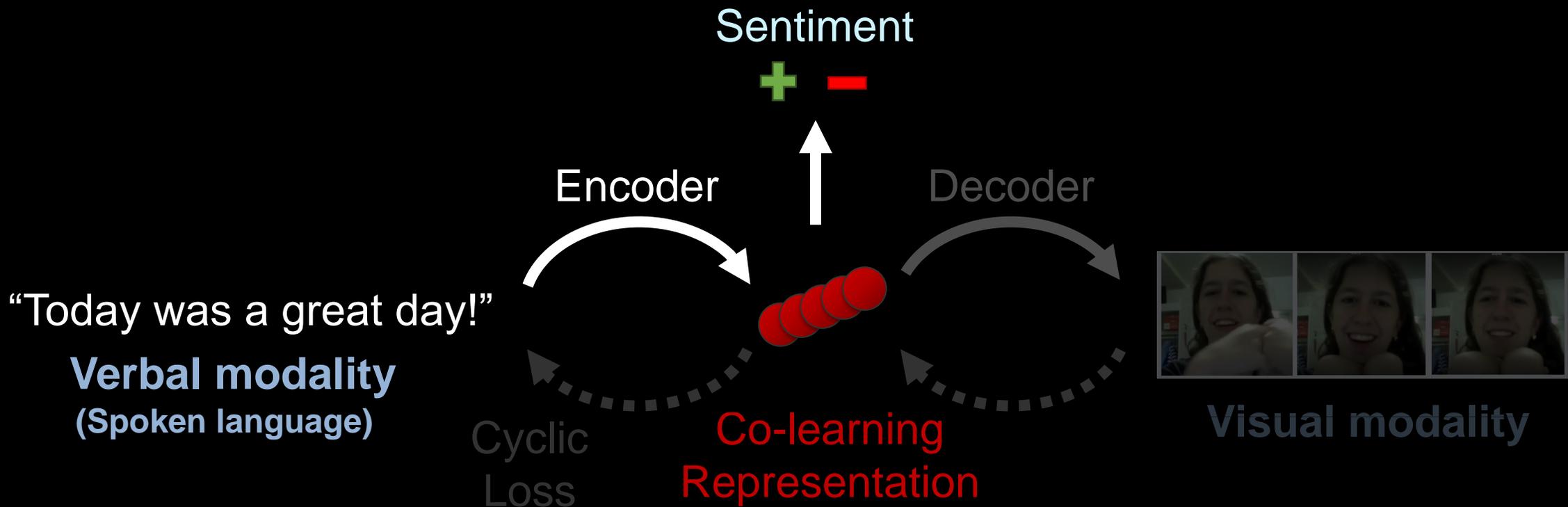
Multimodal AI – Core Challenges

[Survey: TPAMI 2019]



Representations from Cross-modal Translation

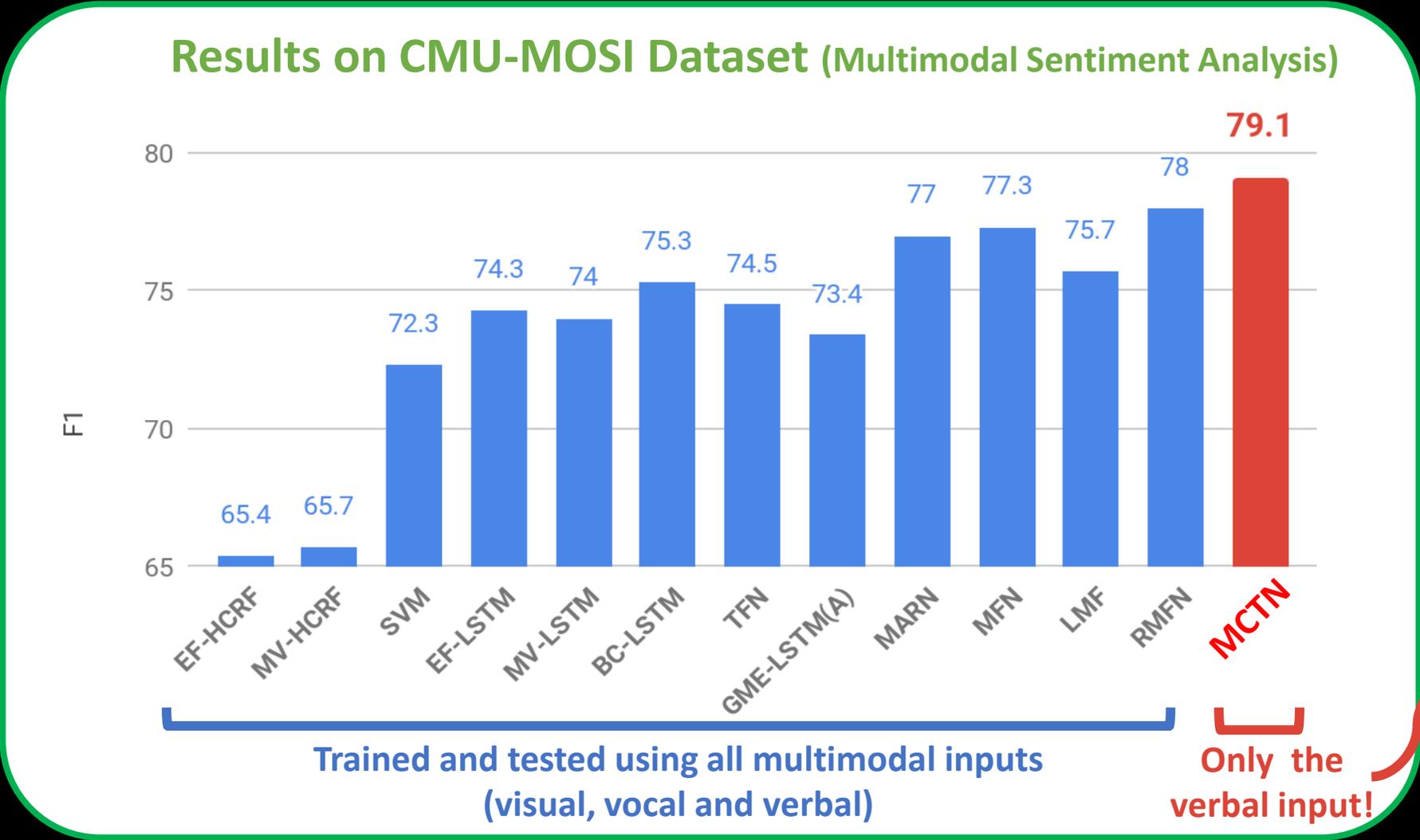
[AAAI 2019]



Multimodal Cyclic Translation Network (MCTN)

Representations from Cross-modal Translation

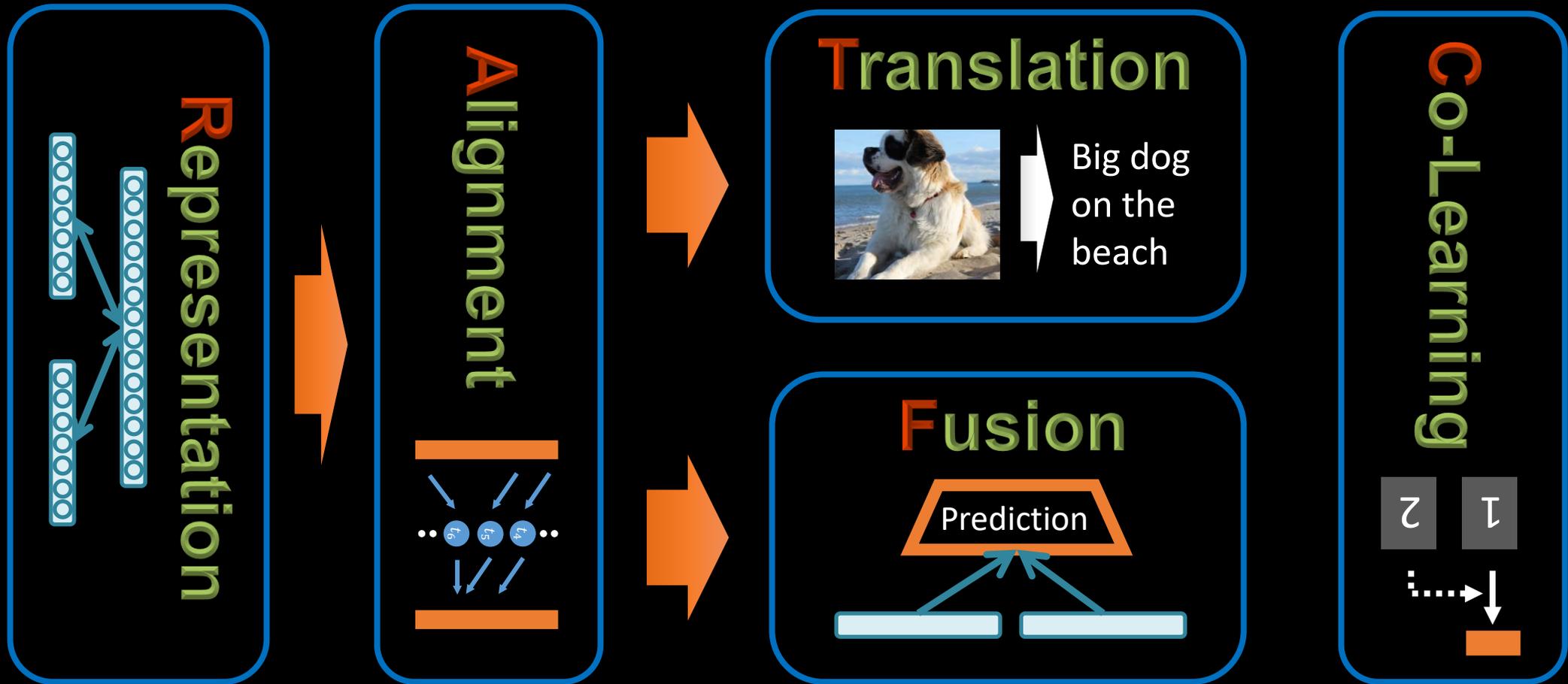
[AAAI 2019]



Trained using multimodal co-learning but tested with only verbal input!

Multimodal AI – Core Challenges

[Survey: TPAMI 2019]



Multimodal Machine Learning – Taxonomy

[Survey: TPAMI 2019]

Representation

Joint

- *Neural networks*
- *Graphical models*
- *Sequential*

Coordinated

- *Similarity*
- *Structured*

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*

Translation

Example-based

- *Retrieval*
- *Combination*

Model-based

- *Grammar-based*
- *Encoder-decoder*
- *Online prediction*

Co-learning

Parallel data

- *Co-training*
- *Transfer learning*

Non-parallel data

- Zero-shot learning*
- Concept grounding*
- Transfer learning*

Hybrid data

- Bridging*

CMU Course on Multimodal Machine Learning

PIAZZA 11-777 Q & A Resources Statistics Manage Class  Louis-Philippe Morency

Carnegie Mellon University - Fall 2017

11-777: Advanced Multimodal Machine Learning

Syllabus  

Course Information **Staff** Resources

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 main technical topics are: (1) multimodal representation learning, including multimodal auto-encoder and deep learning, (2) multimodal component analysis and fusion, including deep canonical correlation analysis and multi-kernel learning, (3) multimodal alignment and multi-stream modeling, including attention models and multimodal recurrent neural networks, and (4) multimodal graphical models, including continuous and fully-connected conditional random fields. 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.

Announcements [show all](#)

Deadline for final report: Monday 12/11 at 11:59pm  

12/07/17 1:59 PM

Hello! Bonjour!

Thank you all for the great presentations this week.

The deadline for the final reports is Monday 12/11 at 11:59pm. Given the grading schedule, we cannot push this deadline later. Any submission after this deadline will be counted as a late submission.

Best,
LP
[View on Piazza](#)

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

MERCI !



<http://multicomp.cs.cmu.edu/>