



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal Machine Learning

Lecture 12.1: Quantification

Paul Liang

** Co-lecturer: Louis-Philippe Morency.
Original course co-developed with Tadas Baltrusaitis.
Spring 2021 edition taught by Yonatan Bisk*

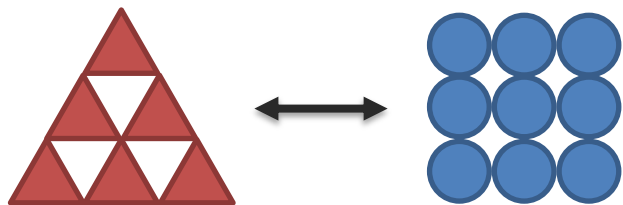
Reading Assignments are Back!

- Four main steps for the reading assignments
 - Monday 8pm: Official start of the assignment
 - Wednesday 8pm: Select your paper
 - **Friday 8pm:** Post your summary
 - **Monday 8pm:** Post your extra comments (5 posts)
- **4 papers:** Latent diffusion, explanations on VQA models, interaction quantification, benchmark quantification

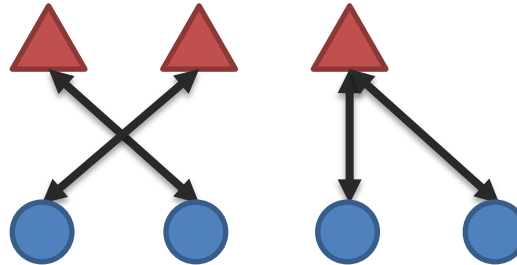
Quantification

Definition: Empirical and theoretical study to better understand heterogeneity, cross-modal interactions, and the multimodal learning process.

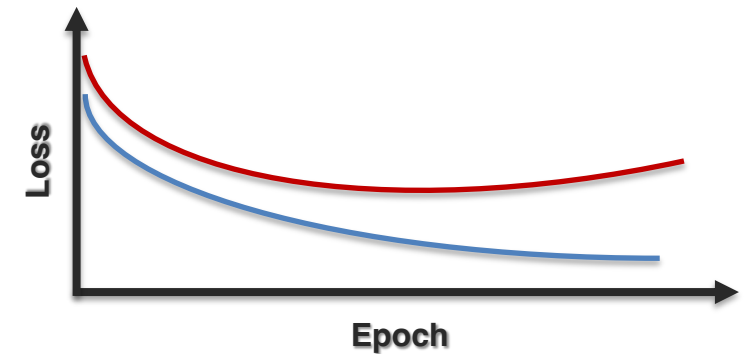
(A) Heterogeneity



(B) Interactions

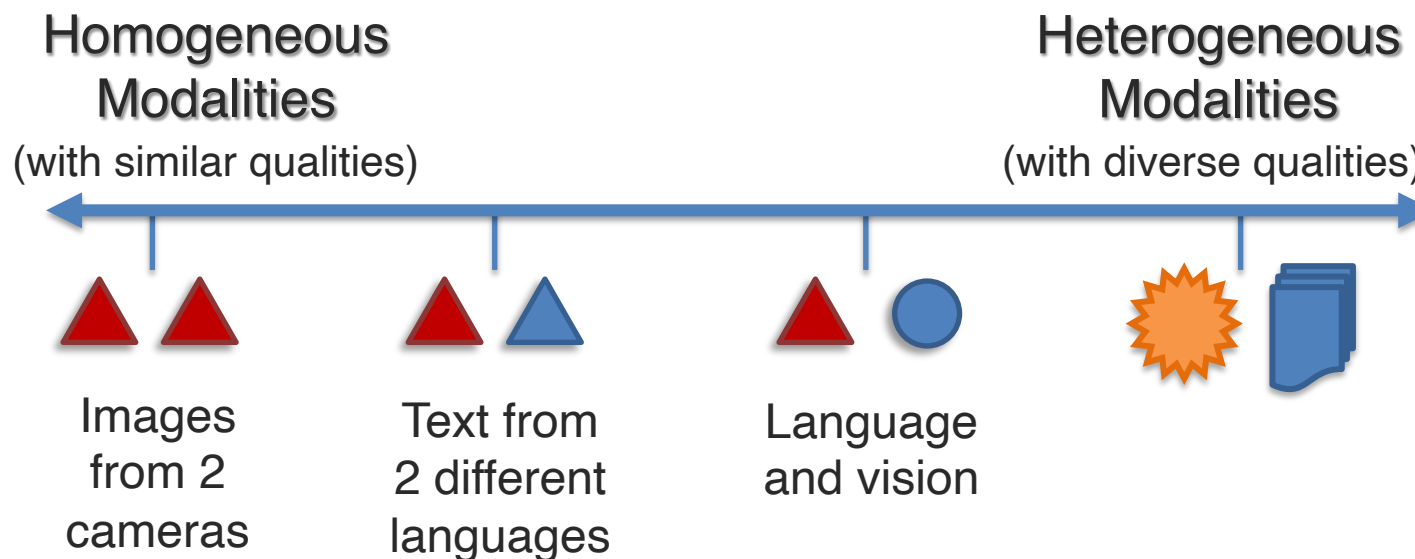
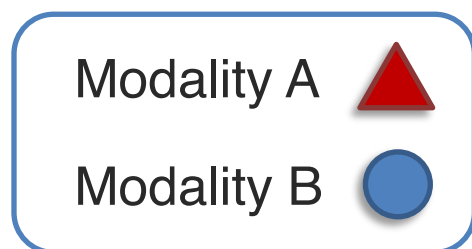


(C) Learning



Sub-Challenge 6a: Heterogeneity

Definition: Quantifying the dimensions of heterogeneity in multimodal datasets and how they subsequently influence modeling and learning.



Examples:

① Element representation

② Element distribution

③ Structure

④ Information

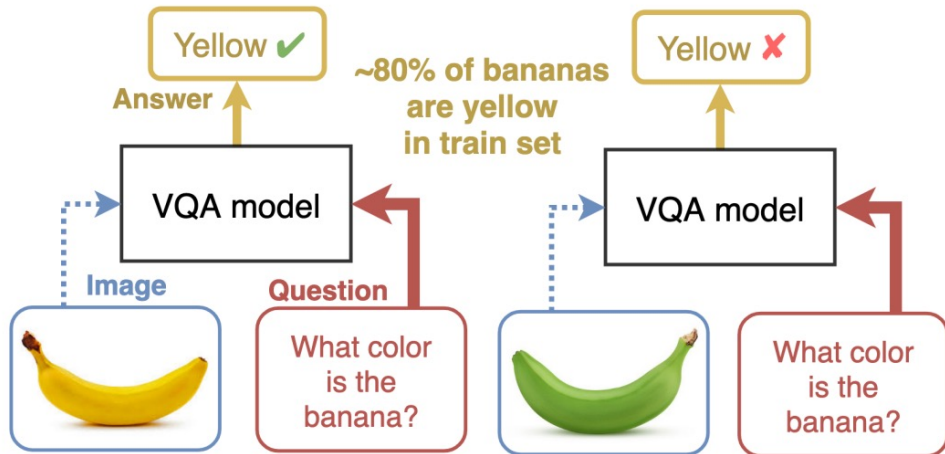
⑤ Noise

⑥ Relevance

Modality Biases

Heterogeneity in information and relevance
Unimodal biases and modality collapse

VQA models answer the question without looking at the image

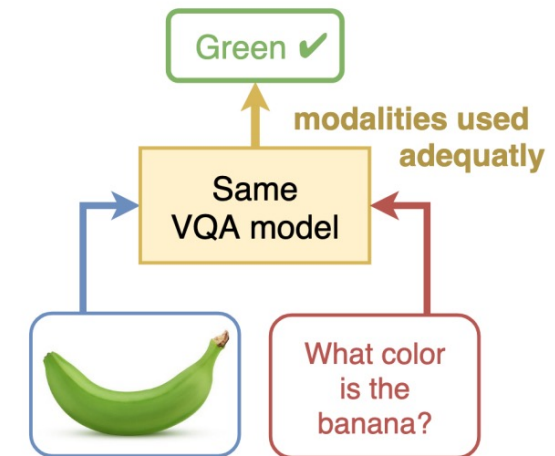


Balancing modalities

Balancing training



Not the case when trained with RUBi



[Wu et al., Characterizing and Overcoming the Greedy Nature of Learning in Multi-modal Deep Neural Networks. ICML 2022]

[Javaloy et al., Mitigating Modality Collapse in Multimodal VAEs via Impartial Optimization. ICML 2022]

[Goyal et al., Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. CVPR 2017]

Modality Biases

Heterogeneity in information and relevance

Fairness and social biases – unimodal social biases

Finding: Image captioning models capture spurious correlations between gender and generated actions

Wrong



Baseline:

*A **man** sitting at a desk with a laptop computer.*

[Hendricks et al., Women also Snowboard: Overcoming Bias in Captioning Models. ECCV 2018]

Modality Biases

Heterogeneity in information and relevance

Fairness and social biases – unimodal social biases

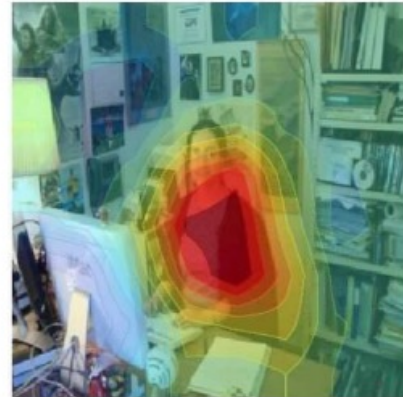
Finding: Image captioning models capture spurious correlations between gender and generated actions

Wrong



Baseline:
*A **man** sitting at a desk with a laptop computer.*

Right for the Right Reasons



Our Model:
*A **woman** sitting in front of a laptop computer.*

[Hendricks et al., Women also Snowboard: Overcoming Bias in Captioning Models. ECCV 2018]

Modality Biases

Heterogeneity in information and relevance

Fairness and social biases – unimodal social biases

Finding: Image captioning models capture spurious correlations between gender and generated actions



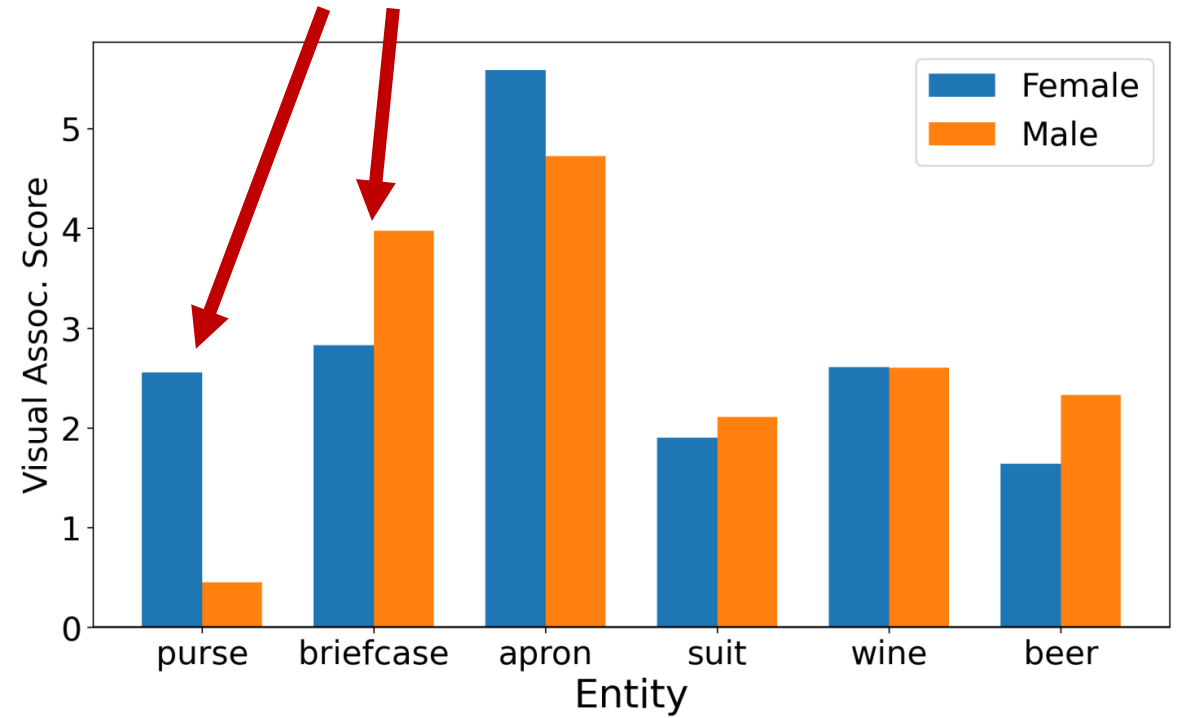
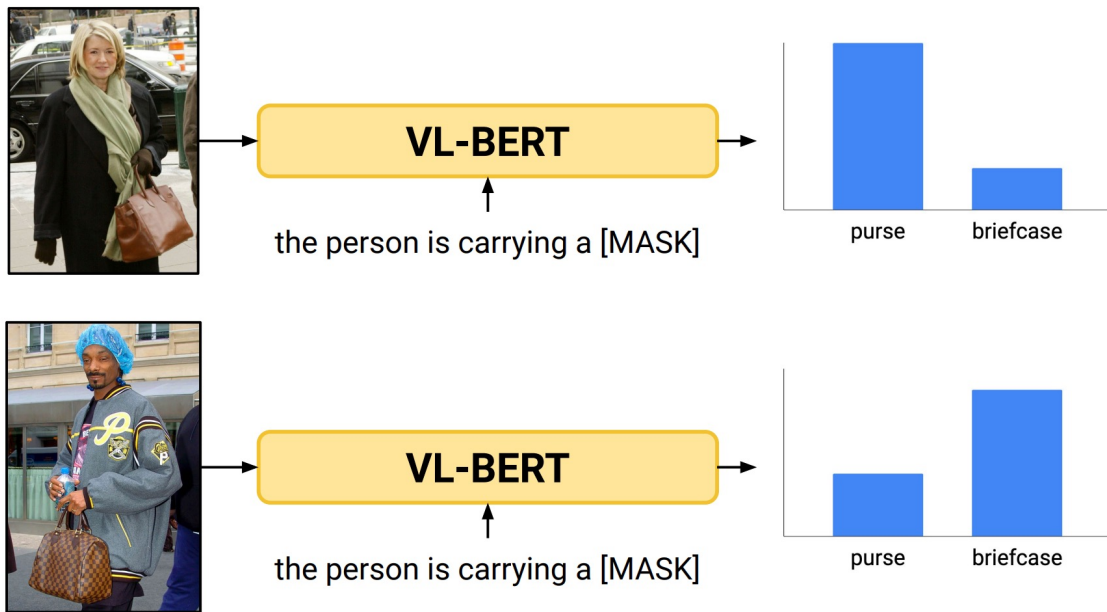
[Hendricks et al., Women also Snowboard: Overcoming Bias in Captioning Models. ECCV 2018]

Modality Biases

Heterogeneity in information and relevance

Fairness and social biases – cross-modal interactions worsen social biases

Visual information makes model more confident in reinforcing gender stereotypes



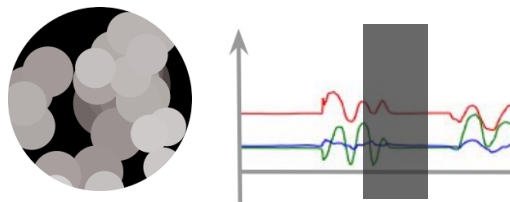
[Srinivasan and Bisk, Worst of Both Worlds: Biases Compound in Pre-trained Vision-and-Language Models. NAACL 2022]

Noise Topologies and Robustness

Heterogeneity in noise

Modality-specific robustness

noise → **nosie**



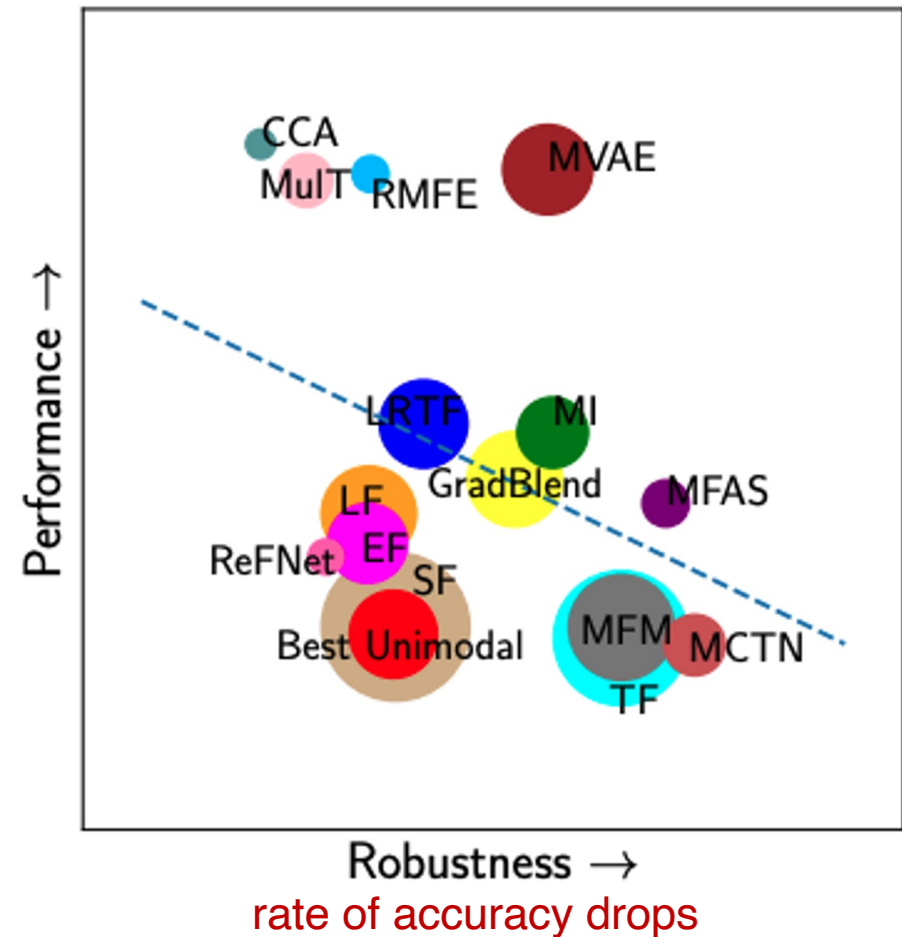
[Belinkov & Bisk, 2018; Subramaniam et al., 2009; Boyat & Joshi, 2015]

Multimodal robustness



[Zadeh et al., 2020]

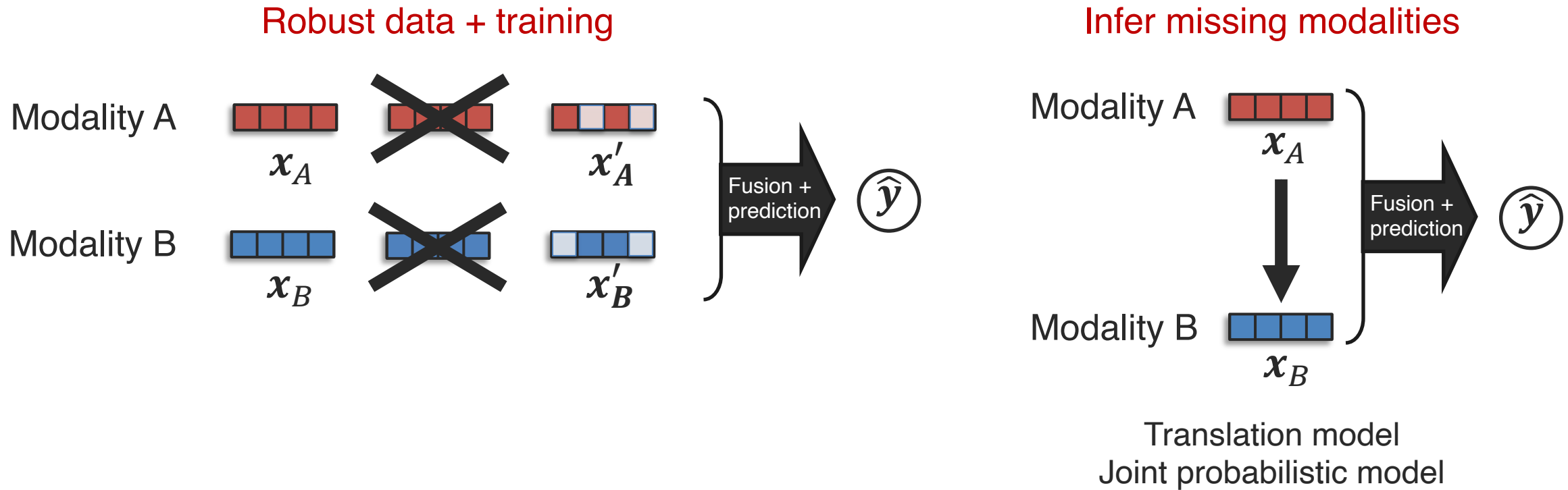
Strong tradeoffs between performance and robustness



[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

Noise Topologies and Robustness

Several approaches towards more robust models



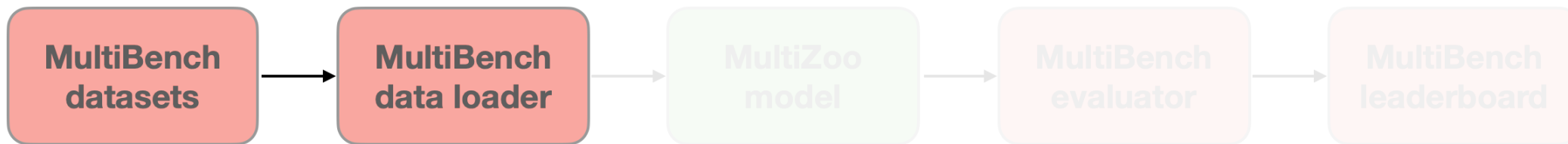
[Ngiam et al., Multimodal Deep Learning. ICML 2011]

[Srivastava and Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines. JMLR 2014]

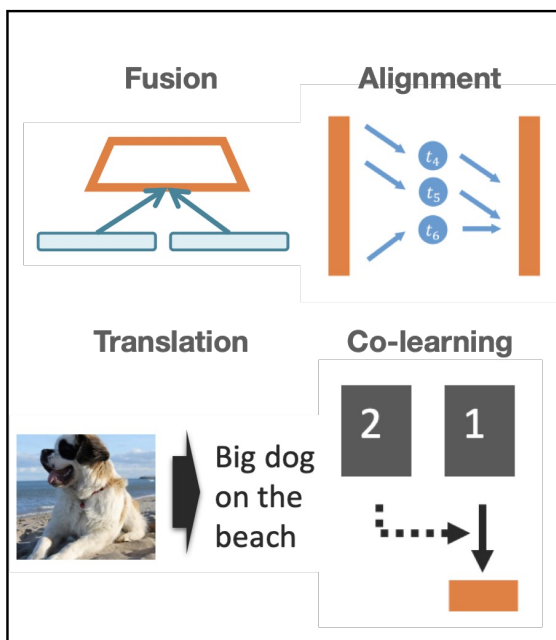
[Tran et al., Missing Modalities Imputation via Cascaded Residual Autoencoder. CVPR 2017]

[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAI 2019]

Going Beyond Language, Vision, and Audio



Challenges



Domains

Affective computing: A video frame showing a woman with the text "And he I don't think he got mad when hah I don't know maybe." and "(frustrated voice)".

Healthcare: A screenshot of a patient's medical record showing various vitals and lab results.

Robotics: A video frame showing a robotic arm with the text "Episode 100" and "21% success rate".

Finance: A line chart showing the stock price of McDonald's Corp. over time.

HCI: A screenshot of a user interface for a mobile application.

Multimedia: A grid of various images and videos from different sources.

Modalities

Language: The text "All I can say is he's a pretty average guy." with a video frame of a person speaking.

Image: A photo of a person in a red shirt.

Video: A video frame of a person in a white shirt.

Audio: A spectrogram showing the frequency components of a speech signal.

Time-series: A line chart showing a time series of data points.

Force sensors: A line graph showing force sensor data over time.

Proprioception: A 3D model of a robotic arm.

Set: A collection of various small images.

Table: A table with columns for SUBJECT_ID, Age, Sex, and Ethnicity.

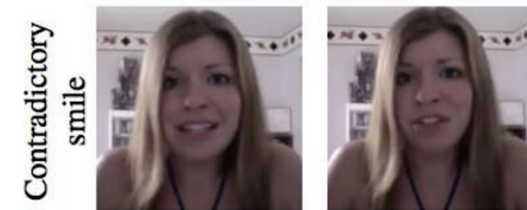
Optical flow: A video frame showing optical flow vectors overlaid on a street scene.

[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

MultiBench: Generalization to Diverse Modalities

Multimodal affect recognition

Language: *And he I don't think he got mad when hah I don't know maybe.* *Too much too fast, I mean we basically just get introduced to this character...* *All I can say is he's a pretty average guy.*



Acoustic: (frustrated voice)

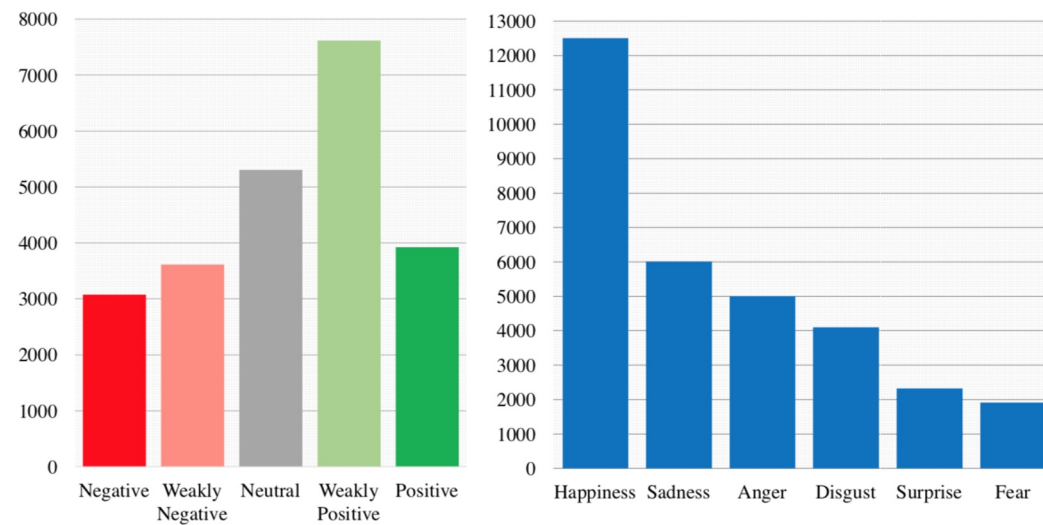
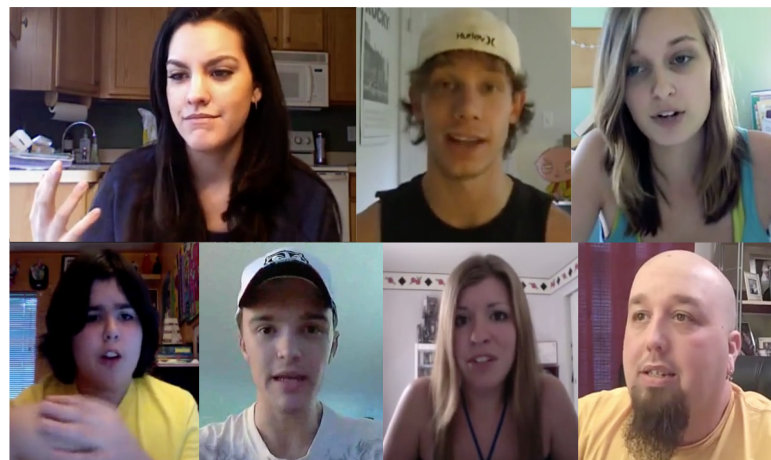
(angry voice)

(disappointed voice)

1,000 speakers

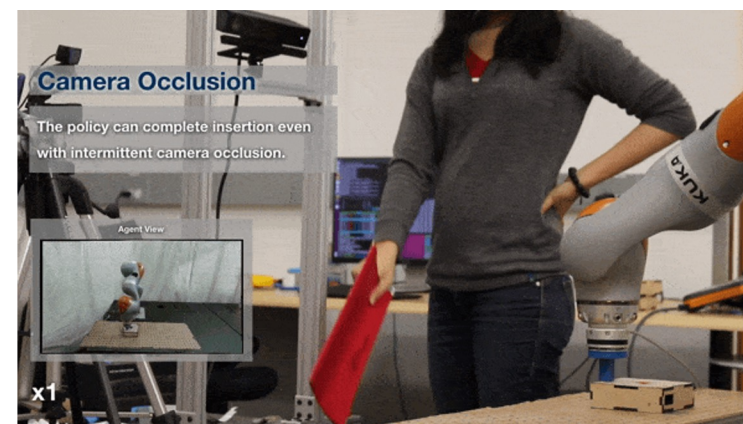
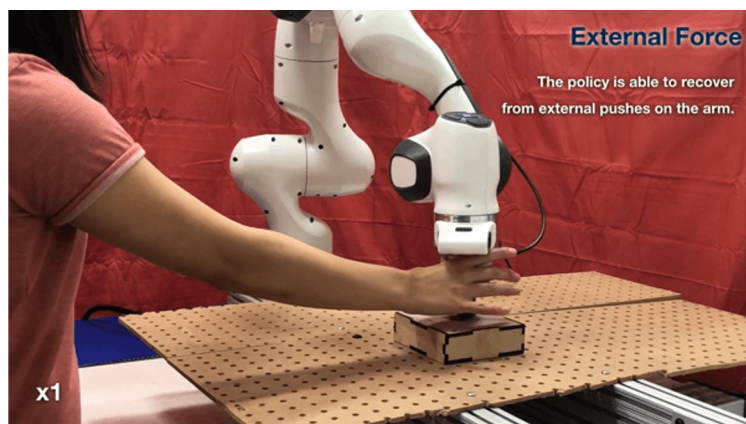
250 topics

Diverse annotations



MultiBench: Generalization to Diverse Modalities

Multisensor fusion in robotics

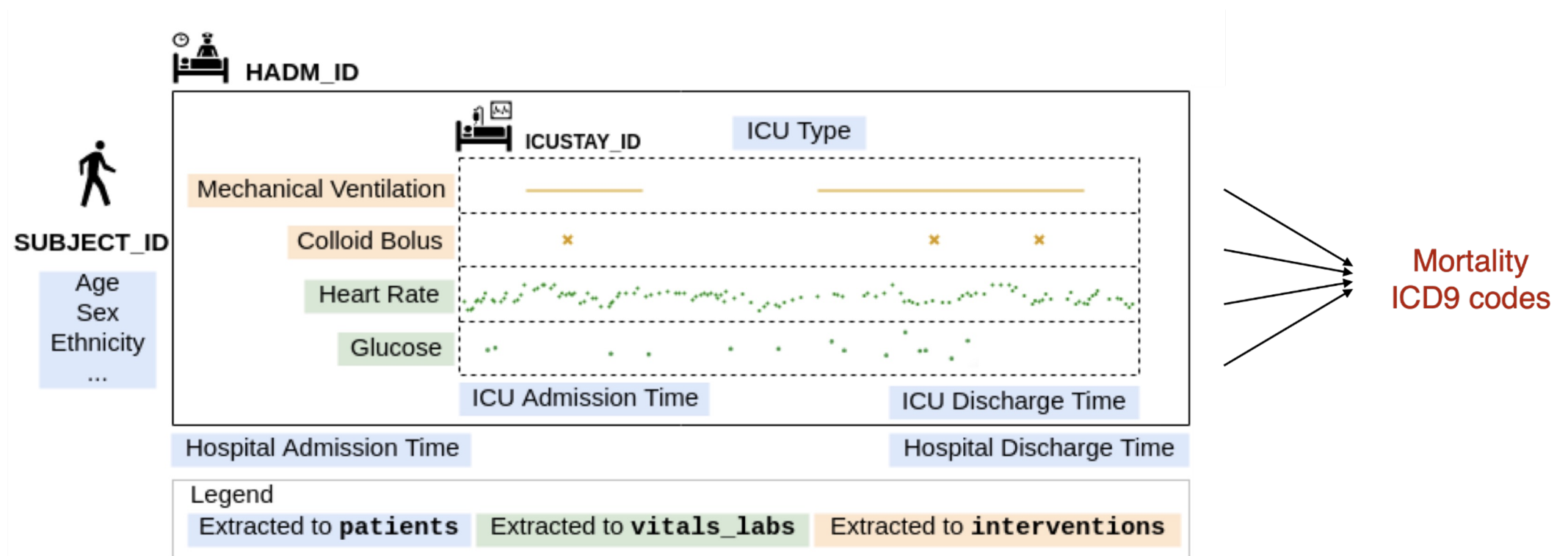


[Lee et al., ICRA 2019]

[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

MultiBench: Generalization to Diverse Modalities

Multimodal learning in healthcare



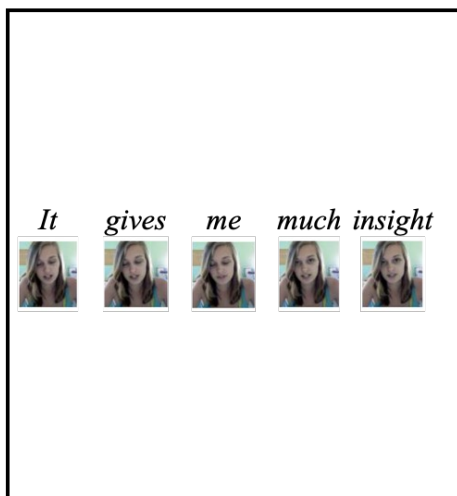
[Johnson et al., Nature 2016]

[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

MultiBench: Generalization to Diverse Modalities

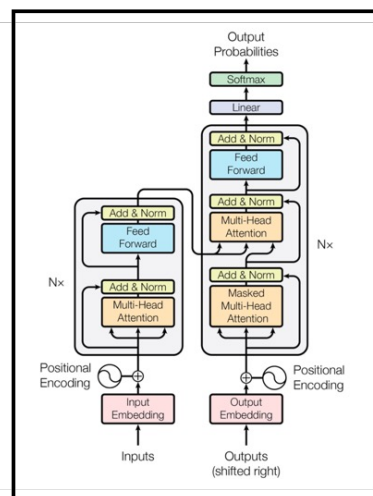


Data preprocessing



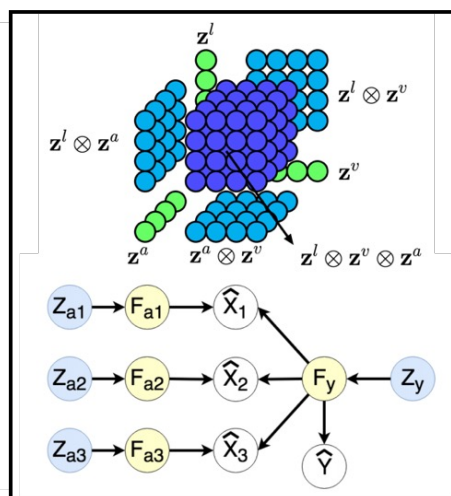
[Chen et al., 2017]

Unimodal models



[Vaswani et al., 2017]

Fusion paradigms



[Zadeh et al., 2017]

Optimization objectives

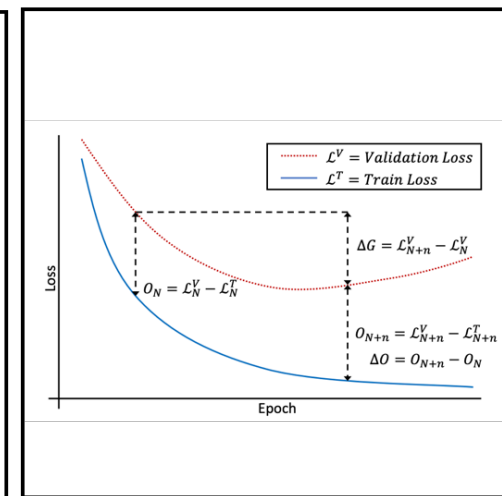
$$\mathcal{L}_{\text{sup}} = -\log p(y|\mathbf{x}_1, \mathbf{x}_2)$$

$$\mathcal{L}_{\text{CCA}} = -\text{corr}(g_1(\mathbf{z}_1), g_2(\mathbf{z}_2))$$

$$\mathcal{L}_{\text{rec}} = \|\mathbf{g}_1(\mathbf{z}_{\text{mm}}) - \mathbf{x}_1\|_2 + \|\mathbf{g}_2(\mathbf{z}_{\text{mm}}) - \mathbf{x}_2\|_2$$

[Sankaran et al., 2021]

Training procedures



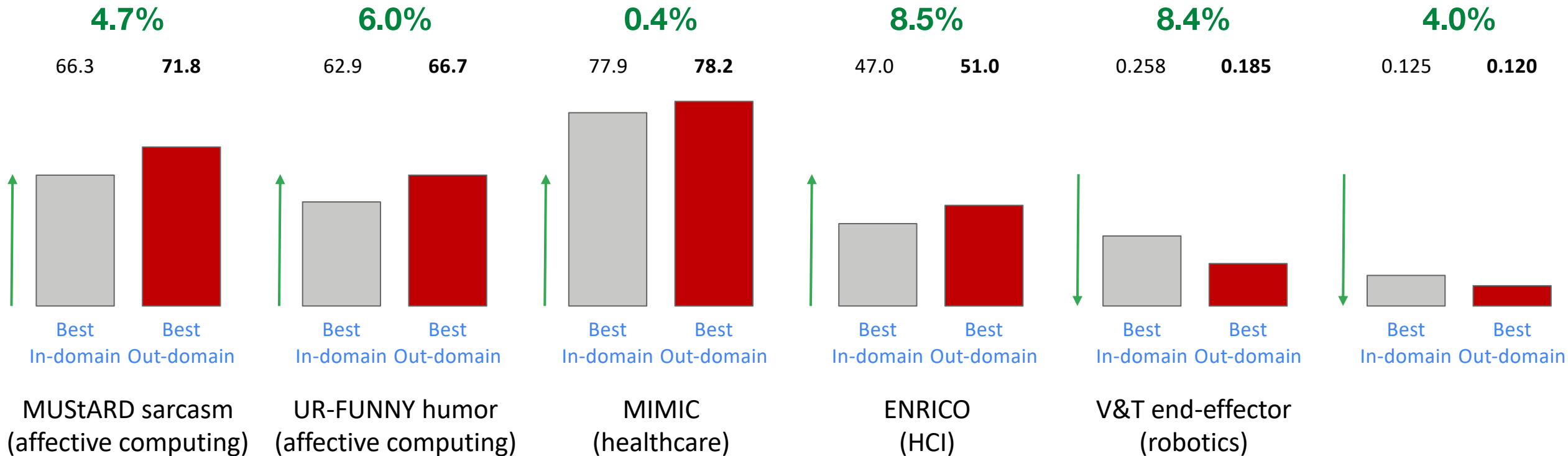
[Wang et al., 2020]

20 recent models, but composable to up to >200 combinations!

[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

MultiBench: Generalization to Diverse Modalities

Benefits of standardization

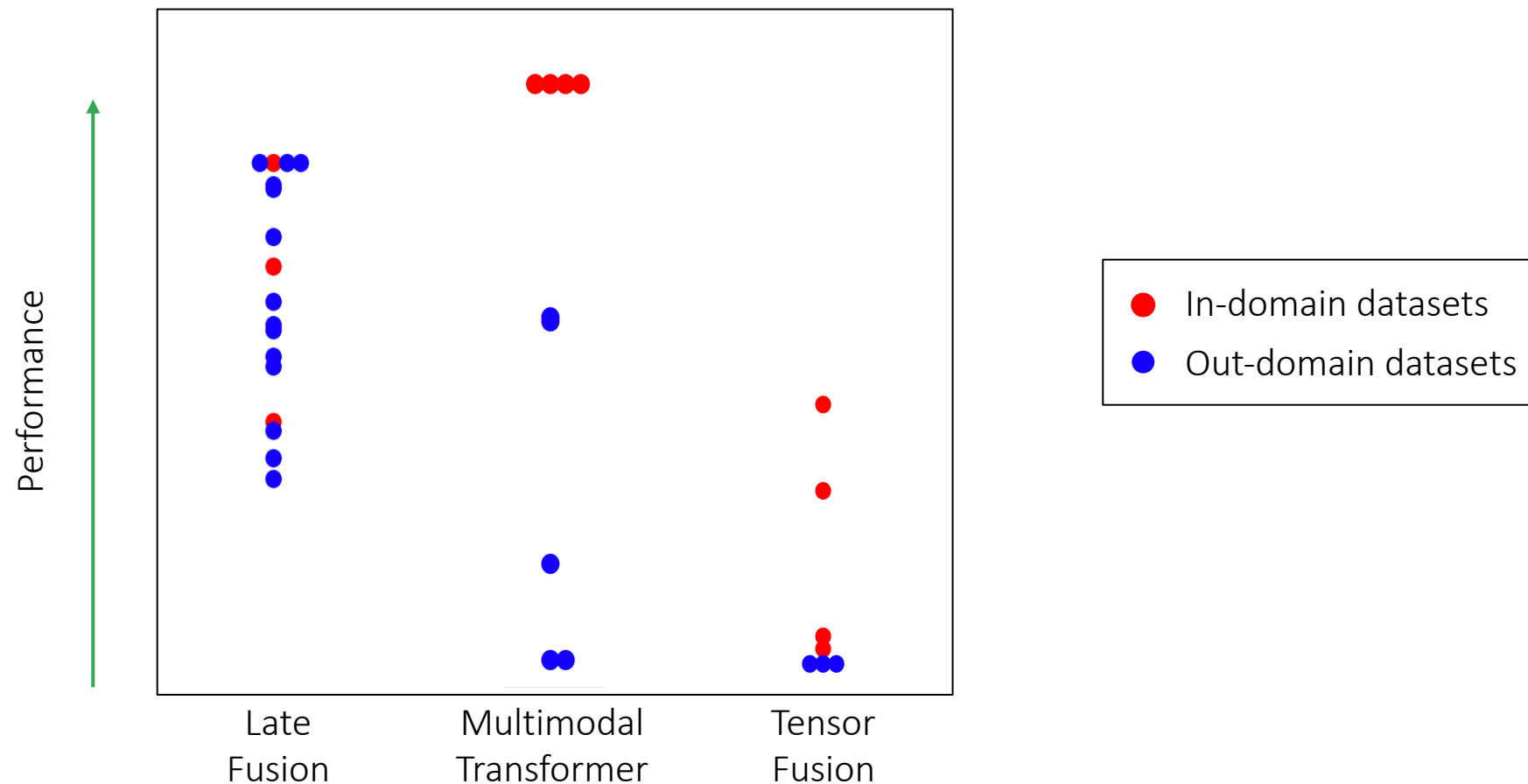


Simply applying methods in other areas improves performance on 9/15 datasets

[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

MultiBench: Generalization to Diverse Modalities

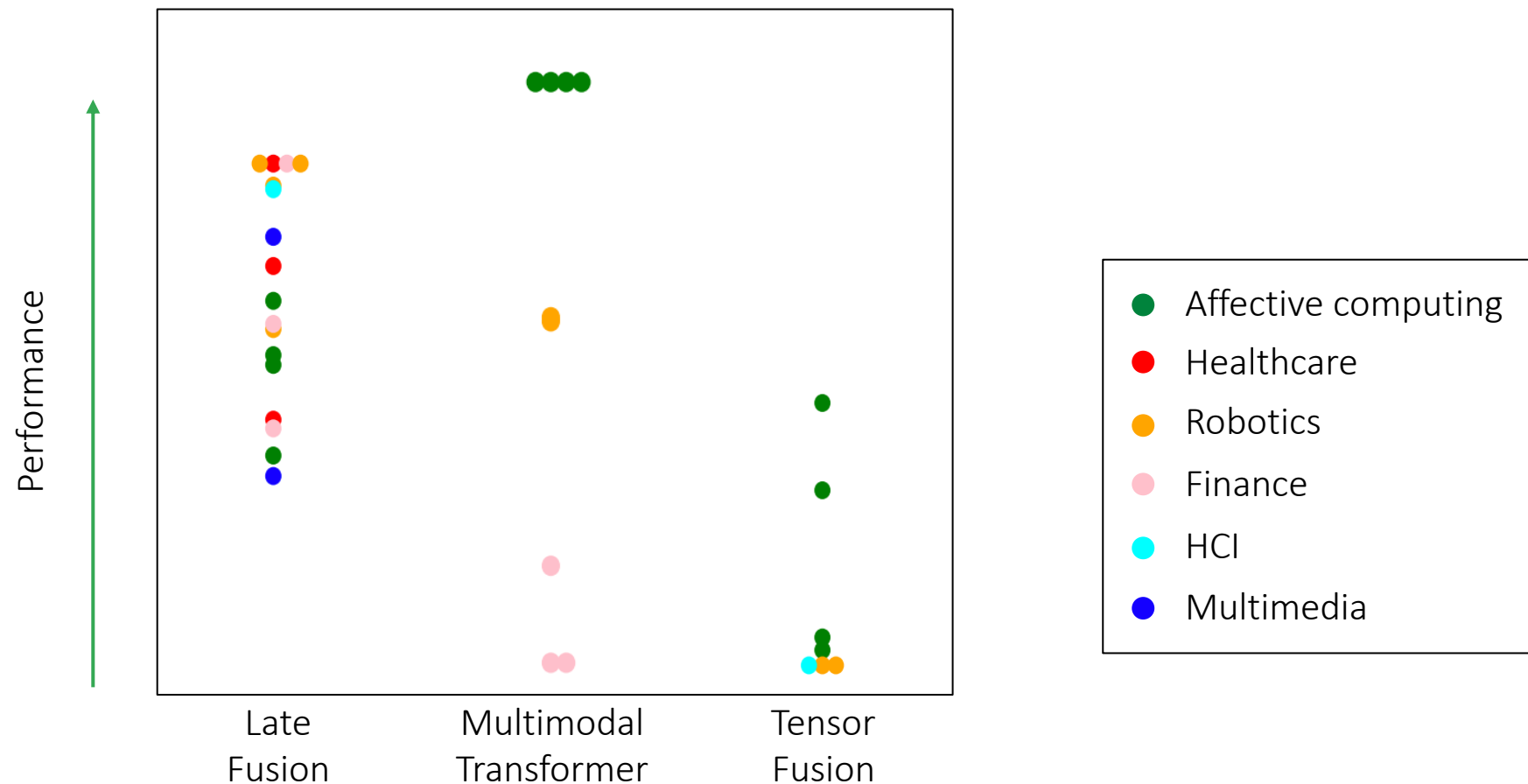
Methods struggle to perform outside of their own domain



[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

MultiBench: Generalization to Diverse Modalities

Generalization across modalities and tasks is difficult!

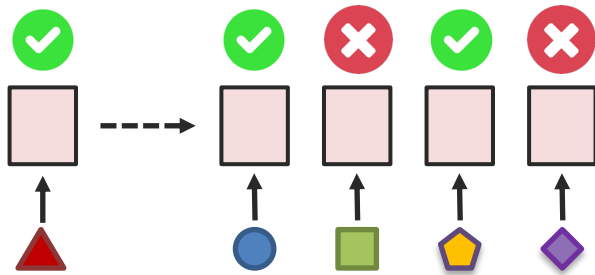


[Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021]

Quantifying Modality Heterogeneity

Information transfer, transfer learning perspective

1a. Estimate modality heterogeneity via transfer



Implicitly captures these:

① Element representation

③ Structure

⑤ Noise

② Element distribution

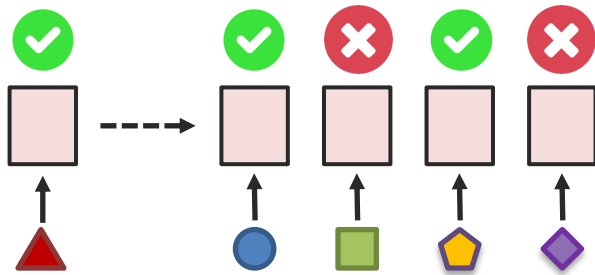
④ Information

⑥ Relevance

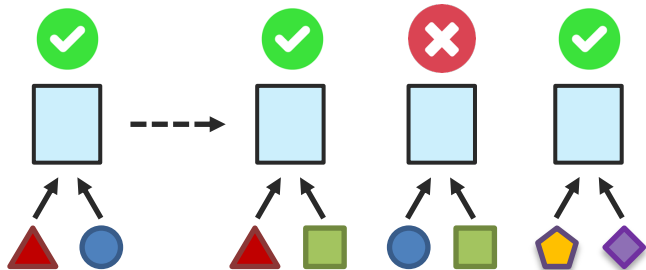
Quantifying Modality Heterogeneity

Information transfer, transfer learning perspective

1a. Estimate modality heterogeneity via transfer



1b. Estimate interaction heterogeneity via transfer



2a. Compute modality heterogeneity matrix

	0				
	1	0			
	3	2	0		
	1	2	3	0	
	5	4	6	3	0











2b. Compute interaction heterogeneity matrix

	0			
	1	0		
	3	2	0	
	1	2	4	0

















Quantifying Modality Heterogeneity

Information transfer, transfer learning perspective

2a. Compute modality heterogeneity matrix

					
	0				
	1	0			
	3	2	0		
	1	2	3	0	
	5	4	6	3	0

2b. Compute interaction heterogeneity matrix

								
		0						
		1	0					
		3	2	0				
		1	2	4	0			

3. Determine parameter clustering

$$\mathbb{U}_1 = \{U_1, U_2, U_4\} \quad \mathbb{C}_1 = \{C_{12}, C_{13}, C_{45}\}$$

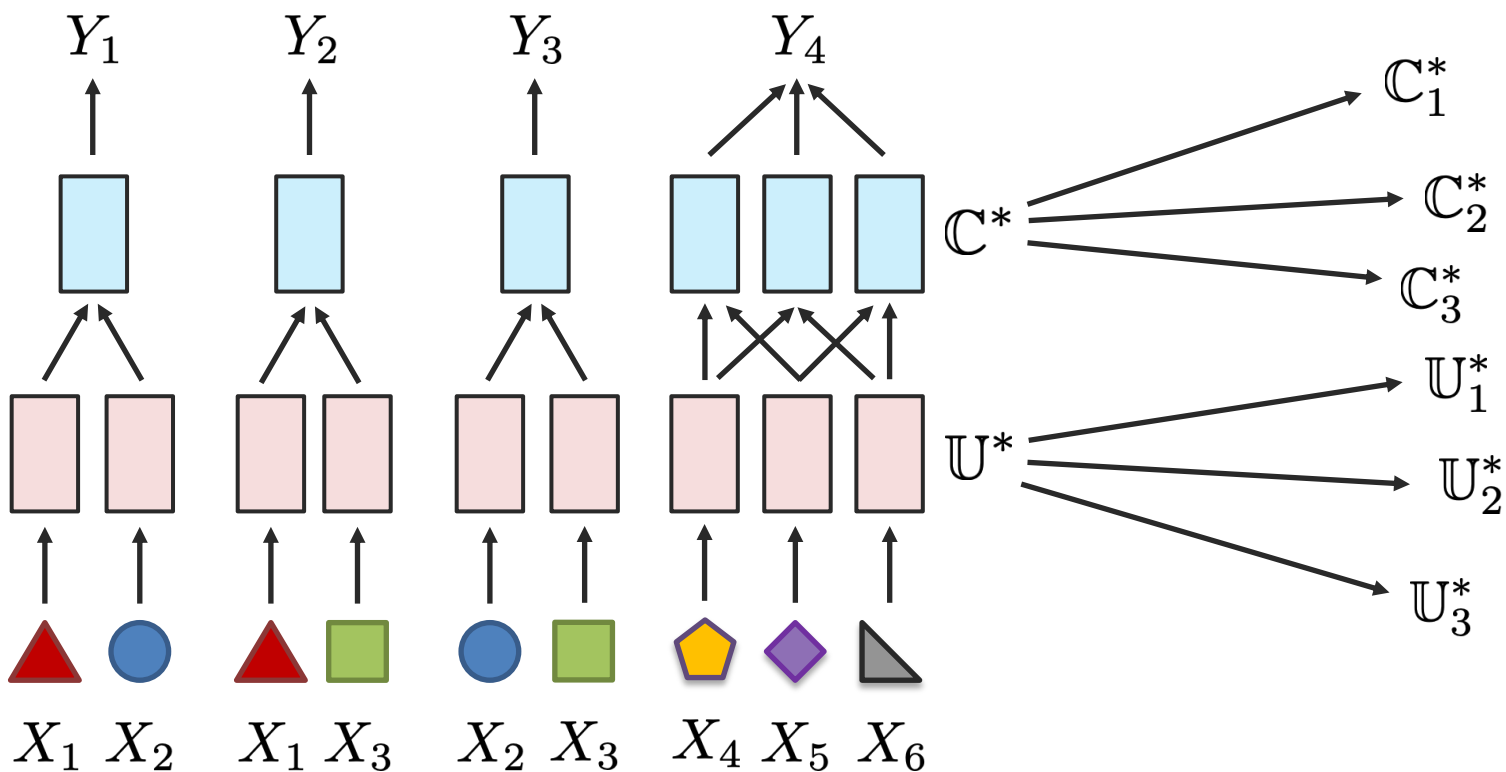
$$\mathbb{U}_2 = \{U_3\} \quad \mathbb{C}_2 = \{C_{23}\}$$

$$\mathbb{U}_3 = \{U_5\}$$

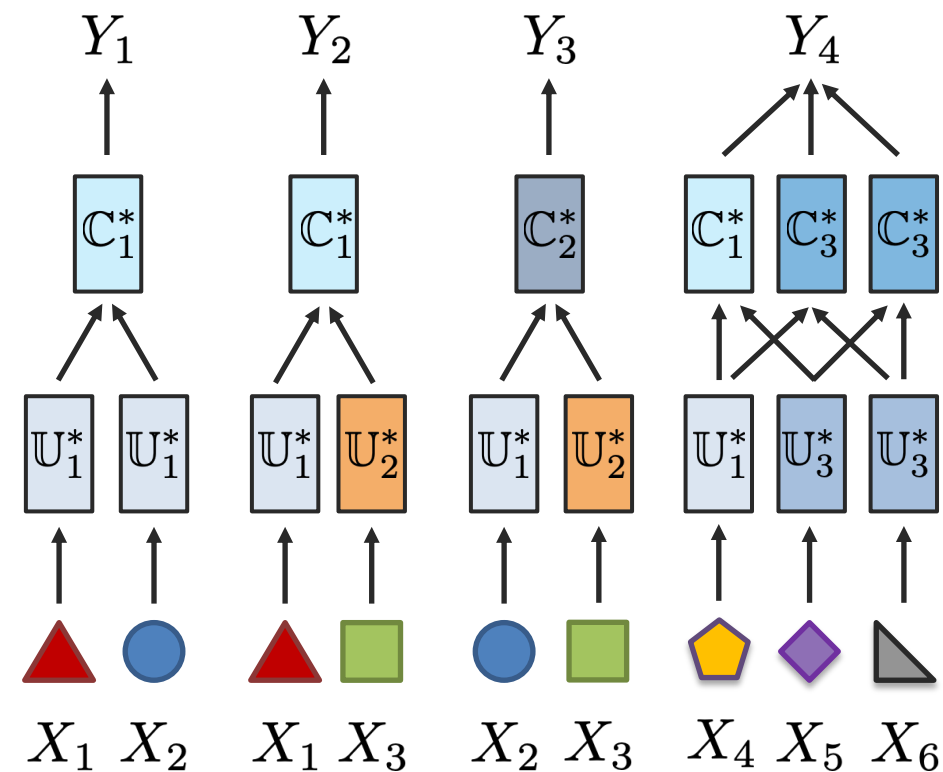
Quantifying Modality Heterogeneity

Information transfer, transfer learning perspective

1. Homogeneous Pre-training



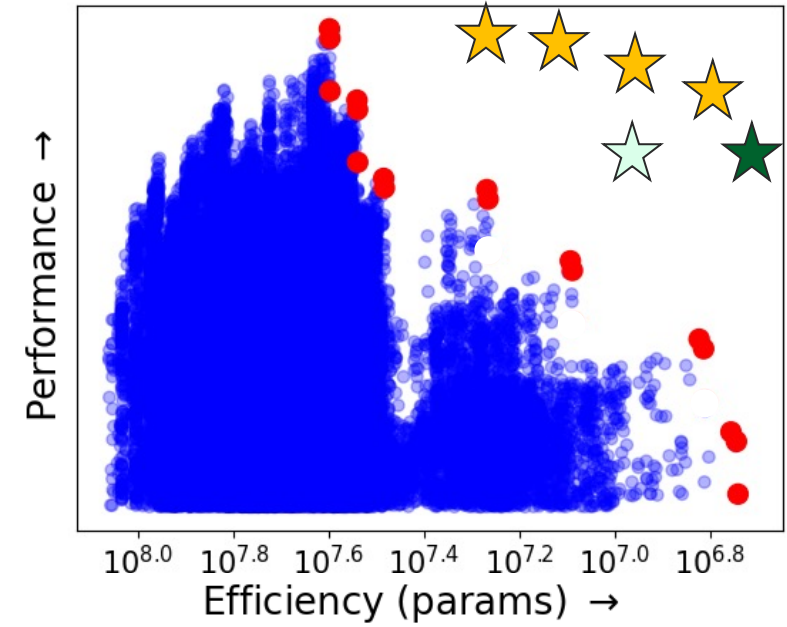
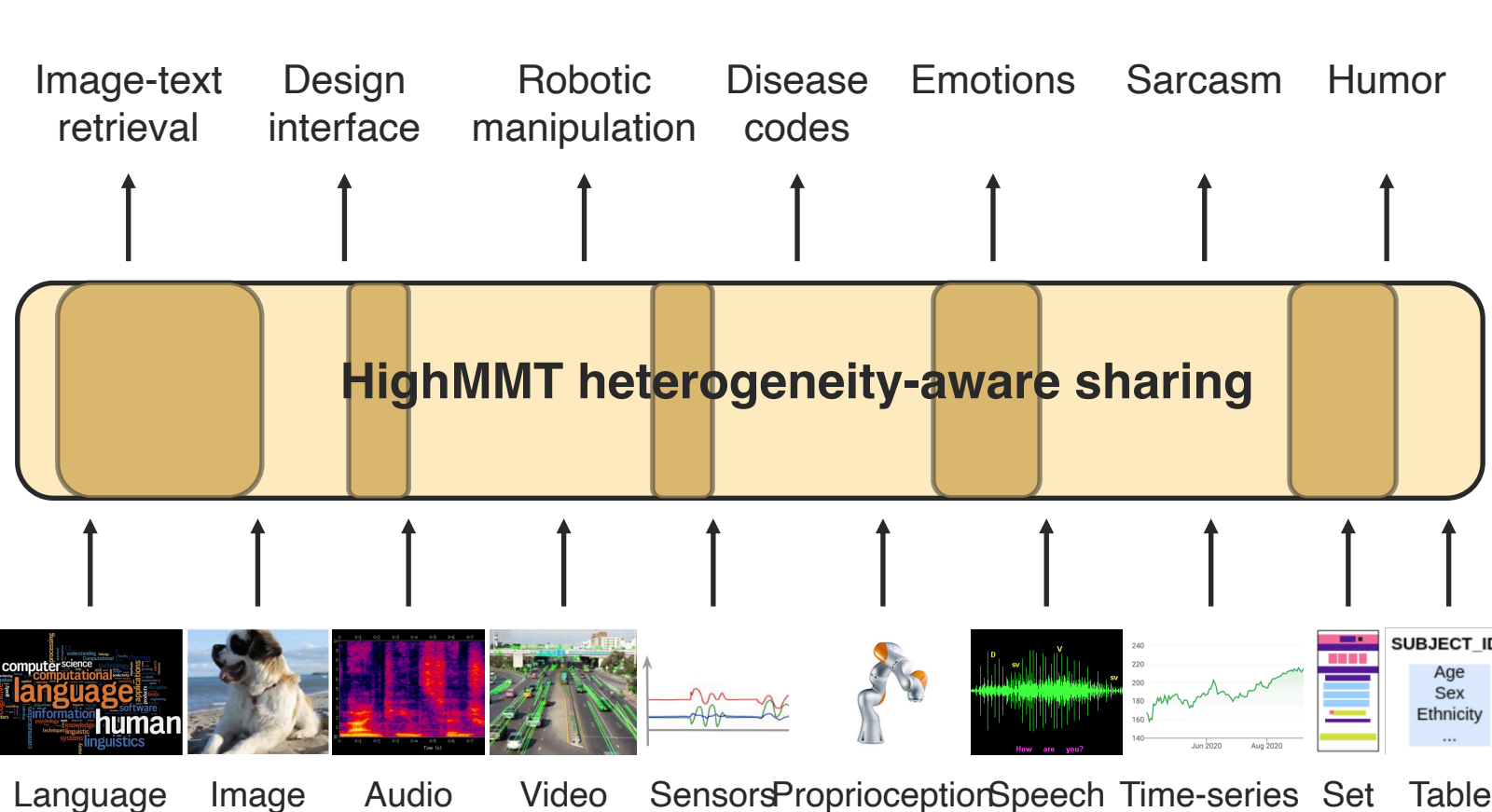
2. Heterogeneity-aware Fine-tuning



[Liang et al., HighMMT: Quantifying Modality and Interaction Heterogeneity for High-Modality Representation Learning. arXiv 2022]

Quantifying Modality Heterogeneity

HighMMT heterogeneity-aware: estimate heterogeneity to determine parameter sharing



- All model combinations (>10,000)
- Pareto front
- HighMMT single-task
- HighMMT multitask
- **HighMMT heterogeneity-aware**

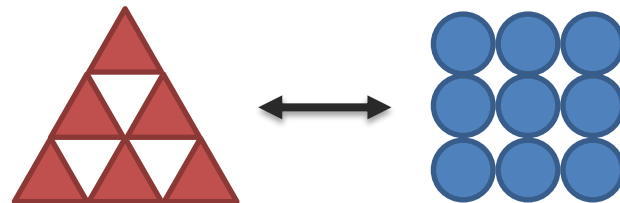
[Liang et al., HighMMT: Towards Modality and Task Generalization for High-Modality Representation Learning. arXiv 2022]

Challenges: Quantifying Heterogeneity

Open
challenges

Open challenges:

- Quantifying and modeling: chicken and egg problem.
- Bottom-up vs top-down, data-driven vs hypothesis-driven.
- Noisy and missing modalities.
- New and understudied modalities.
- Large number of modalities.
- Cases where its unclear which modalities are useful.



Sub-Challenge 6b: Cross-modal Connections

Connected: Shared information that relates modalities



Statistical



Association

Dependency



e.g., correlation,
co-occurrence



e.g., causal,
temporal

Semantic



Correspondence

Relationship



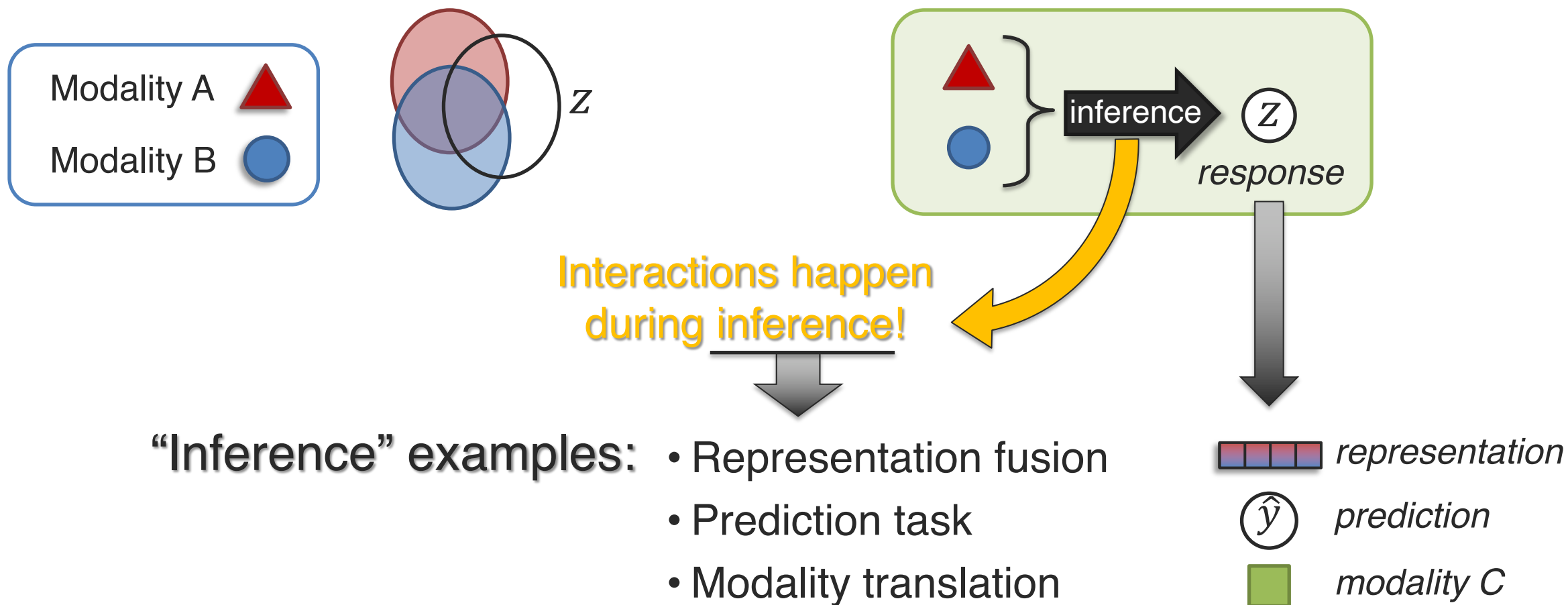
e.g., grounding



e.g., function

Sub-Challenge 6b: Cross-modal Interactions

Interacting: process affecting each modality, creating new response

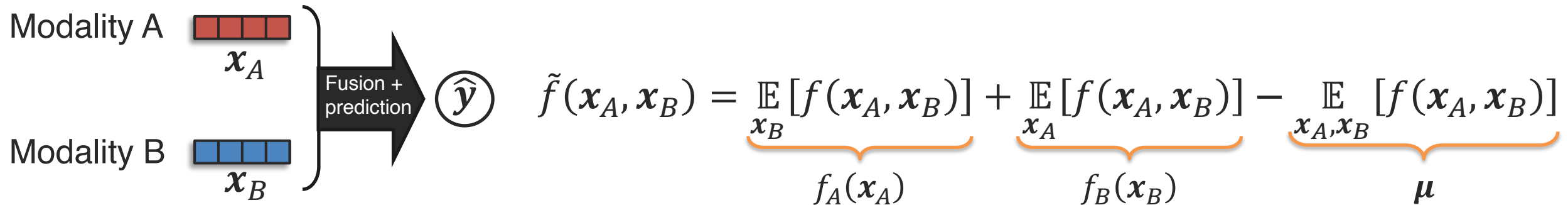


Quantifying Cross-modal Interactions

Identifying overall presence of cross-modal interactions

Statistical non-additive interactions [Friedman & Popescu, 2008, Sorokina et al., 2008]

f exhibits interactions between 2 features x_A and x_B iff f cannot be decomposed into a sum of unimodal subfunctions f_A, f_B such that $f(x_A, x_B) = f_A(x_A) + f_B(x_B)$.



If the additive projection $\tilde{f}(x_A, x_B)$ is equal to nonlinear fusion $f(x_A, x_B)$ then the non-additive interactions are not modeled.

μ measures **overall quantity** of cross-modal interactions on a trained model + dataset.

[Hessel and Lee, Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think!, EMNLP 2020]

Quantifying Cross-modal Interactions

Identifying individual cross-modal interactions

Statistical non-additive interactions [Friedman & Popescu, 2008, Sorokina et al., 2008]

f exhibits interactions between 2 features x_A and x_B iff f cannot be decomposed into a sum of unimodal subfunctions f_A, f_B such that $f(x_A, x_B) = f_A(x_A) + f_B(x_B)$.

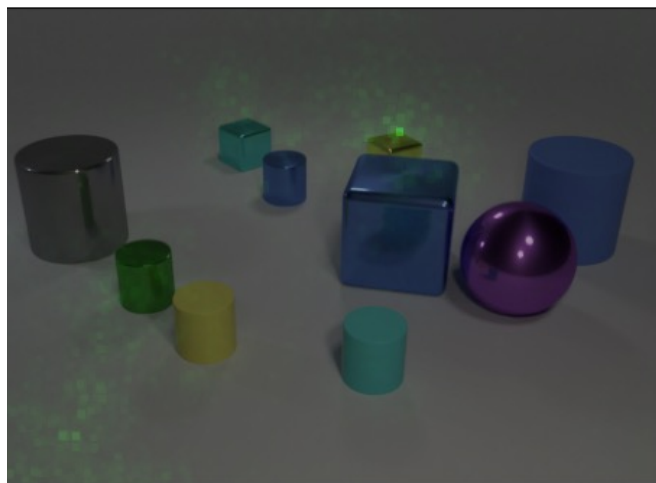
f exhibits interactions between 2 features x_A and x_B iff $\frac{\partial^2 f}{\partial x_A \partial x_B} > 0$.

Natural second-order extension of gradient-based approaches!

Quantifying Cross-modal Interactions

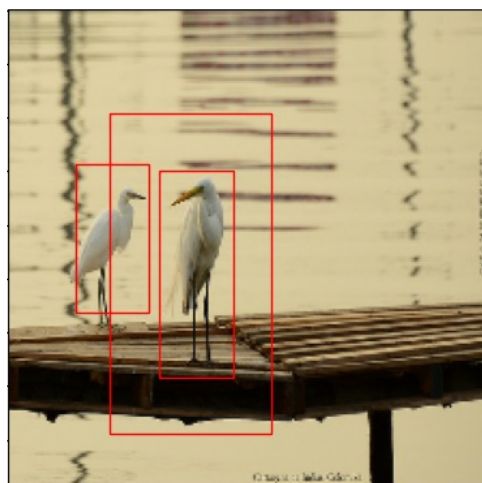
Identifying individual cross-modal interactions

CLEVR



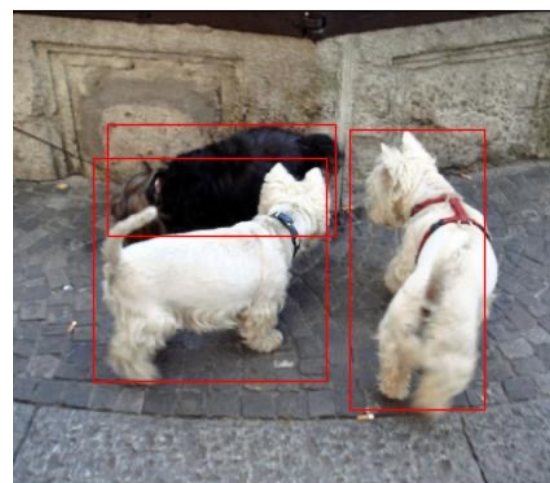
The other small shiny thing that is the same shape as the **tiny yellow shiny object** is what color?

VQA 2.0



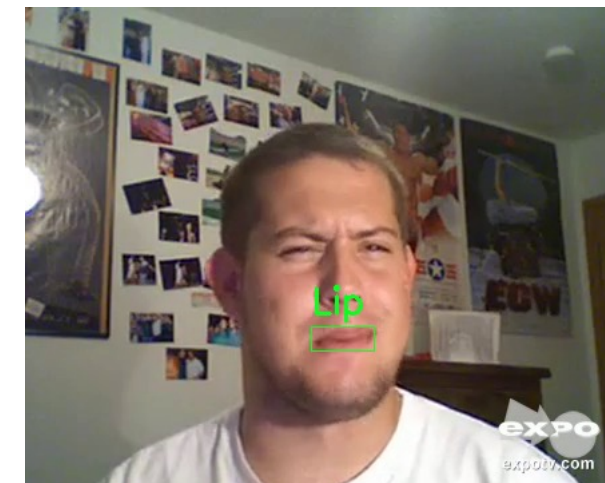
How many **birds**?

Flickr-30k



Three small dogs, two white and one black and white, on a sidewalk.

CMU-MOSEI



Why am I spending my money watching this? (**sigh**) I think I was more **sad**...

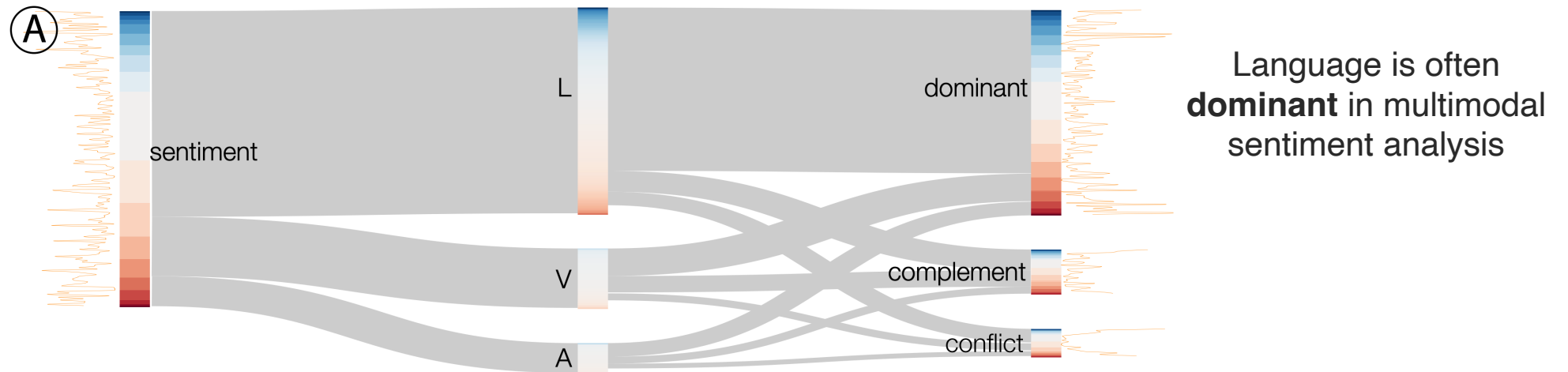
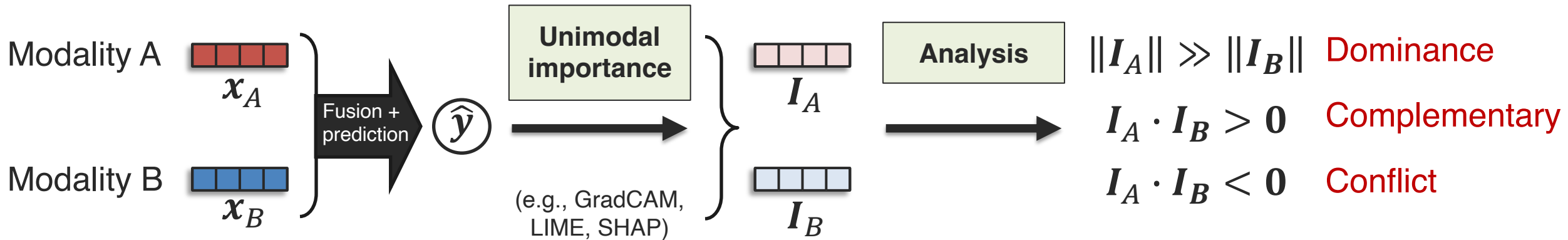
Correspondence

Relationships

[Liang et al., MultiViz: An Analysis Benchmark for Visualizing and Understanding Multimodal Models. arXiv 2022]

Quantifying Cross-modal Interactions

Classification of cross-modal interactions

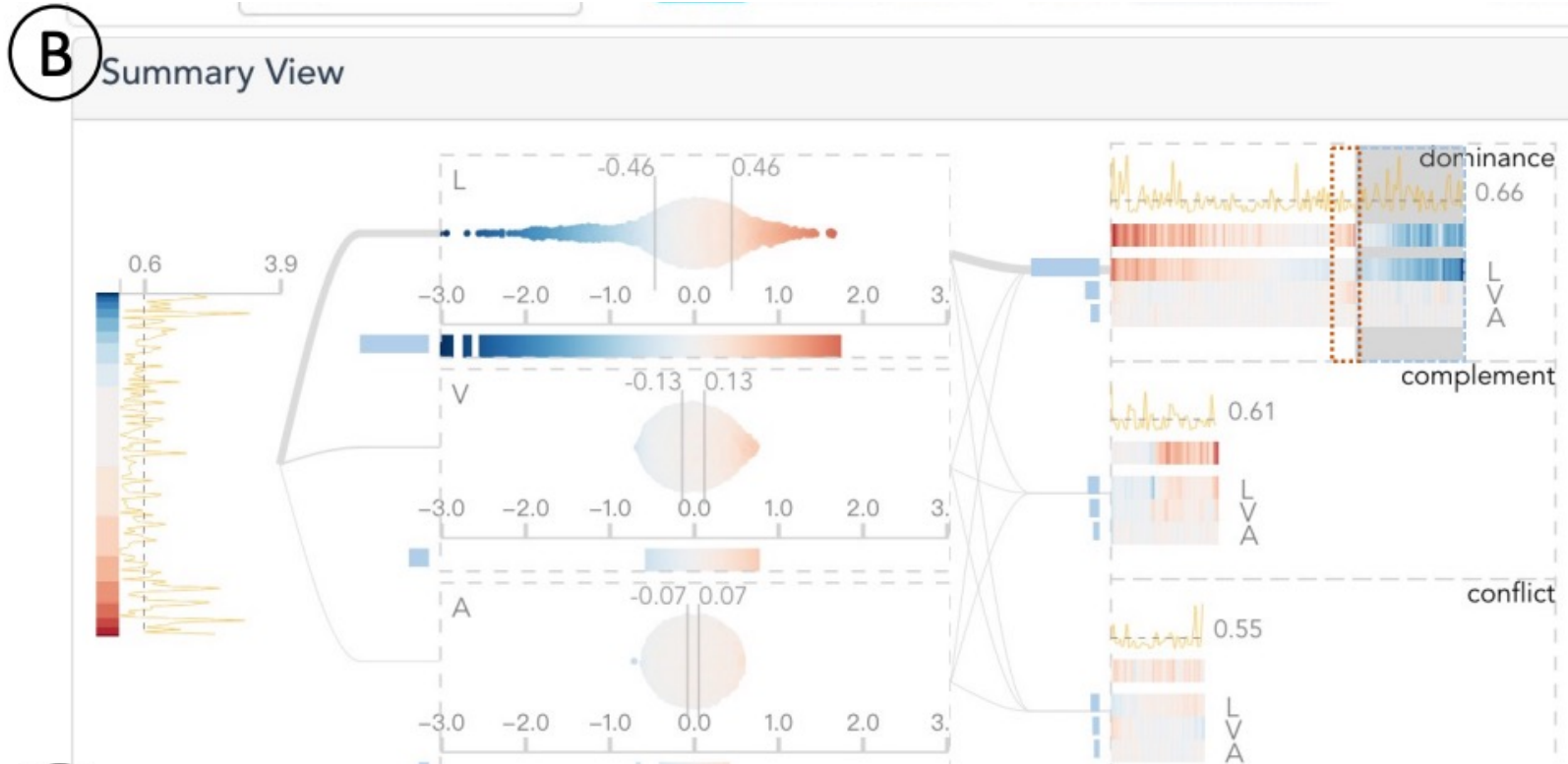


[Wang et al., M2Lens: Visualizing and Explaining Multimodal Models for Sentiment Analysis. IEEE Trans Visualization and Computer Graphics 2021]

Quantifying Cross-modal Interactions

Visualization website

See interactive website: <https://andy-xingbowang.com/m2lens/>

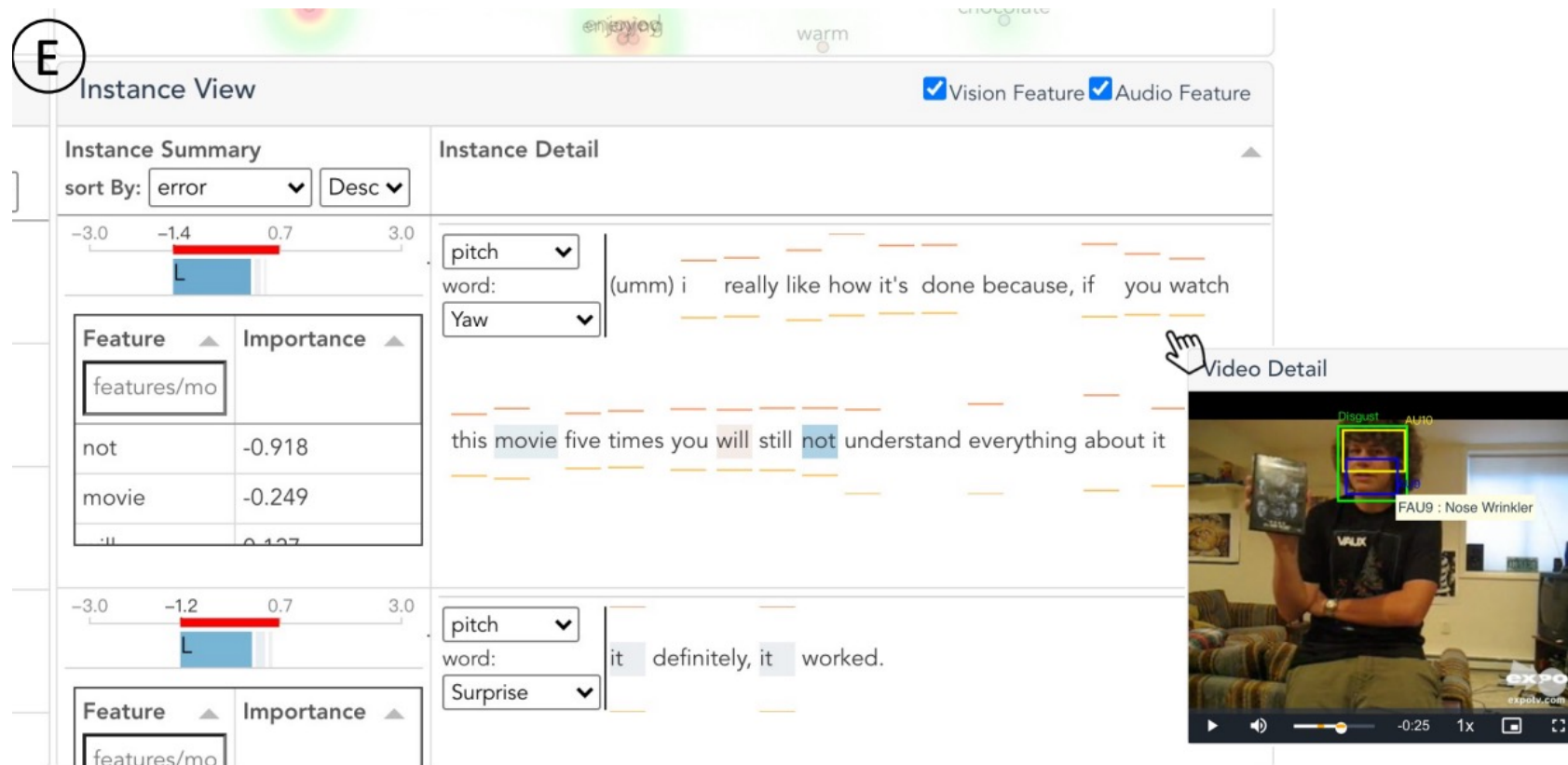


Summary of cross-modal interactions across entire dataset.

Quantifying Cross-modal Interactions

Visualization website

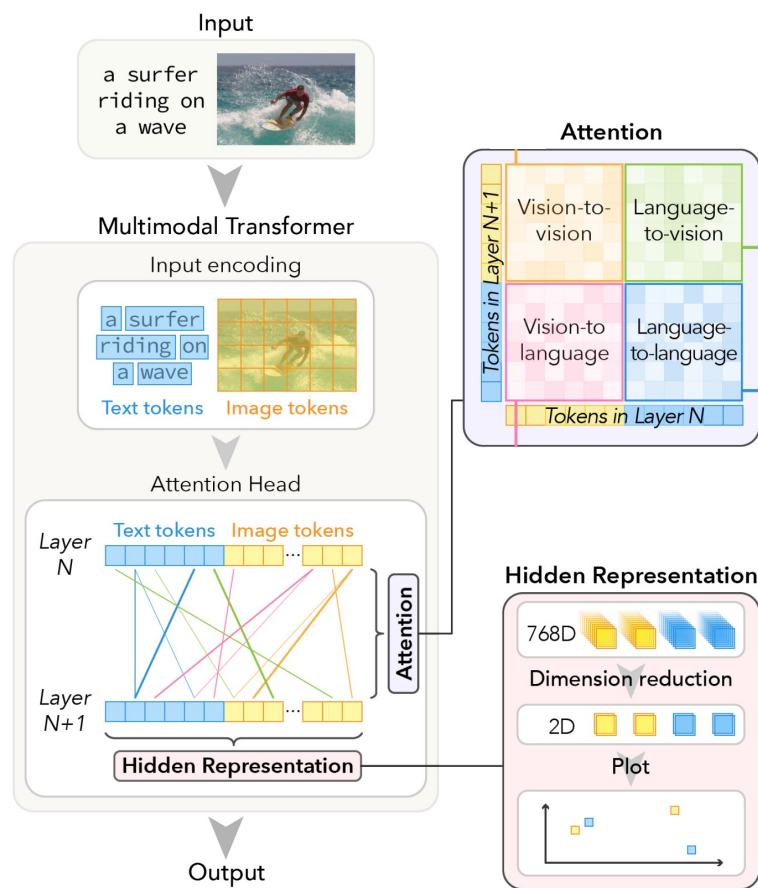
See interactive website: <https://andy-xingbowang.com/m2lens/>



Summary of cross-modal interactions in a single instance.

Quantifying Cross-modal Interactions

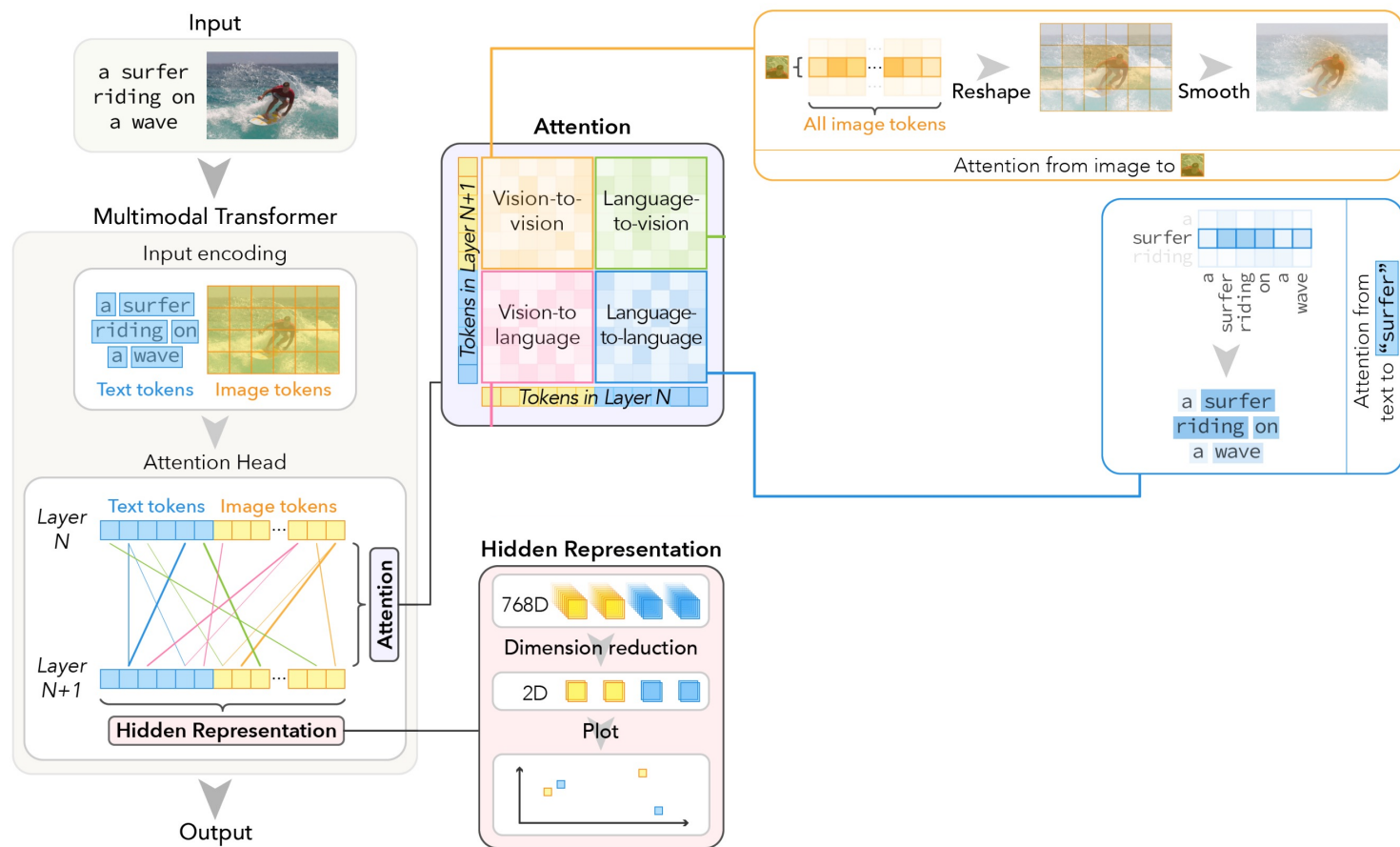
Visualizing multimodal transformers See interactive website: <https://github.com/IntelLabs/VL-InterpreT>



[Aflalo et al., VL-InterpreT: An Interactive Visualization Tool for Interpreting Vision-Language Transformers. CVPR 2022]

Quantifying Cross-modal Interactions

Visualizing multimodal transformers See interactive website: <https://github.com/IntelLabs/VL-InterpreT>



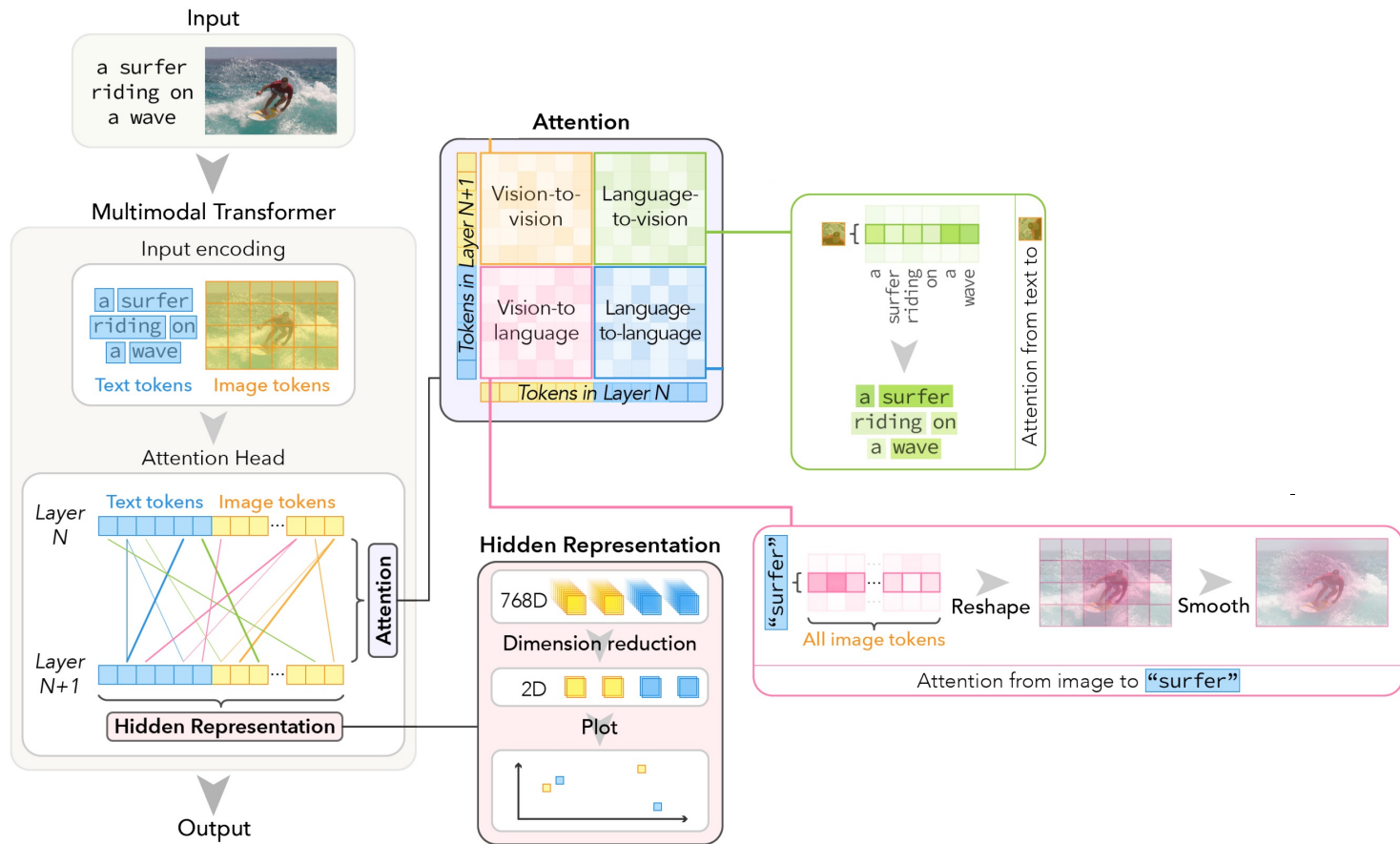
Unimodal image importance

Unimodal text importance

[Aflalo et al., VL-InterpreT: An Interactive Visualization Tool for Interpreting Vision-Language Transformers. CVPR 2022]

Quantifying Cross-modal Interactions

Visualizing multimodal transformers See interactive website: <https://github.com/IntelLabs/VL-InterpreT>



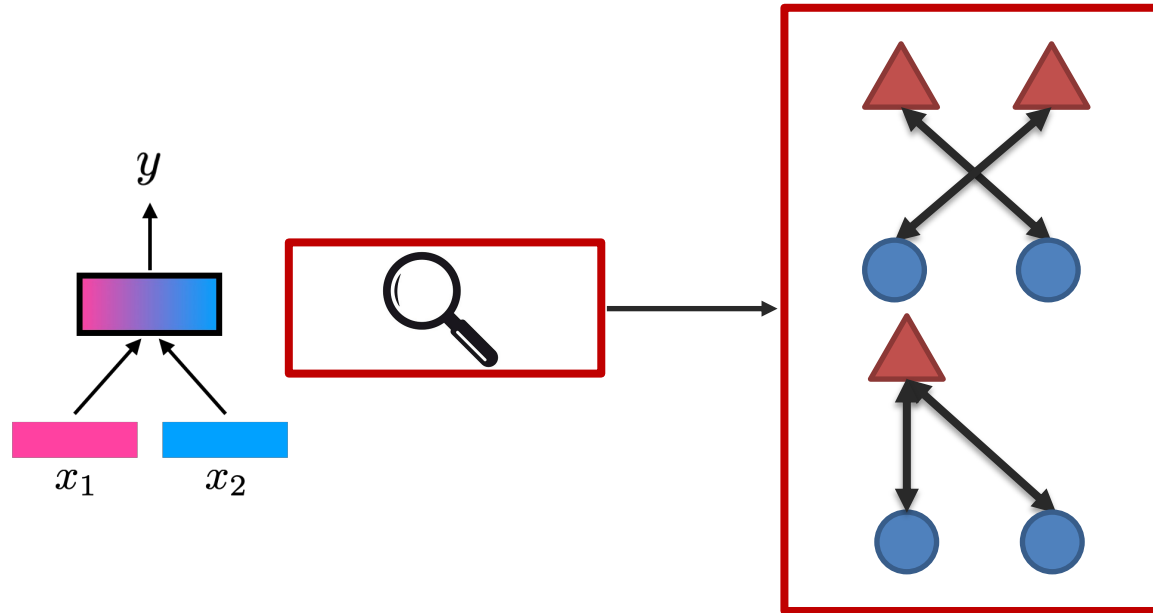
Correspondence and complementary interactions

[Aflalo et al., VL-InterpreT: An Interactive Visualization Tool for Interpreting Vision-Language Transformers. CVPR 2022]

Evaluating Interpretability

How can we evaluate the success of interpreting cross-modal interactions?

Problem: real-world datasets and models do not have cross-modal interactions annotated!



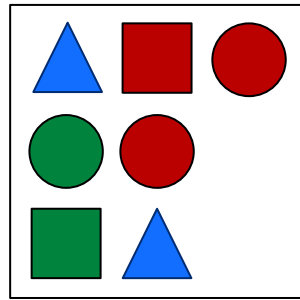
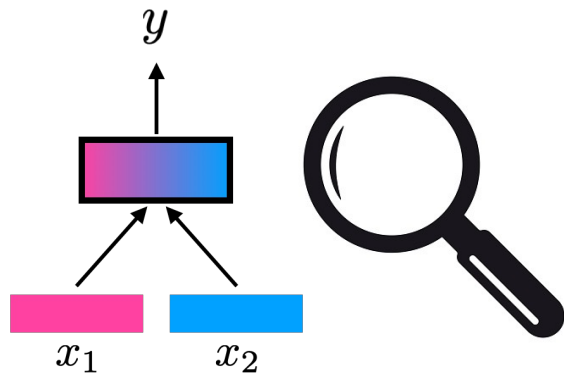
[Liang et al., MultiViz: A Framework for Visualizing and Understanding Multimodal Models. arXiv 2022]

Evaluating Interpretability: A Multi-stage Approach

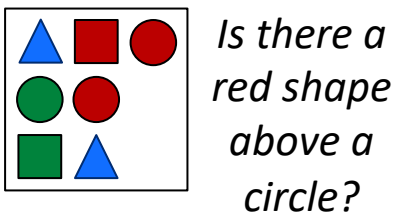
Unimodal importance: Does the model correctly identify keywords in the question?

Yes!

1. Unimodal importance



*Is there a **red shape** above a **circle**?*



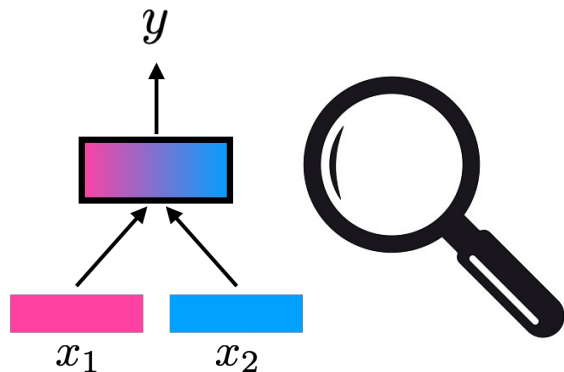
*Is there a **red shape** above a **circle**?*

[Liang et al., MultiViz: An Analysis Benchmark for Visualizing and Understanding Multimodal Models. arXiv 2022]

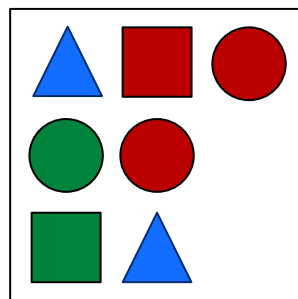
Evaluating Interpretability: A Multi-stage Approach

Cross-modal interactions: Does the model correctly relate the question with the image?

Yes!

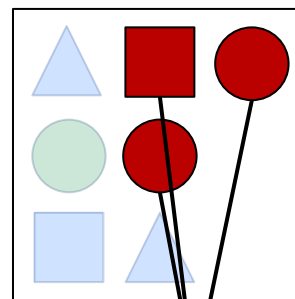


1. Unimodal importance

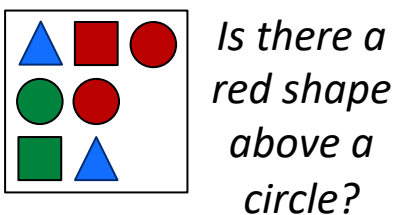
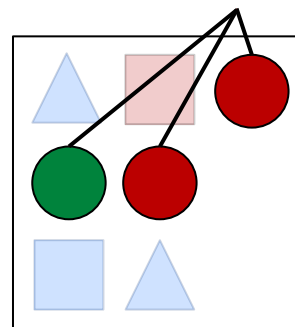


Is there a **red shape** above a **circle**?

2. Cross-modal interactions



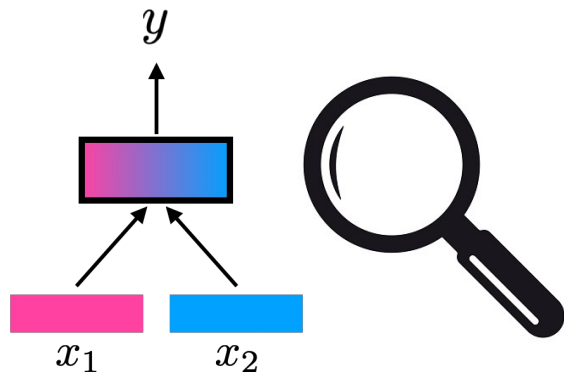
Is there a **red shape** above a **circle**?



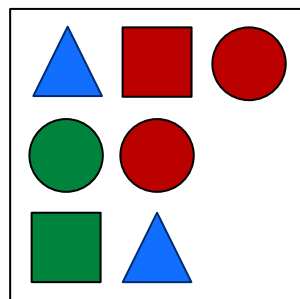
Evaluating Interpretability: A Multi-stage Approach

Multimodal representations: Does the model consistently assign concepts to features?

Yes!

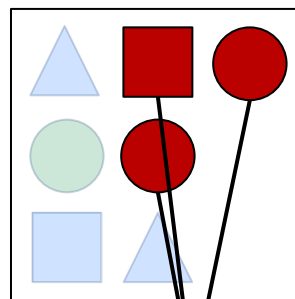


1. Unimodal importance

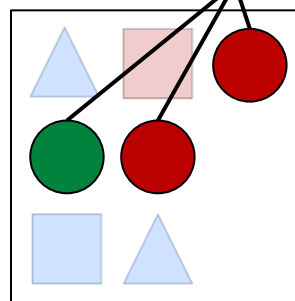


Is there a **red shape** above a circle?

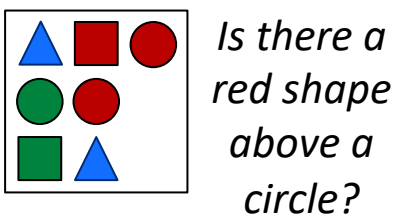
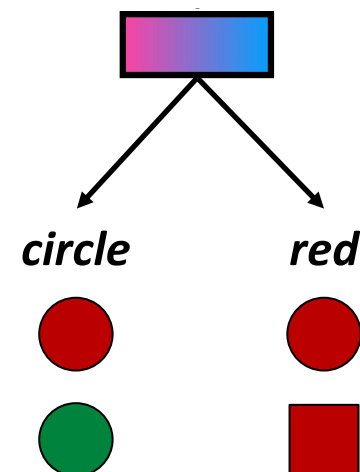
2. Cross-modal interactions



Is there a **red shape** above a circle?



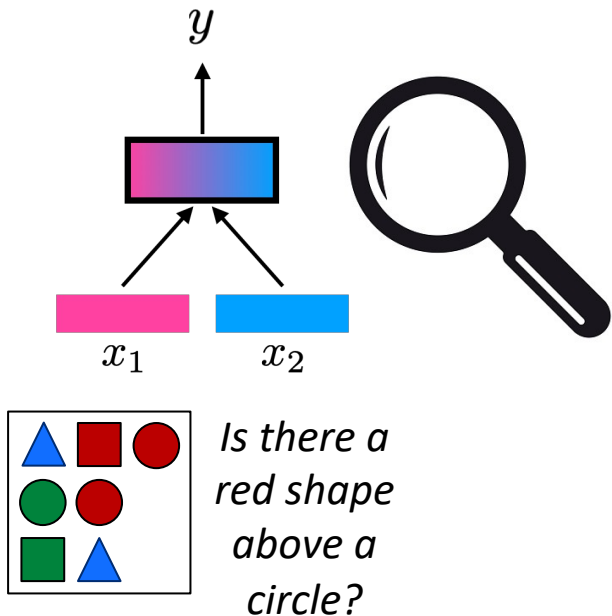
3. Multimodal representations



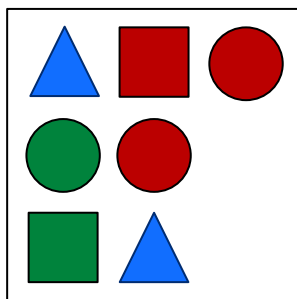
Evaluating Interpretability: A Multi-stage Approach

Multimodal prediction: Does the model correctly compose question and image information?

Yes!

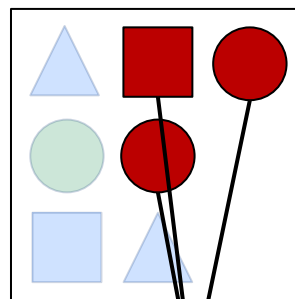


1. Unimodal importance

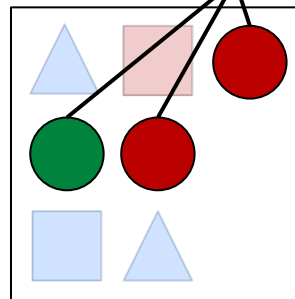


Is there a **red shape** above a circle?

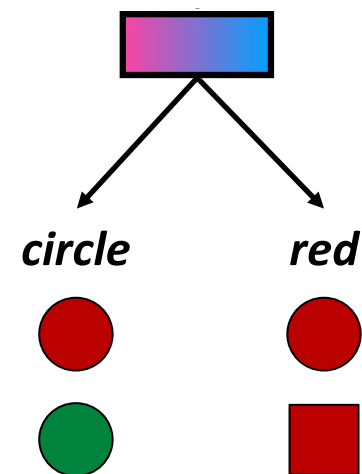
2. Cross-modal interactions



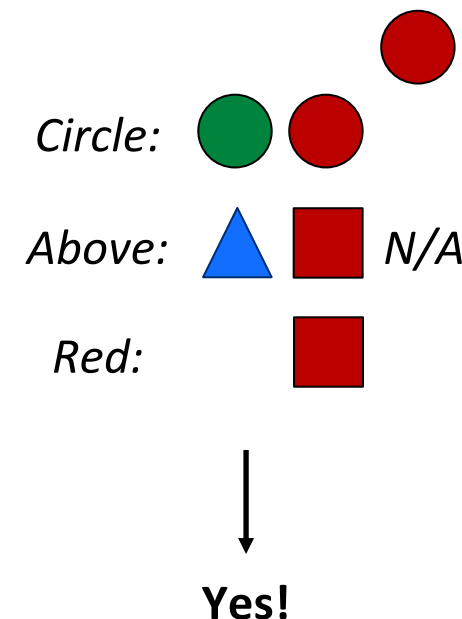
Is there a **red shape** above a circle?



3. Multimodal representations



4. Multimodal prediction



Evaluating Interpretability: A Multi-stage Approach

Identifying individual cross-modal connections

Statistical non-additive interactions [Friedman & Popescu, 2008, Sorokina et al., 2008]

f exhibits interactions between 2 features x_A and x_B iff f cannot be decomposed into a sum of unimodal subfunctions f_A, f_B such that $f(x_A, x_B) = f_A(x_A) + f_B(x_B)$.

f exhibits interactions between 2 features x_A and x_B iff $\frac{\partial^2 f}{\partial x_A \partial x_B} > 0$.

Natural second-order extension of gradient-based approaches!

Evaluating Interpretability: A Multi-stage Approach

How can we understand multimodal representations?

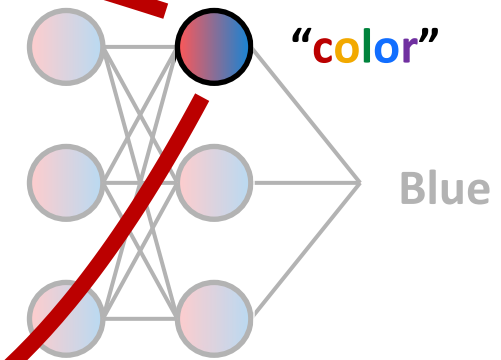


What color is the tie of the second man to the left?

<p><i>What color is the Salisbury Rd sign?</i></p>	<p><i>What color is the building?</i></p>	<p><i>What color are the checkers on the wall?</i></p>
--	---	--

Local analysis

3. Multimodal representations

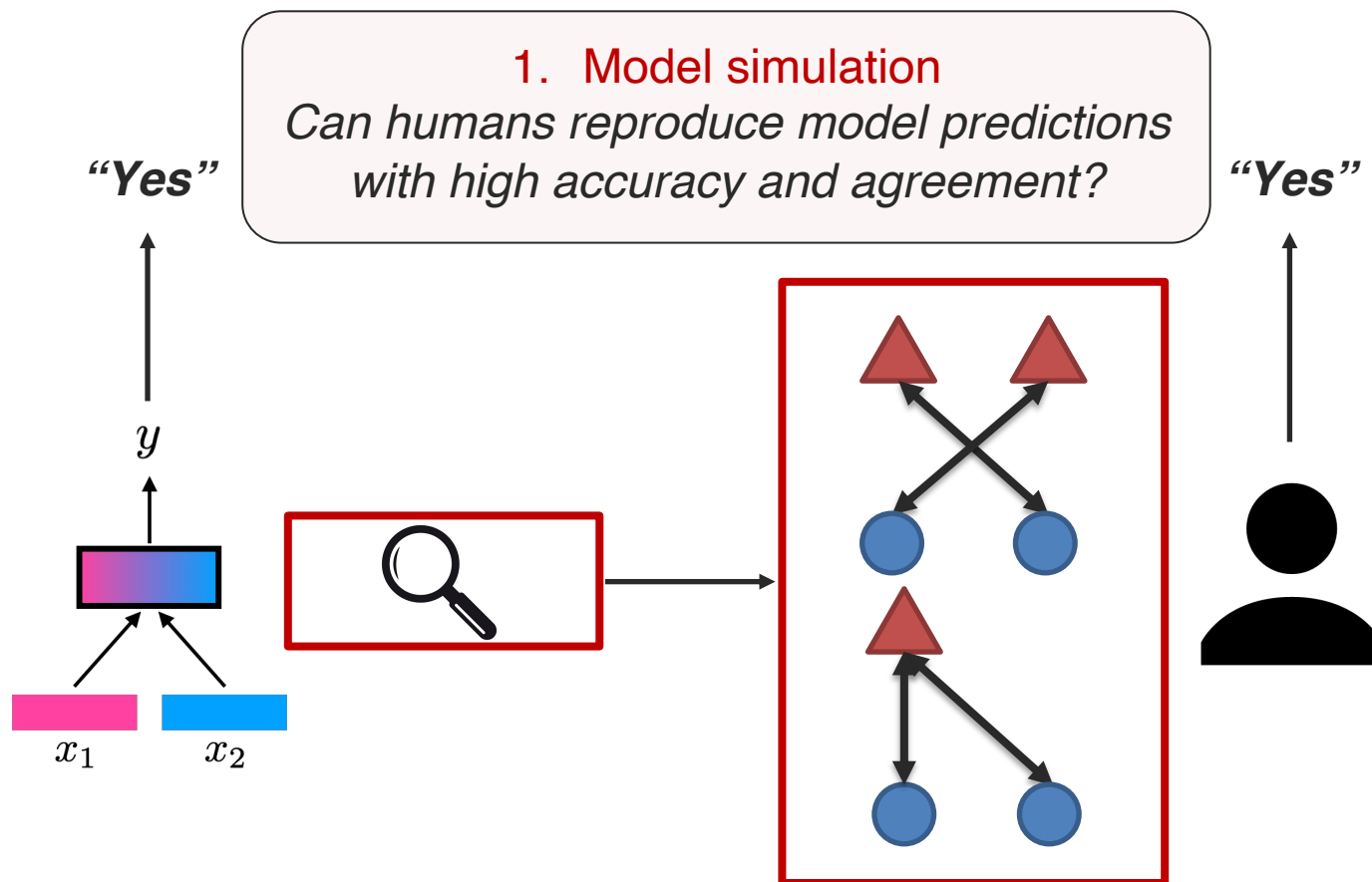


Global analysis

[Liang et al., MultiViz: An Analysis Benchmark for Visualizing and Understanding Multimodal Models. arXiv 2022]

Evaluating Interpretability: MultiViz

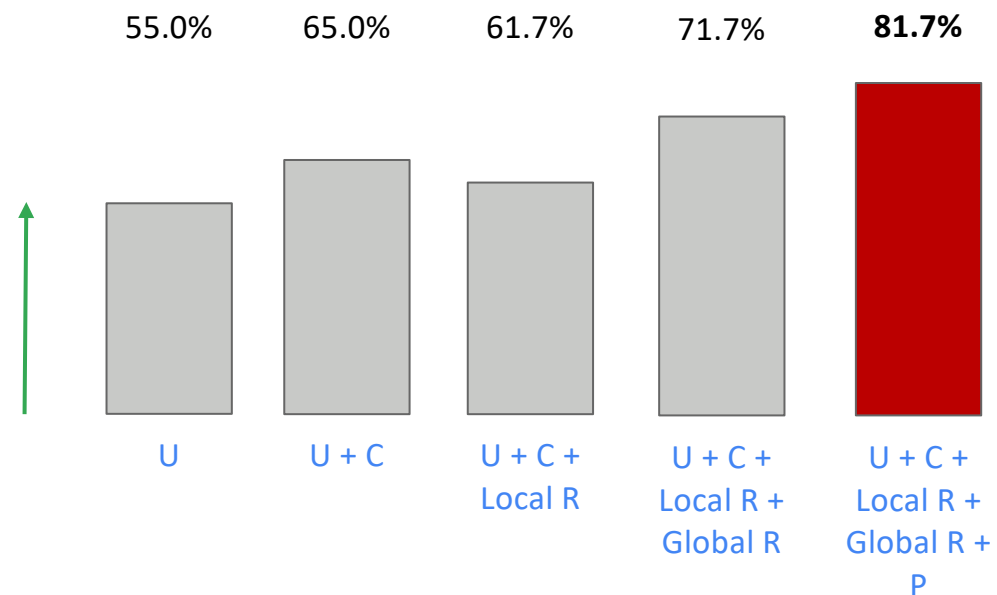
Model simulation



[Liang et al., MultiViz: A Framework for Visualizing and Understanding Multimodal Models. arXiv 2022]

Evaluating Interpretability: MultiViz

Model simulation

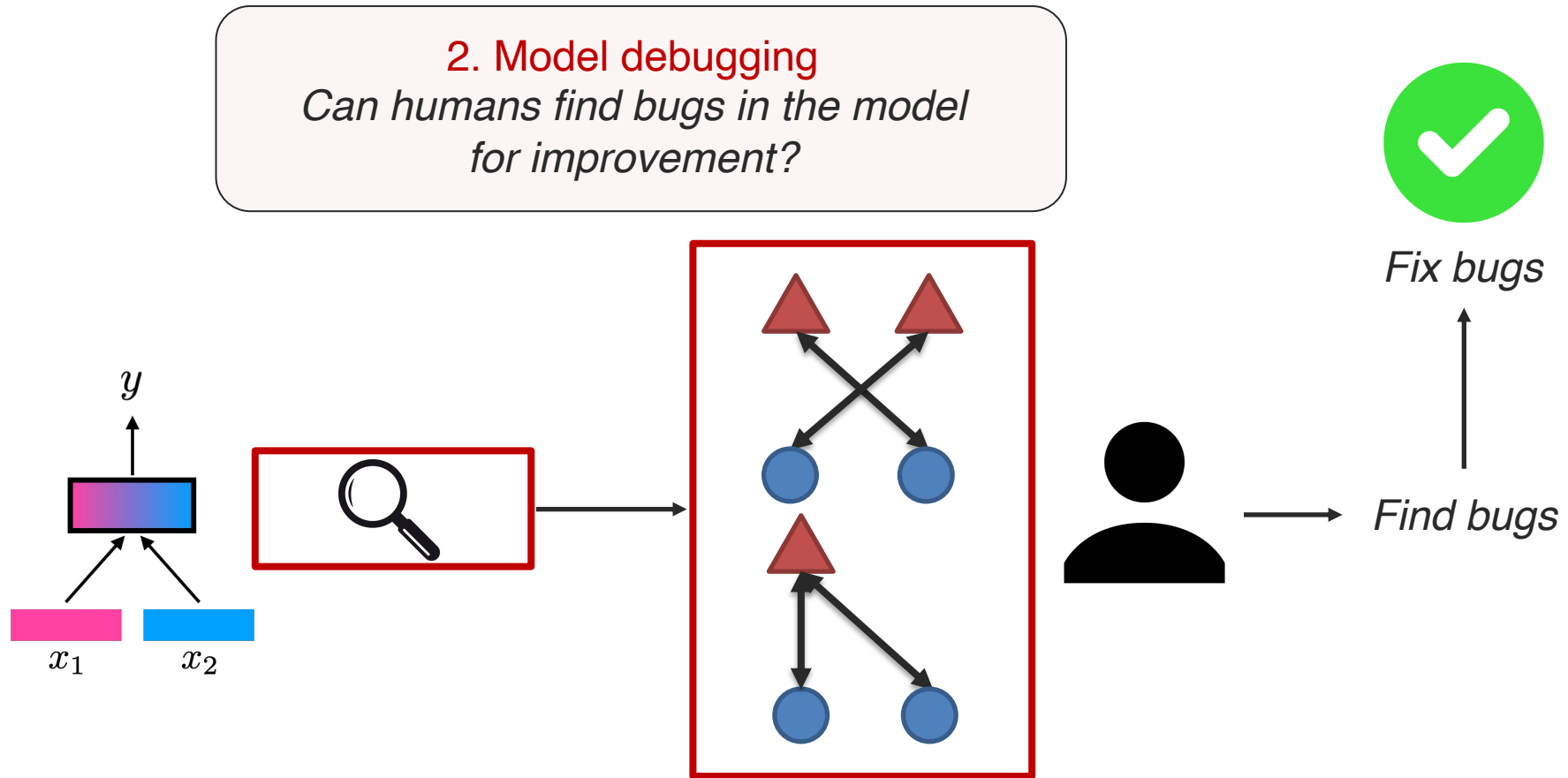


MultiViz stages leads to higher accuracy and agreement
Blind test + reasonable baselines + quantifiable outcome

[Liang et al., MultiViz: A Framework for Visualizing and Understanding Multimodal Models. arXiv 2022]

Evaluating Interpretability: MultiViz

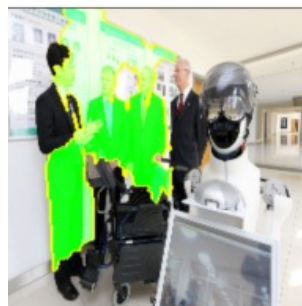
Model error analysis and debugging






[Liang et al., MultiViz: A Framework for Visualizing and Understanding Multimodal Models. arXiv 2022]

Evaluating Interpretability: MultiViz

Model error analysis and debugging

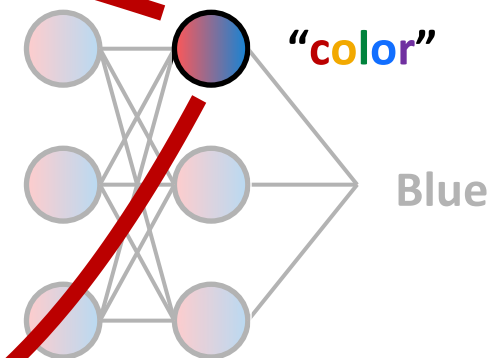


What color is the tie of the second man to the left?

 <p><i>What color is the Salisbury Rd sign?</i></p>	 <p><i>What color is the building?</i></p>	 <p><i>What color are the checkers on the wall?</i></p>
---	---	---

Local analysis

3. Multimodal representations



Global analysis

"Models pick up cross-modal interactions but fail in identifying color!"

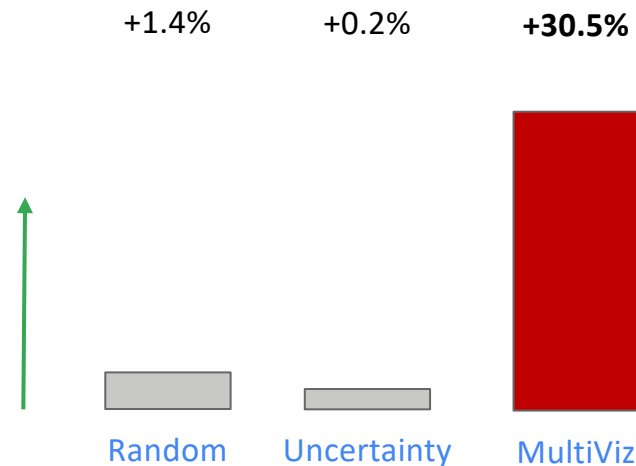
Evaluating Interpretability: MultiViz

Model error analysis and debugging

“Models pick up cross-modal interactions but fail in identifying color!”



Add targeted examples involving color.



Side note: we used this to discover a bug in a popular deep learning code repository.


Transformers

MultiViz enables error analysis and debugging of multimodal models

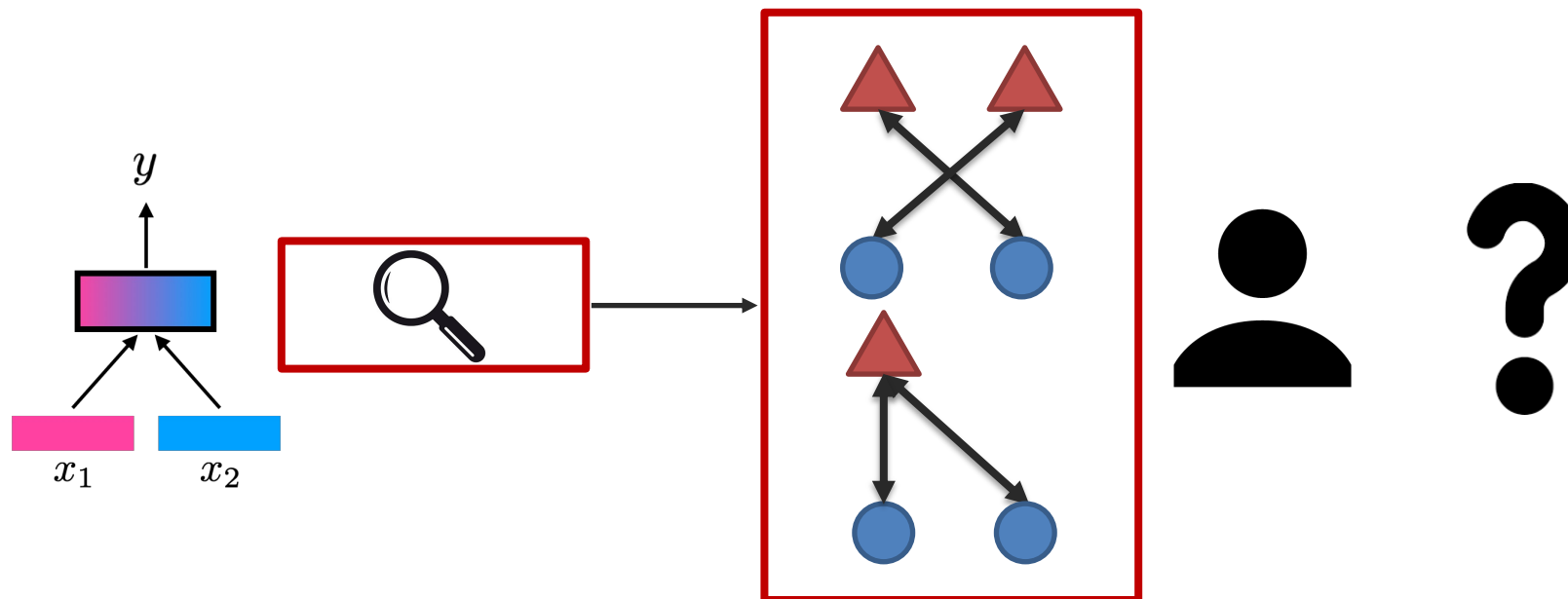
[Liang et al., MultiViz: A Framework for Visualizing and Understanding Multimodal Models. arXiv 2022]

Challenges: Quantifying Multimodal Interactions

Open challenges

Open challenges:

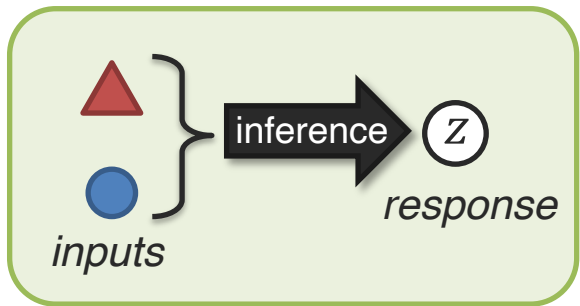
- Faithfulness: do explanations accurately reflect model's internal mechanics?
- Usefulness: unclear if explanations help humans
- Disagreement: different interpretation methods may generate different explanations
- Evaluate: how to best evaluate interpretation methods



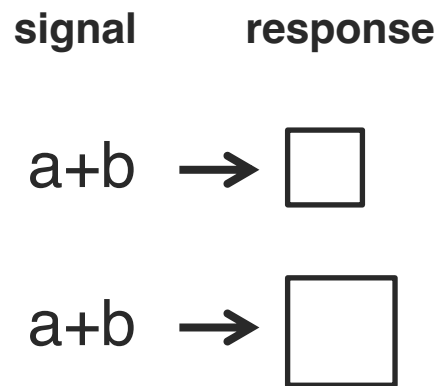
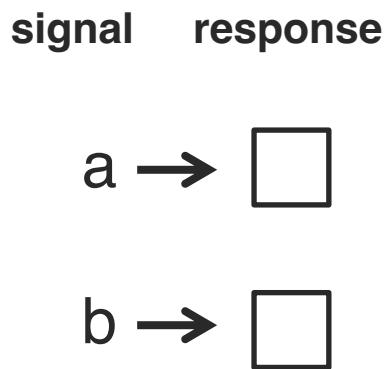
[Chandrasekaran et al., Do explanations make VQA models more predictable to a human? EMNLP 2018]

[Krishna et al., The Disagreement Problem in Explainable Machine Learning: A Practitioner's Perspective. arXiv 2022]

Challenges: Quantifying Multimodal Interactions



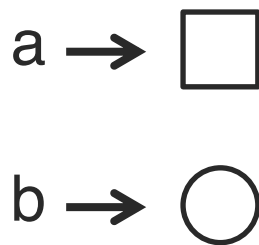
Redundancy



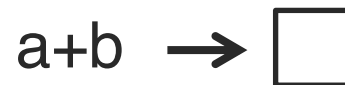
Equivalence

Enhancement

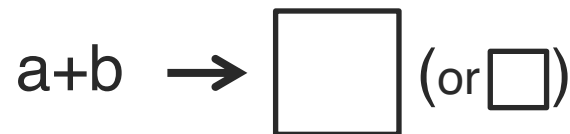
Nonredundancy



Independence



Dominance



Modulation



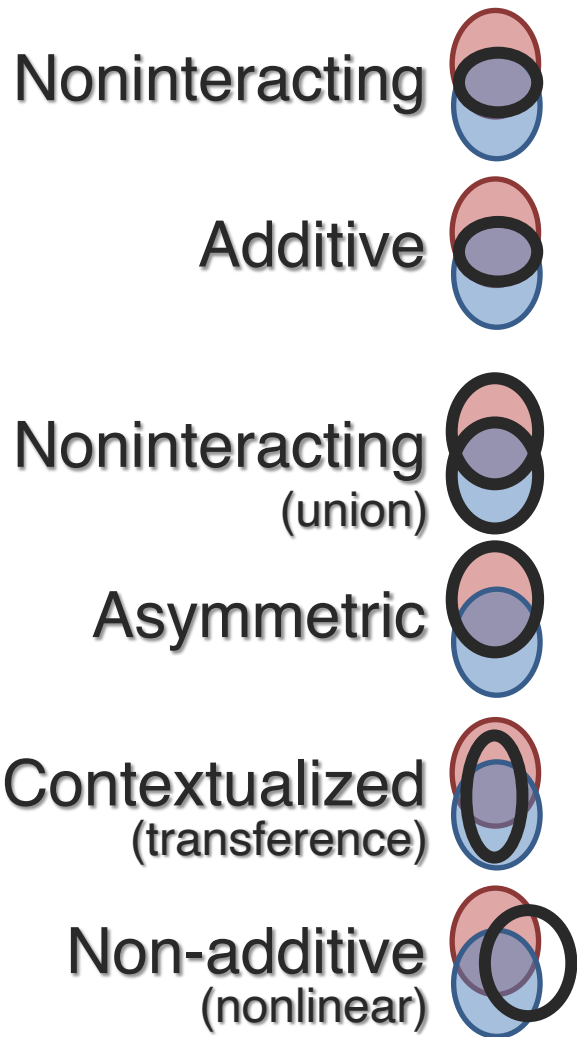
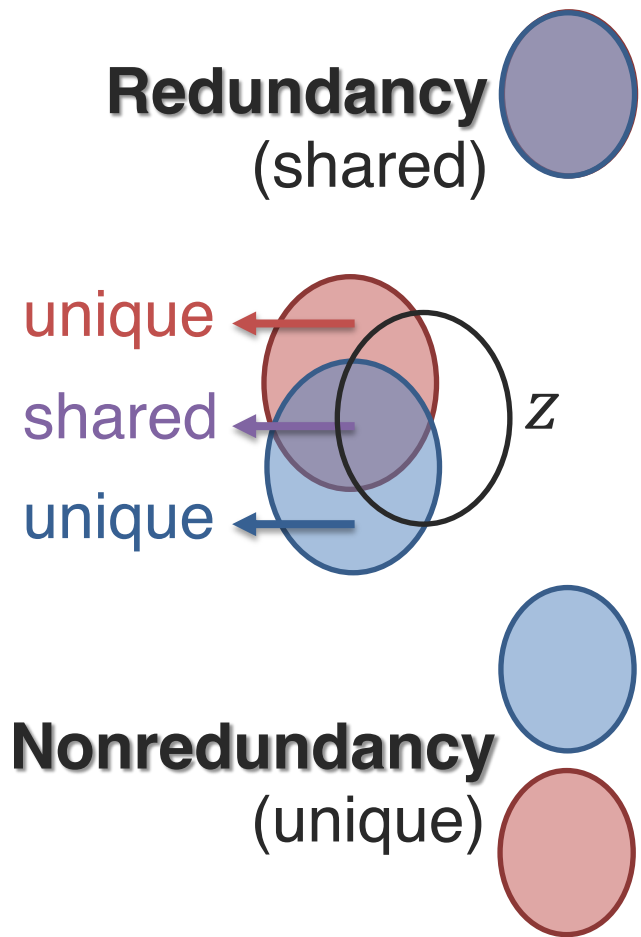
Emergence

Multimodal Communication



Partan and Marler (2005). *Issues in the classification of multimodal communication signals*. *American Naturalist*, 166(2)

Challenges: Quantifying Multimodal Interactions



signal	response	
$a+b$	\rightarrow	Equivalence
$a+b$	\rightarrow	Enhancement
$a+b$	\rightarrow and	Independence
$a+b$	\rightarrow	Dominance
$a+b$	\rightarrow (or)	Modulation
$a+b$	\rightarrow	Emergence

Challenges: Quantifying Multimodal Interactions

Recall error analysis!

Causal, logical interactions beyond additive/multiplicative

Covariant VQA

Target object in question

Q: How many zebras are there in the picture?

A: 2

zebra removed A: 1

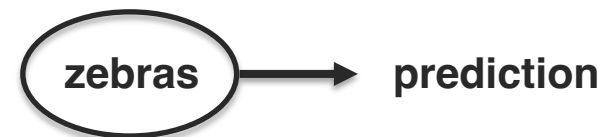


Baselines:

2

2

i.e., treatment variable



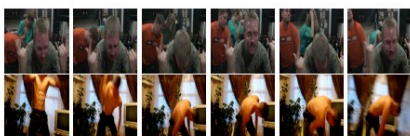
Interventional conditional: $p(y|do(zebras = 1))$

Existing models struggle to adapt to targeted causal interventions.
How can we make them more robust to spurious correlations?

Sub-Challenge 6c: Multimodal Learning Process

Definition: Characterizing the learning and optimization challenges involved when learning from heterogeneous data.

Kinetics dataset



(a) headbanging



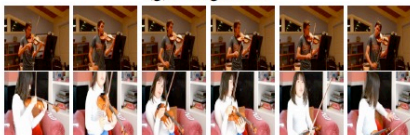
(c) shaking hands



(e) robot dancing



(g) riding a bike



Adding more modalities should always help?

Modalities: RGB (video clips)

A (Audio features)

OF (optical flow - motion)

Dataset	Multi-modal	V@1	Best Uni	V@1	Drop
Kinetics	A + RGB	71.4	RGB	72.6	-1.2
	RGB + OF	71.3	RGB	72.6	-1.3
	A + OF	58.3	OF	62.1	-3.8
	A + RGB + OF	70.0	RGB	72.6	-2.6

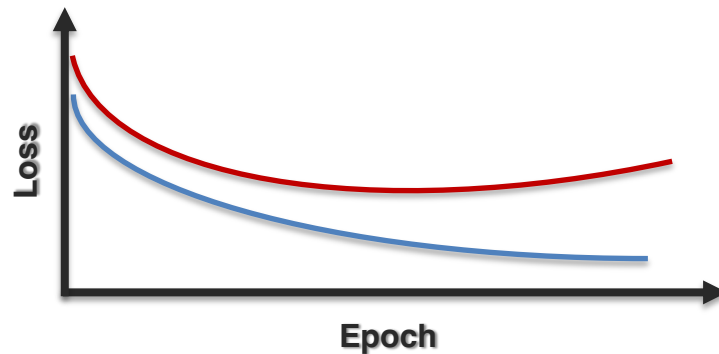
But sometimes multimodal doesn't help! **Why?**

Optimization challenges

Learning and optimization challenges

2 explanations for drop in performance:

1. Multimodal networks are more prone to overfitting due to **increased complexity**
2. Different modalities overfit and generalize at **different rates**



Key idea 1: compute overfitting-to-generalization ratio (OGR)



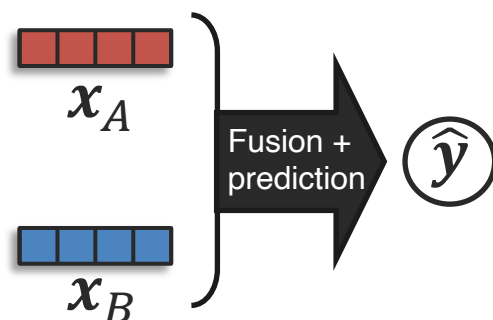
Gap between training and valid loss

OGR wrt each modality tells us how much to train that modality

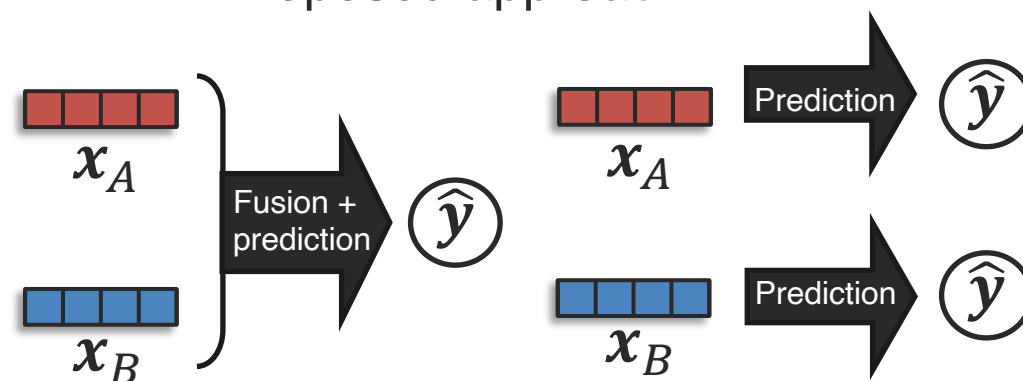
Optimization challenges

Learning and optimization challenges

Conventional approach



Proposed approach



Key idea 2: Simultaneously train unimodal networks to estimate OGR wrt each modality

+ Reweight multimodal loss using unimodal OGR values

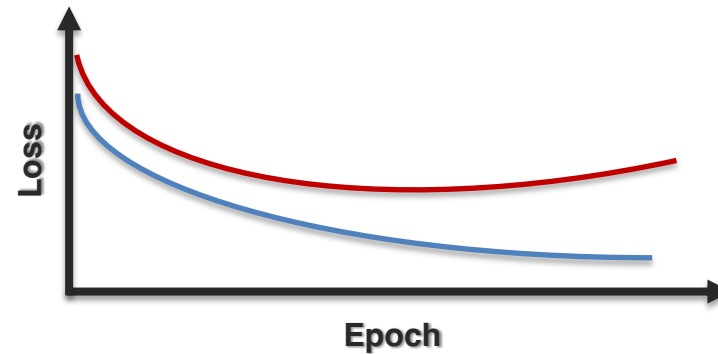
→ Allows to better balance generalization & overfitting rate of different modalities

Challenges

Open
challenges

Open challenges:

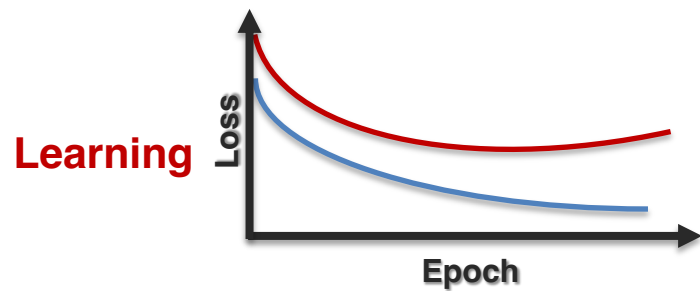
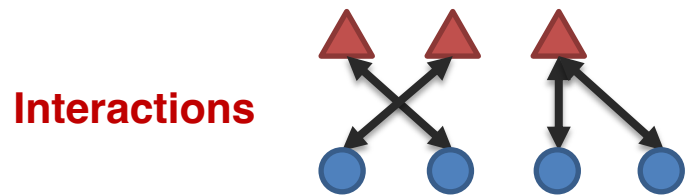
- Learning, generalization, and optimization in high-dimensional settings ($p \gg n$).
- Modality shortcuts and biases.
- Dimensionality reduction, modality selection, approximate inference.
- Reducing time and space complexity, model compression and efficiency.



More Quantification

Dimensions of quantification

Representation Alignment Reasoning Transference Generation



Conclusion

What is a Modality?

Multimodal Behaviors and Signals

Language

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

Acoustic

- **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

Touch

- **Haptics**
- **Motion**

Physiological

- **Skin conductance**
- **Electrocardiogram**

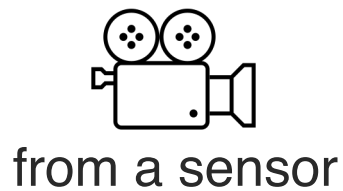
Mobile

- **GPS location**
- **Accelerometer**
- **Light sensors**

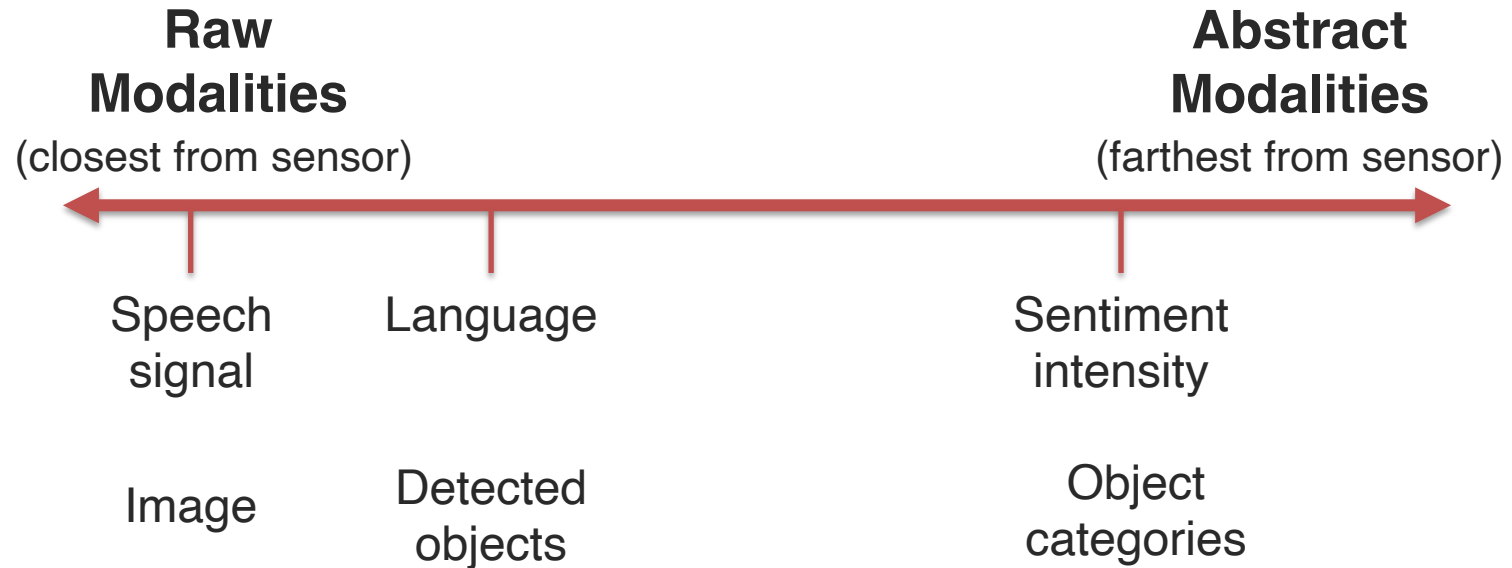
What is a Modality?

Definition

Modality refers to the way in which something expressed or perceived.



Examples:



What is Multimodal?

A dictionary definition...

Multimodal: with multiple modalities

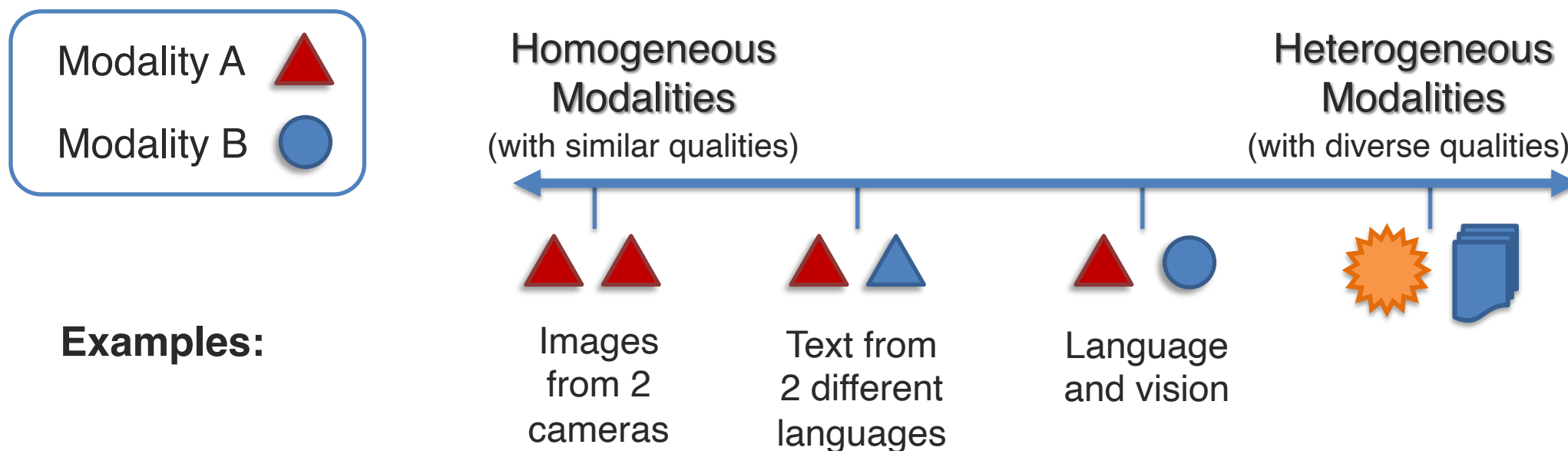
A research-oriented definition...

***Multimodal* is the scientific study of
heterogeneous and interconnected data**

Connected + Interacting

Heterogeneous Modalities

Heterogeneous: Diverse qualities, structures and representations.



Abstract modalities are more likely to be homogeneous

Connected Modalities

Connected: Shared information that relates modalities



Statistical



Association

Dependency



e.g., correlation,
co-occurrence



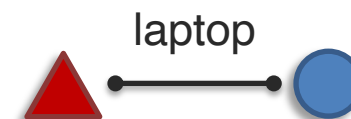
e.g., causal,
temporal

Semantic



Correspondence

Relationship



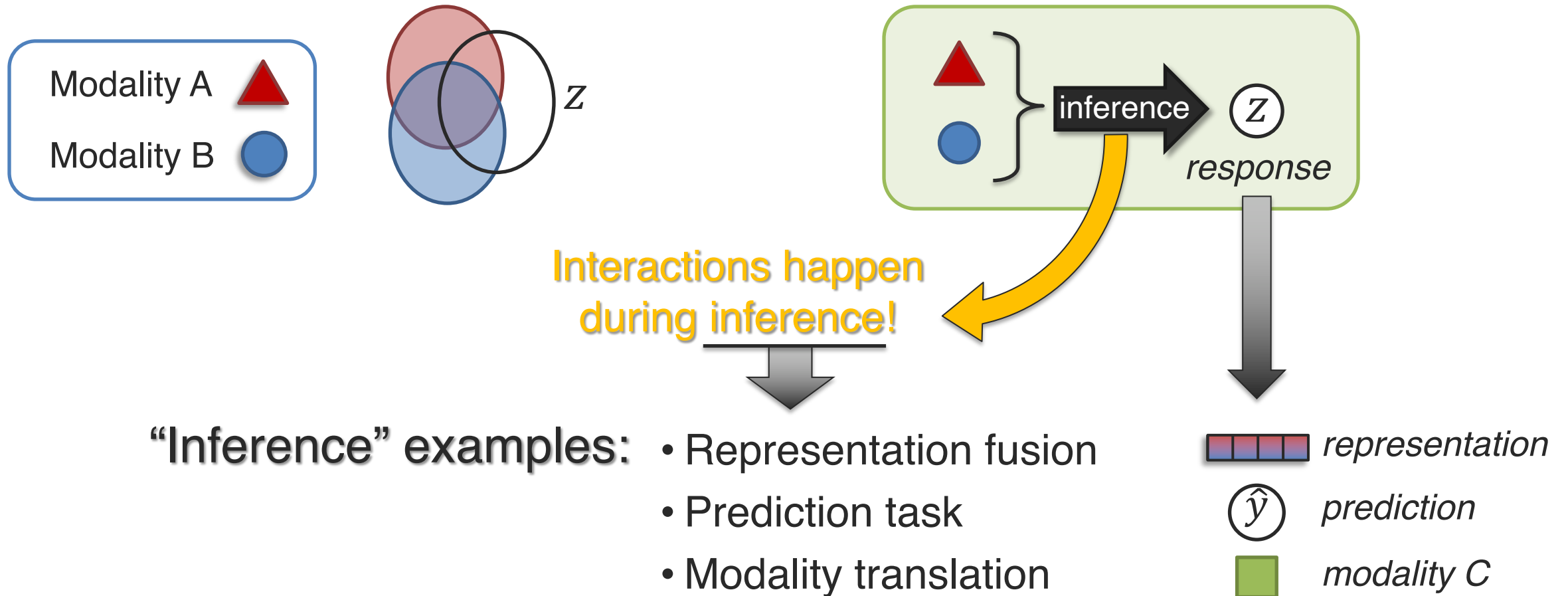
e.g., grounding



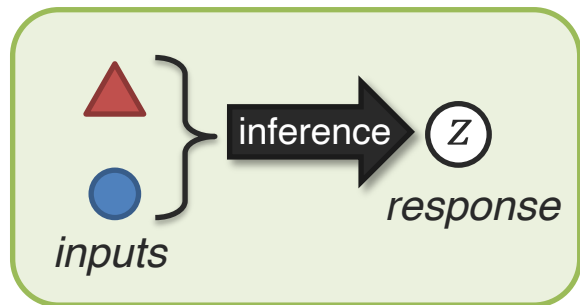
e.g., function

Interacting Modalities

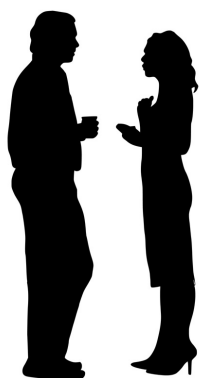
Interacting: process affecting each modality, creating new response



Taxonomy of Interaction Responses – A Behavioral Science View



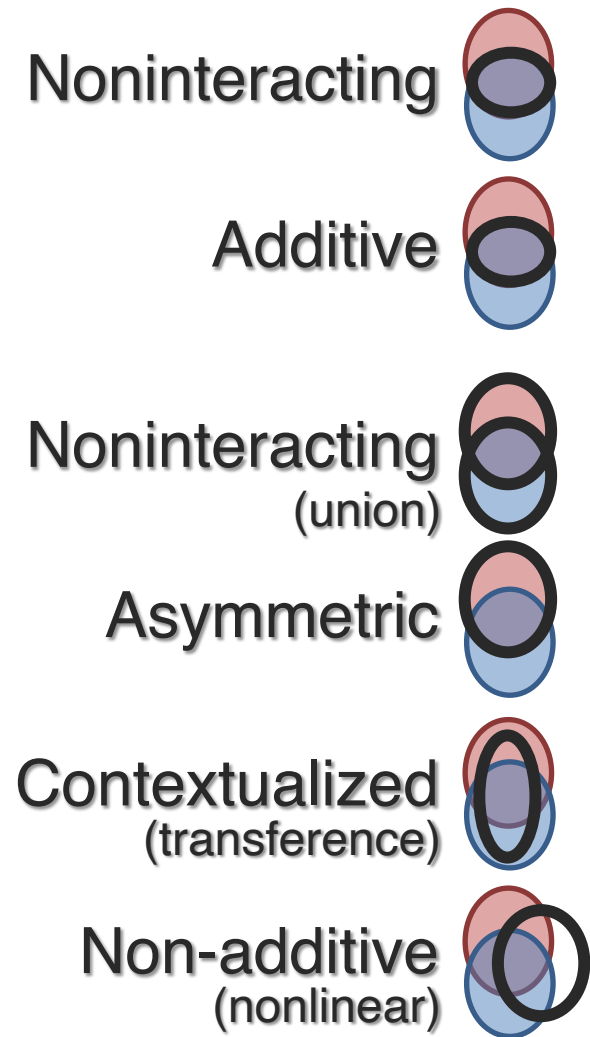
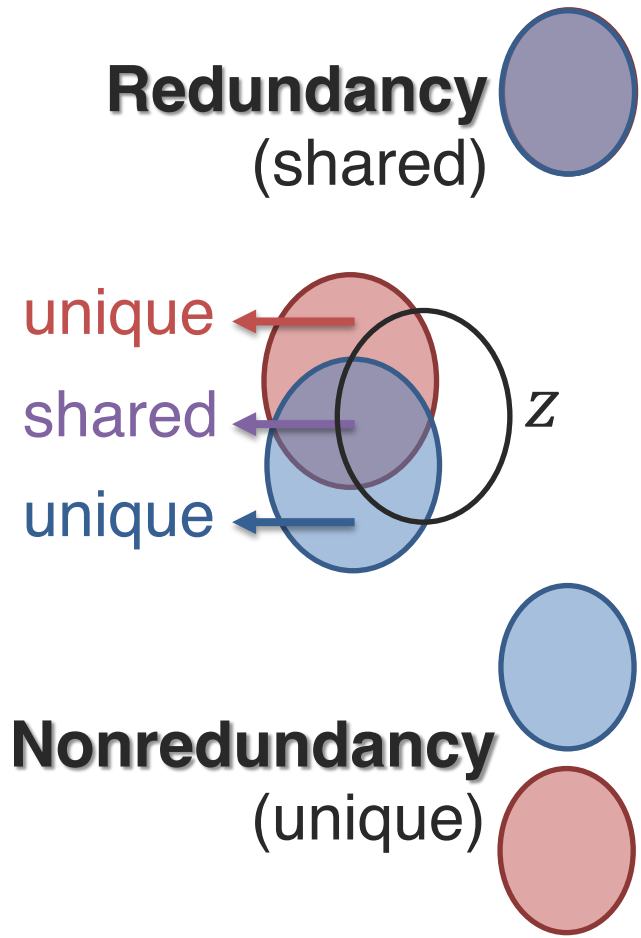
Multimodal Communication



	signal	response	signal	response	
Redundancy	a	→ □	a+b	→ □	Equivalence
	b	→ □	a+b	→ □	Enhancement
Nonredundancy	a+b	→ □ and ○	a+b	→ □ and ○	Independence
	a	→ □	a+b	→ □	Dominance
	b	→ ○	a+b	→ □ (or □)	Modulation
	a+b	→ △	a+b	→ △	Emergence

Partan and Marler (2005). *Issues in the classification of multimodal communication signals*. *American Naturalist*, 166(2)

Interacting Modalities



signal	response	
$a+b$	\rightarrow	Equivalence
$a+b$	\rightarrow	Enhancement
$a+b$	\rightarrow and	Independence
$a+b$	\rightarrow	Dominance
$a+b$	\rightarrow (or)	Modulation
$a+b$	\rightarrow	Emergence

*What is
Multimodal?*



Why is it hard?



What is next?

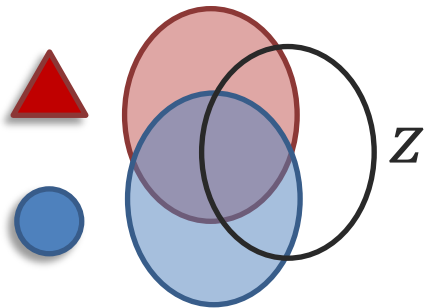
Heterogeneous



Connected

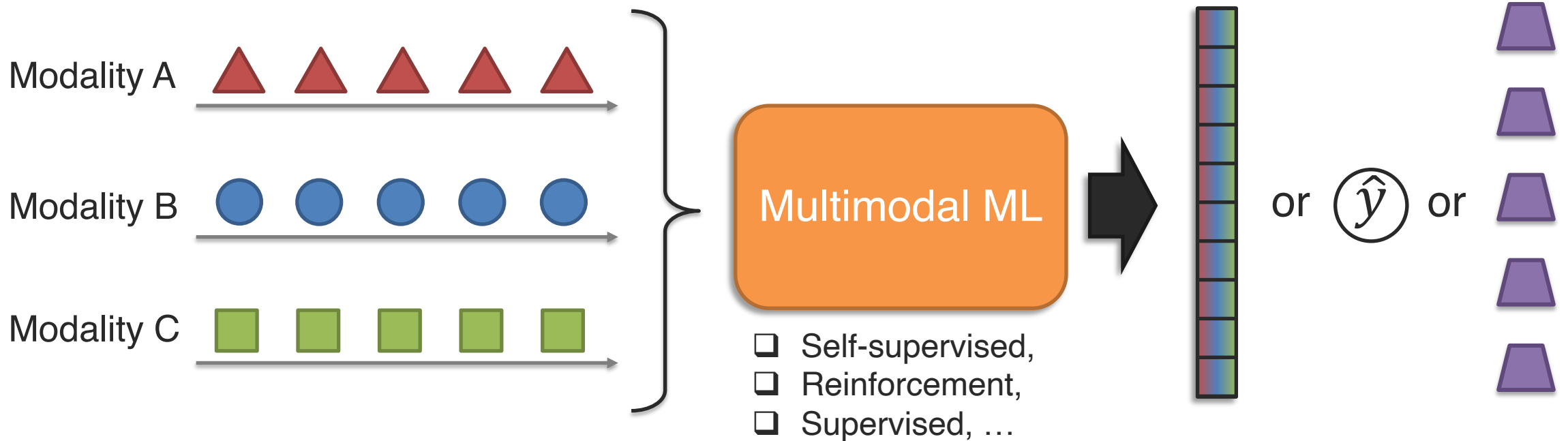


Interacting



**Multimodal is the scientific
study of heterogeneous and
interconnected data 😊**

Multimodal Machine Learning



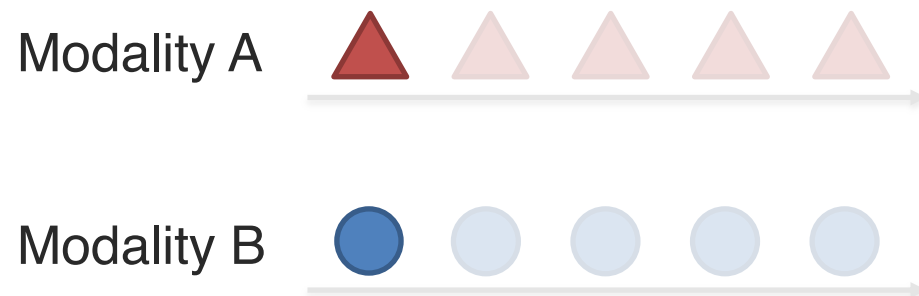
*What are the **core multimodal technical challenges**, understudied in conventional machine learning?*

Challenge 1: Representation

Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

➡ This is a core building block for most multimodal modeling problems!

Individual elements:



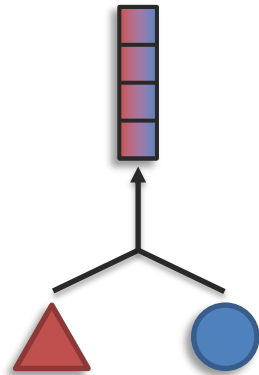
*It can be seen as a “local” representation
or
representation using holistic features*

Challenge 1: Representation

Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

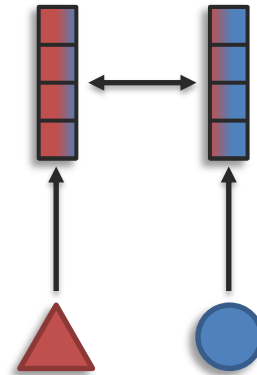
Sub-challenges:

Fusion



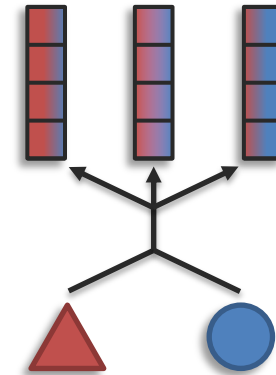
modalities $>$ # representations

Coordination



modalities = # representations

Fission



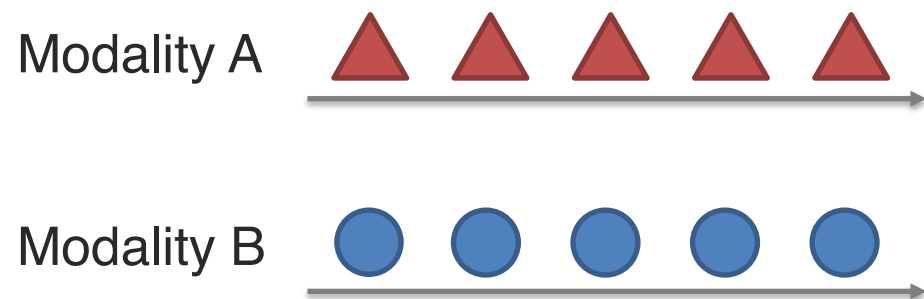
modalities $<$ # representations

Challenge 2: Alignment

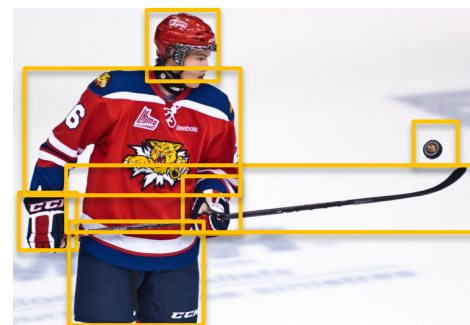
Definition: Identifying and modeling cross-modal connections between all elements of multiple modalities, building from the data structure

➡ Most modalities have internal structure with multiple elements

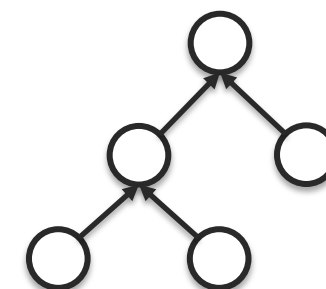
Elements with temporal structure:



Other structured examples:



Spatial



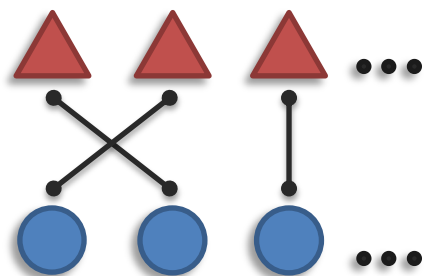
Hierarchical

Challenge 2: Alignment

Definition: Identifying and modeling cross-modal connections between all elements of multiple modalities, building from the data structure

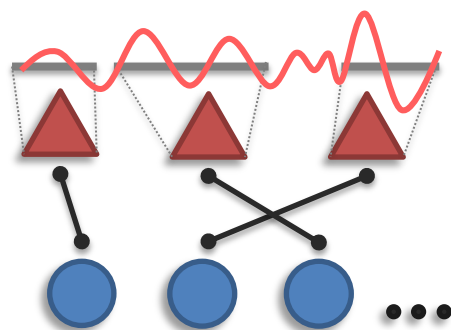
Sub-challenges:

Discrete Alignment



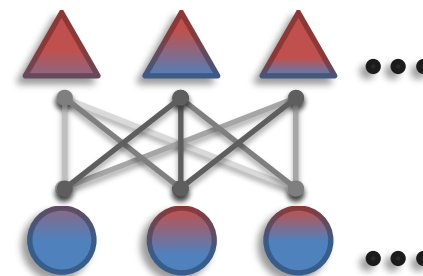
Discrete elements
and connections

Continuous Alignment



Segmentation and
continuous warping

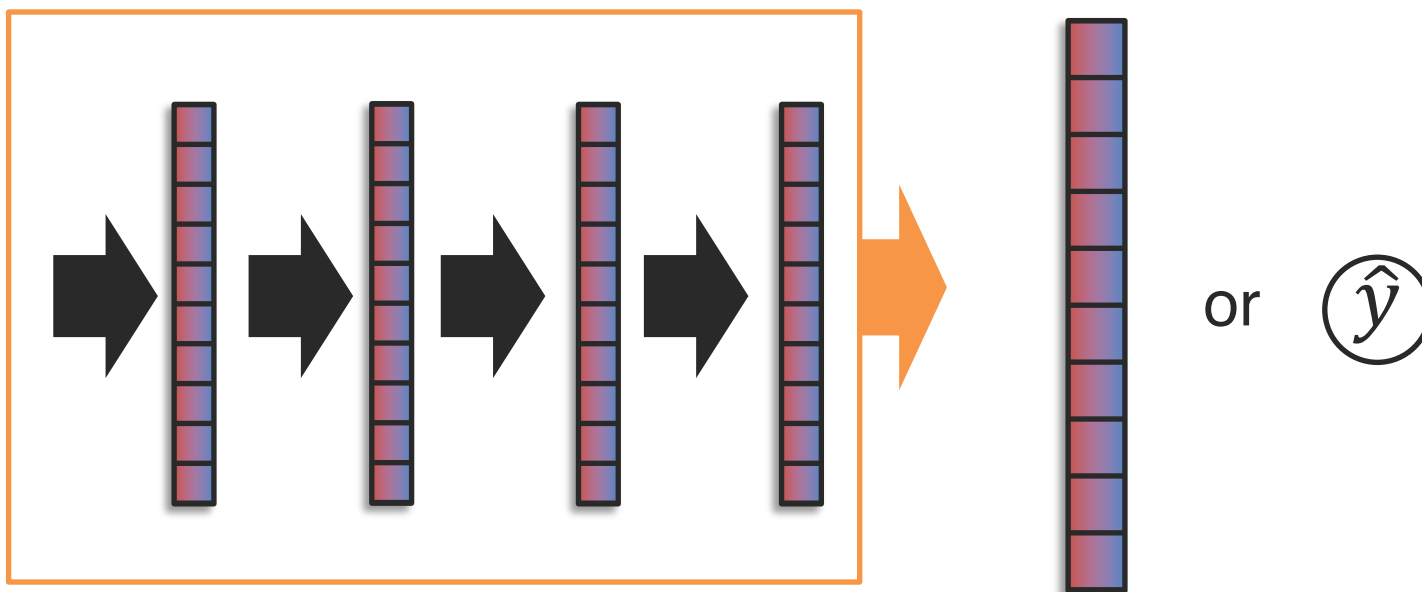
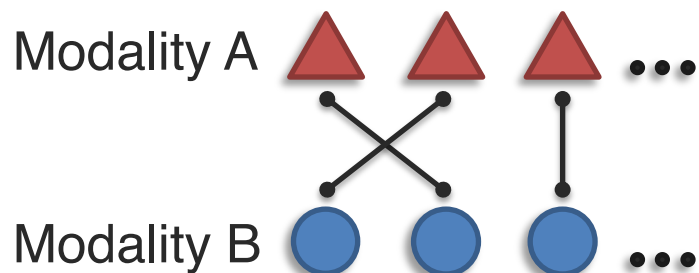
Contextualized Representation



Alignment + representation

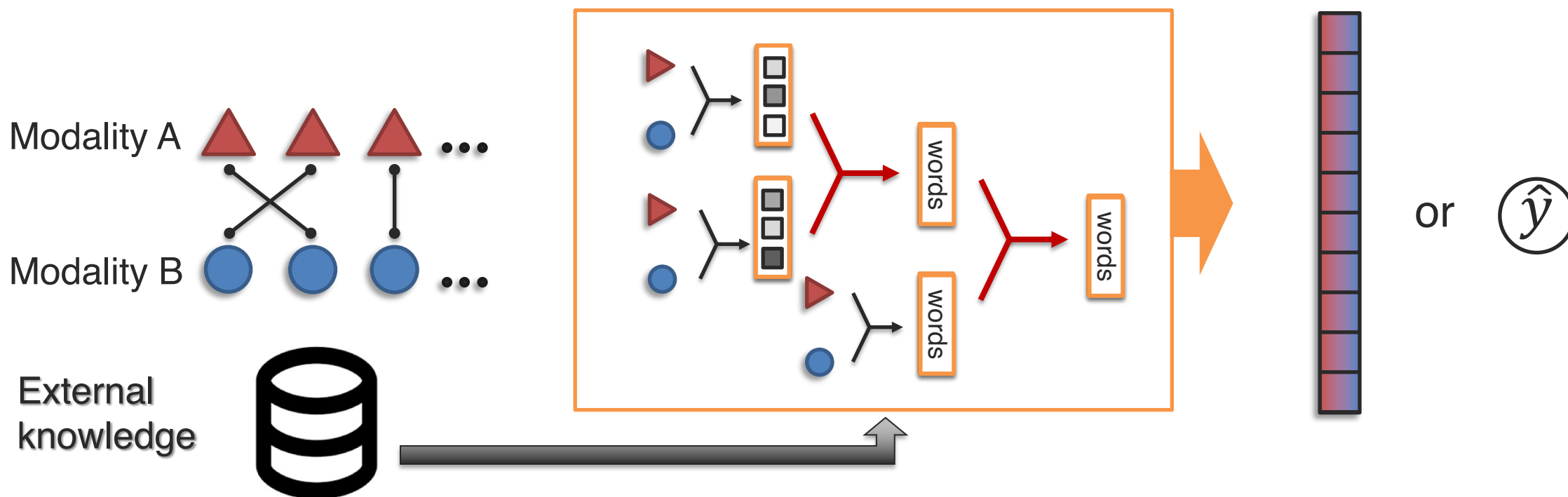
Challenge 3: Reasoning

Definition: Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure



Challenge 3: Reasoning

Definition: Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure

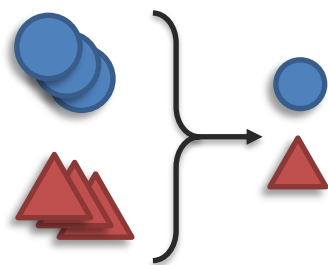


Challenge 4: Generation

Definition: Learning a generative process to produce raw modalities that reflects cross-modal interactions, structure and coherence

Sub-challenges:

Summarization



Reduction



Information:
(content)

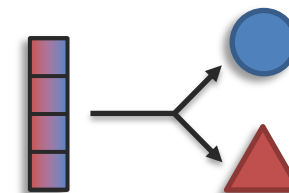
Translation



Maintenance



Creation

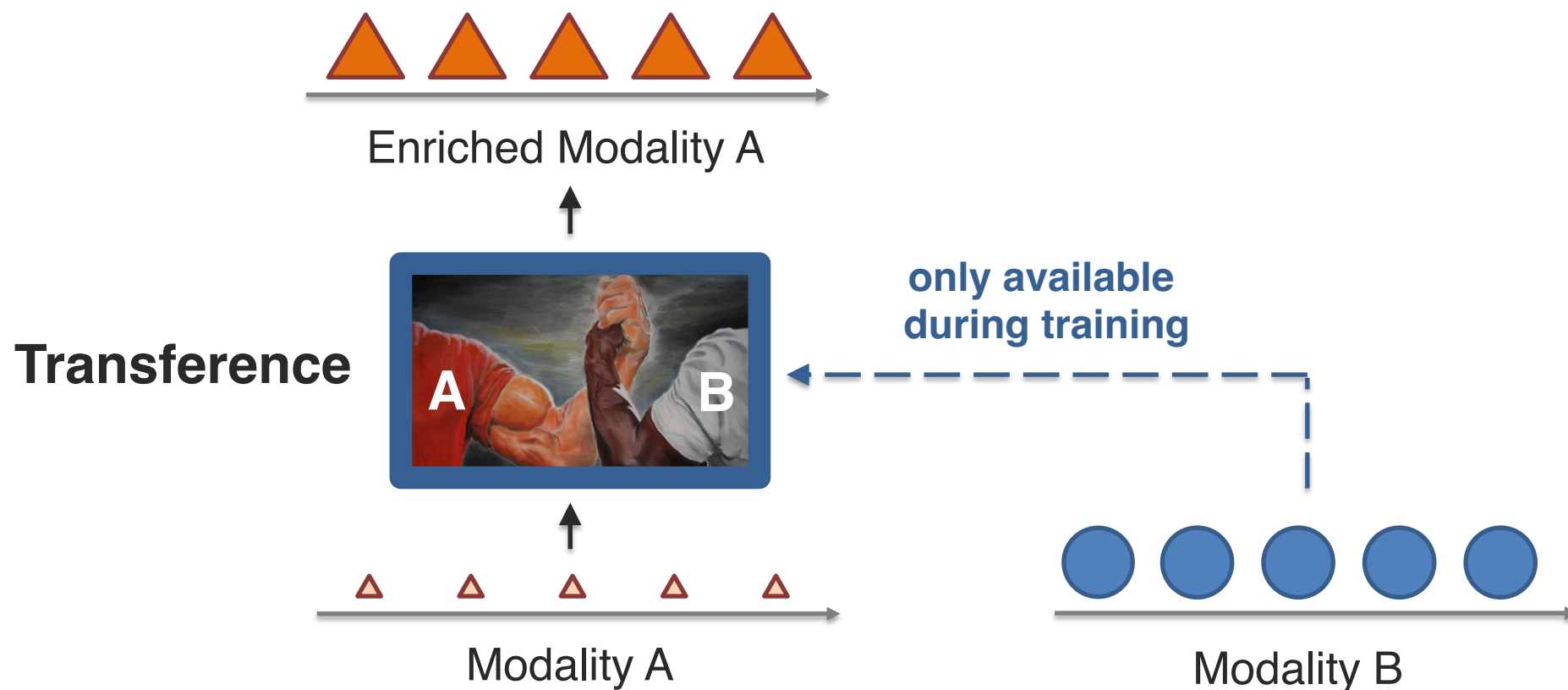


Expansion



Challenge 5: Transference

Definition: Transfer knowledge between modalities, usually to help the target modality which may be noisy or with limited resources

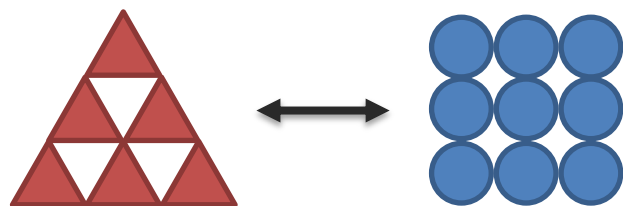


Challenge 6: Quantification

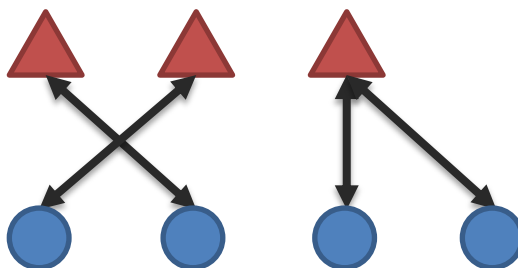
Definition: Empirical and theoretical study to better understand heterogeneity, cross-modal interactions and the multimodal learning process

Sub-challenges:

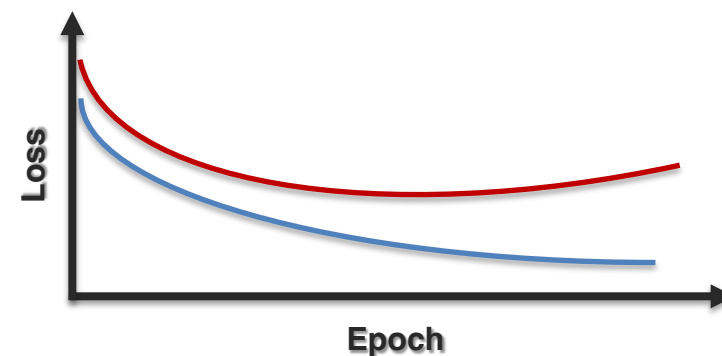
Heterogeneity



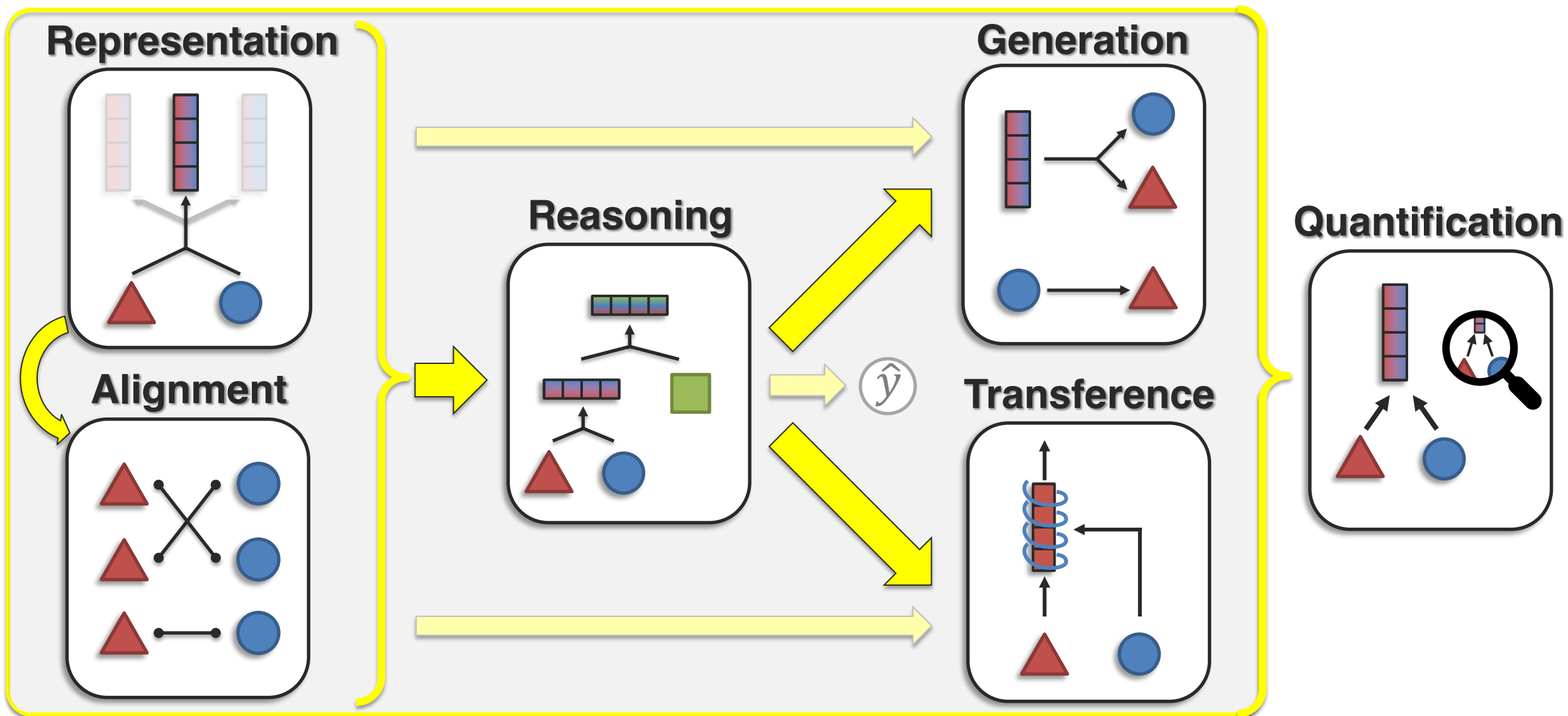
Connections & Interactions



Learning

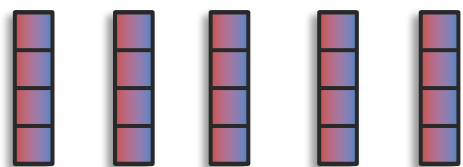


Core Multimodal Challenges



Future Direction: Heterogeneity

Homogeneity



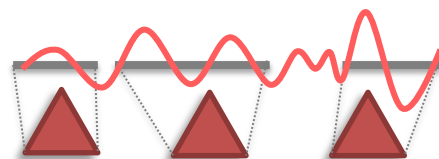
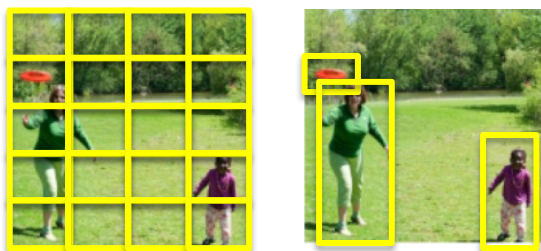
vs

Heterogeneity



Examples:

Arbitrary Tokenization



Beyond Additive Interactions

Causal, logical interactions

Brain-inspired representations

MultiBench

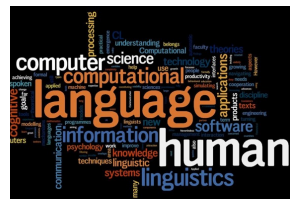
<https://github.com/pliang279/MultiBench>

Future Direction: High-modality

Few modalities



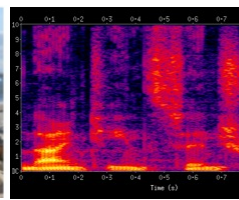
High-modality



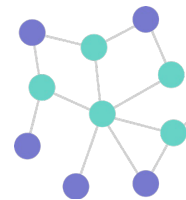
Language



Vision



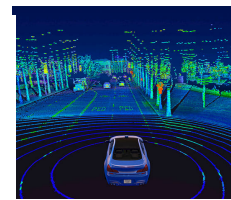
Audio



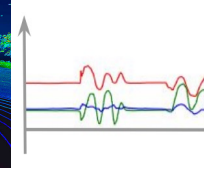
Graphs



Control



LIDAR



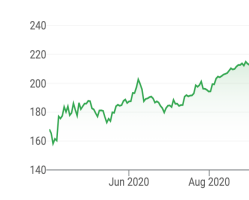
Sensors



Set

SUBJECT_ID
Age
Sex
Ethnicity
...

Table



Financial



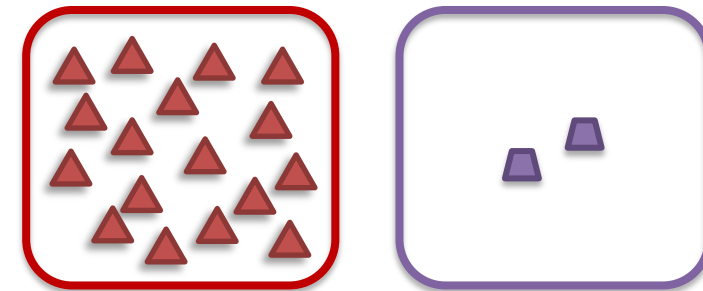
Medical

Examples:

Non-parallel learning

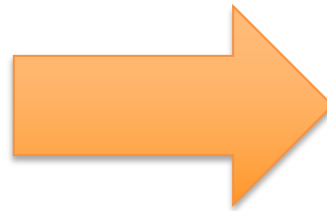
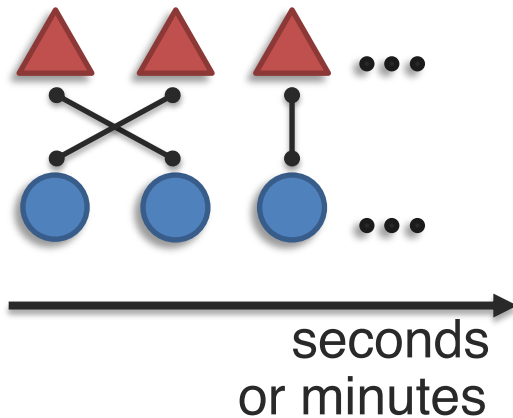


Limited resources



Future Direction: Long-term

Short-term



Long-term



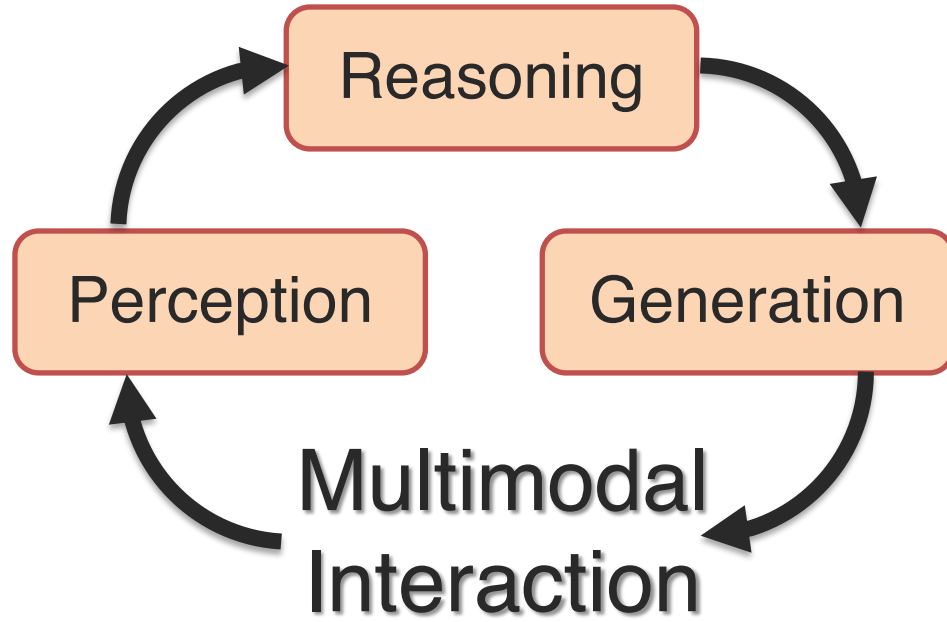
Examples:

Compositionality

Memory

Personalization

Future Direction: Interaction



Social Intelligence



Examples:

Multi-Party

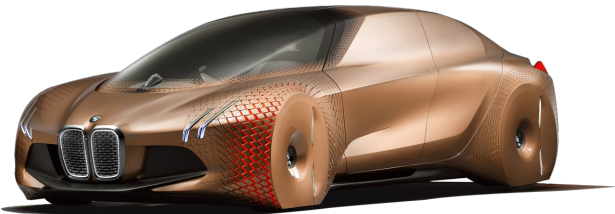
Causality

Ethical

Future Direction: Real-world



Healthcare
Decision Support



Intelligent Interfaces and
Vehicles



Online Learning
and Education

Examples:

Robustness

Fairness

Generalization

What is Multimodal?



Why is it hard?



What is next?

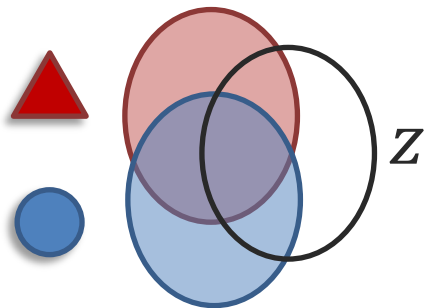
Heterogeneous



Connected



Interacting



Representation

Alignment

Reasoning

Generation

Transference

Quantification

Heterogeneity

High-modality

Long-term

Interaction

Real-world

Advanced Topics in Multimodal ML @ CMU



Advanced Topics in Multimodal Machine Learning

11-877 • Spring 2022 • Carnegie Mellon University

Multimodal machine learning (MML) 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 language, vision, and acoustic. 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 is designed to be a graduate-level course covering recent research papers in multimodal machine learning, including technical challenges with representation, alignment, reasoning, generation, co-learning and quantifications. The main goal of the course is to increase critical thinking skills, knowledge of recent technical achievements, and understanding of future research directions.

- **Time:** Friday 10:10-11:30 am
- **Location:** Virtual for the first 2 weeks (find zoom link in piazza), GHC 5222 thereafter
- **Discussion and Q&A:** [Piazza](#)
- **Assignment submissions:** [Canvas](#) (for registered students only)
- **Contact:** Students should ask all course-related questions on [Piazza](#), where you will also find announcements.



Instructor [Louis-Philippe Morency](#)
Email: morency@cs.cmu.edu



Instructor [Amir Zadeh](#)
Email: abagherz@cs.cmu.edu



Instructor [Paul Liang](#)
Email: pliang@cs.cmu.edu

1/28 Week 2: Cross-modal interactions [synopsis]

- What are the different ways in which modalities can interact with each other in multimodal tasks? Can we formalize a taxonomy of such cross-modal interactions, which will enable us to compare and contrast them more precisely?
- What are the design decisions (aka inductive biases) that can be used when modeling these cross-modal interactions in machine learning models?
- What are the advantages and drawbacks of designing models to capture each type of cross-modal interaction? Consider not just prediction performance, but tradeoffs in time/space complexity, interpretability, etc.
- Given an arbitrary dataset and prediction task, how can we systematically decide what type of cross-modal interactions exist, and how can that inform our modeling decisions?
- Given trained multimodal models, how can we understand or visualize the nature of cross-modal interactions?

- Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think!
- What Does BERT with Vision Look At?
- Multiplicative Interactions and Where to Find Them
- Cooperative Learning for Multi-view Analysis
- Vision-and-Language or Vision-for-Language? On Cross-Modal Influence in Multimodal Transformers
- Seeing past words: Testing the cross-modal capabilities of pretrained V&L models on counting tasks

2/4 Week 3: Multimodal co-learning [synopsis]

- What are the types of cross-modal interactions involved to enable such co-learning scenarios where multimodal training ends up generalizing to unimodal testing?
- What are some design decisions (inductive bias) that could be made to promote transfer of information from one modality to another?
- How do we ensure that during co-learning, only useful information is transferred, and not some undesirable bias? This may become a bigger issue in low-resource settings.
- How can we know if co-learning has succeeded? Or failed? What approaches could we develop to visualize and probe the success of co-learning?
- How can we formally, empirically, or intuitively measure the additional information provided by auxiliary modality? How can we design controlled experiments to test these hypotheses?
- What are the advantages and drawbacks of information transfer during co-learning? Consider not just prediction performance, but also tradeoffs with complexity, interpretability, fairness, etc.

- Multimodal Prototypical Networks for Few-shot Learning
- SMIL: Multimodal Learning with Severely Missing Modality
- Multimodal Co-learning: Challenges, Applications with Datasets, Recent Advances and Future Directions
- Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision
- What Makes Multi-modal Learning Better than Single (Provably)
- Found in Translation: Learning Robust Joint Representations by Cyclic Translations Between Modalities
- Zero-Shot Learning Through Cross-Modal Transfer
- 12-in-1: Multi-Task Vision and Language Representation Learning
- A Survey of Reinforcement Learning Informed by Natural Language

<https://cmu-multicomp-lab.github.io/adv-mmml-course/spring2022/>