

Multimodal Machine Learning: Principles, Challenges, and Open Questions

Paul Pu Liang

Machine Learning Department
Carnegie Mellon University

<https://www.cs.cmu.edu/~pliang/>

pliang@cs.cmu.edu

<https://github.com/pliang279>

 @pliang279

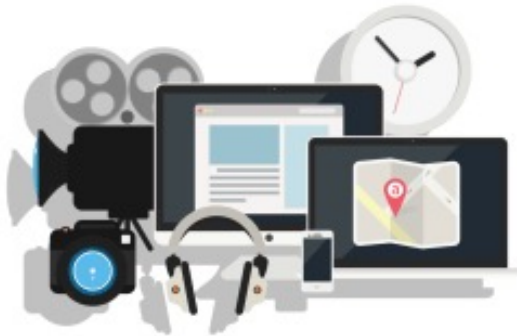
Real-world Artificial Intelligence

Digital intelligence

Multimedia

Image/video description

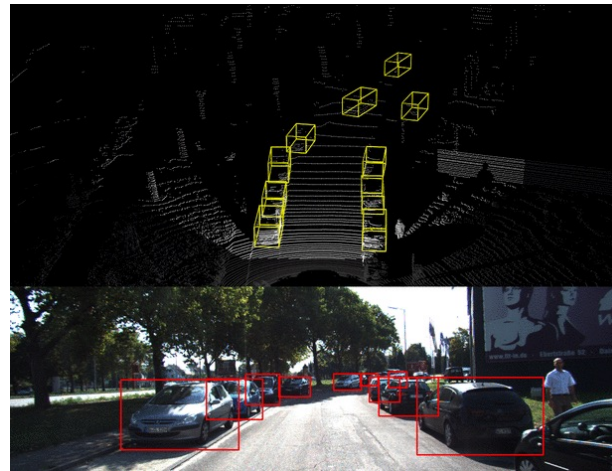
[Rui et al., 1999; Huang 2004]



Physical intelligence

Embodied AI, autonomous driving

[Xu et al., 2017; Szot et al., 2021]

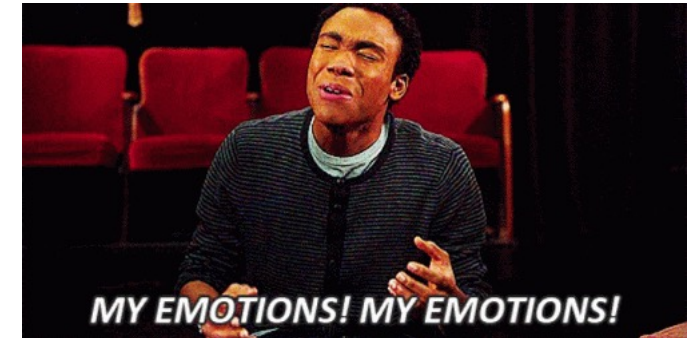


Social intelligence

Affective computing

Human-AI interaction

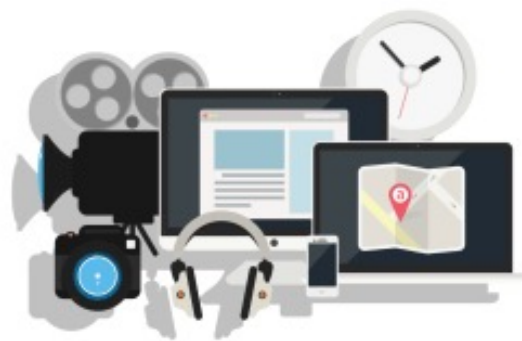
[Picard 1997; Jaimes & Sebe 2007]



Multimodal Artificial Intelligence

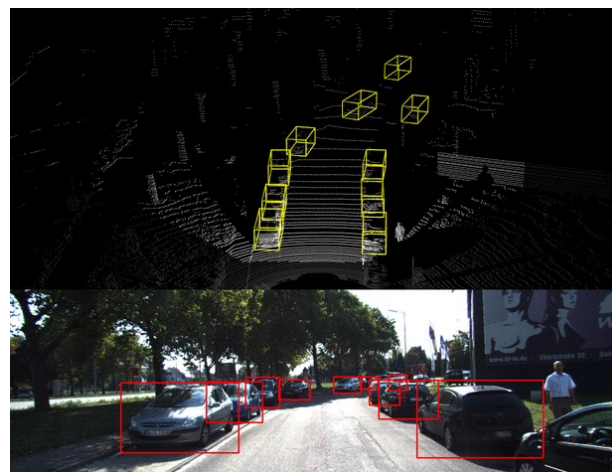
Digital intelligence

Multimedia
Image/video description
[Rui et al., 1999; Huang 2004]



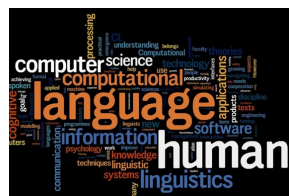
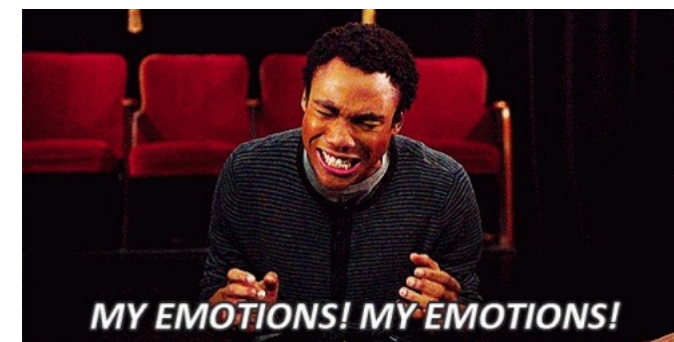
Physical intelligence

Embodied AI, autonomous driving
[Xu et al., 2017; Szot et al., 2021]



Social intelligence

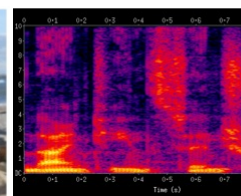
Affective computing
Human-AI interaction
[Picard 1997; Jaimes & Sebe 2007]



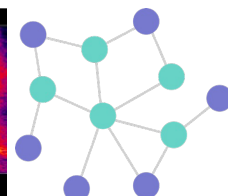
Language
(written)



Image



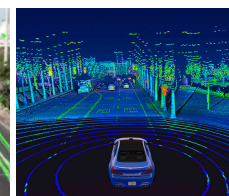
Audio



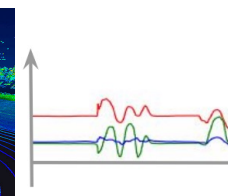
Graphs



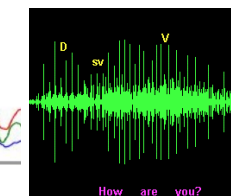
Video



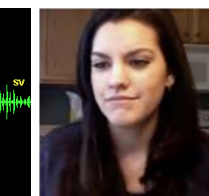
LIDAR



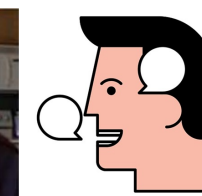
Sensors



Speech



Video
(faces)



Language
(spoken)

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**

Prior Research in 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 ...)
 - ❖ Focus of this talk: last 5 years



Behavioral Study of Multimodal



Language
and gestures

David McNeill

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

McGurk effect



Behavioral Study of Multimodal



Language
and gestures

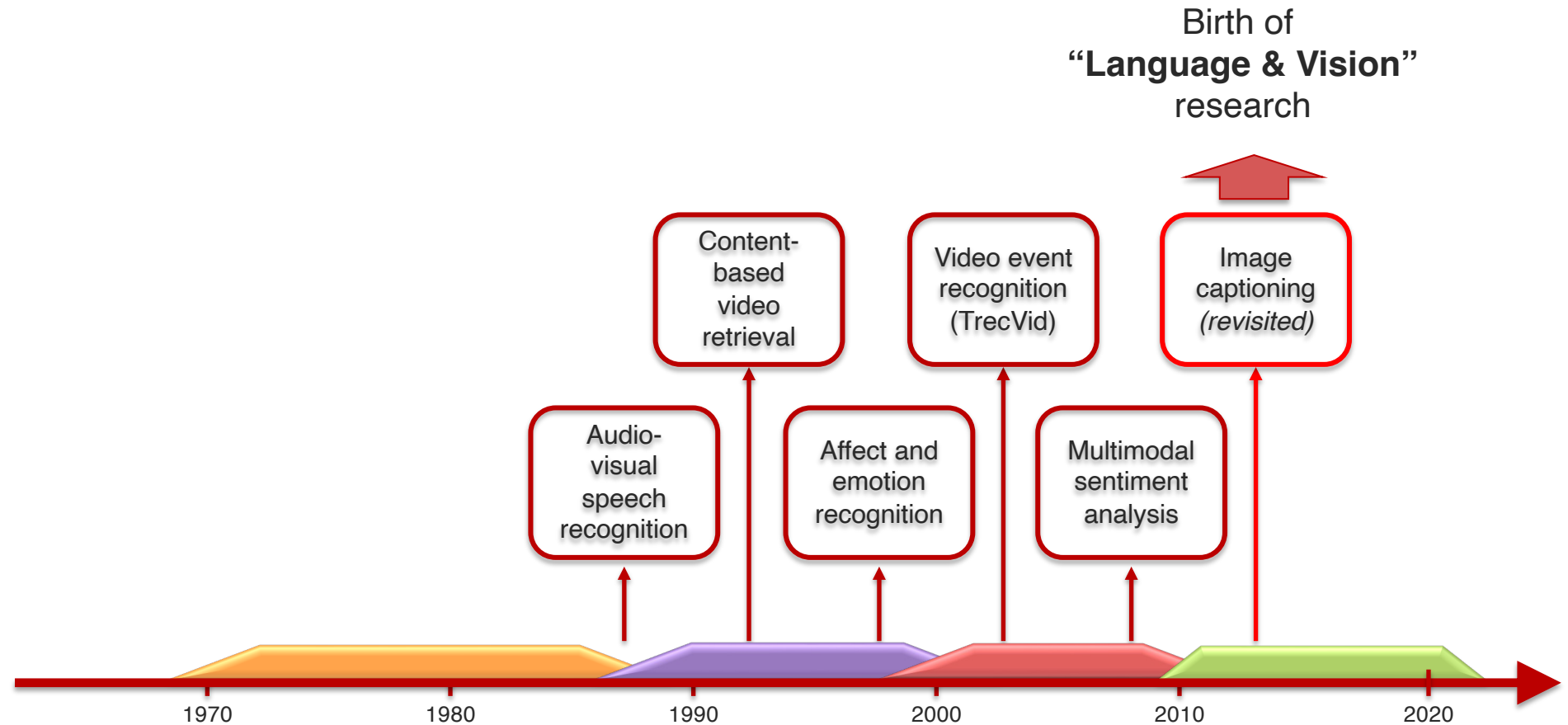
David McNeill

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

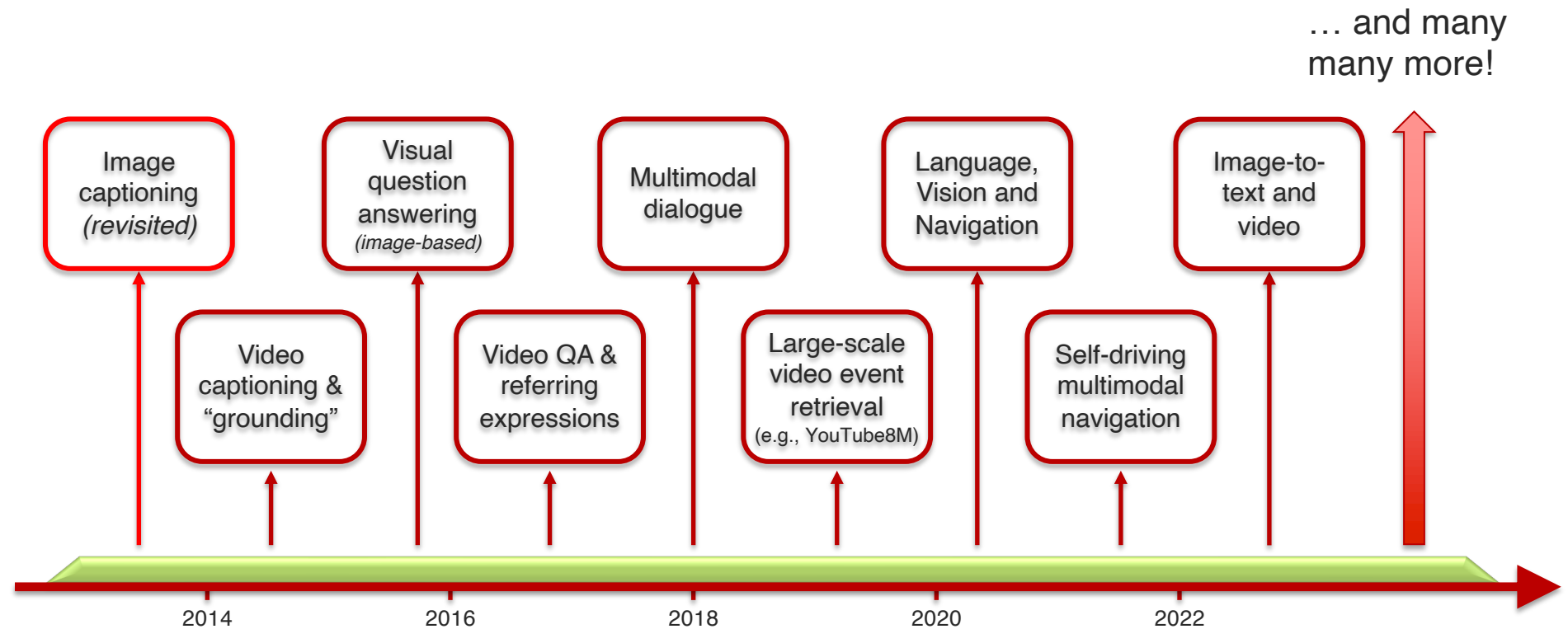
McGurk effect



Multimodal Research Tasks



Multimodal Research Tasks



Multimodal ML – Surveys, Tutorials and Courses

2016

Multimodal Machine Learning: A Survey and Taxonomy

Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency

(Arxiv 2017, IEEE TPAMI journal, February 2019)

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

Tutorials: CVPR 2016, ACL 2016, ICMI 2016, ...

Graduate-level courses:

Multimodal Machine Learning (11th edition)

<https://cmu-multicomp-lab.github.io/mmml-course/fall2020/>

Advanced Topics in Multimodal Machine Learning

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

2022

Foundations and Trends in Multimodal ML

Paul Liang, Amir Zadeh, and Louis-Philippe Morency

- 6 core challenges
- 50+ taxonomic classes
- 600+ referenced papers

Tutorials: CVPR 2022, NAACL 2022, ...

Updated graduate-level course:

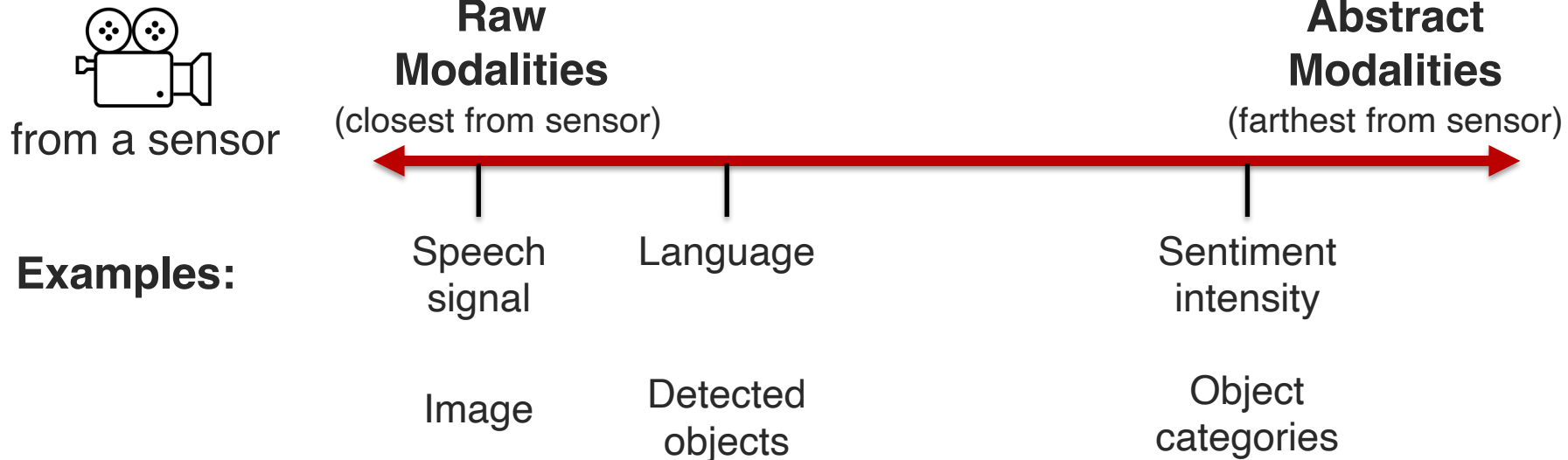
Multimodal Machine Learning (12th edition)

<https://cmu-multicomp-lab.github.io/mmml-course/fall2022/>

What is a Modality?

Modality

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



What is Multimodal?

A dictionary definition...

Multimodal: with multiple modalities

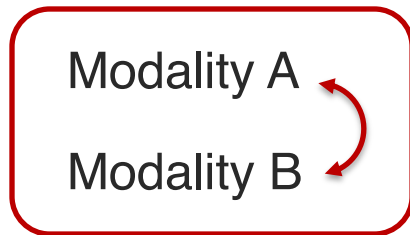
A research-oriented definition...

***Multimodal* is the science of**

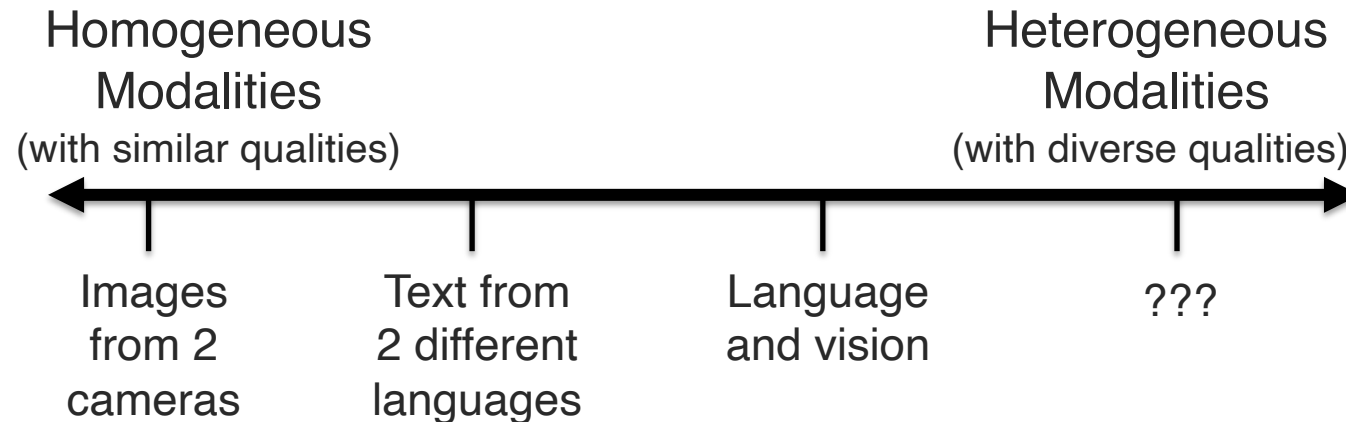
heterogeneous and interconnected data

Heterogeneous Modalities

Information present in different modalities will often show diverse qualities, structures, and representations.



Examples:



Abstract modalities are more likely to be homogeneous

Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



*A teacup on the right of a laptop
in a clean room.*

Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



A *teacup* on the *right* of a *laptop* in a *clean room*.

① **Distribution:** discrete or continuous, support



● {*teacup*, *right*, *laptop*, *clean*, *room*}

Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



*A teacup on the right of a laptop
in a clean room.*

2 **Granularity:** sampling rate and frequency



objects per image



words per minute

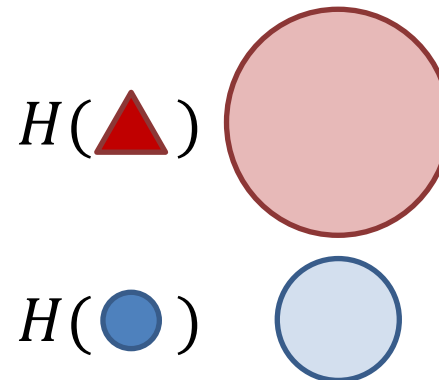
Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



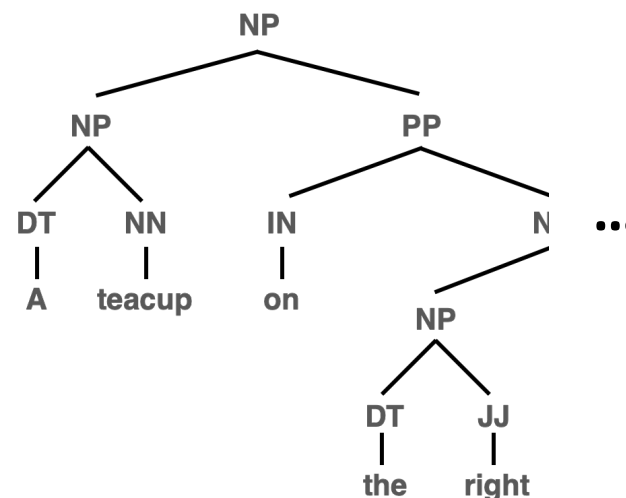
*A teacup on the right of a laptop
in a clean room.*

3 **Information:** entropy and density

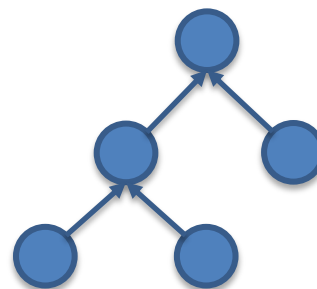
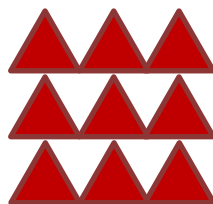


Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



4 **Structure:** static, temporal, spatial, hierarchical



Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



*A teacup on the right of a laptop
in a clean room.*

5 **Noise:** uncertainty, signal-to-noise ratio, missing data



teacup → **teacip**

right → **rihjt**

Dimensions of Heterogeneity

Information present in different modalities will often show diverse qualities, structures, and representations.



*A teacup on the right of a laptop
in a clean room.*

6

Relevance: task relevance, context dependence

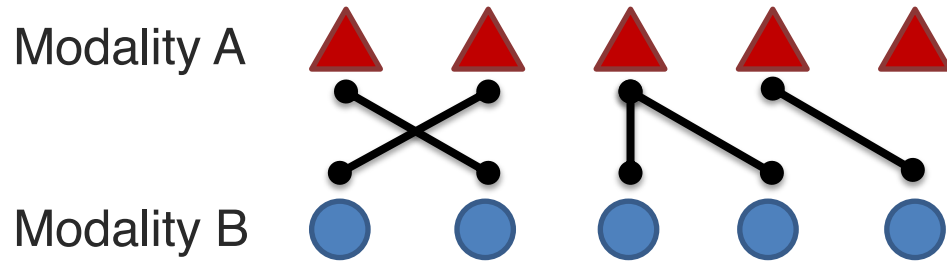


→ recreational
→ living room
→ right-handed

*A teacup on the
right of a laptop
in a clean room.*

→ workspace
→ study room

Interconnected Modalities



① Connections

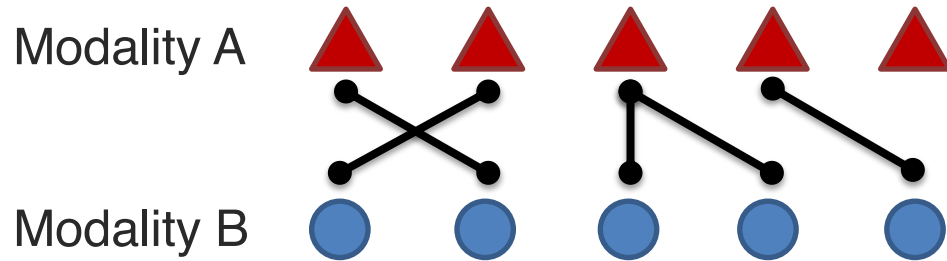
Which elements are connected and why?

② Cross-modal interactions

How are connected elements interacting during inference?

Modalities are often related and share complementary information that interact

Interconnected Modalities



① Connections

*Which elements are connected
and why?*



*A teacup on the right of a laptop
in a clean room.*

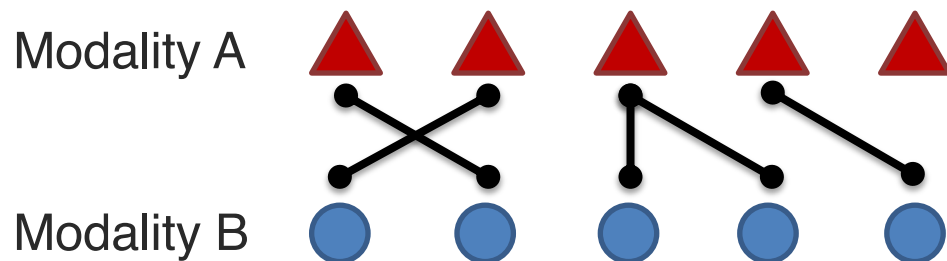


teacup



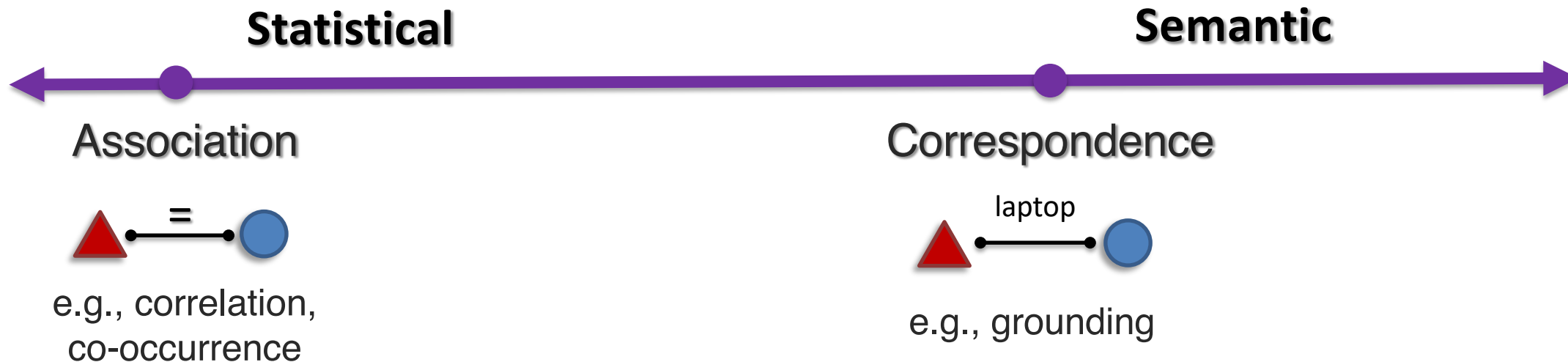
laptop

Interconnected Modalities

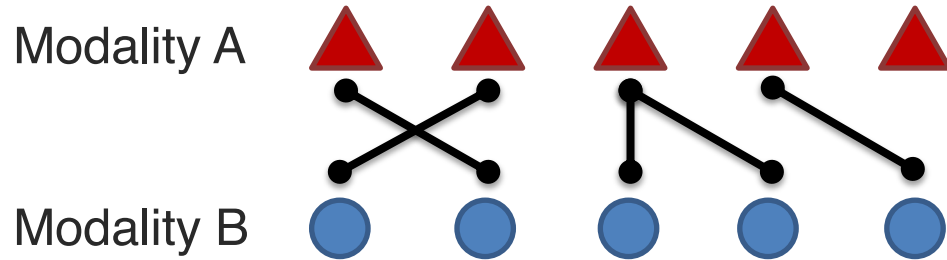


① Connections

*Which elements are connected
and why?*



Interconnected Modalities



① Connections

*Which elements are connected
and why?*



*A teacup on the right of a laptop
in a clean room.*

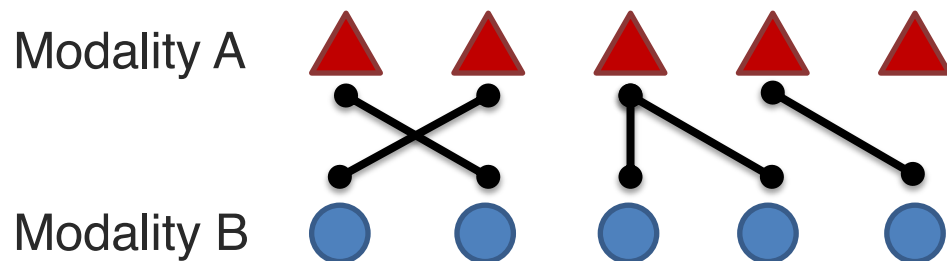


clean



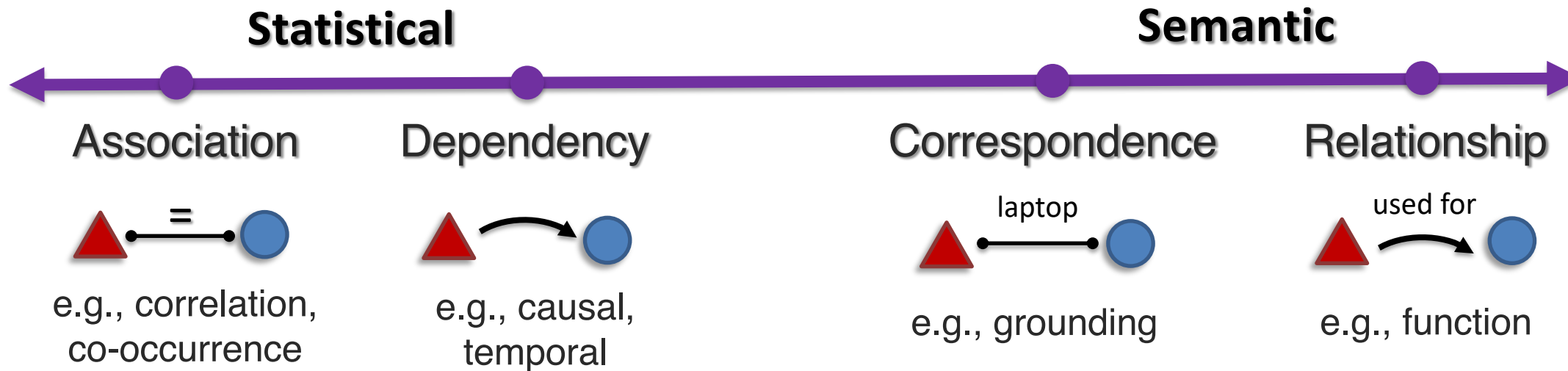
room

Interconnected Modalities

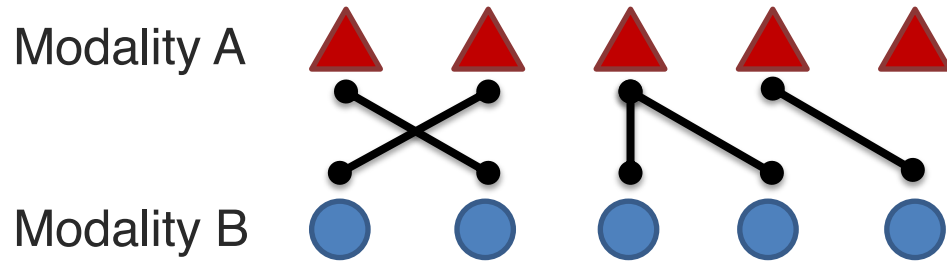


① Connections

*Which elements are connected
and why?*

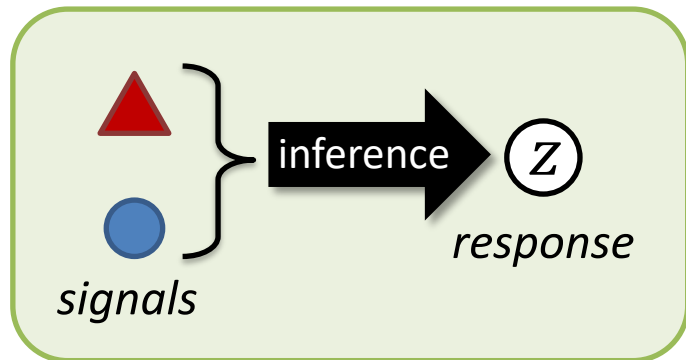


Interconnected Modalities



① Connections

Which elements are connected and why?

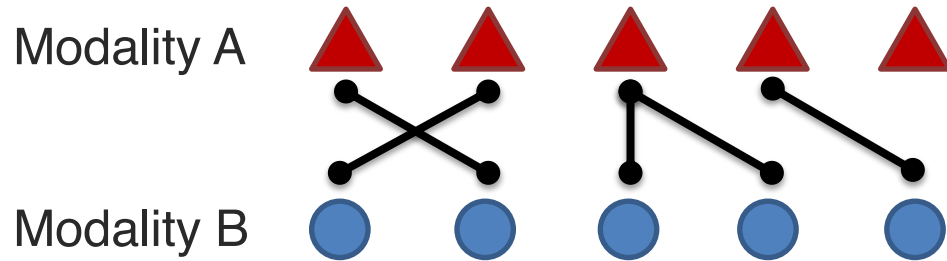


Interactions happen during inference!

② Cross-modal interactions

How are connected elements interacting during inference?

Interconnected Modalities



Is this indoors?

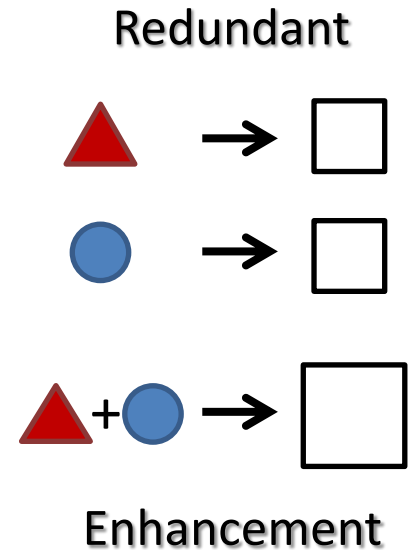
A teacup on the right of a laptop in a clean room.

inference → *Yes!*

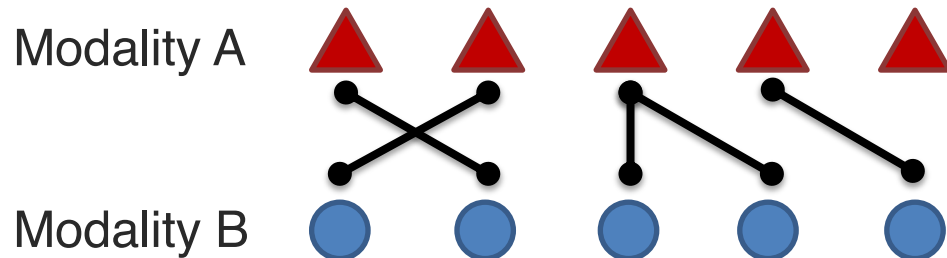
inference → *Yes!*

② Cross-modal interactions

How are connected elements interacting during inference?



Interconnected Modalities



Is this indoors?

A teacup on the right of a laptop in a clean room.

inference

Yes!

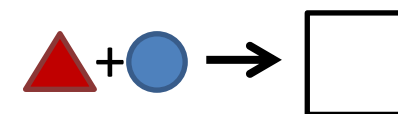
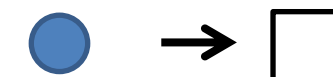
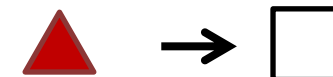
inference

Yes!

② Cross-modal interactions

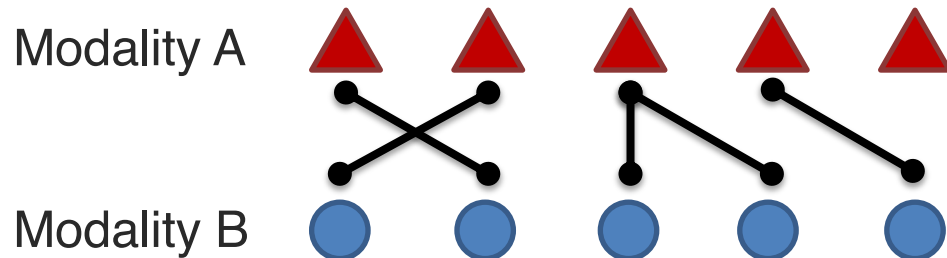
How are connected elements interacting during inference?

Redundant



Enhancement

Interconnected Modalities



*Is this
a living
room?*

*A teacup on the right of a
laptop in a clean room.*



Yes!

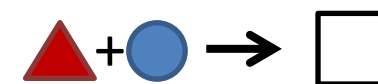
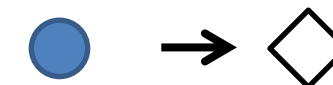
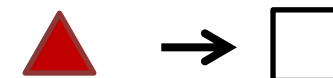


*No, probably
study room.*

② Cross-modal interactions

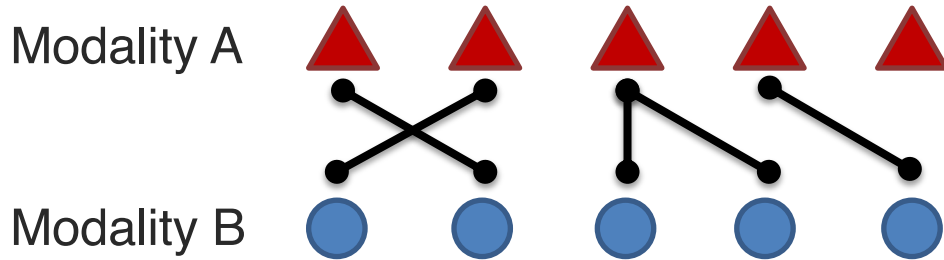
*How are connected elements
interacting during inference?*

Non-redundant



Dominance

Interconnected Modalities



*Is this
a living
room?*

*A teacup on the right of a
laptop in a clean room.*

inference

Yes!

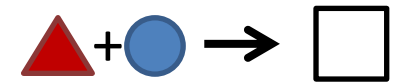
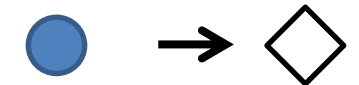
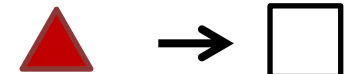
inference

*No, probably
study room.*

② Cross-modal interactions

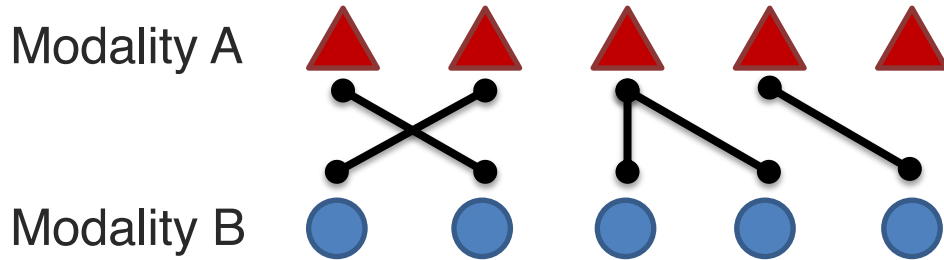
*How are connected elements
interacting during inference?*

Non-redundant



Dominance

Interconnected Modalities



***Should I
work here?***

*A teacup on the right of a
laptop in a clean room.*



*Maybe? Comfy
sofa but table's
too small.*

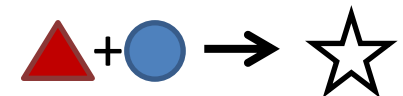
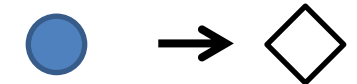
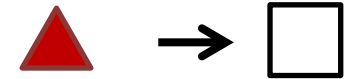


*Maybe? Clean
and there's tea.*

② Cross-modal interactions

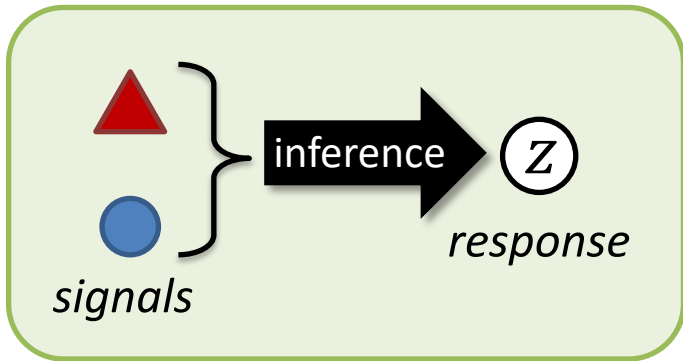
*How are connected elements
interacting during inference?*

Non-redundant

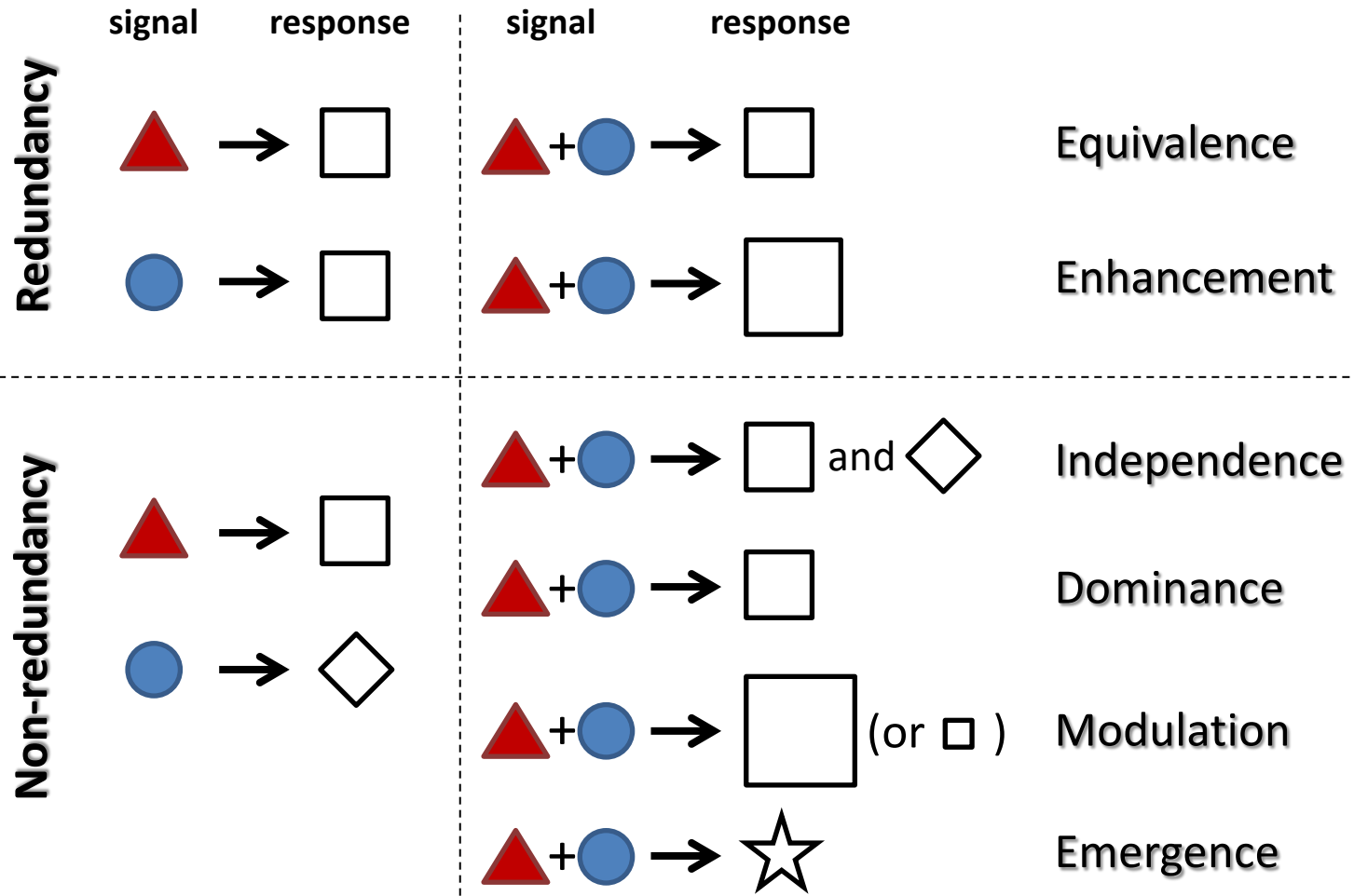


Emergence

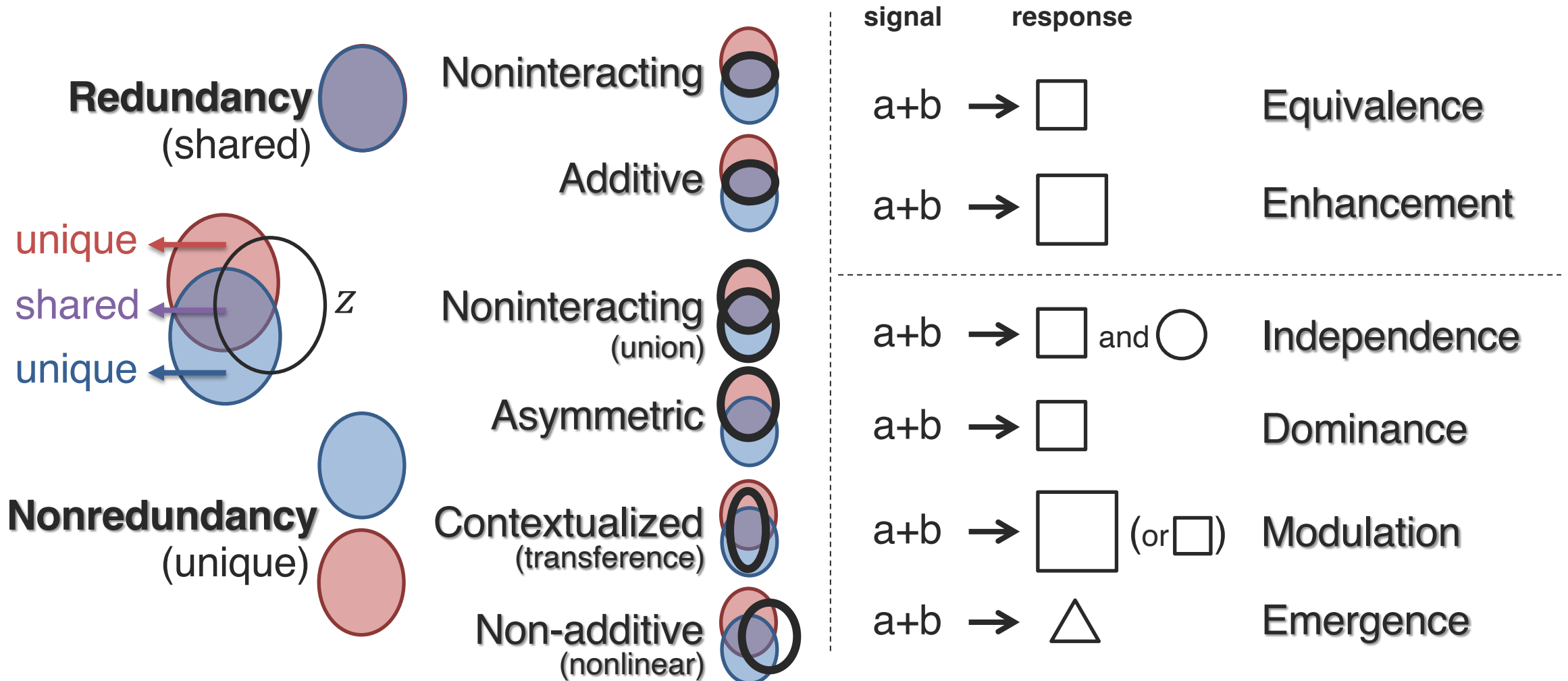
Cross-modal Interactions – A Behavioral Science View



Multimodal Communication



Cross-modal Interaction Mechanics



*What is
Multimodal?*



Why is it hard?



What is next?

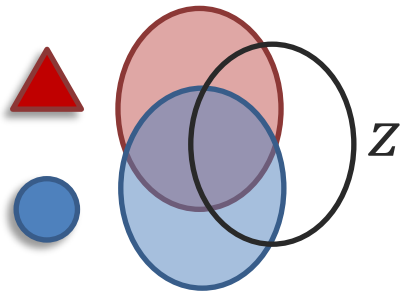
Heterogeneous



Connected

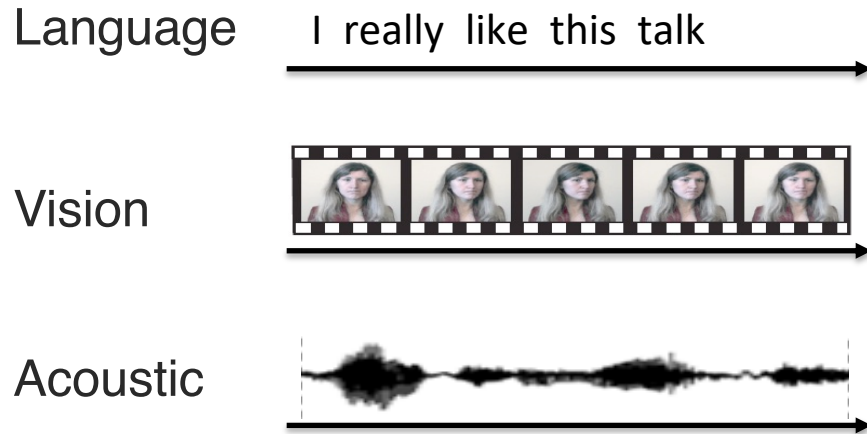


Interacting

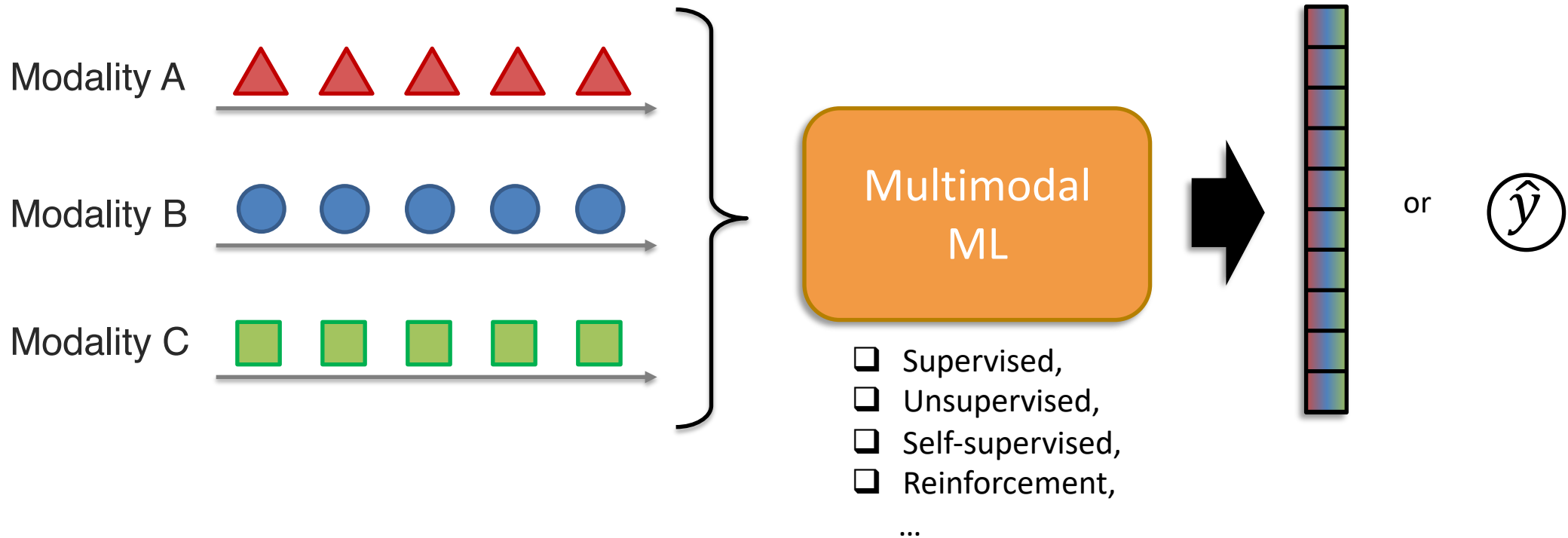


**Multimodal is the scientific
study of heterogeneous and
interconnected data 😊**

Multimodal Machine Learning



Multimodal Machine Learning



Multimodal Machine Learning

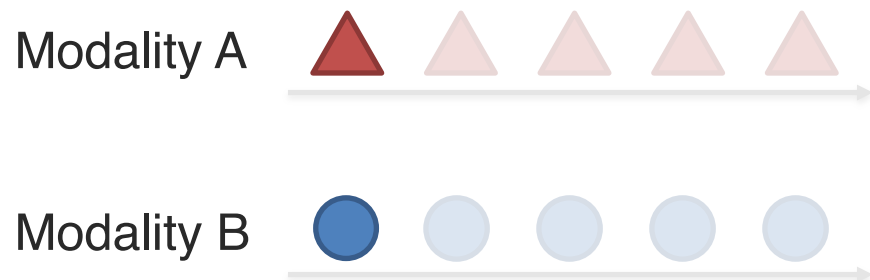
*What are the **core multimodal technical challenges**,
understudied in conventional machine learning?*

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:

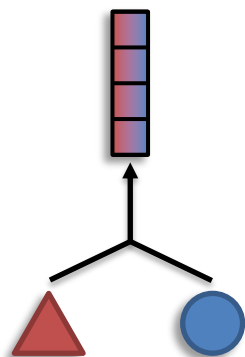


Challenge 1: Representation

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

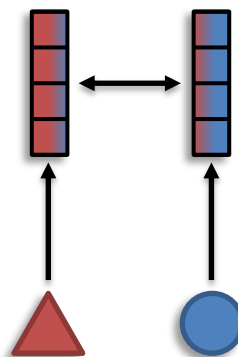
Sub-challenges:

Fusion



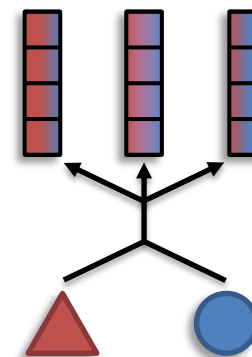
modalities $>$ # representations

Coordination



modalities = # representations

Fission



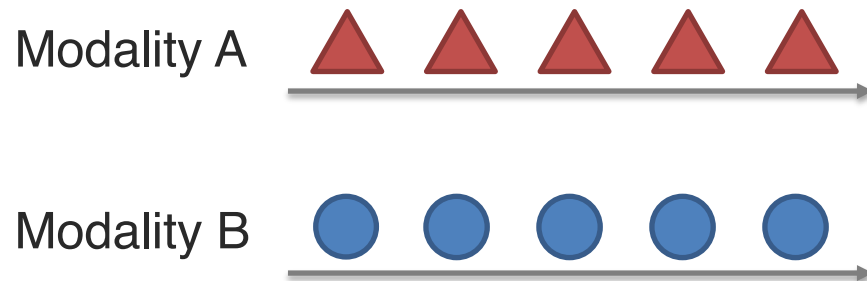
modalities $<$ # representations

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**

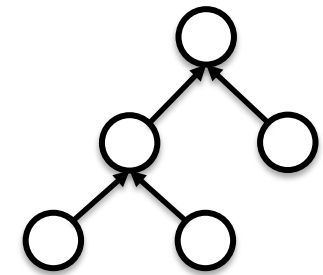
Elements with temporal structure:



Other structured examples:



Spatial



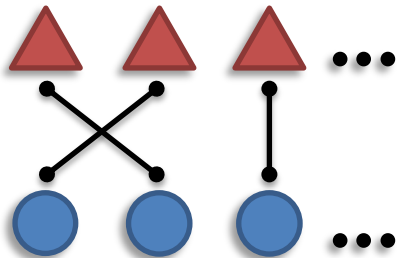
Hierarchical

Challenge 2: Alignment

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

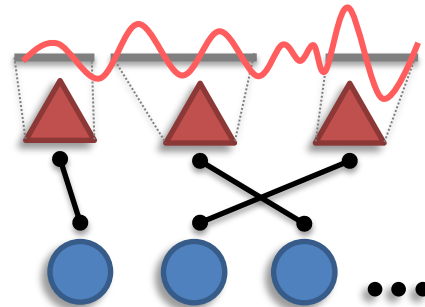
Sub-challenges:

Discrete connections



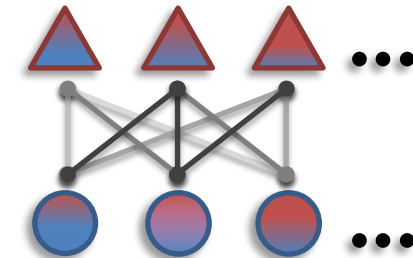
Explicit alignment
(e.g., grounding)

Continuous alignment



Granularity of
individual elements

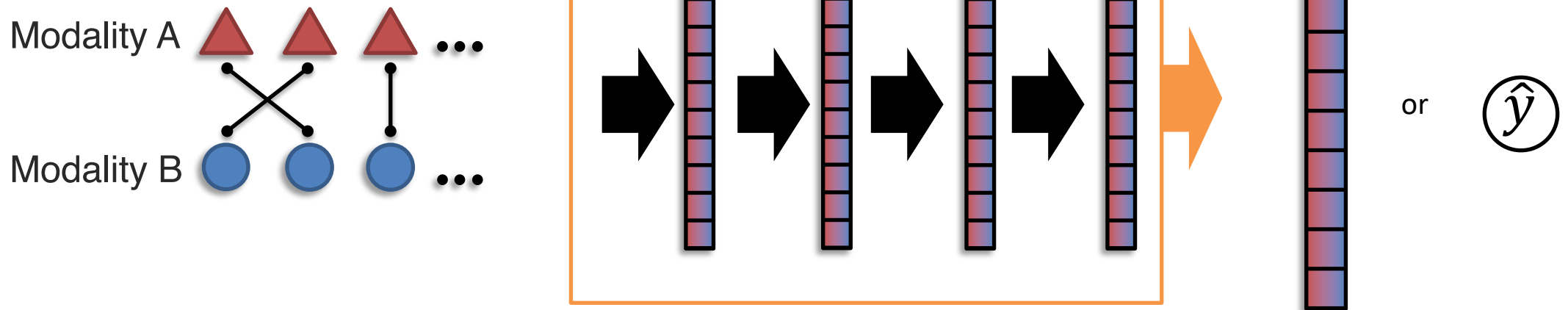
Contextualized representation



Implicit alignment
+ representation

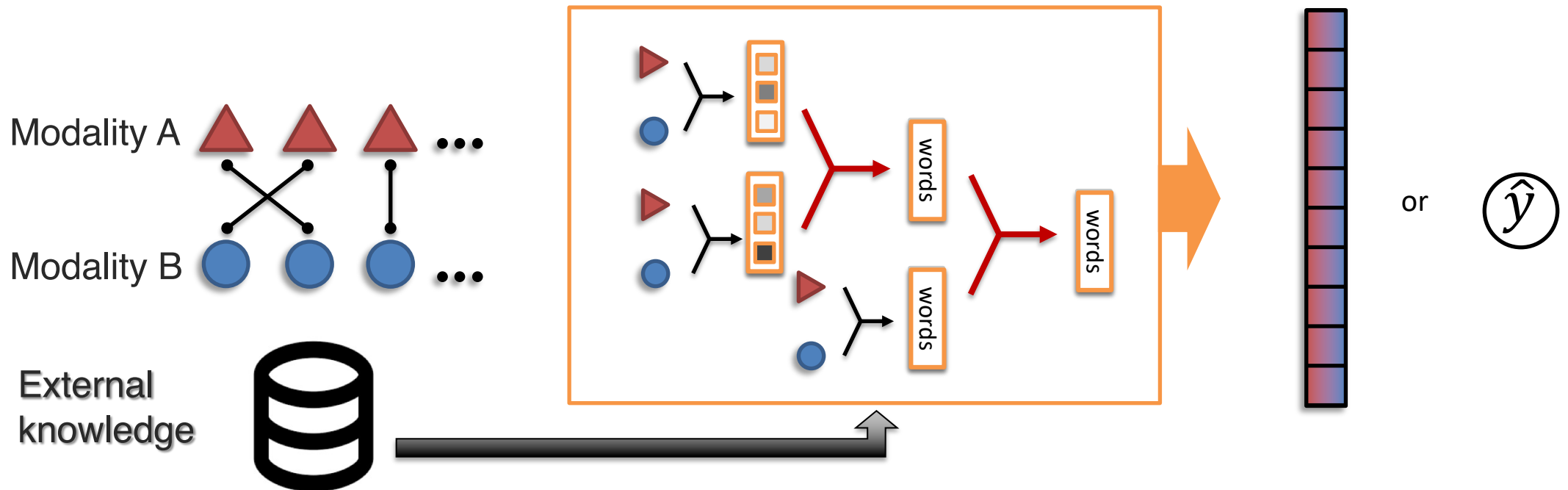
Challenge 3: Reasoning

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.

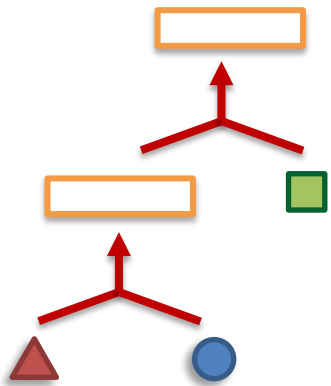


Challenge 3: Reasoning

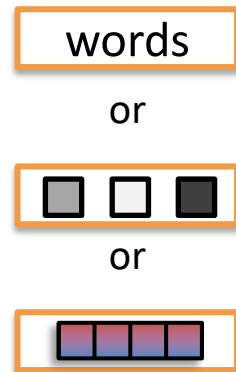
Definition: Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.

Sub-challenges:

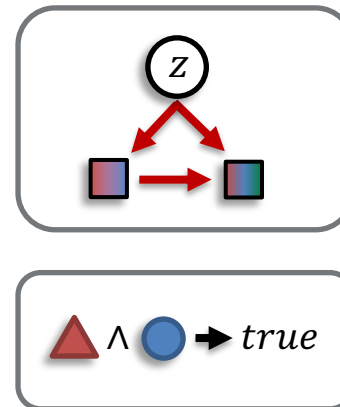
Structure modeling



Intermediate concepts



Inference paradigm



External knowledge

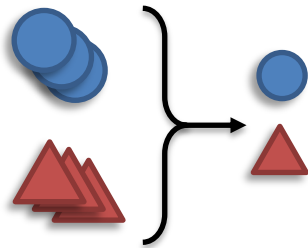


Challenge 4: Generation

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

Sub-challenges:

Summarization



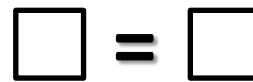
Reduction



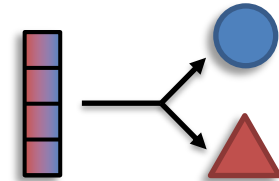
Translation



Maintenance



Creation



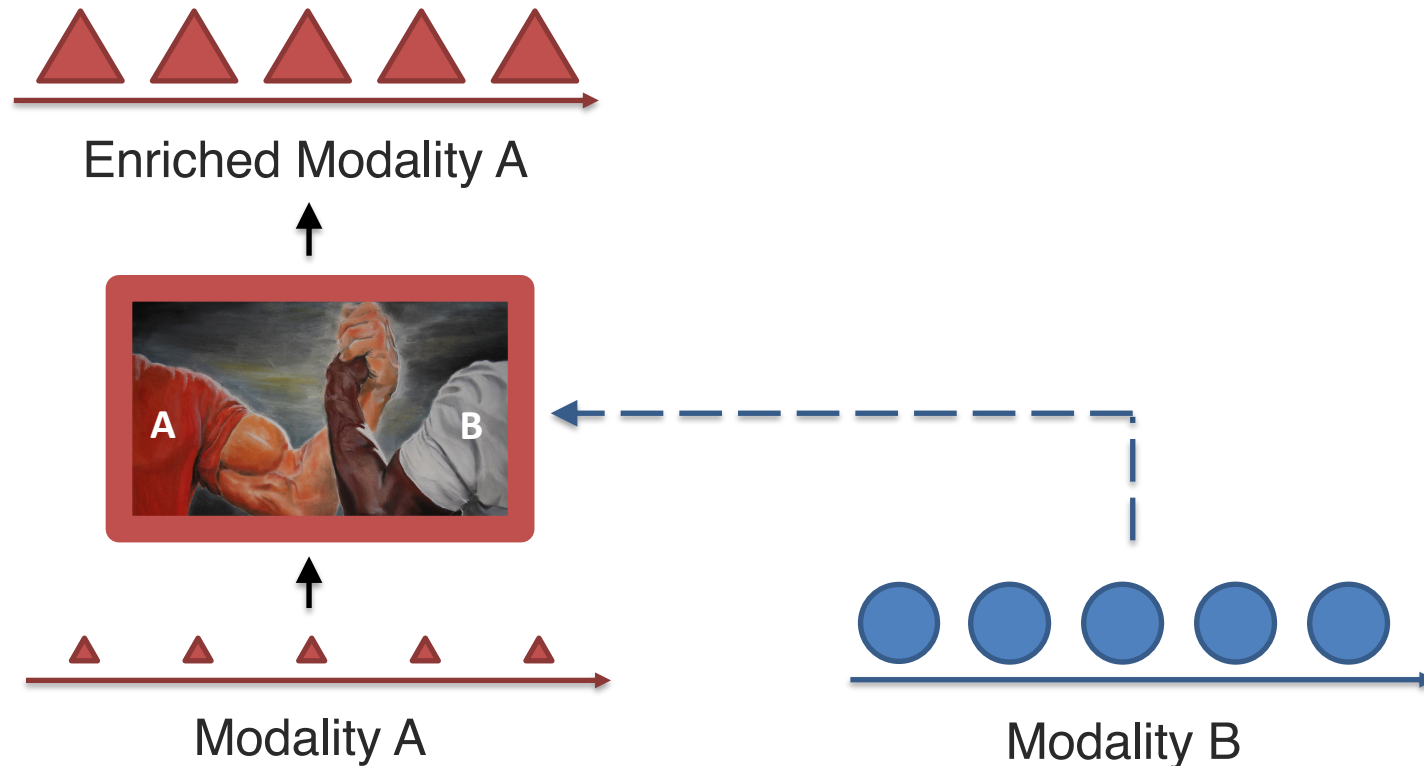
Expansion



Information:
(content)

Challenge 5: Transference

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

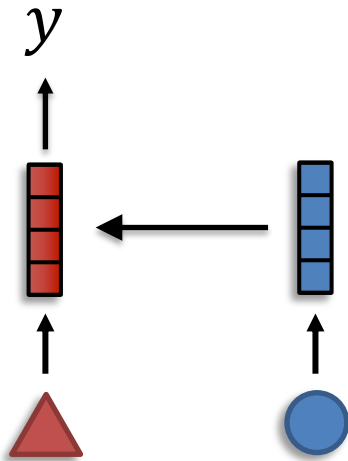


Challenge 5: Transference

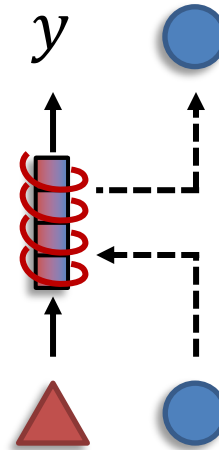
Definition: Transfer knowledge between modalities, usually to help the target modality which may be noisy or with limited resources.

Sub-challenges:

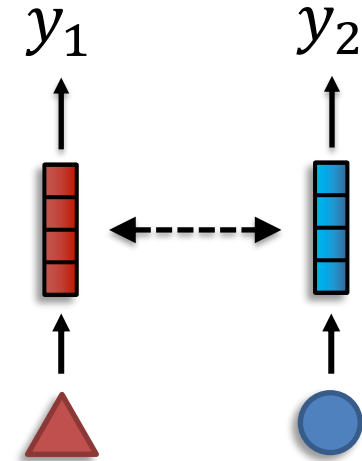
Transfer



Co-learning



Model Induction

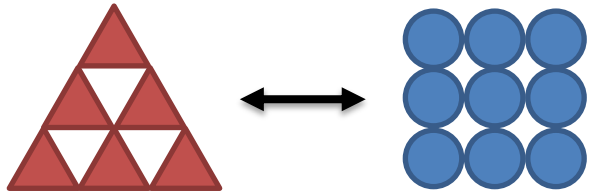


Challenge 6: Quantification

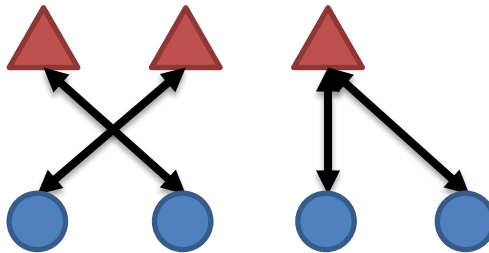
Definition: Empirical and theoretical study to better understand heterogeneity, cross-modal interactions, and the multimodal learning process.

Sub-challenges:

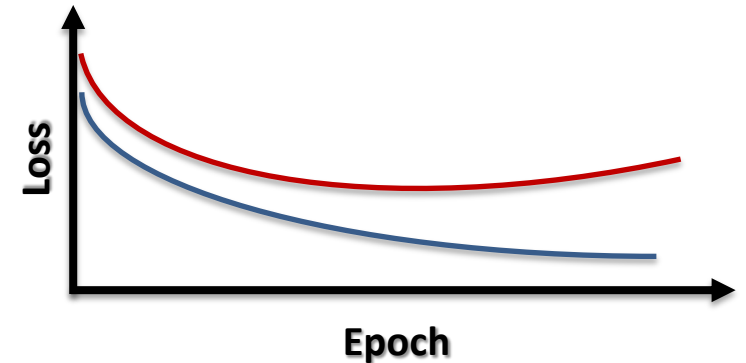
Heterogeneity



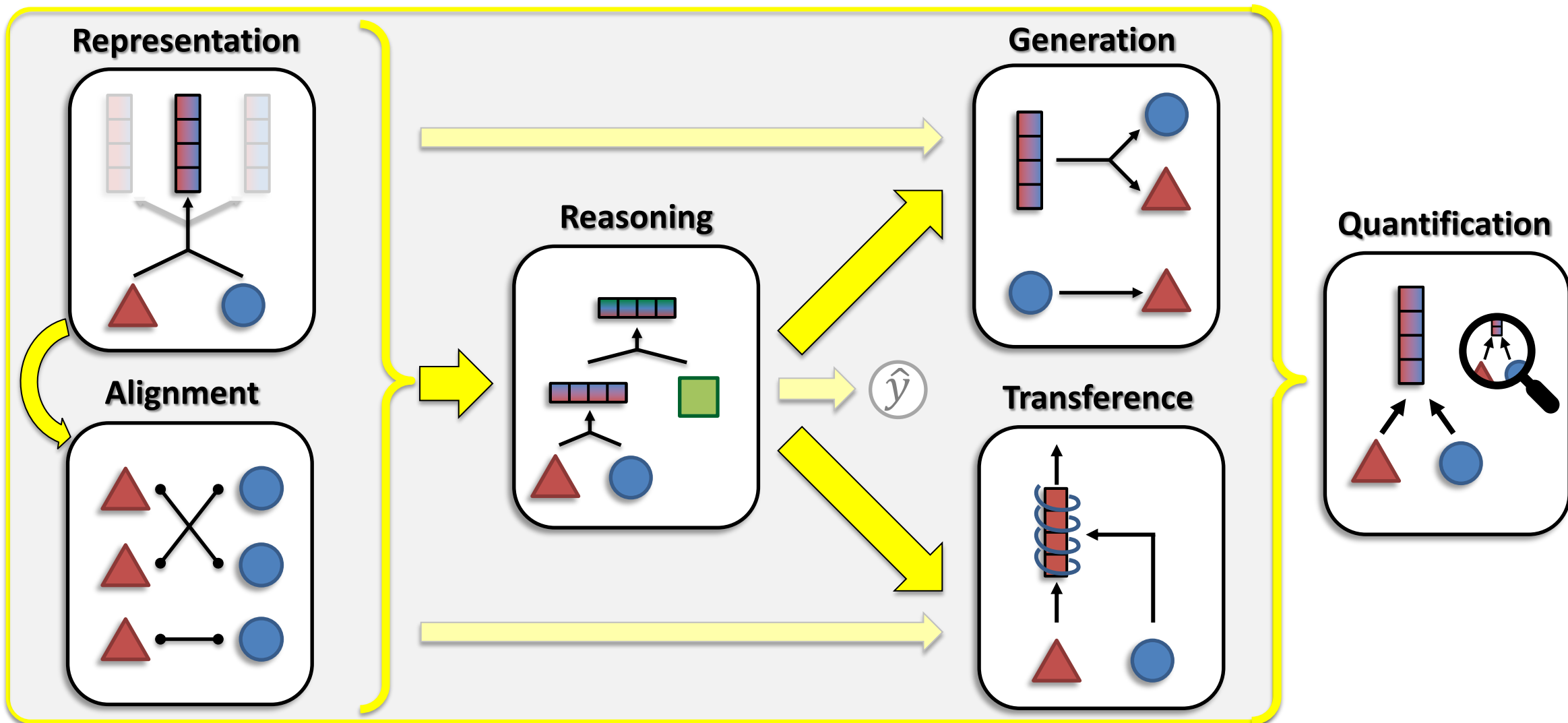
Interactions



Learning



Core Multimodal Challenges



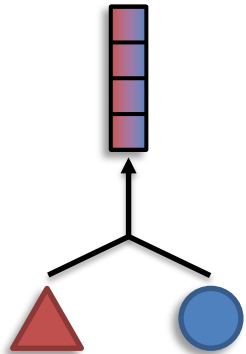
Challenge 1: Representation

Challenge 1: Representation

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

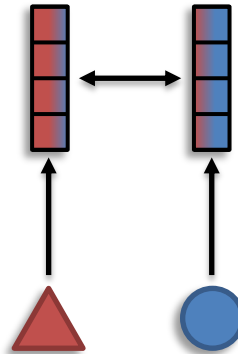
Sub-challenges:

Fusion



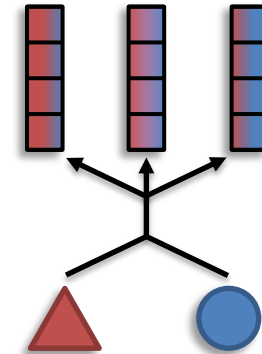
modalities $>$ # representations

Coordination



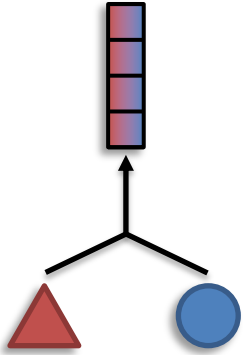
modalities = # representations

Fission



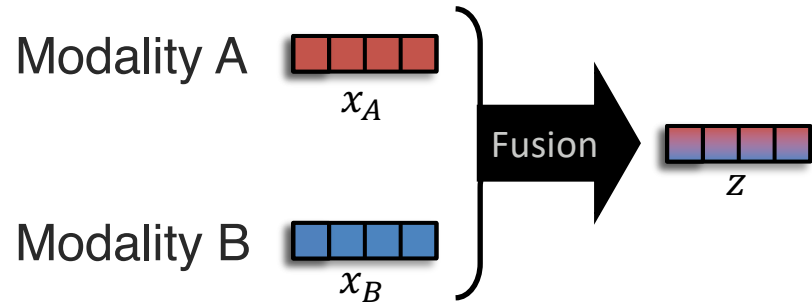
modalities $<$ # representations

Sub-Challenge 1a: Representation Fusion



Definition: Learn a joint representation that models cross-modal interactions between individual elements of different modalities.

Concepts for Representation Fusion



Goal: Model *cross-modal interactions* between the multimodal elements

→ Let's study the univariate case first
↳ (only 1-dimensional features)

Linear regression:

$$Z = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

constant Additive terms Multiplicative term error

① Additive interaction:

$$Z = w_1 x_A + w_2 x_B + \epsilon$$

② Multiplicative interaction:

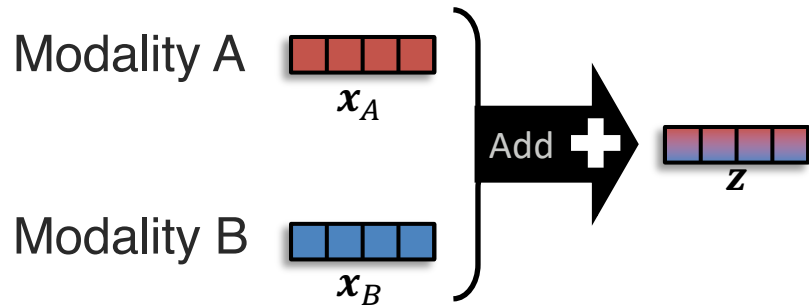
$$Z = w_3 (x_A \times x_B) + \epsilon$$

③ Additive and multiplicative interactions:

$$Z = w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

Additive Fusion Back to multivariate case!

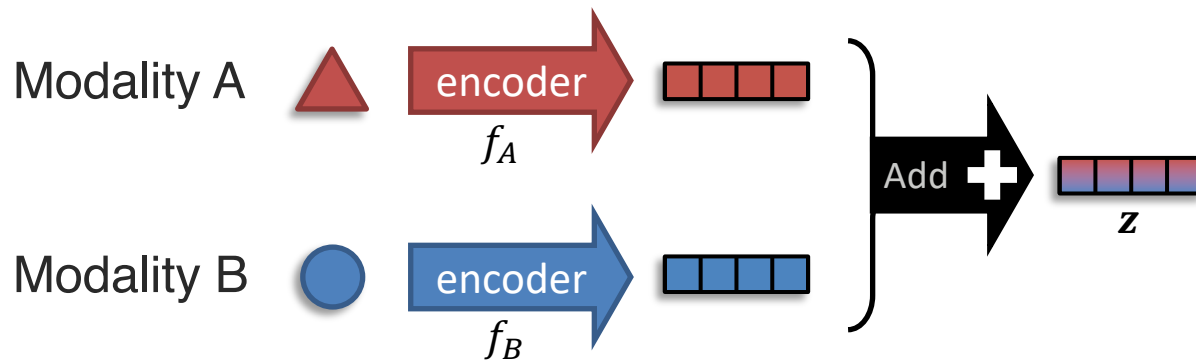
 (multi-dimensional features)



Additive fusion:

$$z = w_1 x_A + w_2 x_B$$

With unimodal encoders:

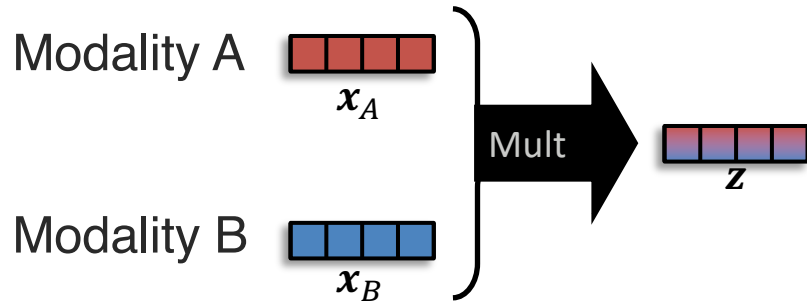


Additive fusion:

$$z = f_A(\triangle) + f_B(\circ)$$

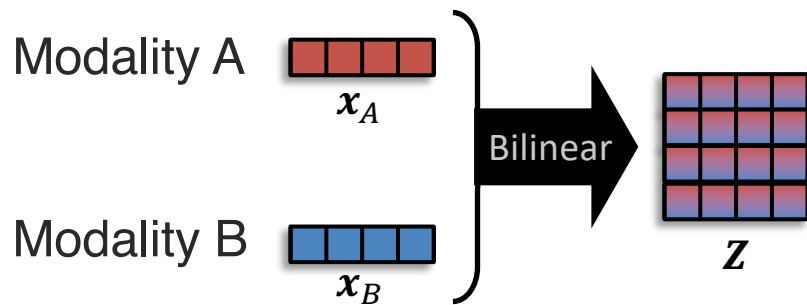
 It could be seen as an ensemble approach
(late fusion)

Multiplicative Fusion



Multiplicative fusion:

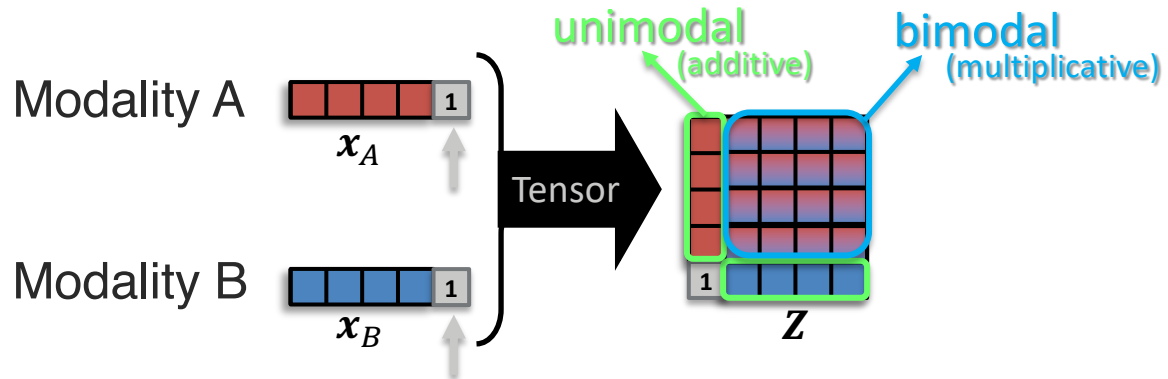
$$z = w(x_A \times x_B)$$



Bilinear Fusion:

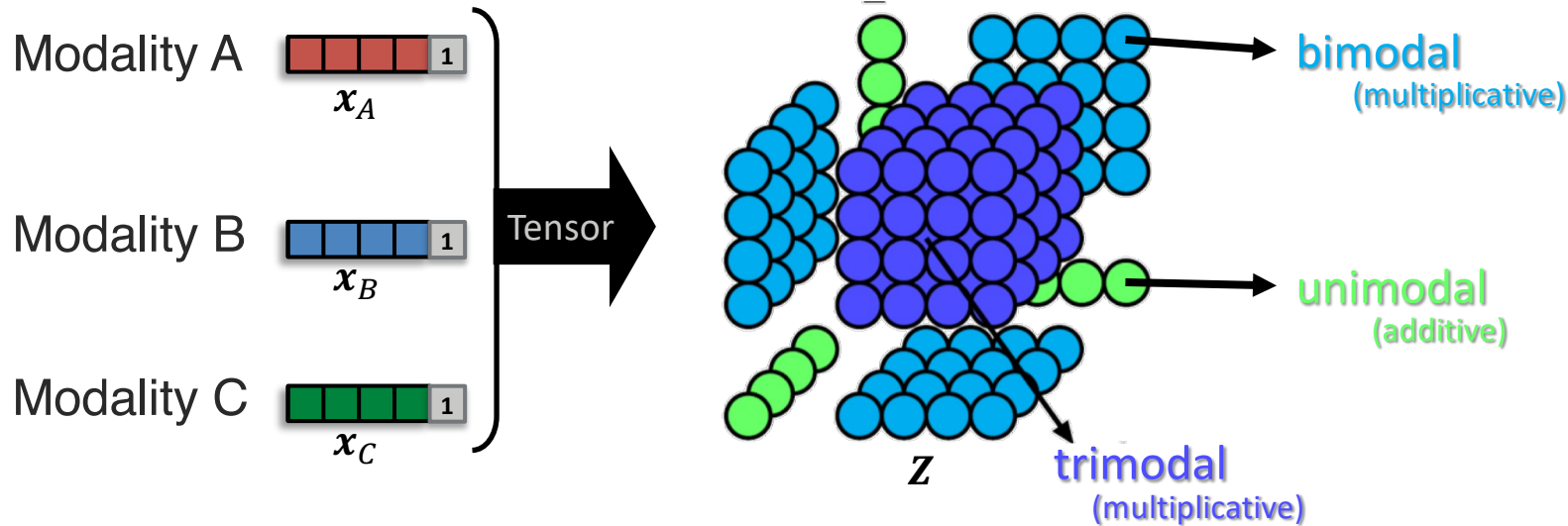
$$Z = w(x_A^T \cdot x_B)$$

Tensor Fusion



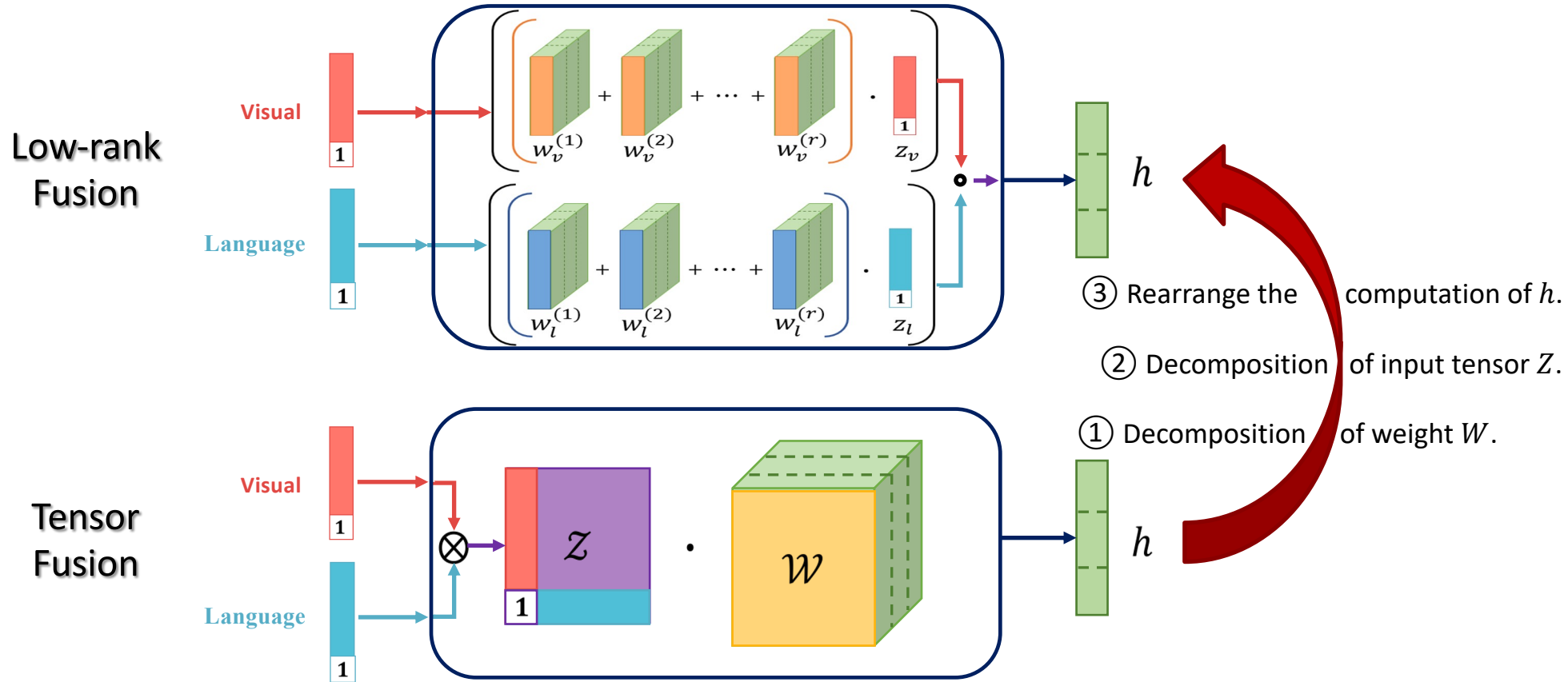
Tensor Fusion (bimodal):

$$Z = w([\mathbf{x}_A \ 1]^T \cdot [\mathbf{x}_B \ 1])$$

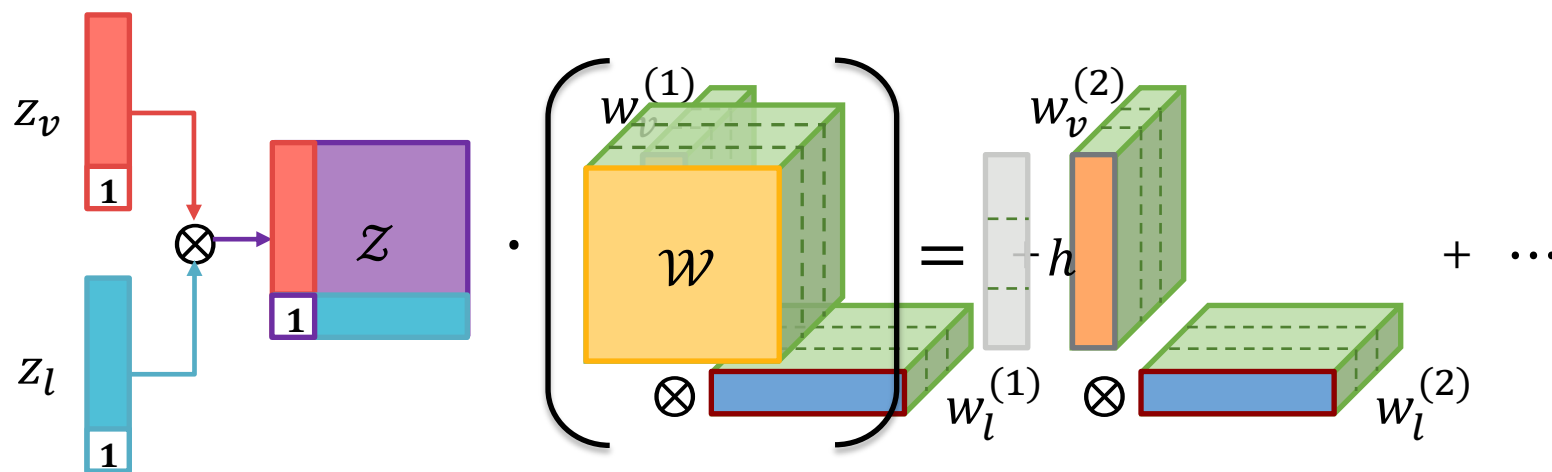


... but the weight matrix may end up quite large!

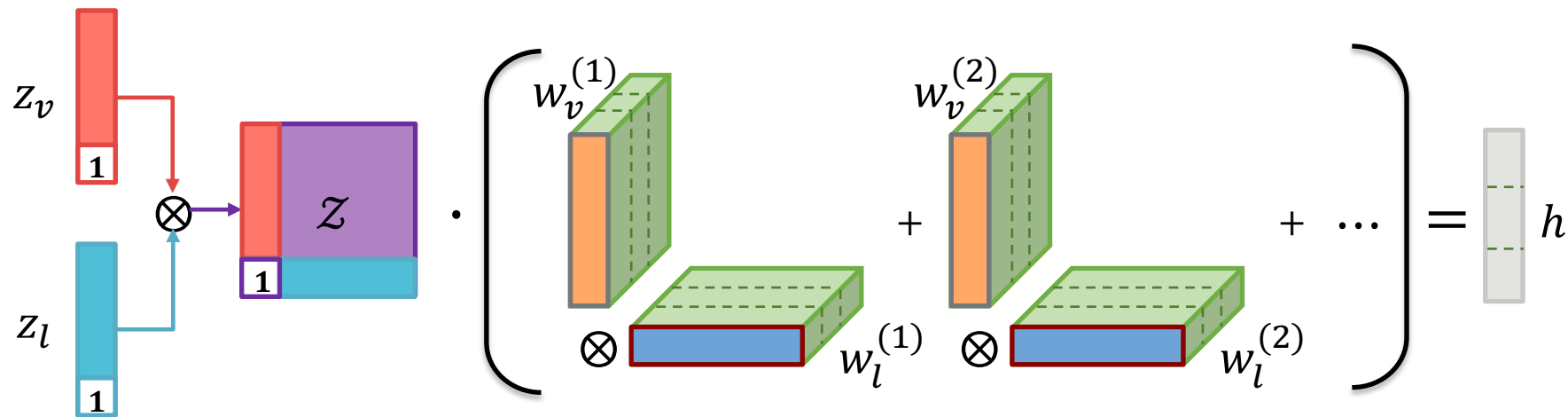
Low-rank Fusion



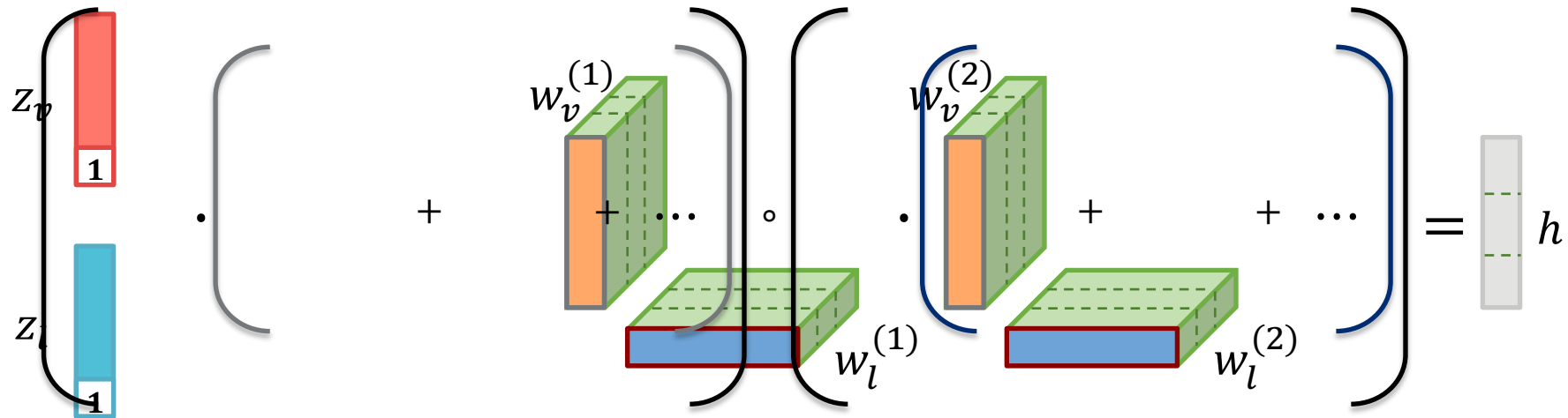
Low-rank Fusion



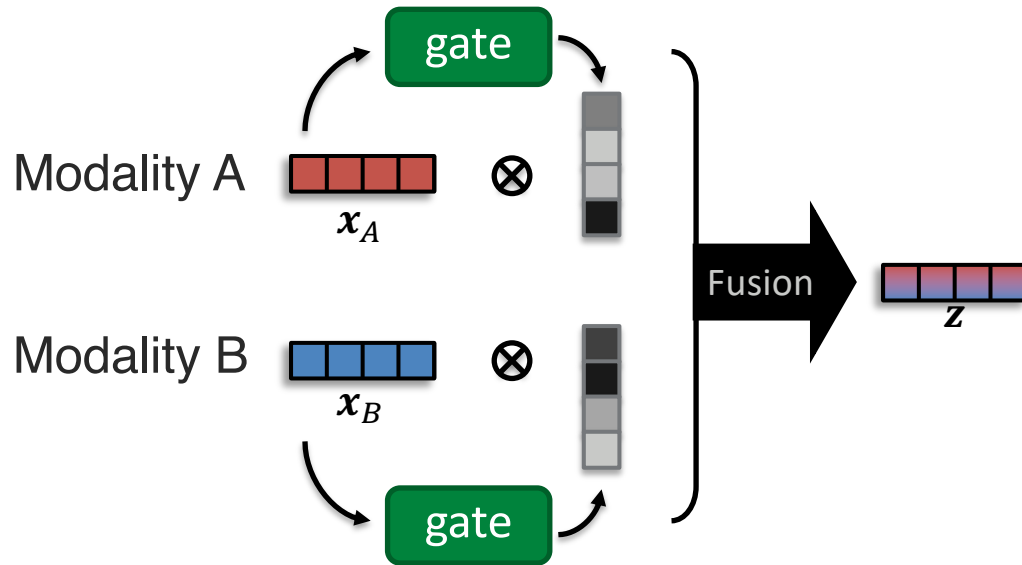
Low-rank Fusion



Low-rank Fusion



Gated Fusion



Example with additive fusion:

$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

→ g_A and g_B can be seen as attention functions

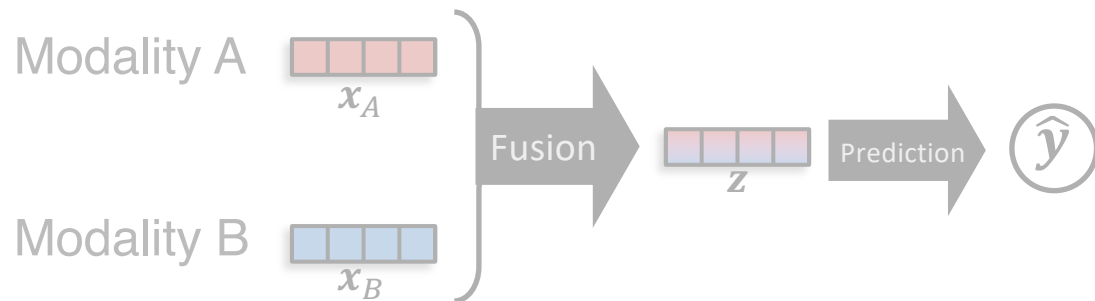
Linear: $x_A w_1 \cdot (x_B W_2)^T$

Nonlinear: $f_A(x_A) \cdot (f_B(x_B))^T$

Kernel: $k(x_A, x_B)$

- Linear
- Polynomial
- Exponential
- RBF

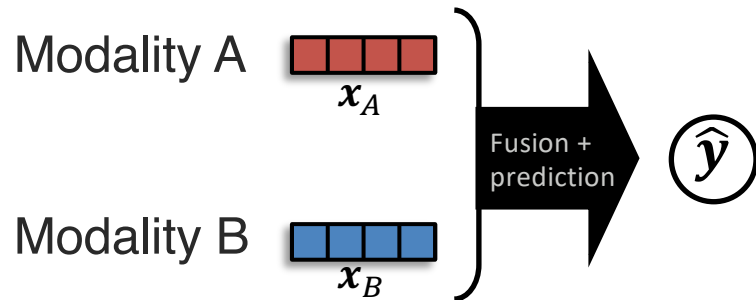
Nonlinear Fusion



Nonlinear fusion:

$$\hat{y} = f(x_A, x_B) \in \mathbb{R}^d$$

where f could be a multi-layer perceptron or any nonlinear model

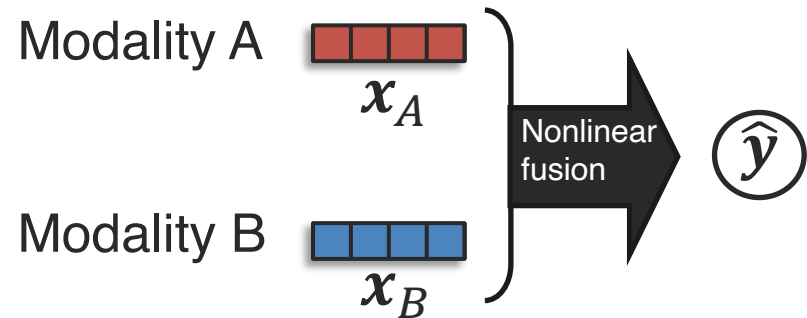


→ This could be seen as *early fusion*:

$$\hat{y} = f([x_A, x_B])$$

... but will our neural network learn the nonlinear interactions?

Measuring Non-Additive Interactions



Nonlinear fusion:

$$\hat{\mathbf{y}} = f(\mathbf{x}_A, \mathbf{x}_B)$$

Projection?

Additive fusion:

$$\hat{\mathbf{y}}' = f_A(\mathbf{x}_A) + f_B(\mathbf{x}_B)$$

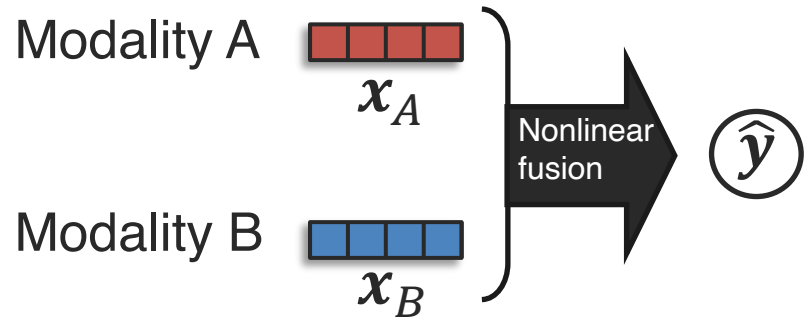
Projection from nonlinear to additive (using EMAP):

$$\tilde{f}(\mathbf{x}_A, \mathbf{x}_B) = \underbrace{\mathbb{E}_{\mathbf{x}_B} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_A(\mathbf{x}_A)} + \underbrace{\mathbb{E}_{\mathbf{x}_A} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_B(\mathbf{x}_B)}$$

Modality A + Modality B

Additive fusion
(approximation)

Measuring Non-Additive Interactions



Nonlinear fusion:

$$\hat{\mathbf{y}} = f(\mathbf{x}_A, \mathbf{x}_B)$$

EMAP projection

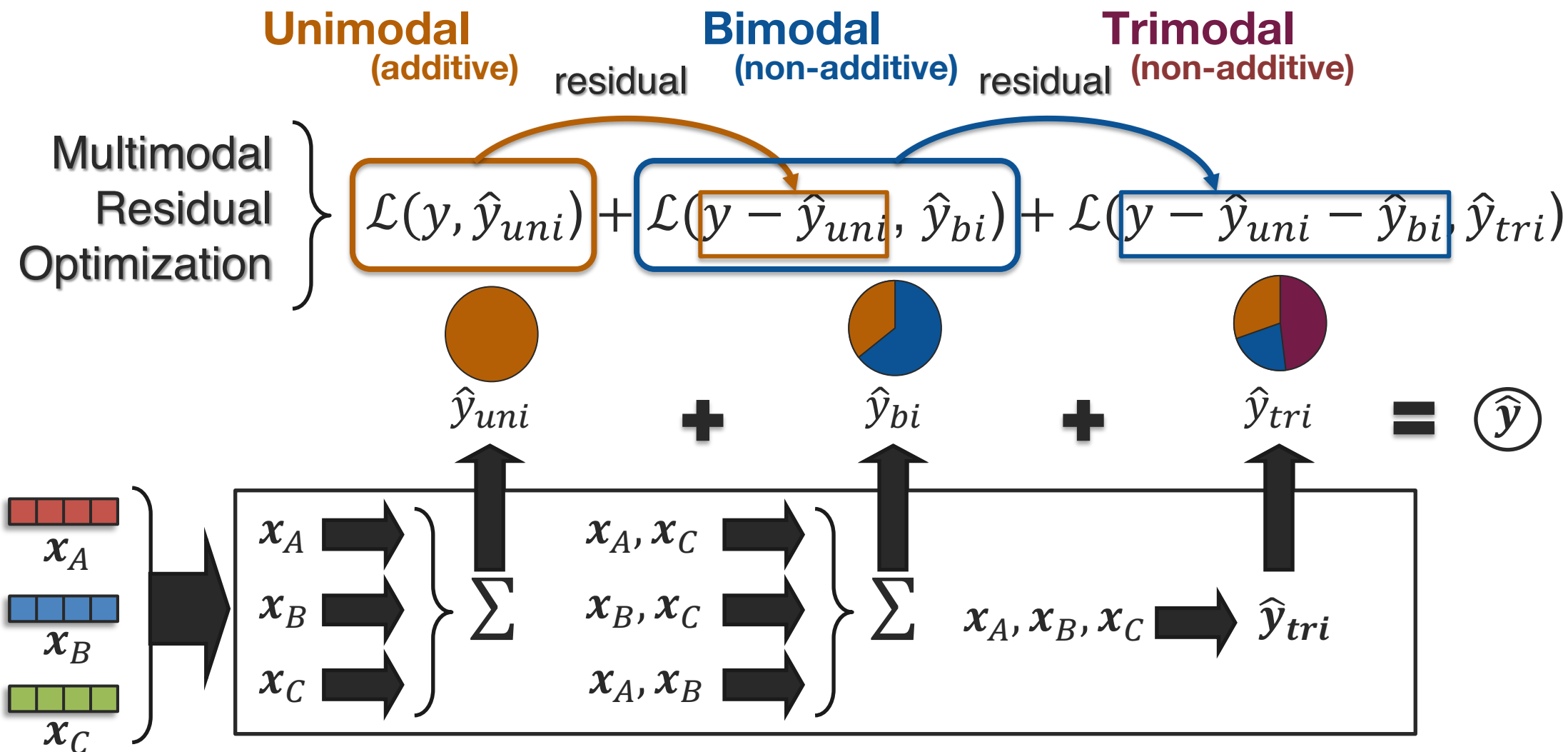
Additive fusion:

$$\hat{\mathbf{y}}' = \hat{f}_A(\mathbf{x}_A) + \hat{f}_B(\mathbf{x}_B) + \hat{\boldsymbol{\mu}}$$

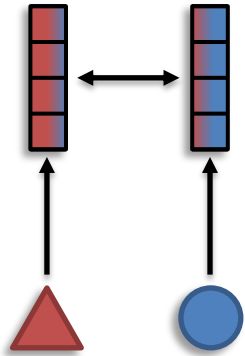
		I-INT	I-SEM	I-CTX	T-VIS	R-POP	T-ST1	T-ST2	
Nonlinear	Neural Network	90.4	69.2	78.5	51.1	63.5	71.1	79.9	
Polynomial	Polykernel SVM	91.3	74.4	81.5	50.8	–	72.1	80.9	
Nonlinear	FT LXMERT	83.0	68.5	76.3	53.0	63.0	66.4	78.6	
Nonlinear	\hookrightarrow + Linear Logits	89.9	73.0	80.7	53.4	64.1	75.5	80.3	
Additive	Linear Model	90.4	72.8	80.9	51.3	63.7	75.6	76.1	Always a good baseline!
	Best Model	91.3	74.4	81.5	53.4	64.2	75.5	80.9	
Additive	\hookrightarrow + EMAP	91.1	74.2	81.3	51.0	64.1	75.9	80.7	Differences are small!!!

Non-Additive Interactions

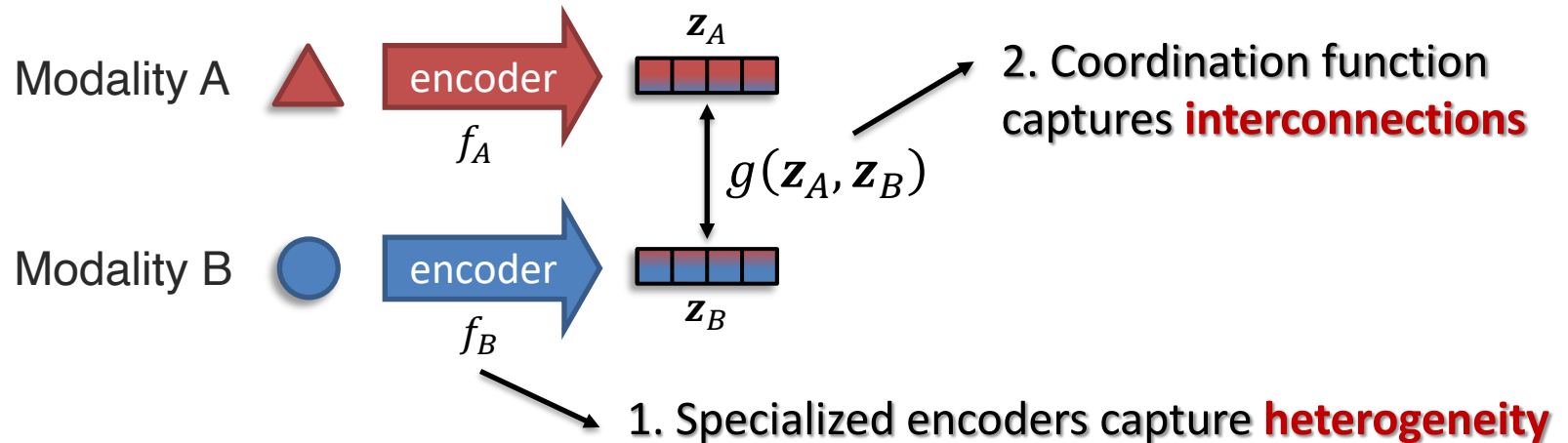
Idea: prioritize simpler interactions



Sub-Challenge 1b: Representation Coordination



Definition: Learn multimodal contextualized representations coordinated through their interconnections.

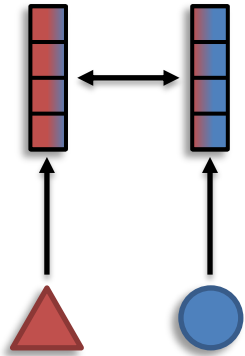


Learning with coordination function:

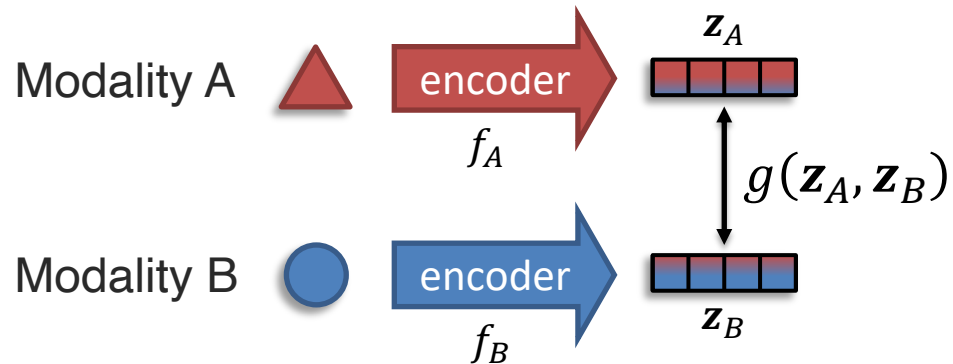
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

Coordinated Representations



Definition: Learn multimodal contextualized representations coordinated through their interconnections.



Learning with coordination function:

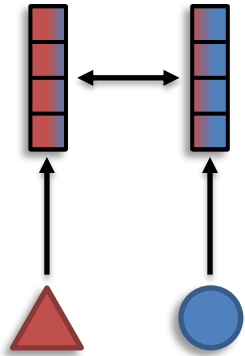
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

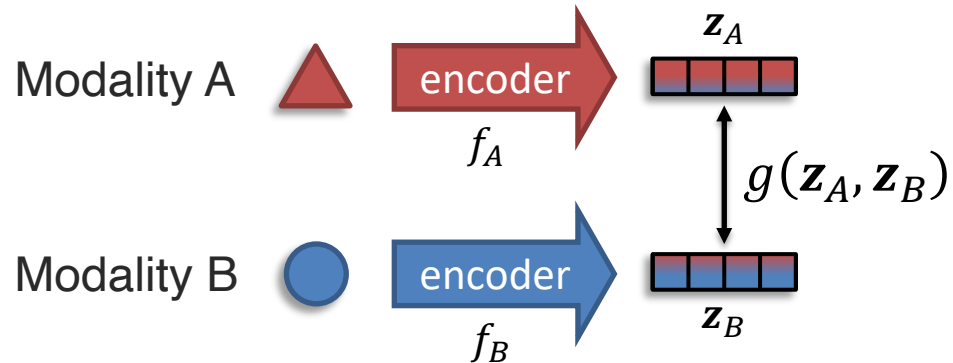
① Cosine similarity:

$$g(\mathbf{z}_A, \mathbf{z}_B) = \frac{\mathbf{z}_A \cdot \mathbf{z}_B}{\|\mathbf{z}_A\| \|\mathbf{z}_B\|}$$

Coordinated Representations



Definition: Learn multimodal contextualized representations coordinated through their interconnections.



Learning with coordination function:

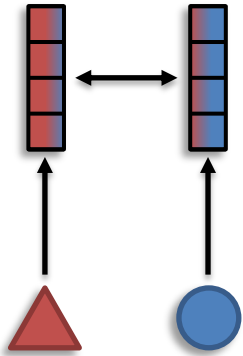
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

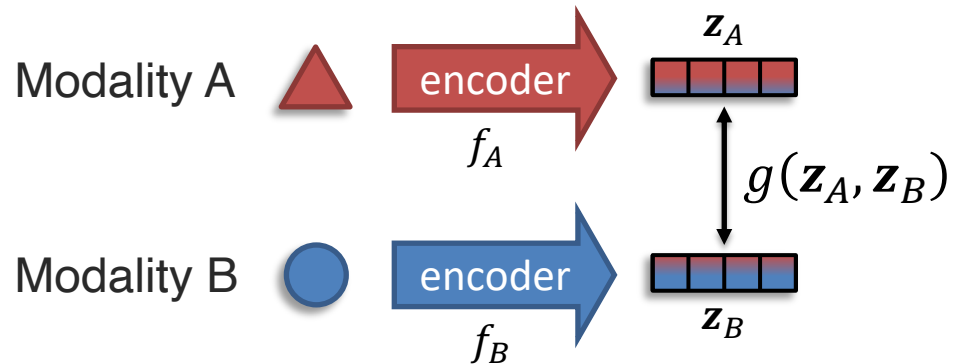
② Kernel similarity functions:

$$g(\mathbf{z}_A, \mathbf{z}_B) = k(\mathbf{z}_A, \mathbf{z}_B) \left\{ \begin{array}{l} \bullet \text{ Linear} \\ \bullet \text{ Polynomial} \\ \bullet \text{ Exponential} \\ \bullet \text{ RBF} \end{array} \right.$$

Coordinated Representations



Definition: Learn multimodally-contextualized representations coordinated through their cross-modal connections.



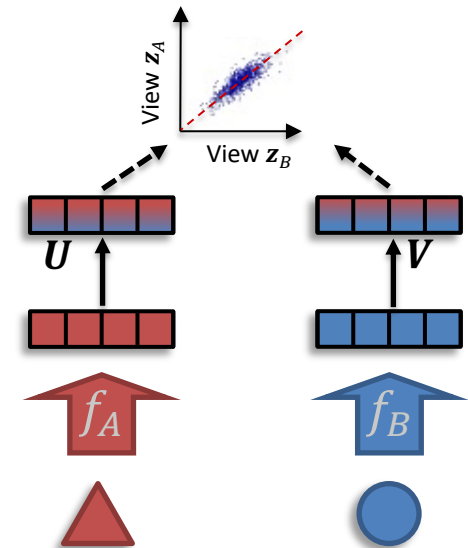
Learning with coordination function:

$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

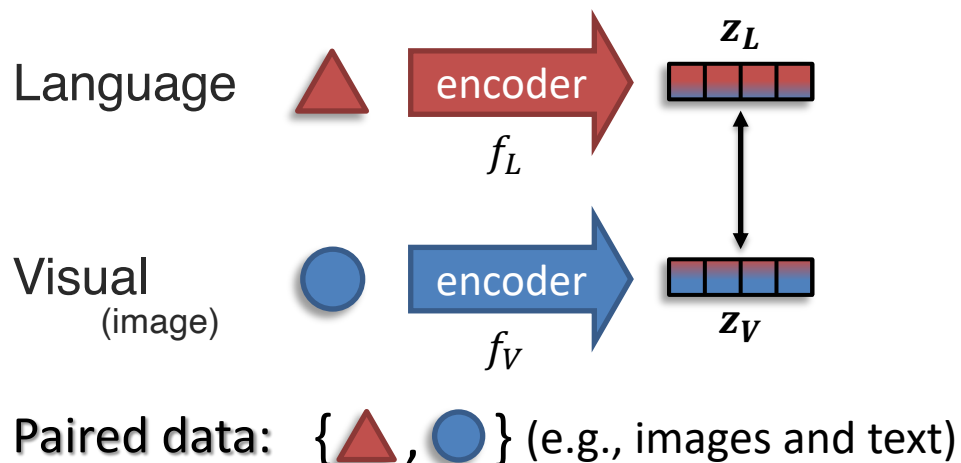
with model parameters θ_g , θ_{f_A} and θ_{f_B}

③ Canonical Correlation Analysis (CCA):

$$\operatorname{argmax}_{V, U, f_A, f_B} \operatorname{corr}(z_A, z_B)$$



Coordination with Contrastive Learning



Blue car



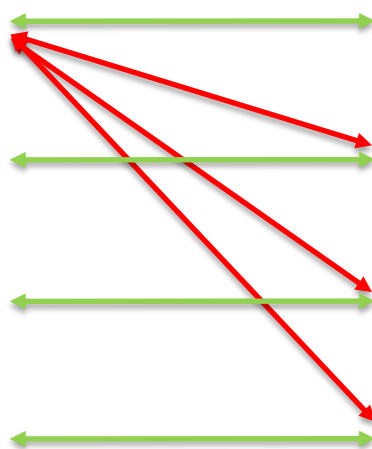
Yellow bus



Airplane



Bowl of cats



Contrastive loss:

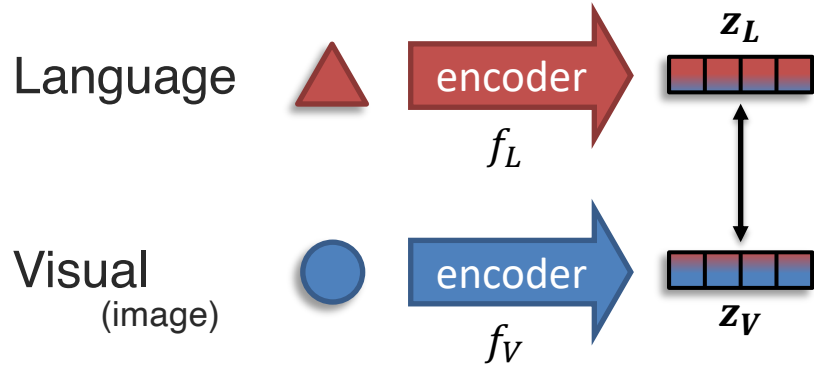
 brings **positive pairs** closer and pushes **negative pairs** apart

Simple contrastive loss:

$$\max\{0, \alpha + \underbrace{g(z_A, z_B^+)}_{\text{positive pairs}} - \underbrace{g(z_A, z_B^-)}_{\text{negative pair}}\}$$

Coordination function (e.g., cosine similarity)

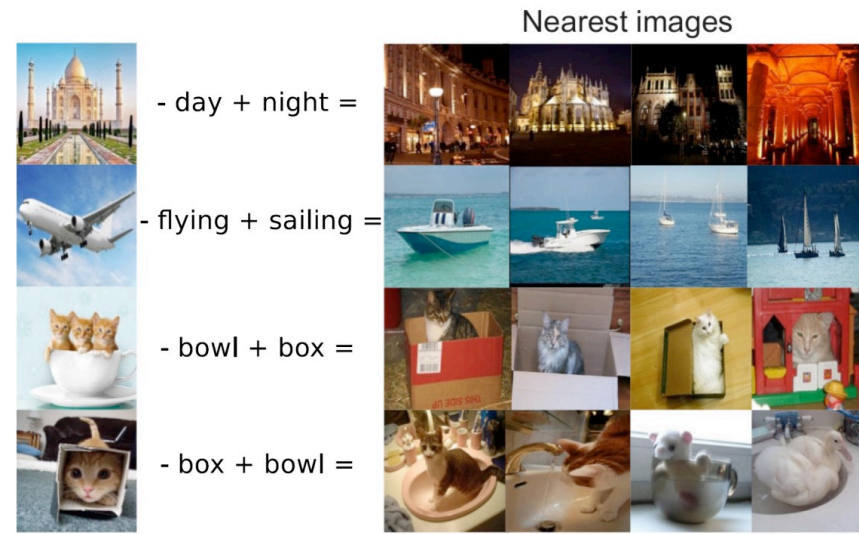
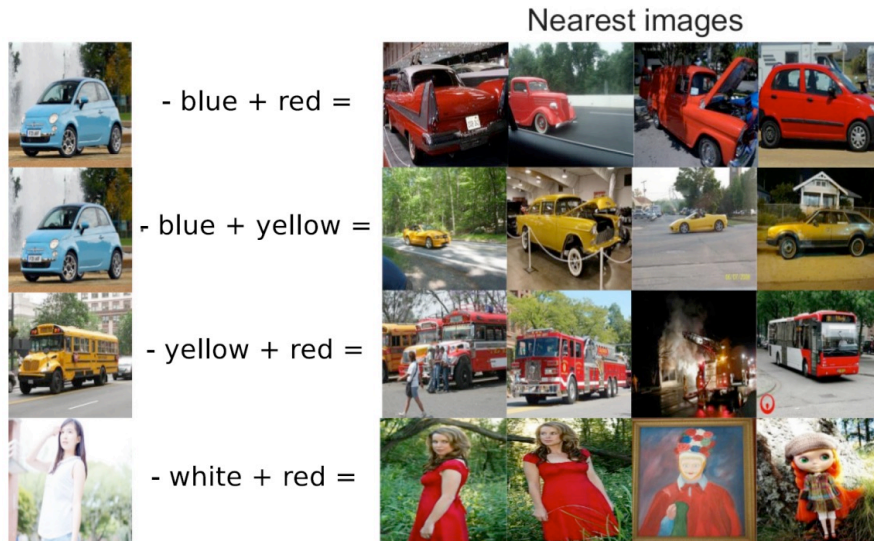
Visual-Semantic Representations



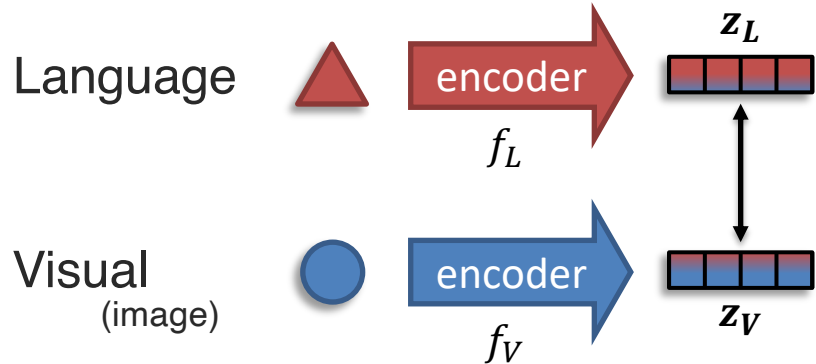
Two contrastive loss terms:

$$\max\{0, \alpha + \text{sim}(z_L, z_V^+) - \text{sim}(z_L, z_V^-)\}$$

$$+ \max\{0, \alpha + \text{sim}(z_V, z_L^+) - \text{sim}(z_V, z_L^-)\}$$



Contrastive Language Image Pretraining



Popular contrastive loss: InfoNCE

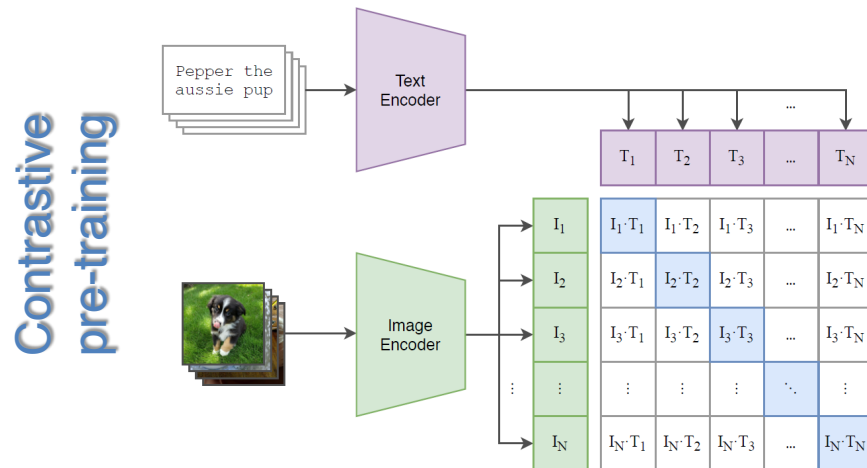
$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(z_A^i, z_B^i)}{\sum_{j=1}^N \text{sim}(z_A^i, z_B^j)}$$

Similarity function can be cosine similarity

positive pairs

negative pairs and positive pairs

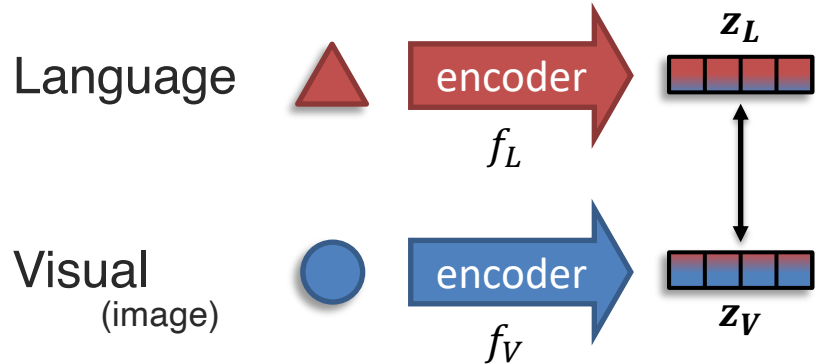
Positive and negative pairs:



\Rightarrow CLIP encoders (f_L and f_V) are great for language-vision tasks

\Rightarrow z_L and z_V are coordinated but not identical representation spaces

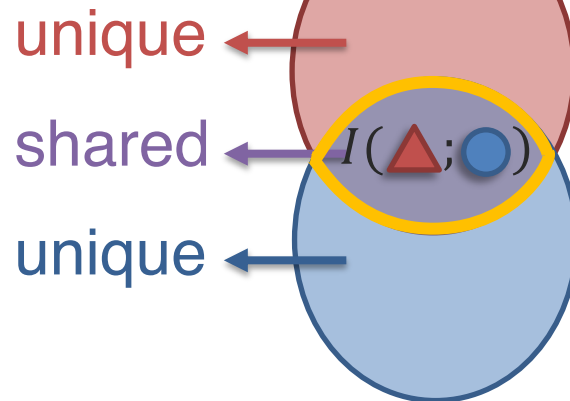
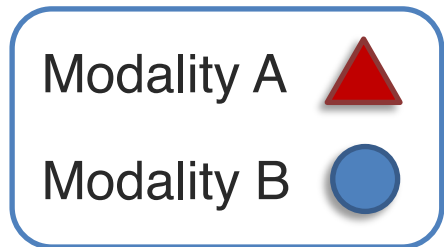
Contrastive Learning and Connected Modalities



Popular contrastive loss: InfoNCE

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(\mathbf{z}_A^i, \mathbf{z}_B^i)}{\sum_{j=1}^N \text{sim}(\mathbf{z}_A^i, \mathbf{z}_B^j)}$$

Connected modalities:



Mutual information $I(X; Y)$

$$\mathbb{E}_{X,Y} \left[\log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} \right]$$

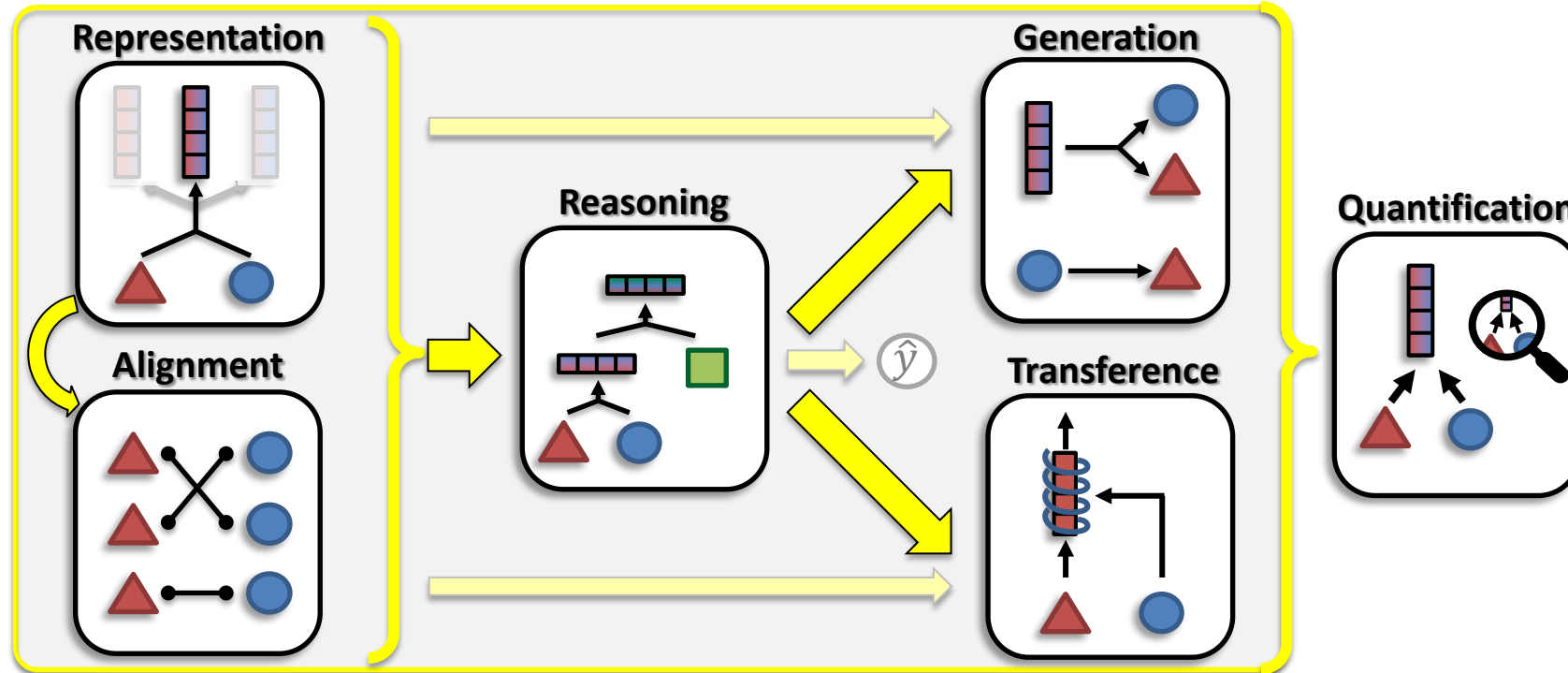


CLIP focuses on shared connections

Summary

Resources: <https://cs.cmu.edu/~pliang/>

Multimodal is the science of **heterogeneous** and **interconnected** data.



What is Multimodal?



Why is it hard?



What is next?

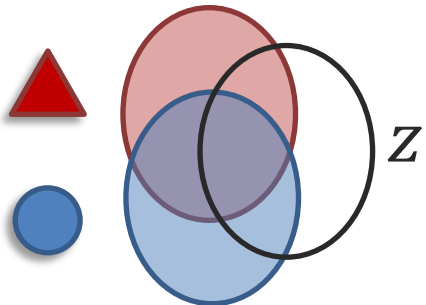
Heterogeneous



Connected



Interacting



Representation

Alignment

Reasoning

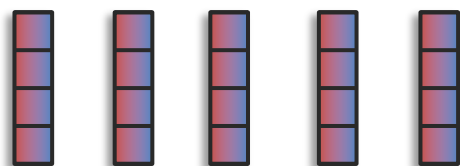
Generation

Transference

Quantification

Future Direction: Heterogeneity & Interactions

Homogeneity



vs

Heterogeneity



Challenges:

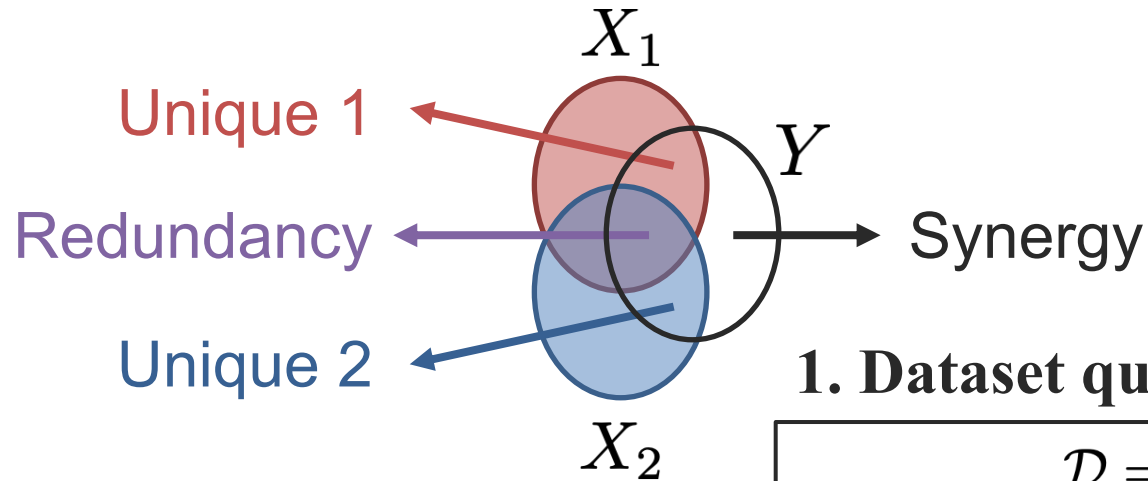
Arbitrary tokenization



Beyond differentiable interactions

Causal, logical, brain-inspired
Theoretical foundations of interactions

Quantifying Interactions



1. Dataset quantification:

$$\mathcal{D} = \{(x_1, x_2, y)\} \longrightarrow \{R, U_1, U_2, S\}_{\mathcal{D}}$$

2. Model quantification:

$$f(\mathcal{D}) = \{(x_1, x_2, \hat{y} = f(x_1, x_2))\} \longrightarrow \{R, U_1, U_2, S\}_{f(\mathcal{D})}$$

$$\{R, U_1, U_2, S\}_{f(\mathcal{D}_1)}, \dots, \{R, U_1, U_2, S\}_{f(\mathcal{D}_k)} \longrightarrow \{R, U_1, U_2, S\}_f$$

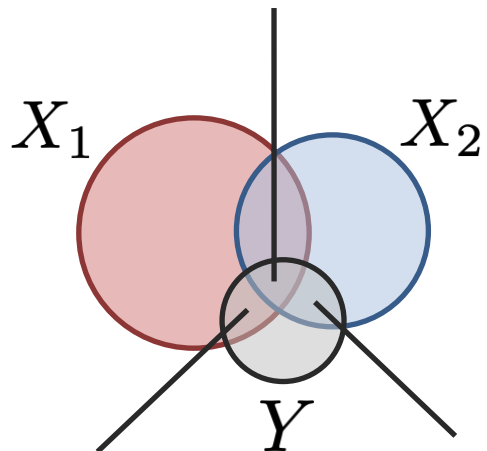
3. Model selection:

$$\{R, U_1, U_2, S\}_{\mathcal{D}} \longleftrightarrow \{R, U_1, U_2, S\}_f$$

Quantifying Interactions

Classical Information Theory

$$R = I(X_1; X_2; Y)$$

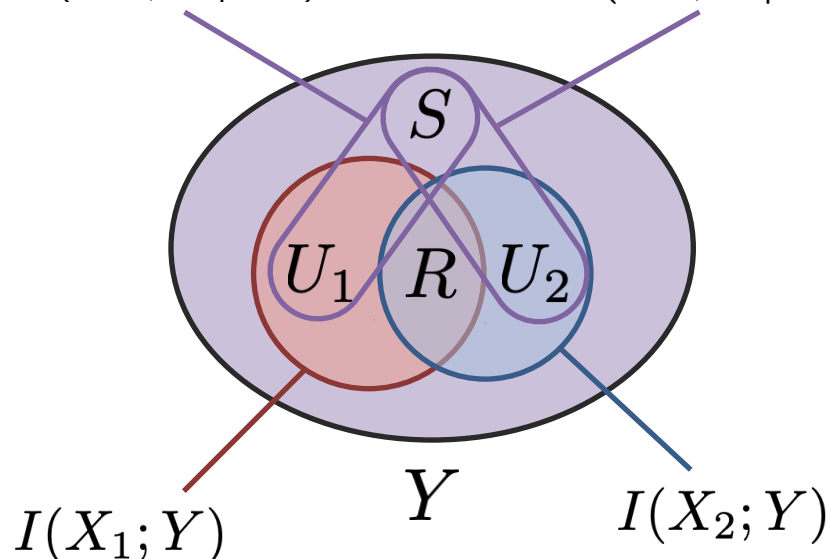


$$U_1 = I(X_1; Y|X_2) \quad U_2 = I(X_2; Y|X_1)$$

Partial Information Decomposition

$$I(X_1; Y|X_2)$$

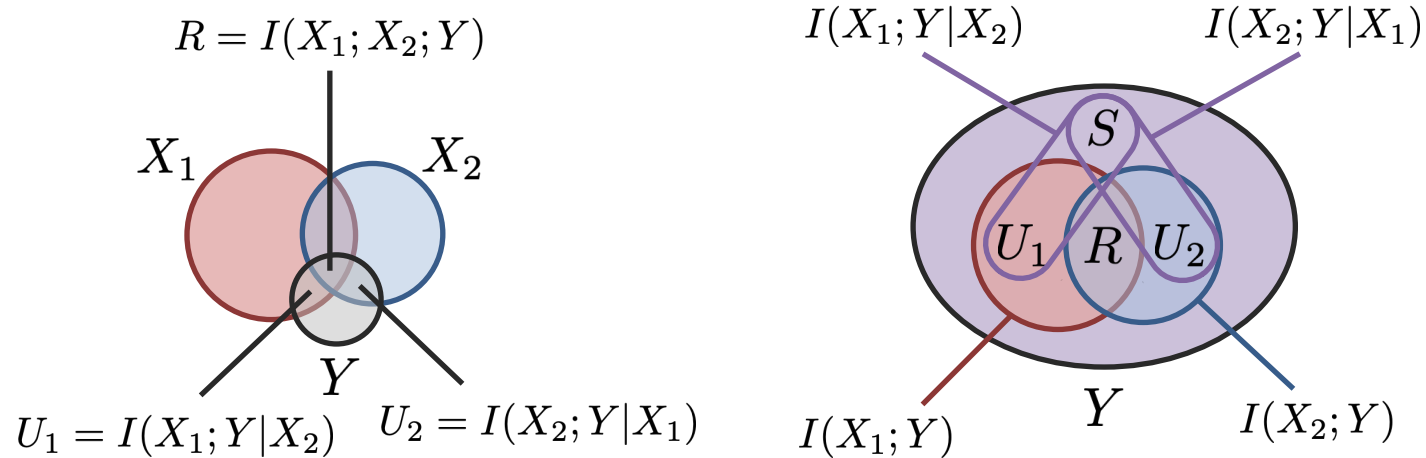
$$I(X_2; Y|X_1)$$



$$R - S = I(X_1; X_2; Y)$$

$$R + U_1 + U_2 + S = I(X_1, X_2; Y)$$

Quantifying Interactions



$$\Delta_p = \{q \in \Delta : q(x_i, y) = p(x_i, y) \forall y, x_i, i \in \{1, 2\}\}$$

Marginal-matching
distributions

$$R = \max_{q \in \Delta_p} I_q(X_1; X_2; Y)$$

$$U_2 = \max_{q \in \Delta_p} I_q(X_2; Y | X_1)$$

$$U_1 = \max_{q \in \Delta_p} I_q(X_1; Y | X_2)$$

$$S = I_p(X_1, X_2; Y) - \min_{q \in \Delta_p} I_q(X_1, X_2; Y)$$

Quantifying Interactions

$$q^* = \arg \max_{q \in \Delta_p} H_q(Y | X_1, X_2)$$

If X_1, X_2, Y have small and discrete support,
exact solution via convex programming with linear constraints.

Else, neural network estimator.

$$\Delta_p = \{q \in \Delta : q(x_i, y) = p(x_i, y) \forall y, x_i, i \in \{1, 2\}\}$$

Marginal-matching
distributions

$$R = \max_{q \in \Delta_p} I_q(X_1; X_2; Y)$$

$$U_2 = \max_{q \in \Delta_p} I_q(X_2; Y | X_1)$$

$$U_1 = \max_{q \in \Delta_p} I_q(X_1; Y | X_2)$$

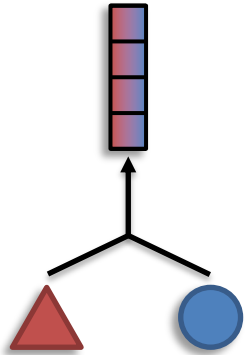
$$S = I_p(X_1, X_2; Y) - \min_{q \in \Delta_p} I_q(X_1, X_2; Y)$$

Representation Models

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

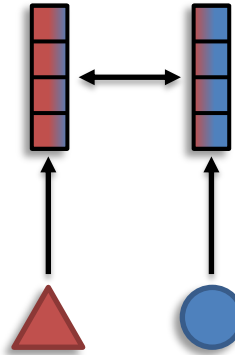
Sub-challenges:

Fusion



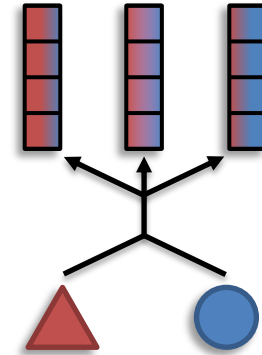
modalities $>$ # representations

Coordination



modalities = # representations

Fission

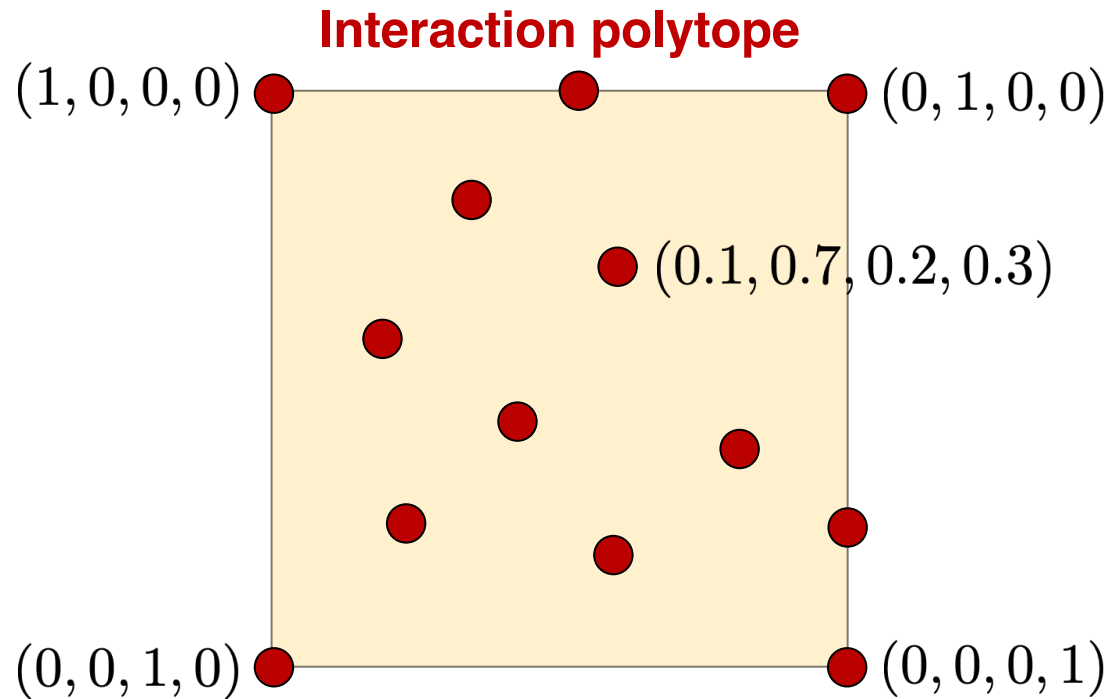


modalities $<$ # representations

Model Selection

1. Dataset quantification:

$$\mathcal{D} = \{(x_1, x_2, y)\} \longrightarrow \{R, U_1, U_2, S\}_{\mathcal{D}} \bullet$$

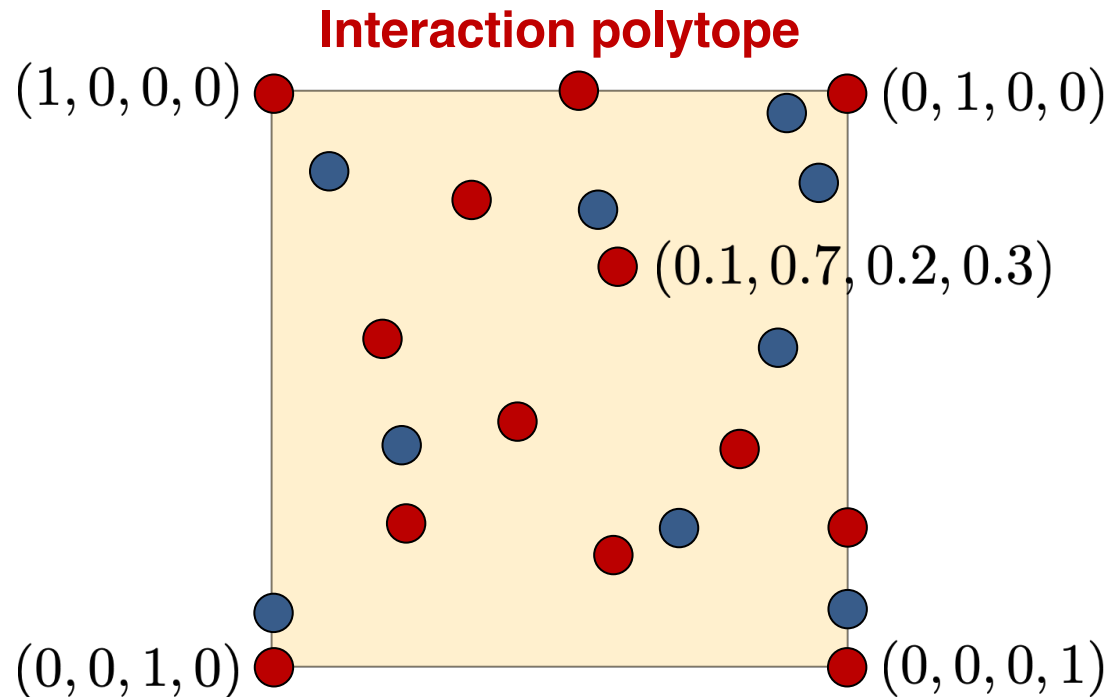


Model Selection

2. Model quantification:

$$f(\mathcal{D}) = \{(x_1, x_2, \hat{y} = f(x_1, x_2))\} \xrightarrow{\text{red arrow}} \{R, U_1, U_2, S\}_{f(\mathcal{D})}$$

$$\{R, U_1, U_2, S\}_{f(\mathcal{D}_1)}, \dots, \{R, U_1, U_2, S\}_{f(\mathcal{D}_k)} \xrightarrow{\text{black arrow}} \{R, U_1, U_2, S\}_f \quad \bullet$$

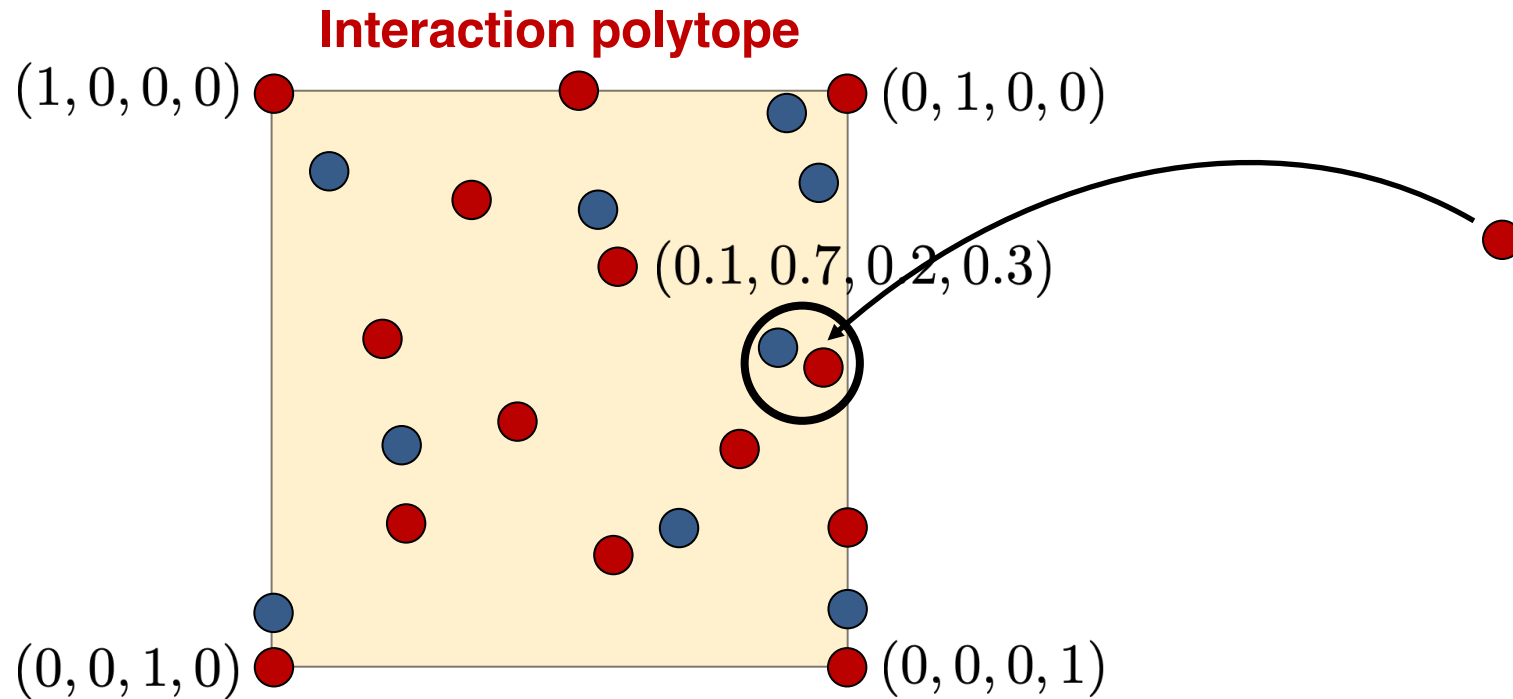


Model Selection

3. Model selection:

$$\{R, U_1, U_2, S\}_{\mathcal{D}} \longleftrightarrow \{R, U_1, U_2, S\}_f$$

**Selects models with
>96% performance**



Future Direction: High-modality

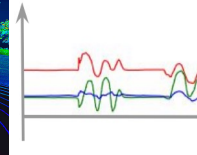
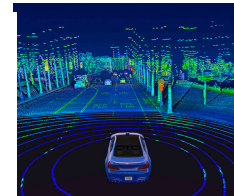
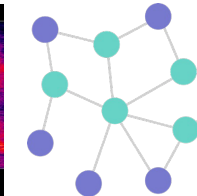
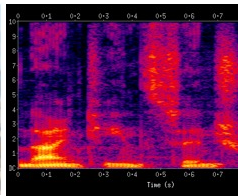
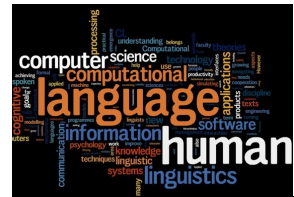
MultiBench

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

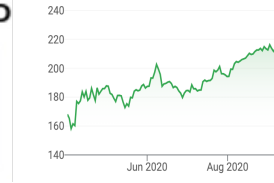
Few modalities



High-modality



SUBJECT_ID
Age
Sex
Ethnicity
...



Language

Vision

Audio

Graphs

Control

LIDAR

Sensors

Set

Table

Financial

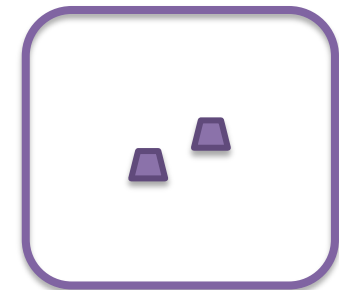
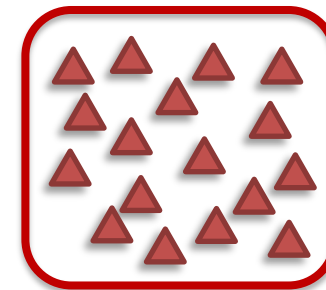
Medical

Challenges:

Non-parallel learning



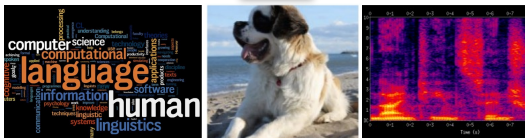
Limited resources



High-Modality Learning

How can we transfer knowledge across multiple tasks, each over a different subset of modalities?

Video
classification



Language

Video

Audio

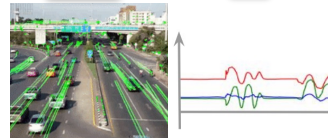
Sentiment,
emotions



Audio

Video

Robot
dynamics



Video

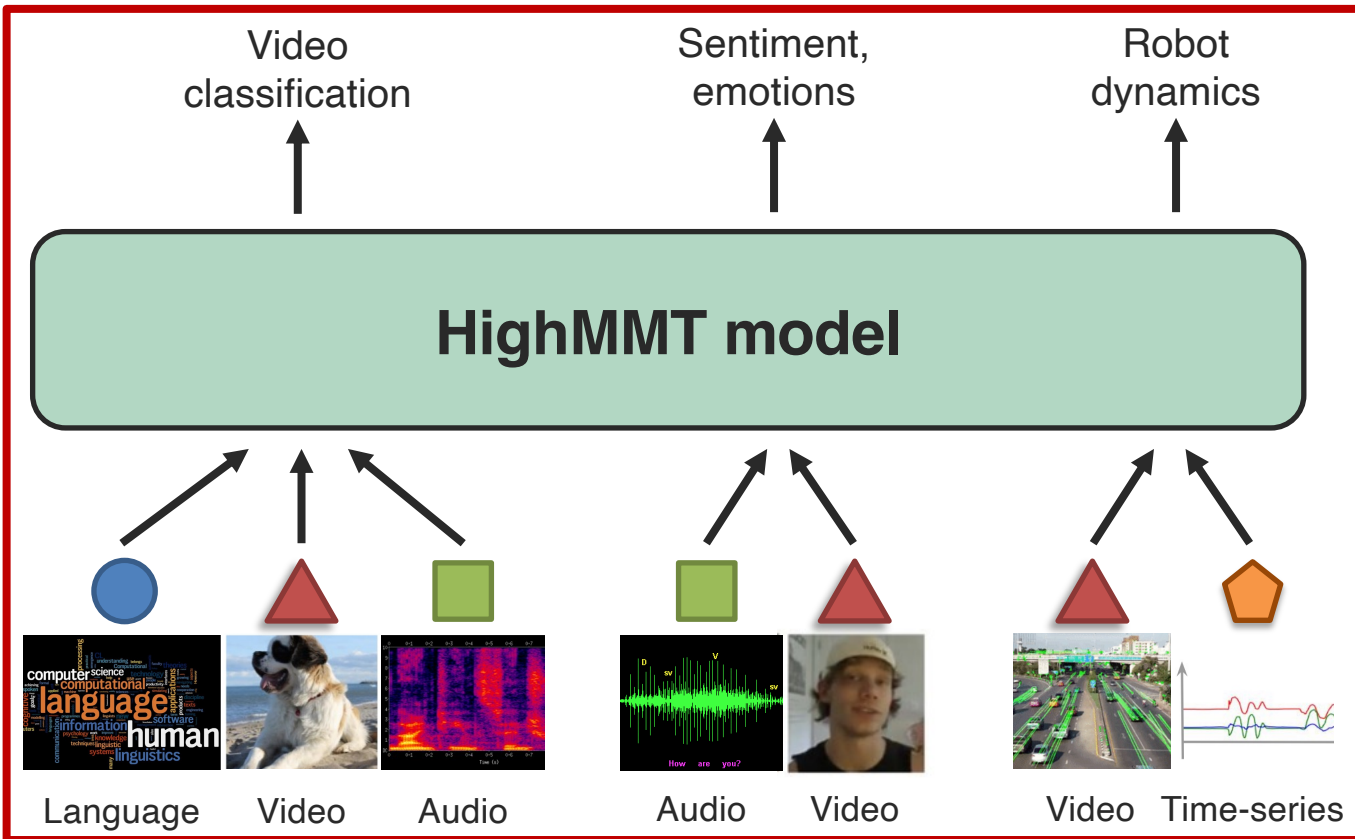
Time-series

Generalization across modalities and tasks
Important if some tasks are low-resource

HighMMT

Transfer across partially observable modalities

HighMMT: unified model + parameter sharing + multitask and transfer learning



Non-parallel multitask learning

Task-specific classifiers

Same model architecture!

Shared multimodal model

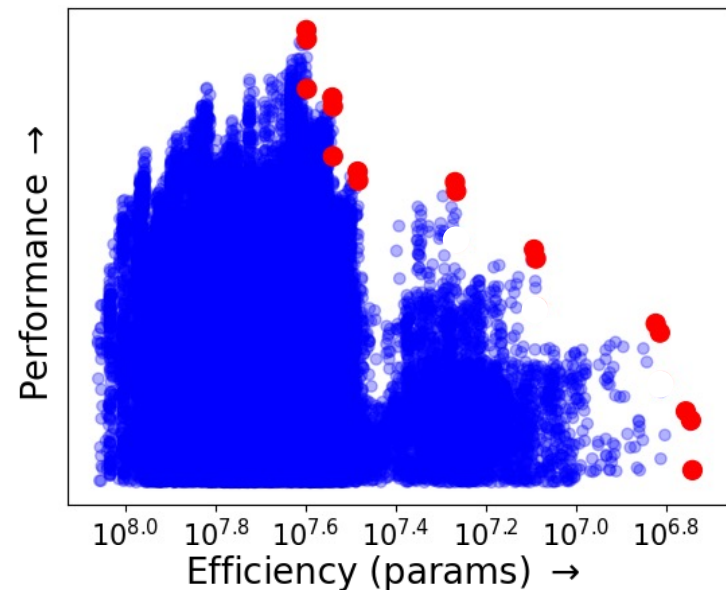
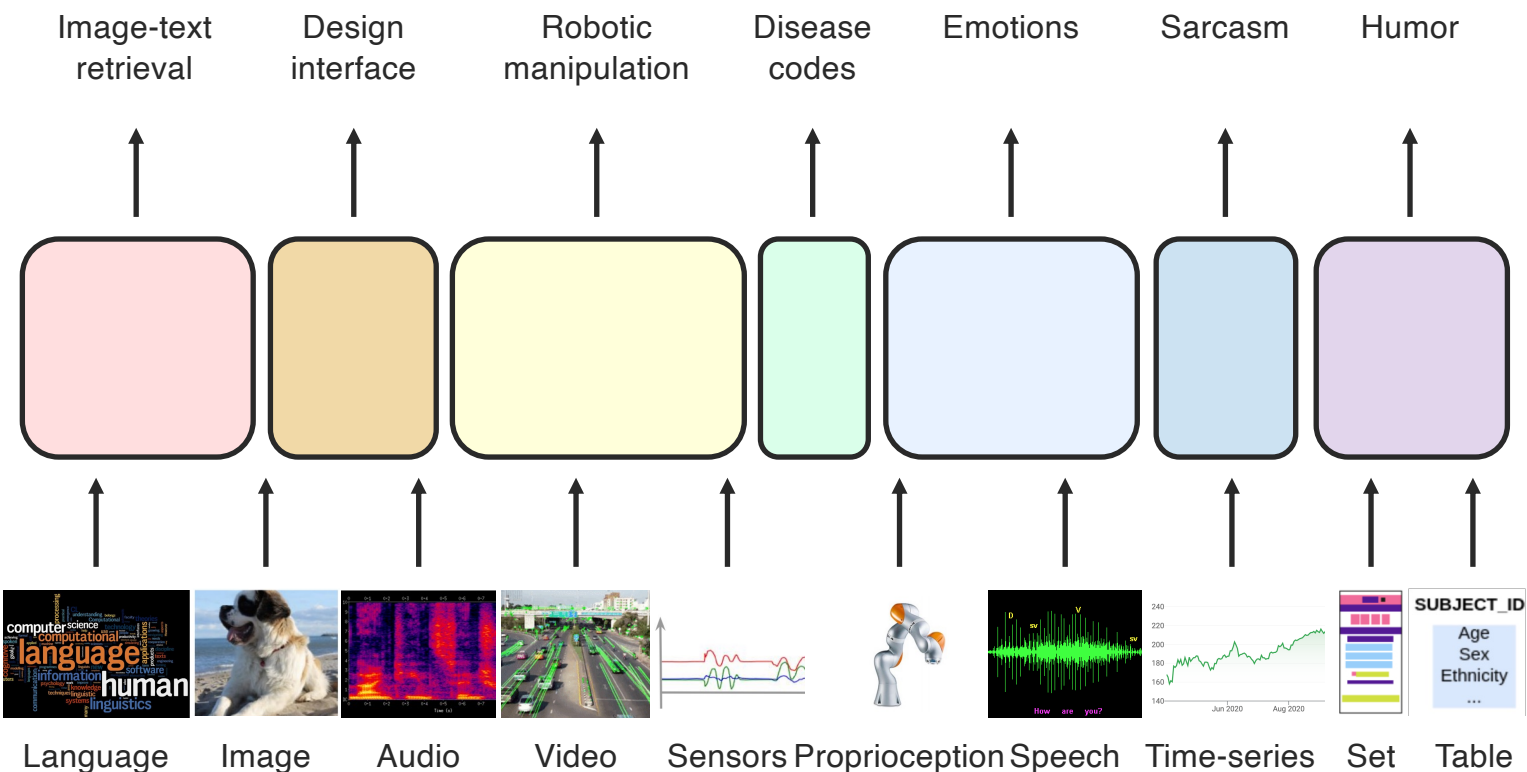
Same parameters!

Modality-specific embeddings

Standardized input sequence

HighMMT

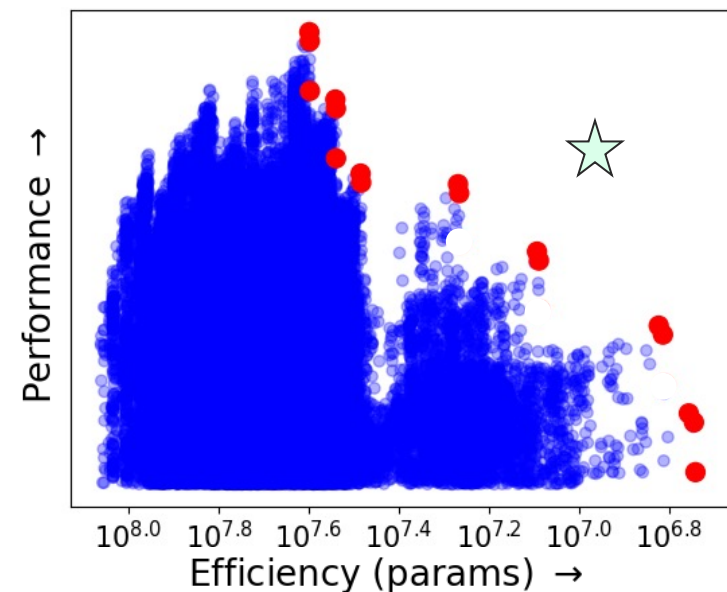
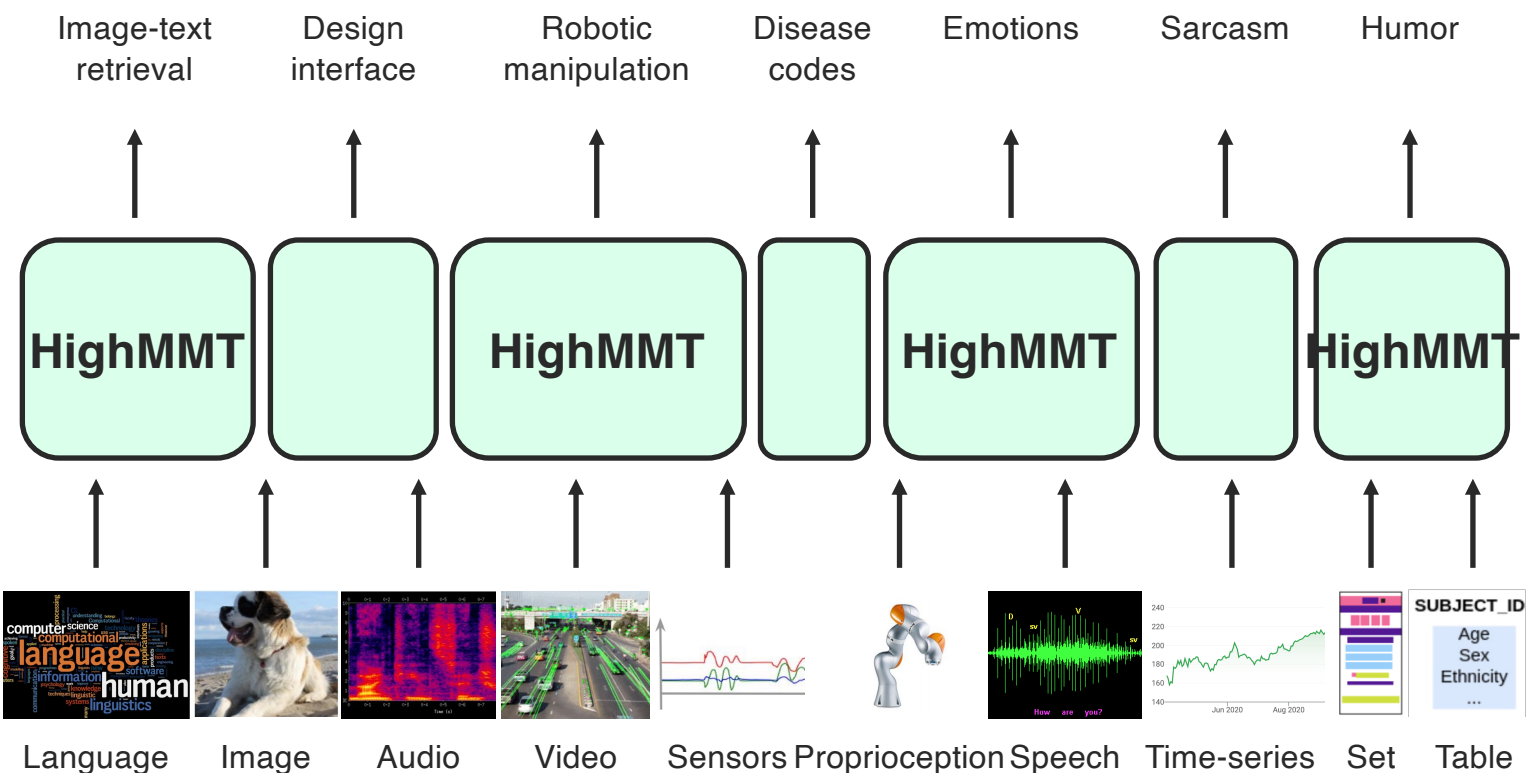
Traditional approaches: different model + different parameters



- All model combinations (>10,000)
- Pareto front

HighMMT

Traditional approaches: different model + different parameters



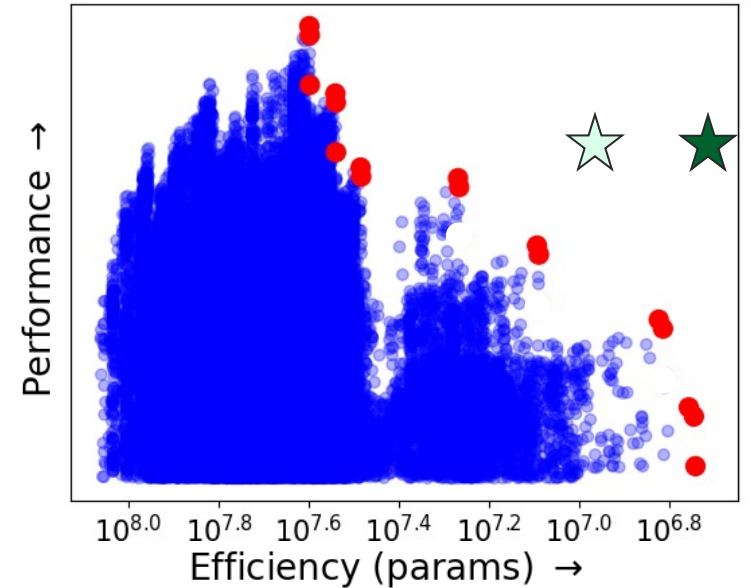
- All model combinations (>10,000)
- Pareto front
- HighMMT single-task

HighMMT

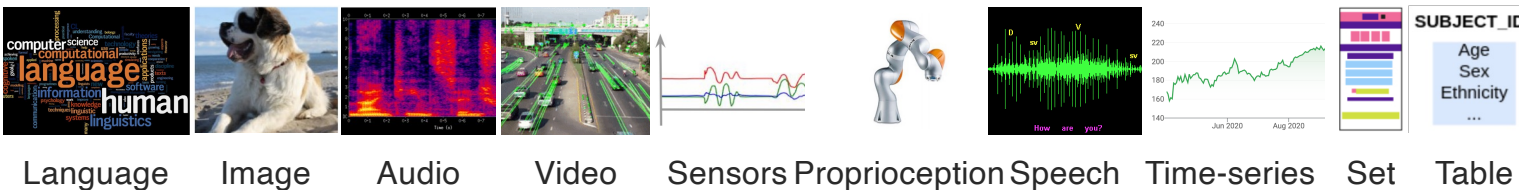
Traditional approaches: different model + different parameters

Image-text retrieval Design interface Robotic manipulation Disease codes Emotions Sarcasm Humor

HighMMT multitask model



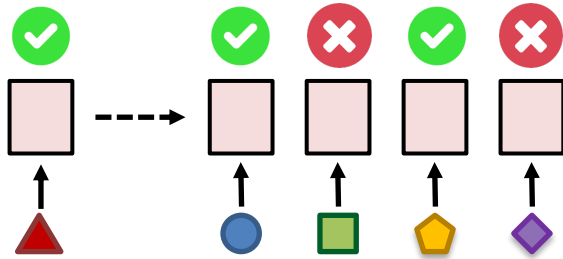
- All model combinations (>10,000)
- Pareto front
- ☆ HighMMT single-task
- ☆ HighMMT multitask



Quantifying Modality Heterogeneity

Information transfer, transfer learning perspective

1a. Estimate modality heterogeneity via transfer



Implicitly captures these:

① Element representation

② Element distribution

③ Structure

④ Information

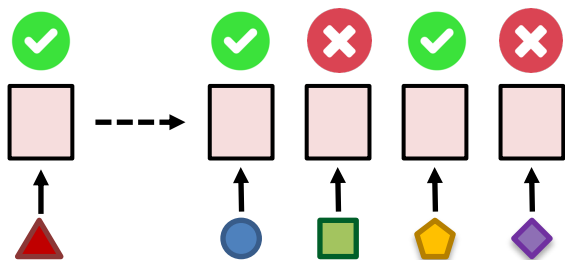
⑤ Noise

⑥ Relevance

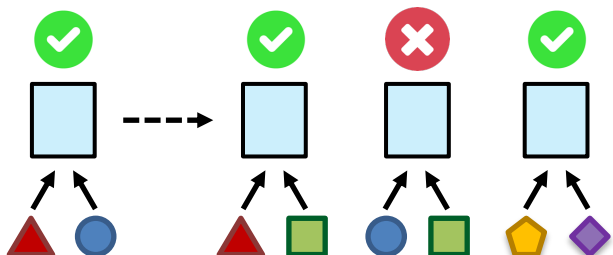
Quantifying Modality Heterogeneity

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

$$U_1 = \{U_1, U_2, U_4\} \quad C_1 = \{C_{12}, C_{13}, C_{45}\}$$

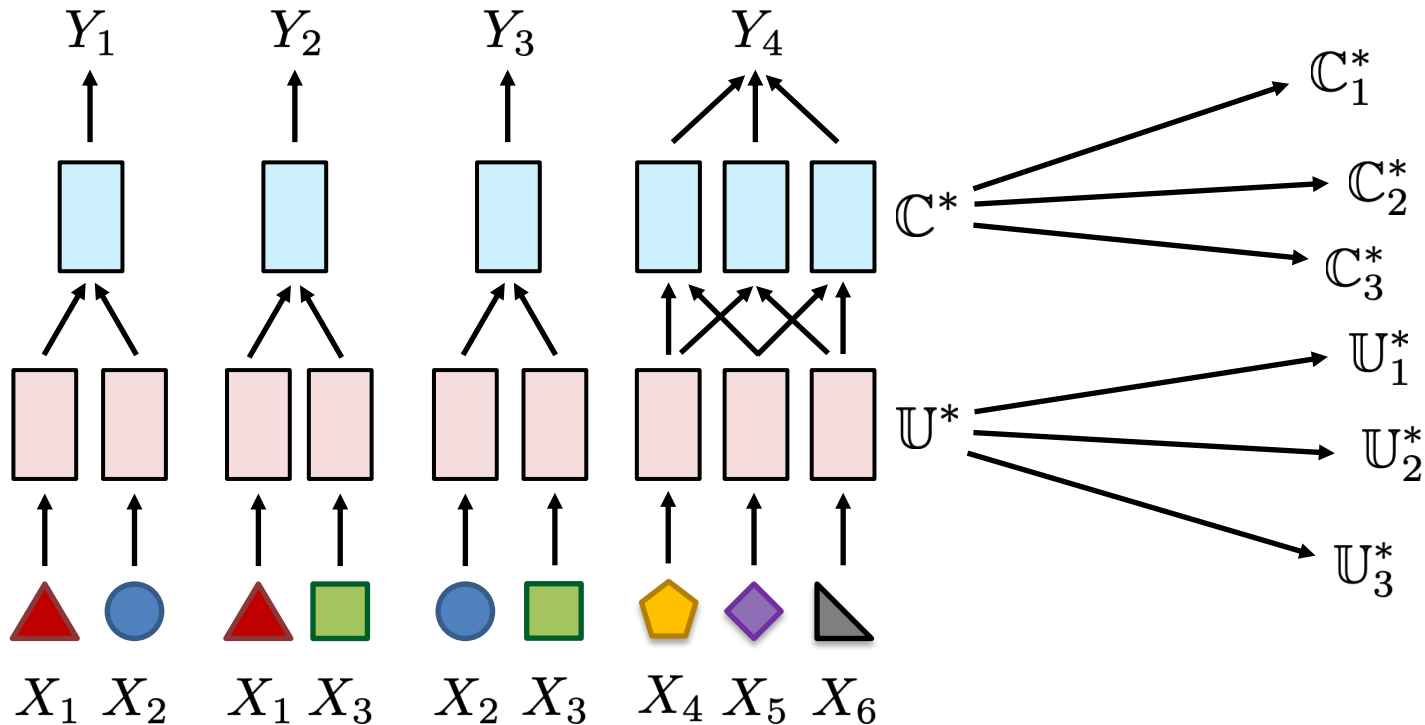
$$U_2 = \{U_3\} \quad C_2 = \{C_{23}\}$$

$$U_3 = \{U_5\}$$

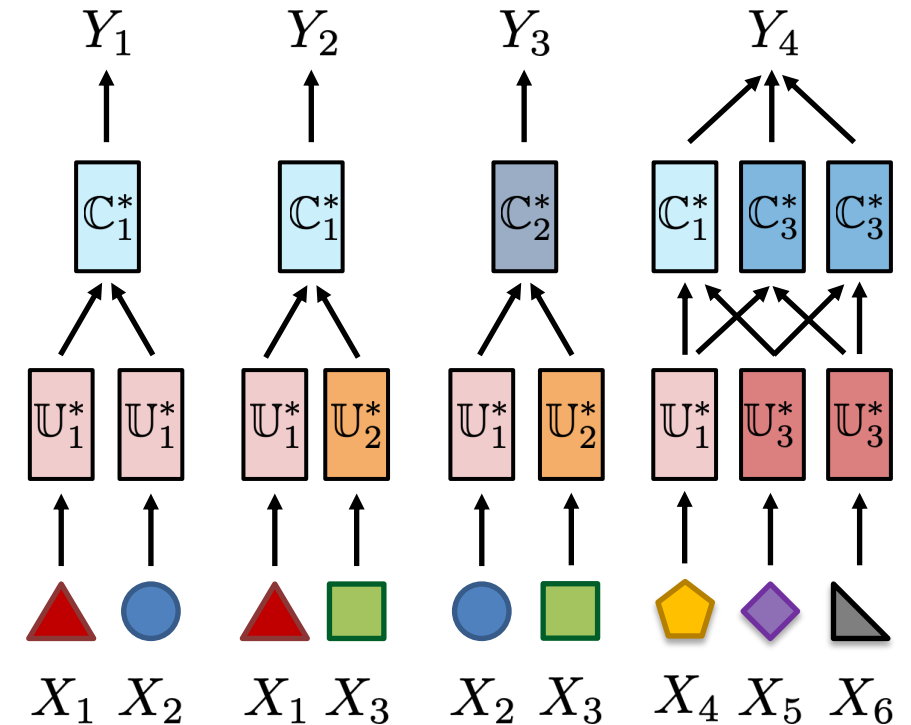
Quantifying Modality Heterogeneity

Information transfer, transfer learning perspective

1. Homogeneous Pre-training

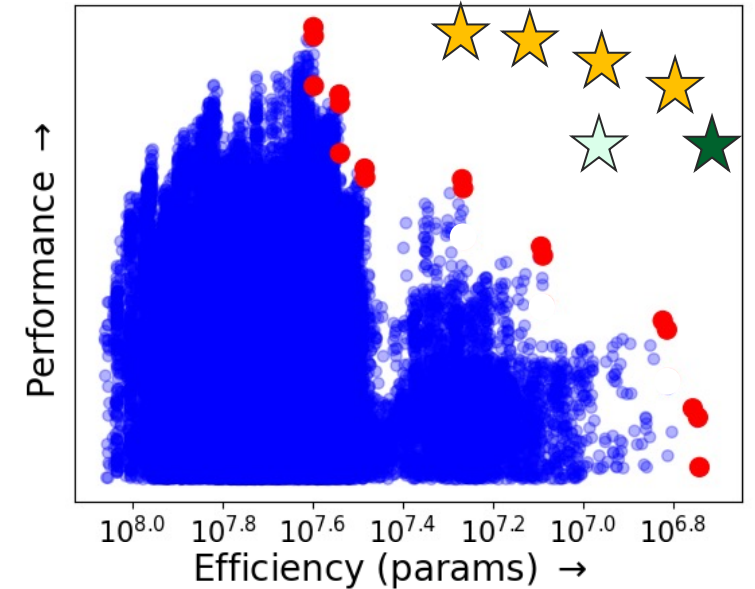
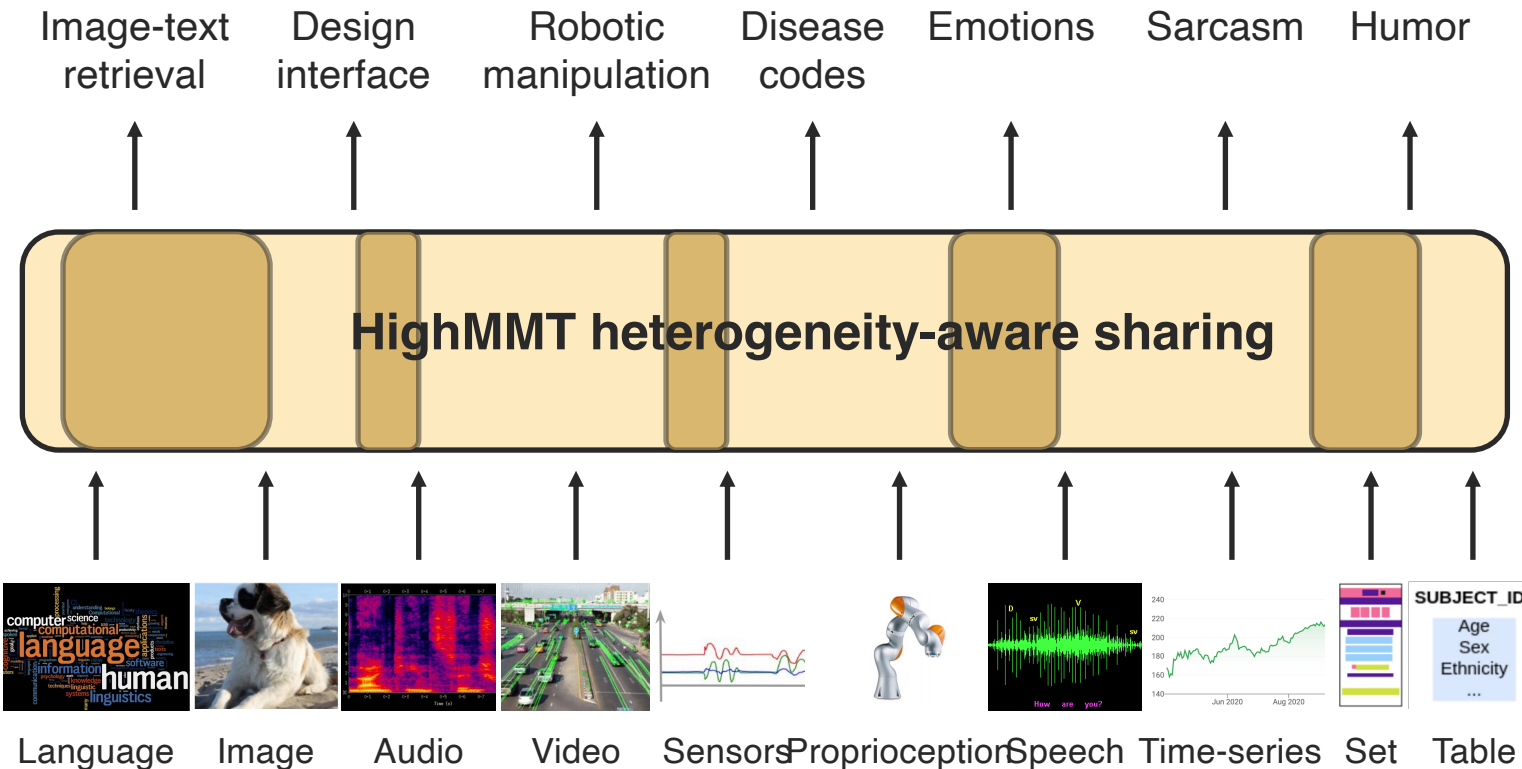


2. Heterogeneity-aware Fine-tuning



HighMMT + Quantifying Modality Heterogeneity

HighMMT heterogeneity-aware sharing: estimate heterogeneity to determine parameter sharing

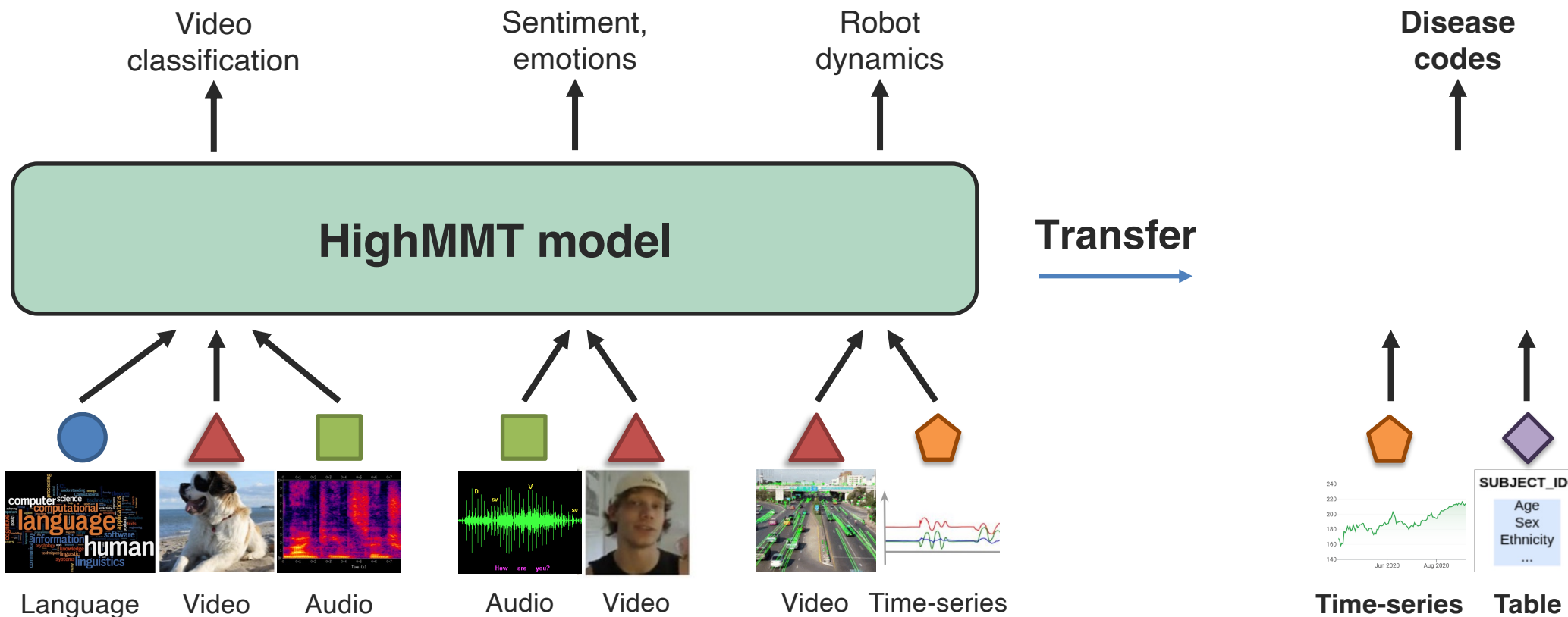


- All model combinations (>10,000)
- Pareto front
- HighMMT single-task
- HighMMT multitask
- **HighMMT heterogeneity-aware**

HighMMT

Transfer across partially observable modalities

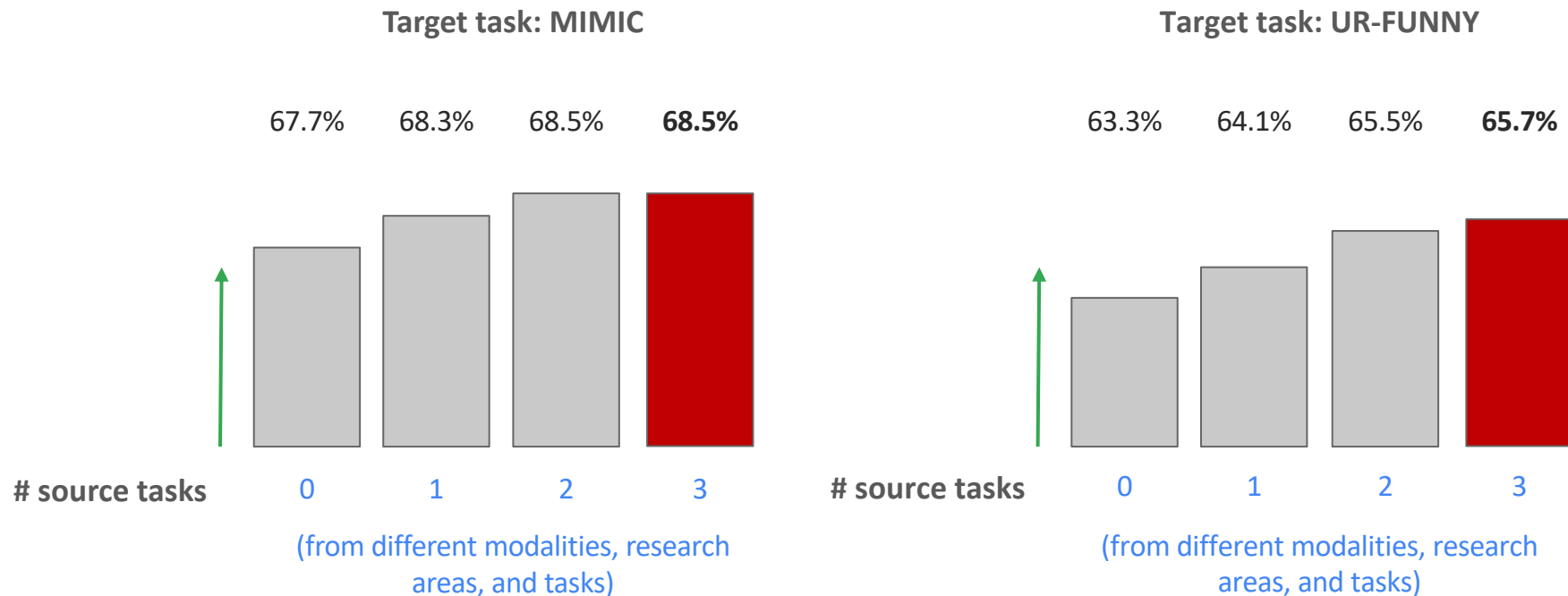
HighMMT: unified model + parameter sharing + multitask and transfer learning



HighMMT

Transfer across partially observable modalities

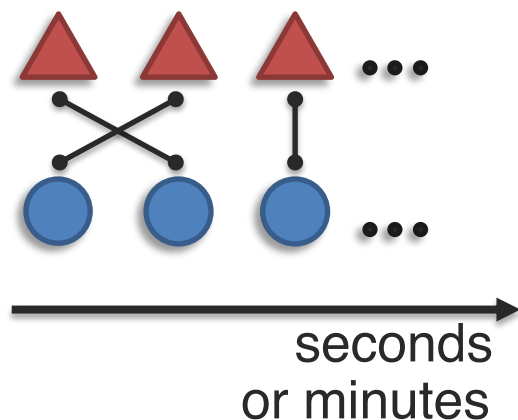
HighMMT: unified model + parameter sharing + multitask and transfer learning



Achieves both multitask and transfer capabilities across modalities and tasks

Future Direction: Long-term

Short-term



Long-term



Challenges:

Compositionality

Memory

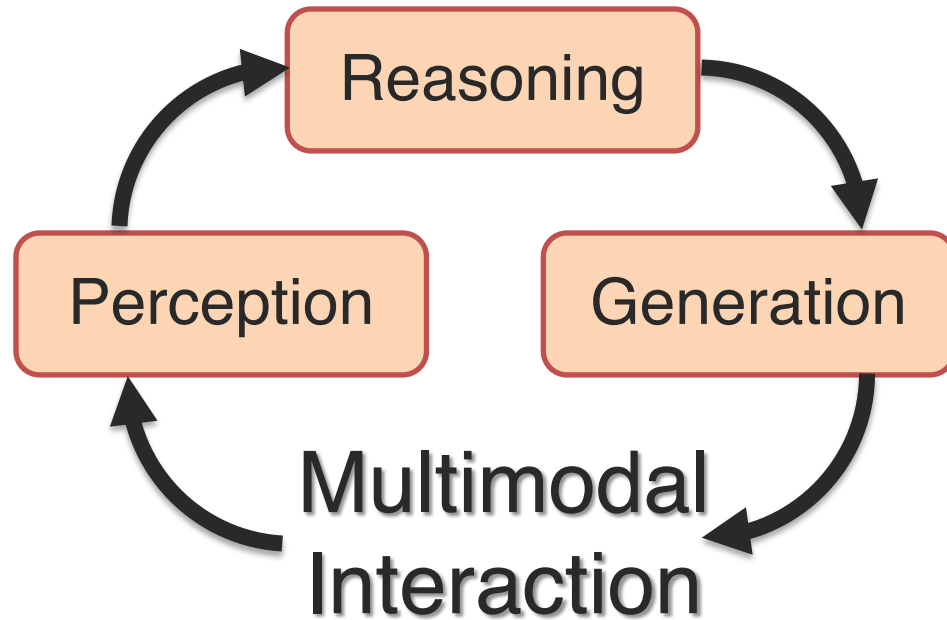
Personalization

Future Direction: Interaction

Social-IQ

<https://www.thesocialiq.com/>

Social Intelligence



Challenges:

Multi-Party

Generation

Ethics

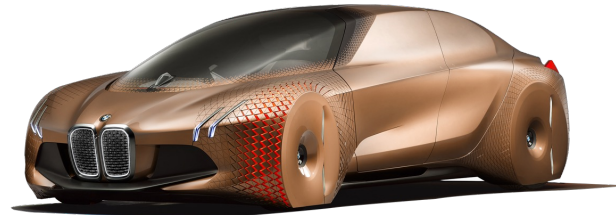
Future Direction: Real-world

MultiViz

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



Healthcare
Decision Support



Intelligent Interfaces and
Vehicles



Online Learning
and Education

Challenges:

Robustness

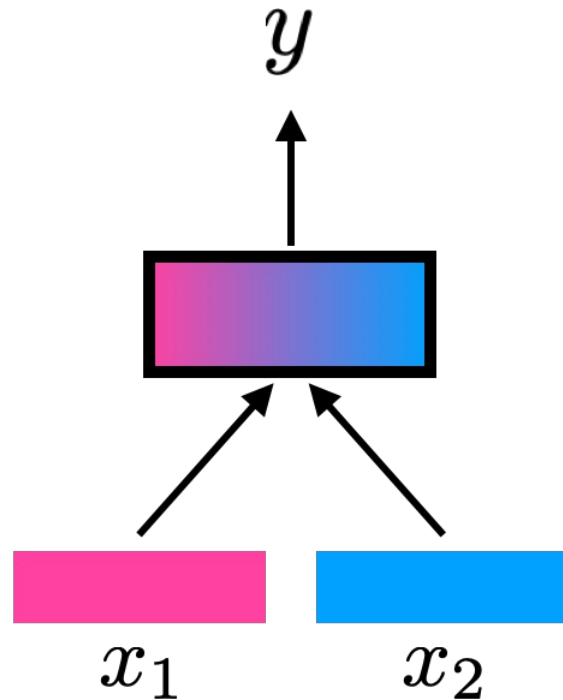
Fairness

Generalization

Interpretation

Real-World Quantification

How can we understand the modeling of **heterogeneity** and **interconnections** and gain insights for safe real-world deployment?

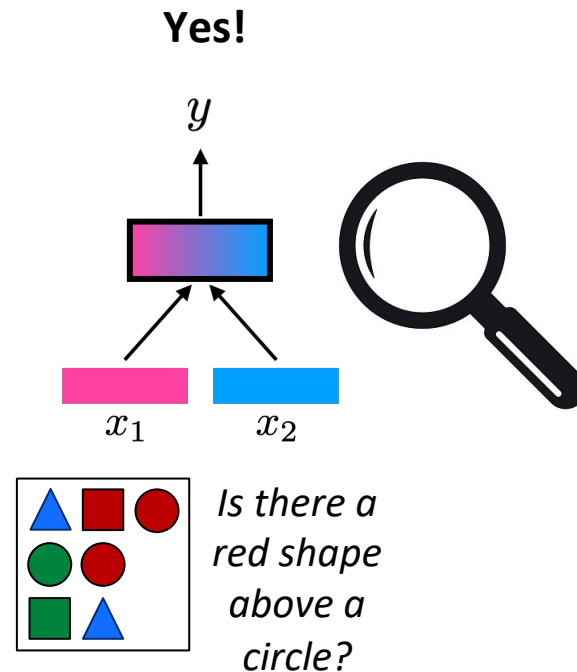


Internal mechanics



MultiViz: Interpreting Internal Mechanics

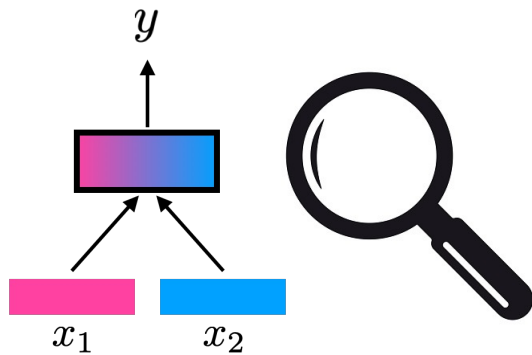
How can we understand the modeling of **heterogeneity** and **interconnections** and gain insights for safe real-world deployment?



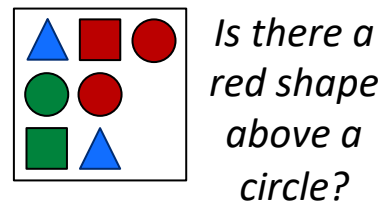
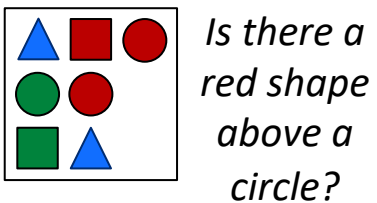
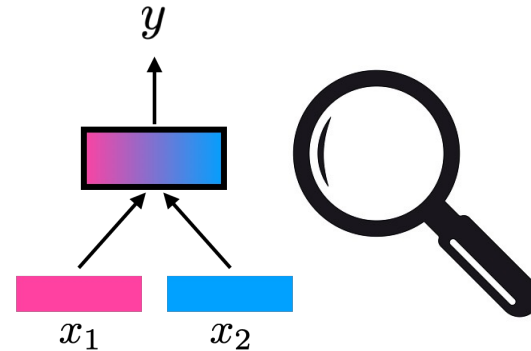
MultiViz: Interpreting Internal Mechanics

How can we understand the modeling of **heterogeneity** and **interconnections** and gain insights for safe real-world deployment?

Yes!



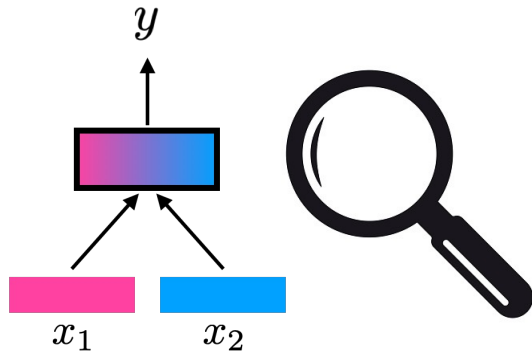
Yes!



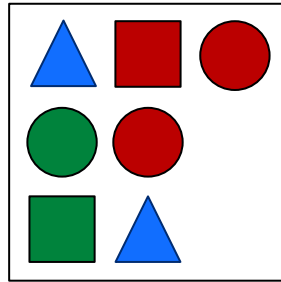
MultiViz: Interpreting Internal Mechanics

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

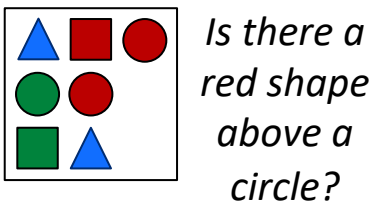
Yes!



1. Unimodal
importance



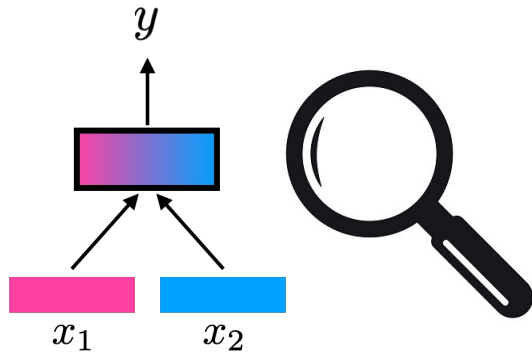
*Is there a **red shape**
above a **circle**?*



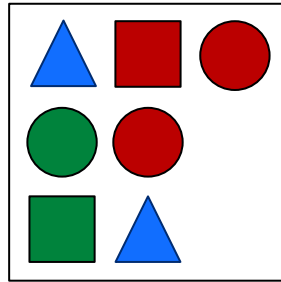
MultiViz: Interpreting Internal Mechanics

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

Yes!

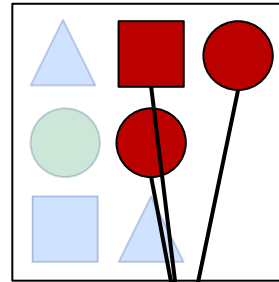


1. Unimodal importance

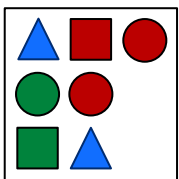
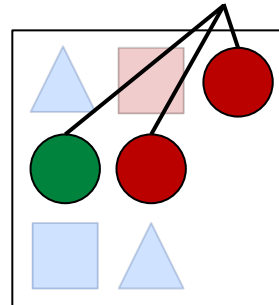


Is there a **red shape** above a **circle**?

2. Cross-modal interactions



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

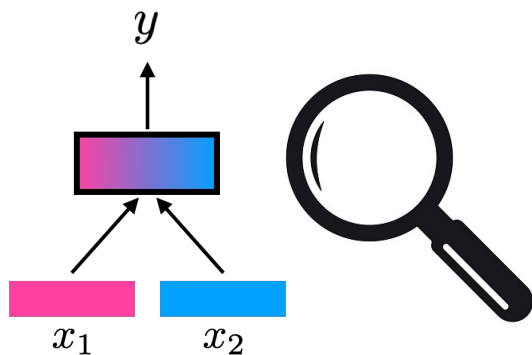


Is there a **red shape** above a **circle**?

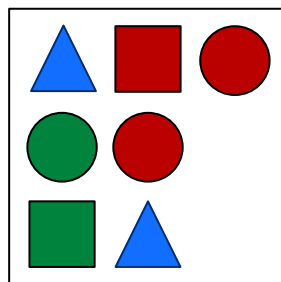
MultiViz: Interpreting Internal Mechanics

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

Yes!

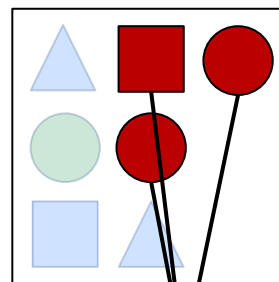


1. Unimodal importance



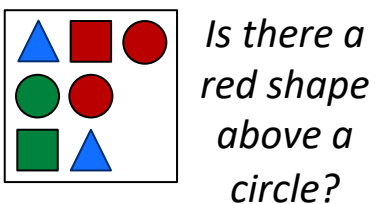
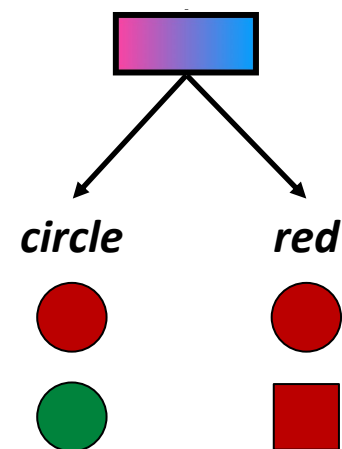
Is there a **red shape** above a **circle**?

2. Cross-modal interactions



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

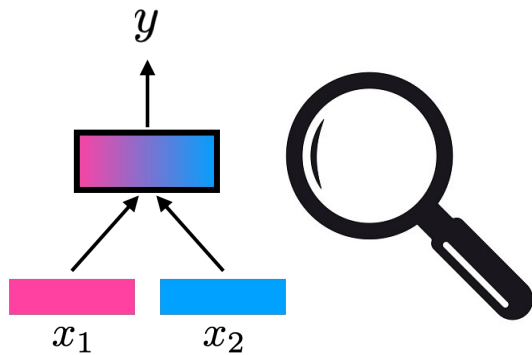
3. Multimodal representations



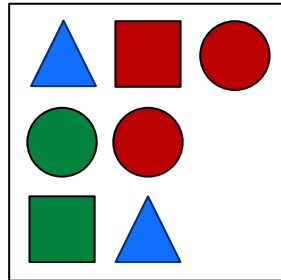
MultiViz: Interpreting Internal Mechanics

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

Yes!

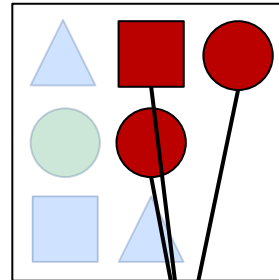


1. Unimodal importance



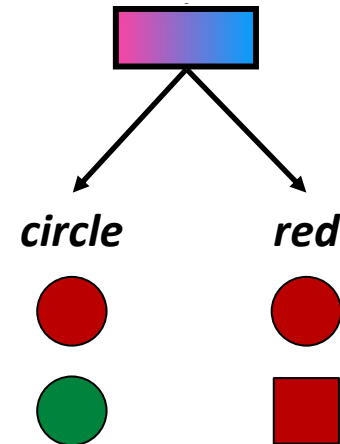
Is there a **red shape** above a **circle**?

2. Cross-modal interactions

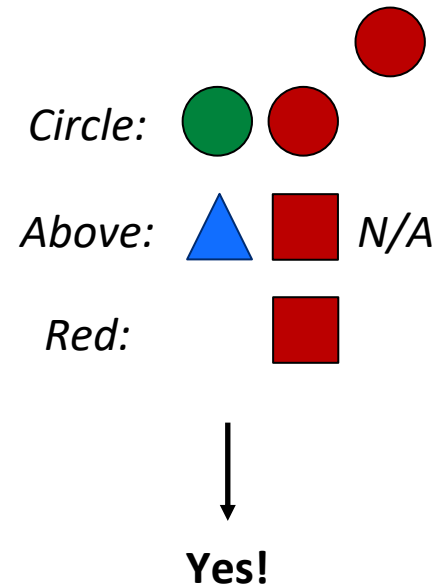


Is there a **red shape** above a **circle**?

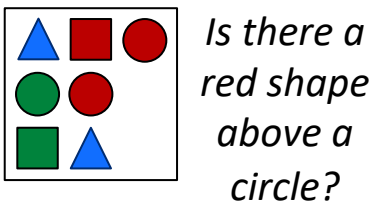
3. Multimodal representations



4. Multimodal prediction



Yes!



MultiViz: Interpreting Internal Mechanics

How can we interpret cross-modal interactions in multimodal models?

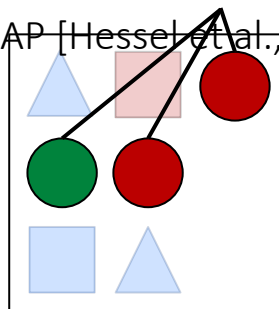
2. Cross-modal
Statistical non-additive interactions [Friedman & Popescu, 2008, Sorokina et al., 2008]
Interactions

f exhibits interactions between 2 features x_1 and x_2 iff
 f cannot be decomposed into a sum of unimodal subfunctions g_1, g_2
such that $f(x_1, x_2) = g_1(x_1) + g_2(x_2)$.

f exhibits interactions between 2 features x_1 and x_2 iff $\frac{\partial^2 f}{\partial x_1 \partial x_2} > 0$.

Is there a **red** shape
above a **circle**?
Natural second-order extension of gradient-based approaches!

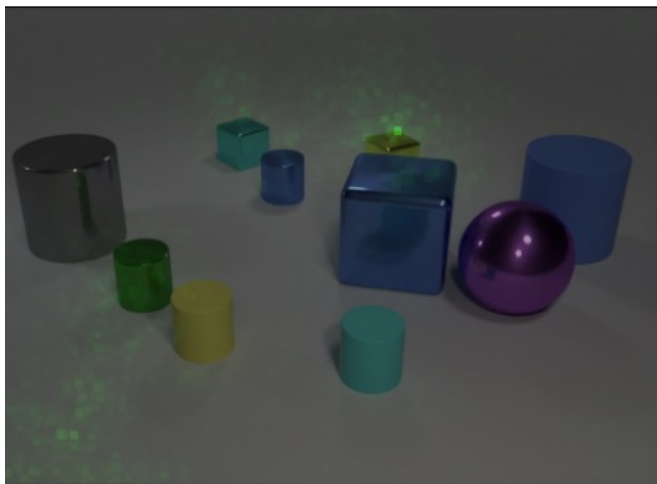
Also related: EMAP [Hessel et al., 2020], DIME [Lyu et al., 2022]



MultiViz: Interpreting Internal Mechanics

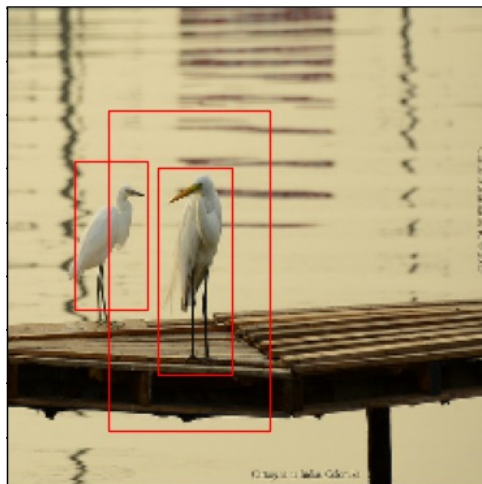
How can we interpret cross-modal interactions in multimodal models?

CLEVR



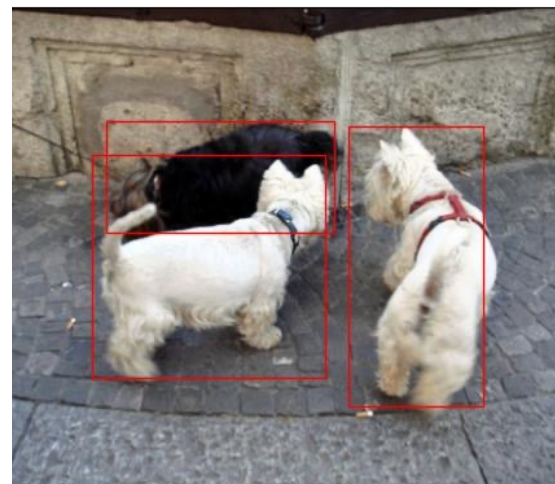
The other small shiny thing that is the same shape as the **tiny yellow shiny object** is what color?

VQA 2.0



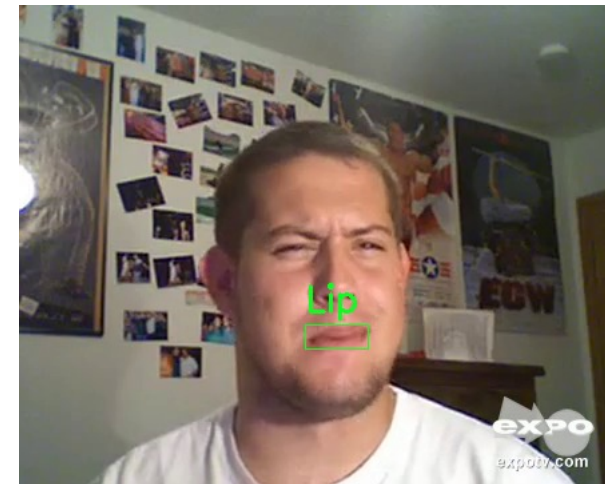
How many **birds**?

Flickr-30k



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

CMU-MOSEI



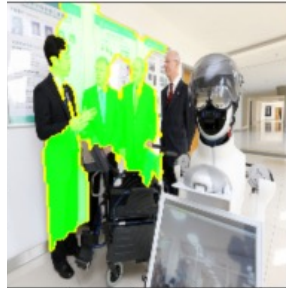
Why am I spending my money watching this? (**sigh**) I think I was more **sad**...

Correspondences

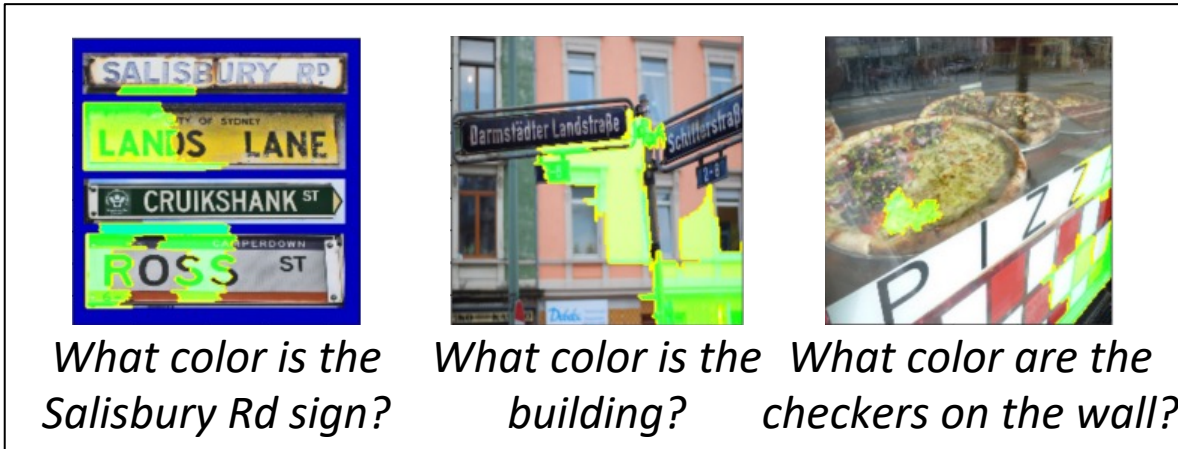
Relationships

MultiViz: Interpreting Internal Mechanics

How can we understand multimodal representations?



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



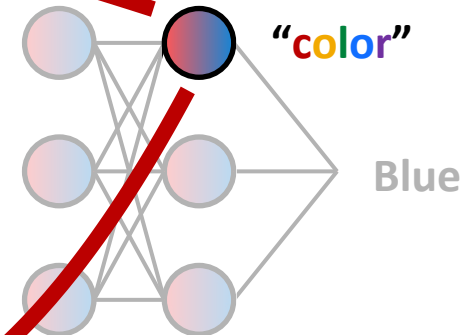
What color is the Salisbury Rd sign?

What color is the building?

What color are the checkers on the wall?

Local analysis

3. Multimodal representations

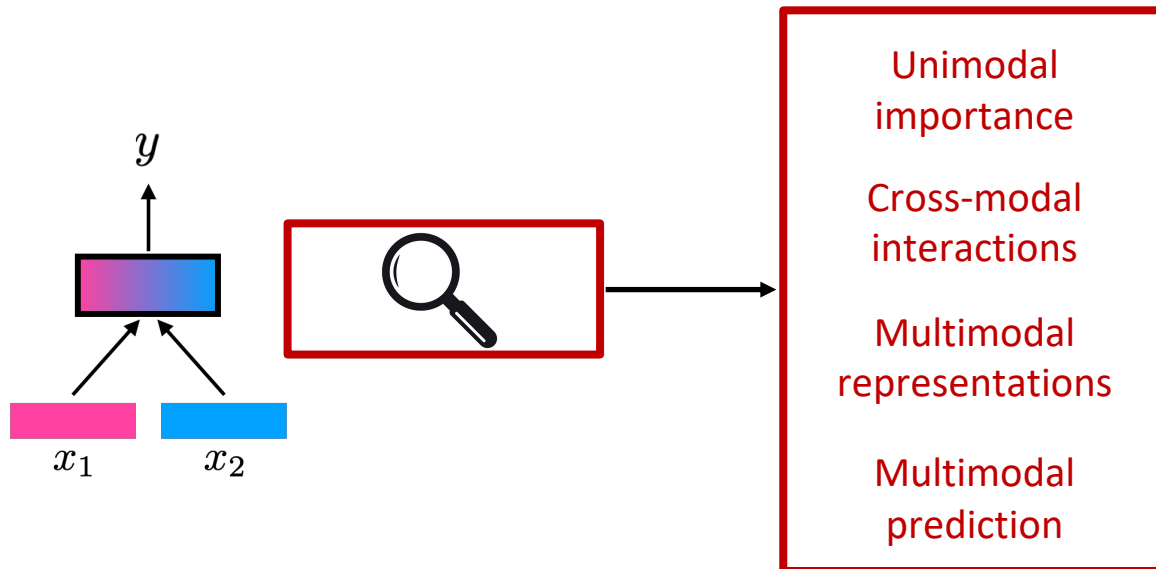


Global analysis

MultiViz: Interpreting Internal Mechanics

How can we evaluate the success of interpreting internal mechanics?

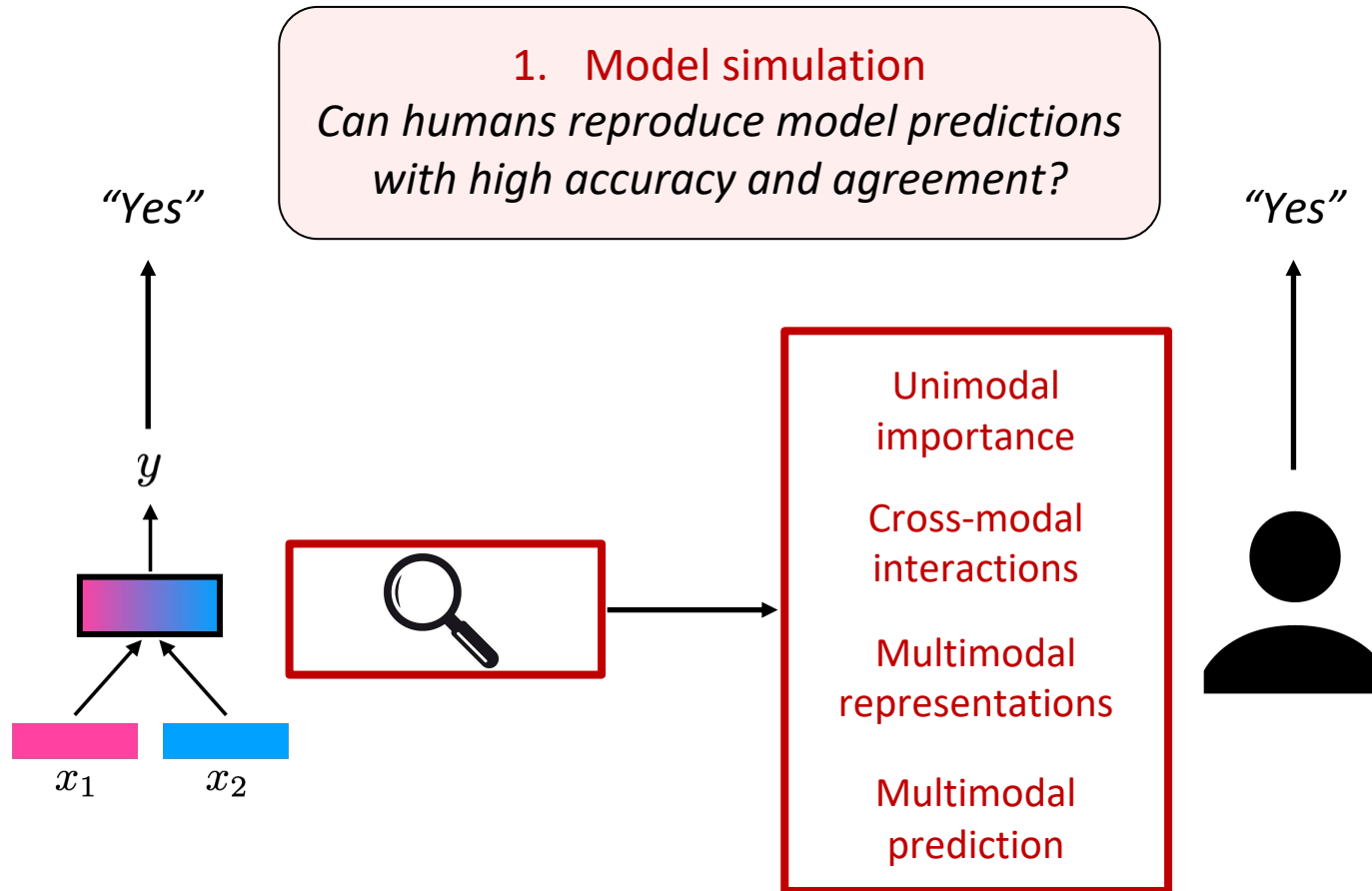
Problem: real-world datasets and models do not have unimodal importance, cross-modal interactions, representations annotated!



MultiViz: Interpreting Internal Mechanics

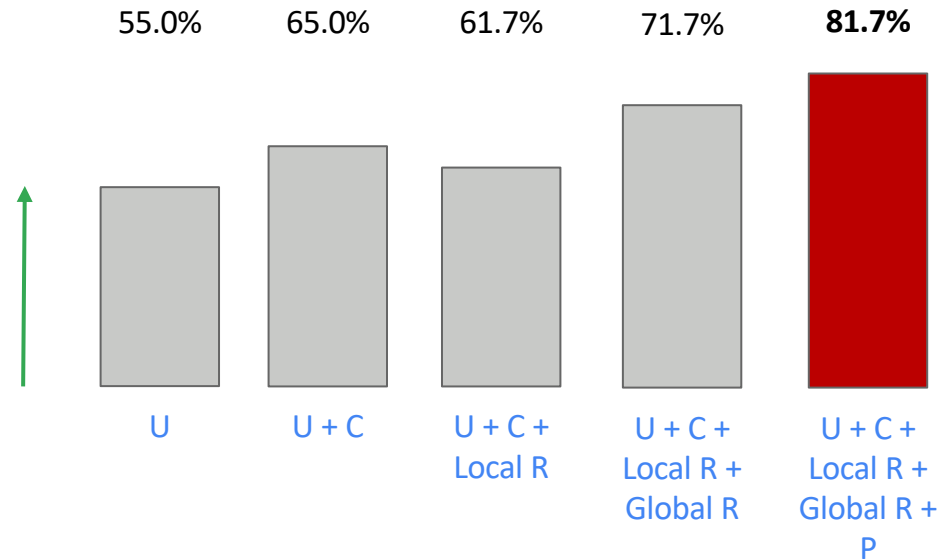
Open
challenges

How can we evaluate the success of interpreting internal mechanics?



MultiViz: Interpreting Internal Mechanics

How can we evaluate the success of interpreting internal mechanics?



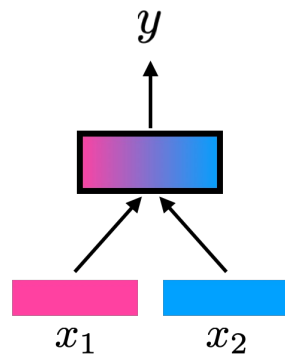
MultiViz stages leads to higher accuracy and agreement

MultiViz: Interpreting Internal Mechanics

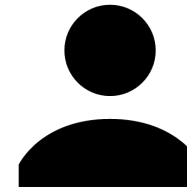
How can we evaluate the success of interpreting internal mechanics?



2. Model debugging
Can humans find bugs in the model
for improvement?



Unimodal importance
Cross-modal interactions
Multimodal representations
Multimodal prediction



Find bugs



Fix bugs

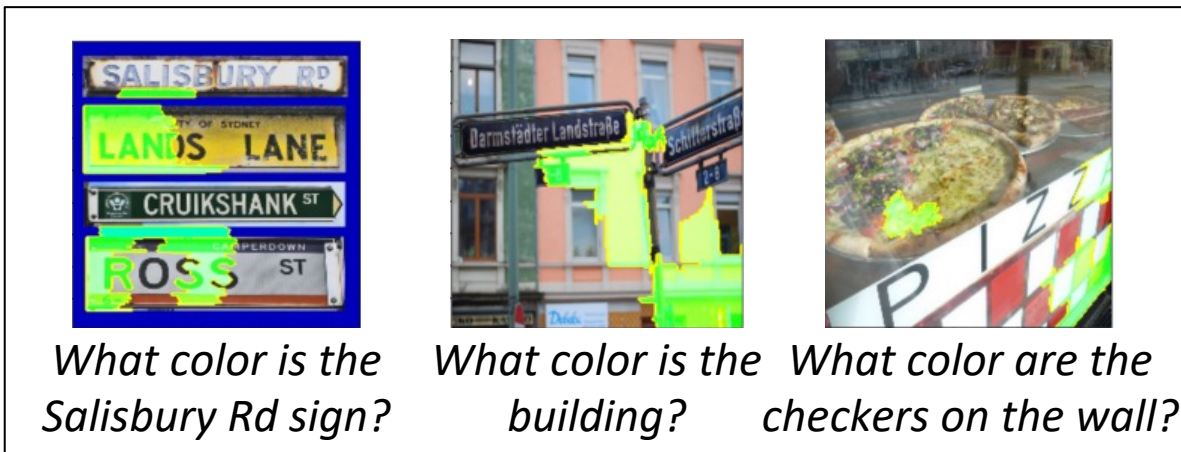


MultiViz: Interpreting Internal Mechanics

How can we understand multimodal representations?



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



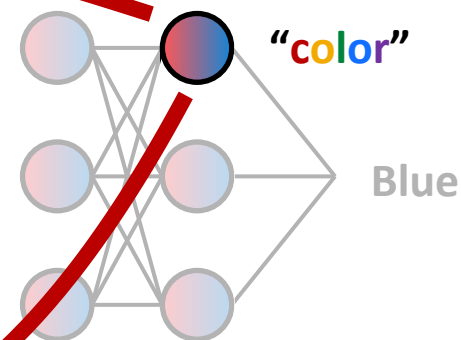
What color is the Salisbury Rd sign?

What color is the building?

What color are the checkers on the wall?

Local analysis

3. Multimodal representations



Global analysis

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

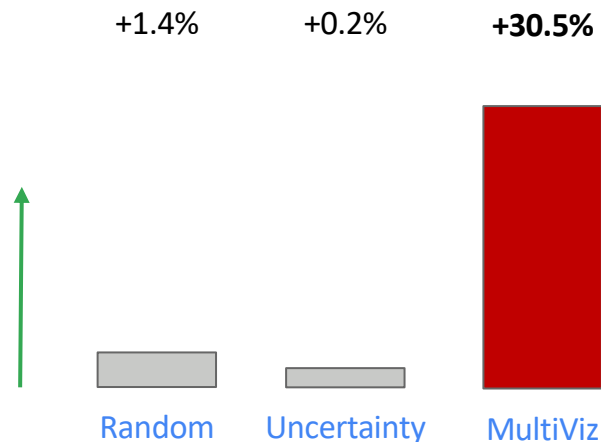
MultiViz: Interpreting Internal Mechanics

How can we evaluate the success of interpreting internal mechanics?

“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.



MultiViz enables error analysis and debugging of multimodal models

What is Multimodal?

Heterogeneous



Connected



Interacting



Why is it hard?

Representation

Alignment

Reasoning

Generation

Transference

Quantification



What is next?

Heterogeneity

High-modality

Long-term

Interaction

Real-world

<https://cmu-multicomp-lab.github.io/mmml-course/fall2022/>

[Liang, Zadeh, and Morency. Foundations and Trends on Multimodal Machine Learning. arXiv 2022]

<https://www.cs.cmu.edu/~pliang/>

pliang@cs.cmu.edu

 [@pliang279](https://twitter.com/pliang279)