# Multimodal Machine Learning

## Lecture 12.2: Quantification

Paul Liang

# Quantification

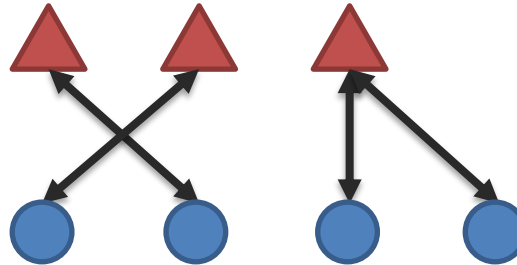**Definition:** Empirical and theoretical study to better understand heterogeneity, cross-modal interactions, and the multimodal learning process.
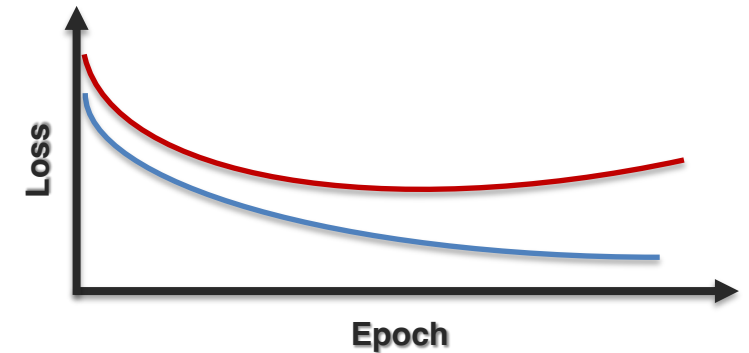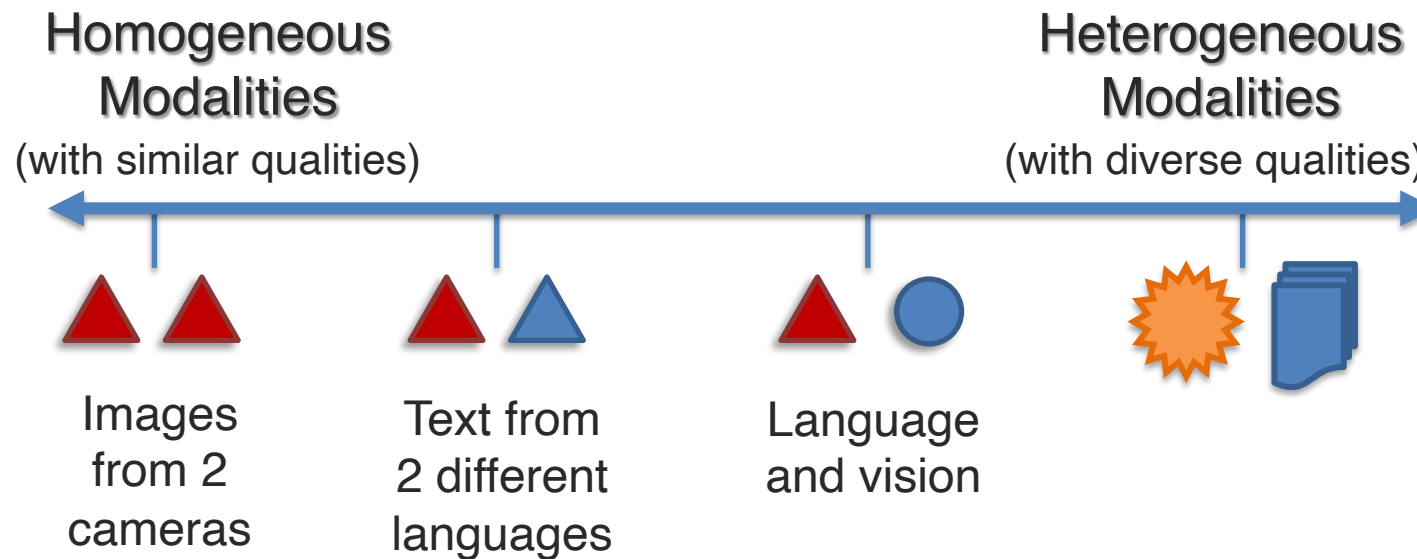


Ⓐ **Heterogeneity**

Ⓑ **Interactions**

Ⓒ **Learning**

# Sub-Challenge 6a: Heterogeneity

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

Modality A ▲

Modality B ●

Homogeneous Modalities
(with similar qualities)

Heterogeneous Modalities
(with diverse qualities)

**Examples:**

Images from 2 cameras

Text from 2 different languages

Language and vision

① **Element representation**

② **Element distribution**

③ **Structure**

④ **Information**

⑤ **Noise**

⑥ **Relevance**

# Distribution heterogeneity

**Inspired by distributed learning**



[Ye et al., Heterogeneous Federated Learning: State-of-the-art and Research Challenges, 2023]
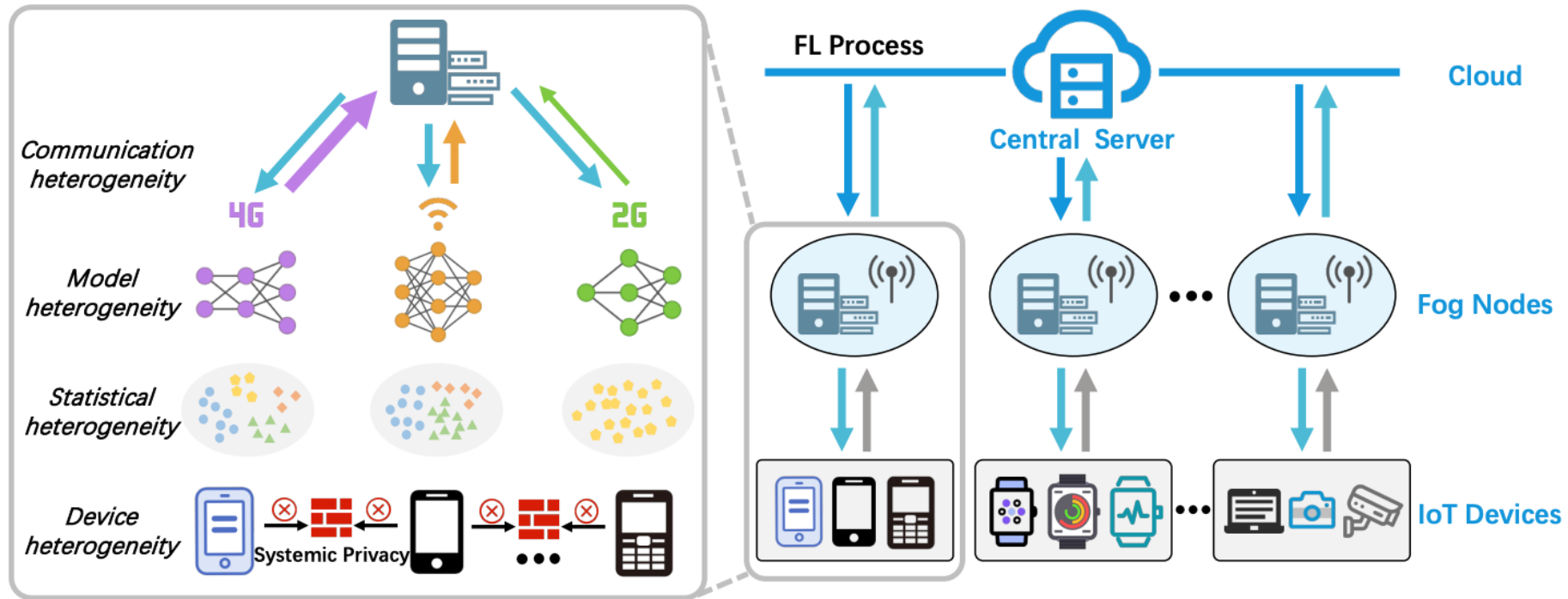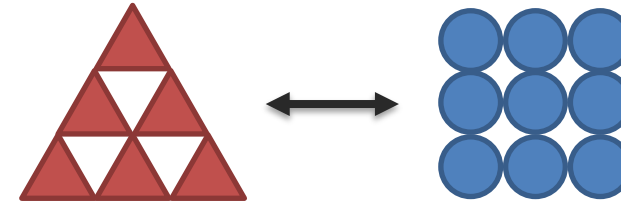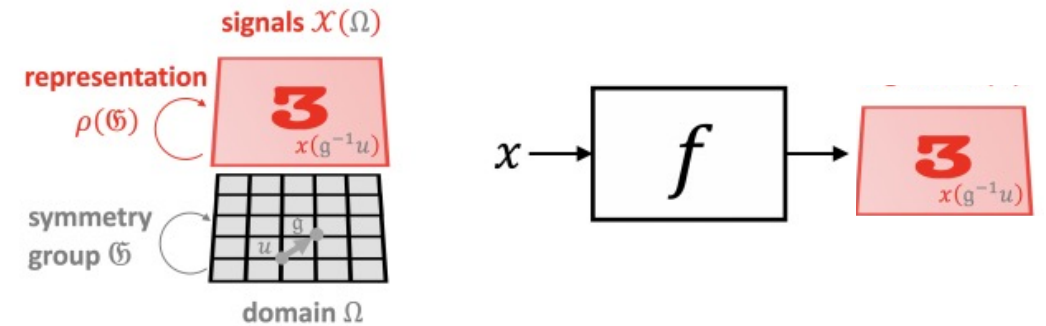
# Structure heterogeneity

**Inspired by structure learning**



A function $f : \mathcal{X}(\Omega) \to \mathcal{Y}$ is $\mathfrak{G}$-*invariant* if $f(\rho(\mathfrak{g})x) = f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega)$, i.e., its output is unaffected by the group action on the input.

A function $f : \mathcal{X}(\Omega) \to \mathcal{X}(\Omega)$ is $\mathfrak{G}$-*equivariant* if $f(\rho(\mathfrak{g})x) = \rho(\mathfrak{g})f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$, i.e., group action on the input affects the output in the same way.



Derives deep learning architectures on grids, graphs, sets, etc.

[Bronstein et al., Geometric Deep Learning Grids, Groups, Graphs, Geodesics, and Gauges, 2021]

# Modality Biases

**Heterogeneity in information and relevance**
Unimodal biases and modality collapse
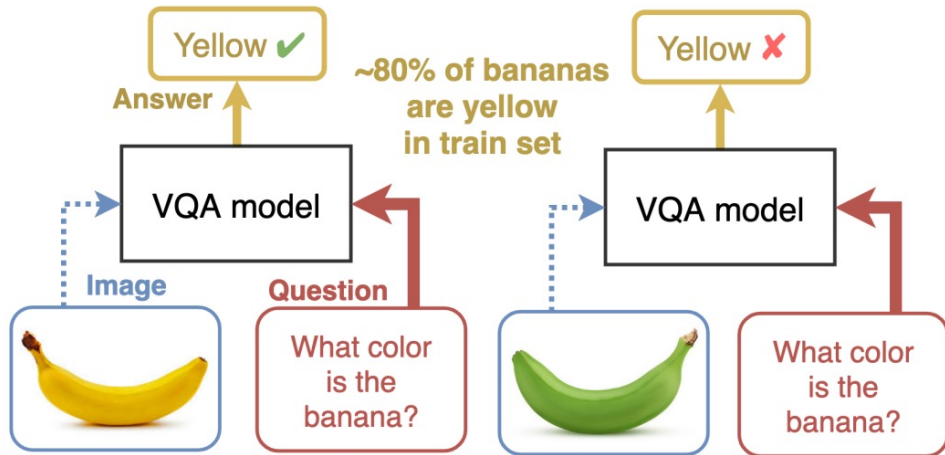


Balancing modalities

Balancing training

Who is wearing glasses?
man          woman

Where is the child sitting?
fridge          arms

Is the umbrella upside down?
yes          no

How many children are in the bed?
2          1

**Not the case when trained with RUBi**

VQA models answer the question without looking at the image

Yellow ✔
Answer
~80% of bananas are yellow in train set
Yellow ✘
VQA model          VQA model
Image          Question
What color is the banana?          What color is the banana?

Green ✔
modalities used adequatly
Same VQA model
What color is the banana?

[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



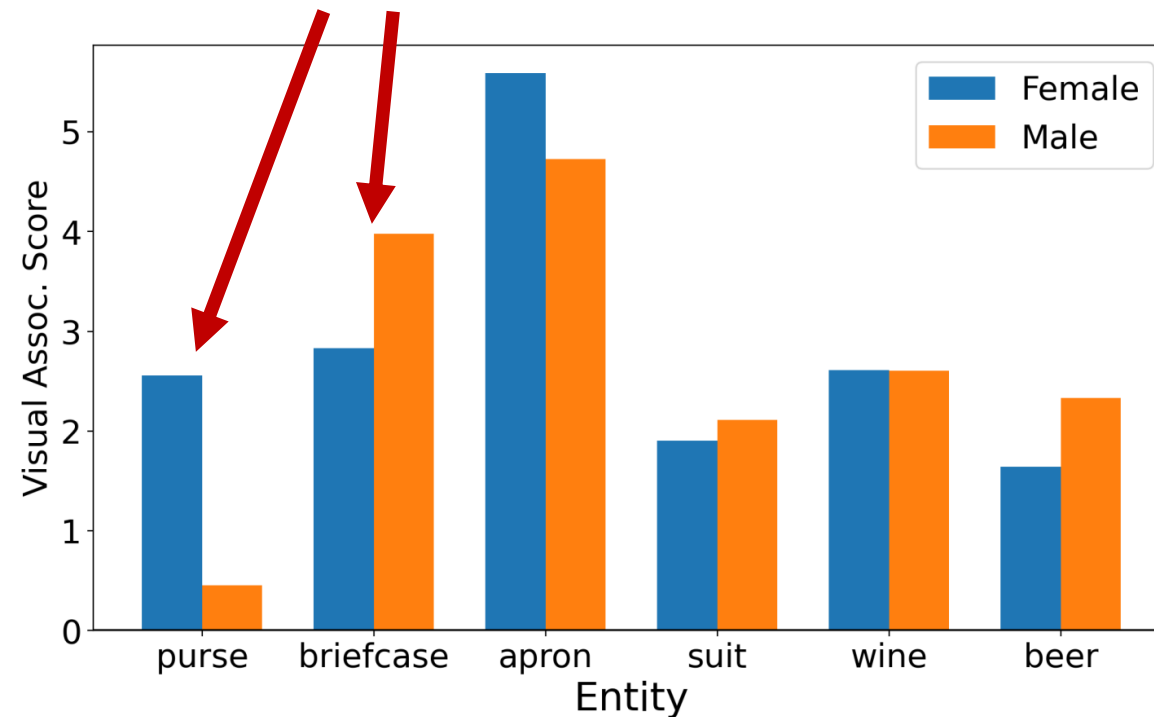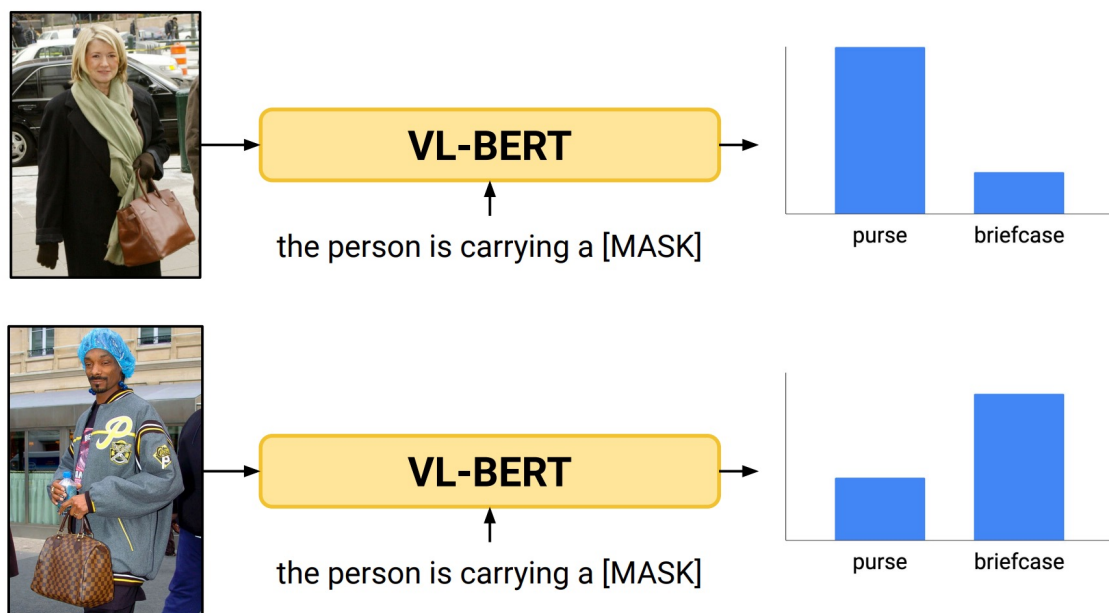| Wrong | Right for the Right Reasons | Right for the Wrong Reasons | Right for the Right Reasons |
|---|---|---|---|
| Baseline: *A **man** sitting at a desk with a laptop computer.* | Our Model: *A **woman** sitting in front of a laptop computer.* | Baseline: *A **man** holding a tennis racquet on a tennis court.* | Our Model: *A **man** holding a tennis racquet on a tennis court.* |

[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



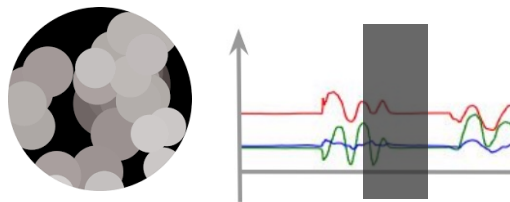Visual information makes model more confident in reinforcing gender stereotypes

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

# Noise Topologies and Robustness
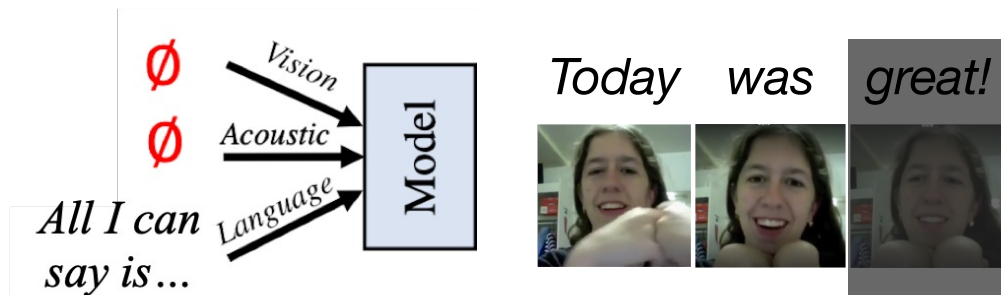
## Heterogeneity in noise

### Modality-specific robustness

noise → nosie

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

### Multimodal robustness

Vision, Acoustic, Language → Model

All I can say is…

Today was great!

[Zadeh et al., 2020]

## Strong tradeoffs between performance and robustness



Performance →

CCA  MulT  RMFE  MVAE

LRTF  MI  GradBlend  MFAS

LF  EF  ReFNet  SF

Best Unimodal  MFM  MCTN  TF

Robustness →

rate of accuracy drops

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

# Noise Topologies and Robustness

**Several approaches towards more robust models**

Robust data + training



Infer missing modalities



Translation model
Joint probabilistic model

[Ngiam et al., Multimodal Deep Learning. ICML 2011]
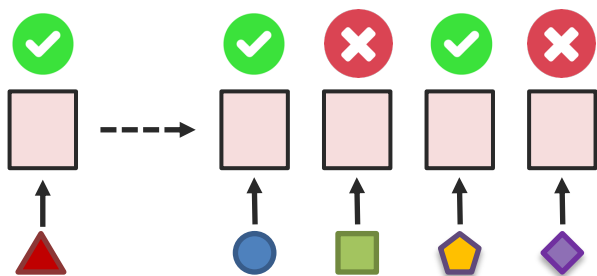[Srivastava and Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines. JMLR 2014]
[Tran et al., Missing Modalities Imputation via Cascaded Residual Autoencoder. CVPR 2017]
[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAAI 2019]

# Quantifying Heterogeneity via Transfer

## Information transfer, transfer learning perspective

**1a. Estimate modality heterogeneity via transfer**



In practice, efficient by pre-trained models and few-shot transfer

## Implicitly captures these:

1. **Element representation**
2. **Element distribution**
3. **Structure**
4. **Information**
5. **Noise**
6. **Relevance**

[Liang et al., HighMMT: Quantifying Modality & Interaction Heterogeneity for High-Modality Learning. TMLR 2022]

# Heterogeneity-aware Fusion

## Information transfer, transfer learning perspective

**1a. Estimate modality heterogeneity via transfer**



**(Implicitly captures heterogeneity)**

**2a. Compute modality heterogeneity matrix**



**3. Determine parameter clustering**

$$\mathbb{U}_1 = \{U_1, U_2, U_4\}$$
$$\mathbb{U}_2 = \{U_3\}$$
$$\mathbb{U}_3 = \{U_5\}$$

$$\mathbb{C}_1 = \{C_{12}, C_{13}, C_{45}\}$$
$$\mathbb{C}_2 = \{C_{23}\}$$

**1b. Estimate interaction heterogeneity via transfer**



**2b. Compute interaction heterogeneity matrix**



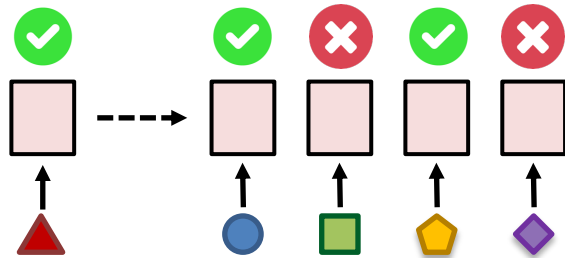[Zamir et al., Taskonomy: Disentangling Task Transfer Learning. CVPR 2018]
[Liang et al., HighMMT: Quantifying Modality & Interaction Heterogeneity for High-Modality Learning. TMLR 2022]
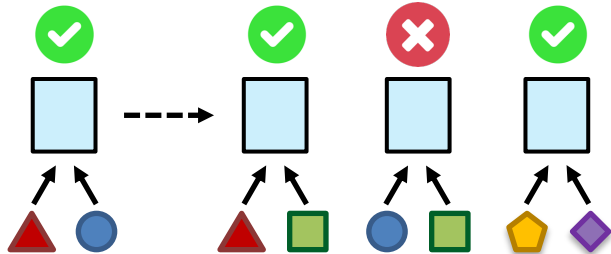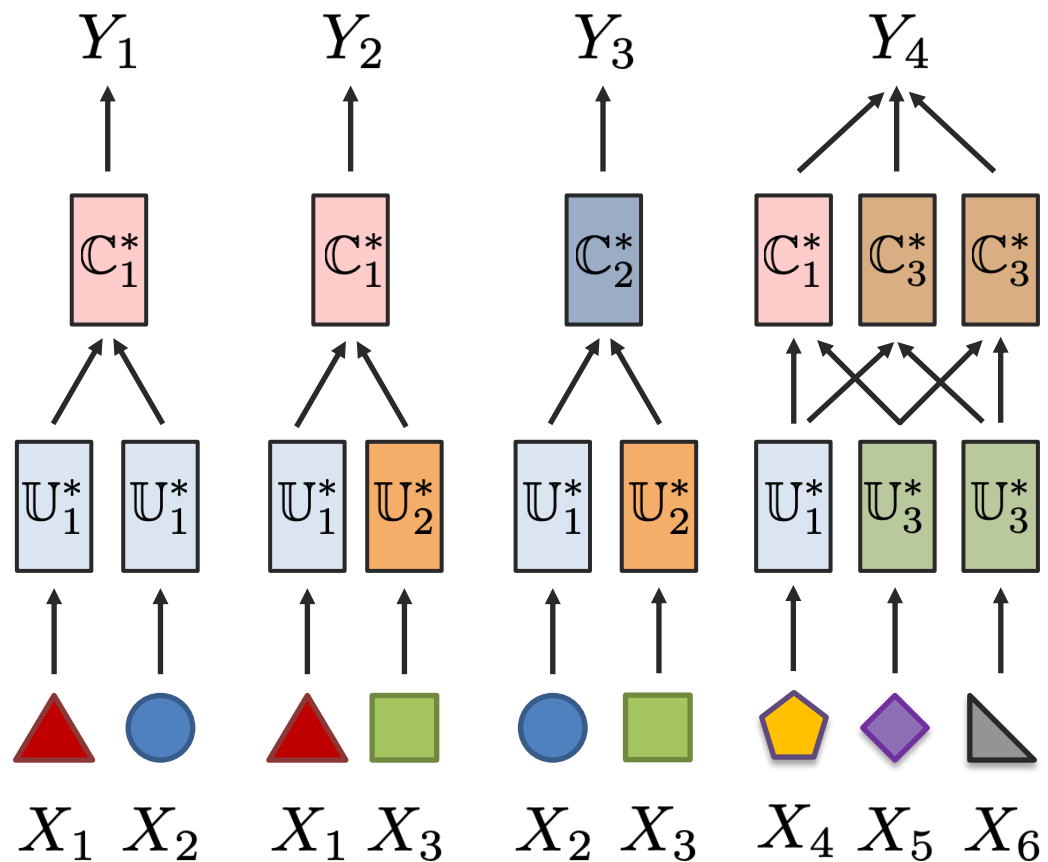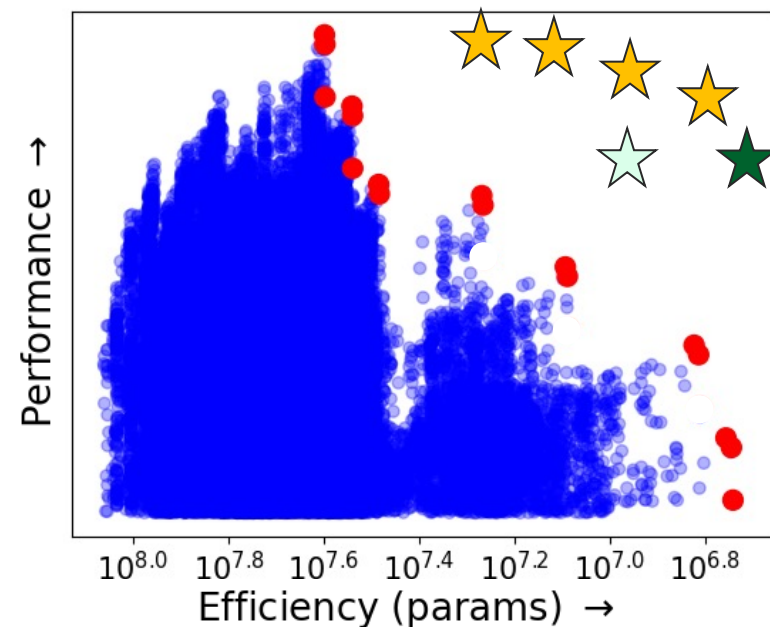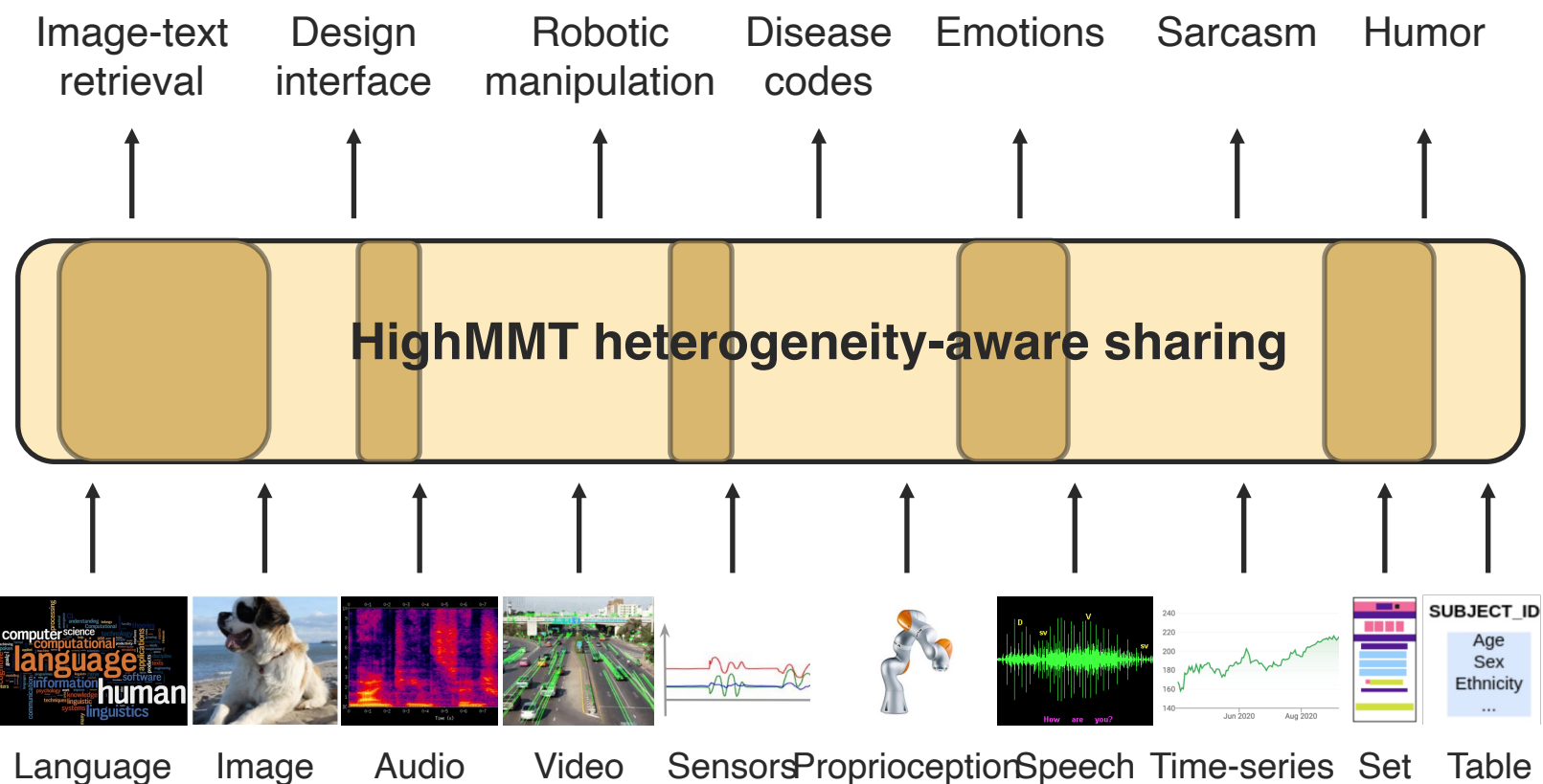
# Heterogeneity-aware Fusion

**Information transfer, transfer learning perspective**



[Liang et al., HighMMT: Quantifying Modality & Interaction Heterogeneity for High-Modality Learning. TMLR 2022]

# Quantifying Modality Heterogeneity

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



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

**HighMMT heterogeneity-aware sharing**

Language | Image | Audio | Video | Sensors | Proprioception | Speech | Time-series | Set | Table

- Performance ↑
- Efficiency (params) →
- $10^{8.0}$ $10^{7.8}$ $10^{7.6}$ $10^{7.4}$ $10^{7.2}$ $10^{7.0}$ $10^{6.8}$

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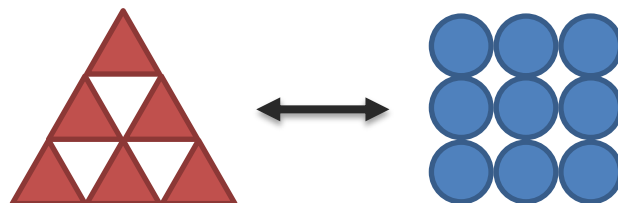- ● **HighMMT heterogeneity-aware**

[Liang et al., HighMMT: Quantifying Modality & Interaction Heterogeneity for High-Modality Learning. TMLR 2022]

# Challenges: Quantifying Heterogeneity

**Open challenges:**
- Noisy and missing modalities.
- New and understudied modalities.
- Large number of modalities.
- Cases where its unclear which modalities are useful – active selection

- Related fields: federated learning, active learning, distributed systems, structure & invariances

# Sub-Challenge 6b: Cross-modal Connections

Connected: Shared information that relates modalities



Modality A ▲
Modality B ●

unique
unique

stronger
weaker
unconnected

## Statistical

**Association**

e.g., correlation, co-occurrence

**Dependency**

e.g., causal, temporal

## Semantic

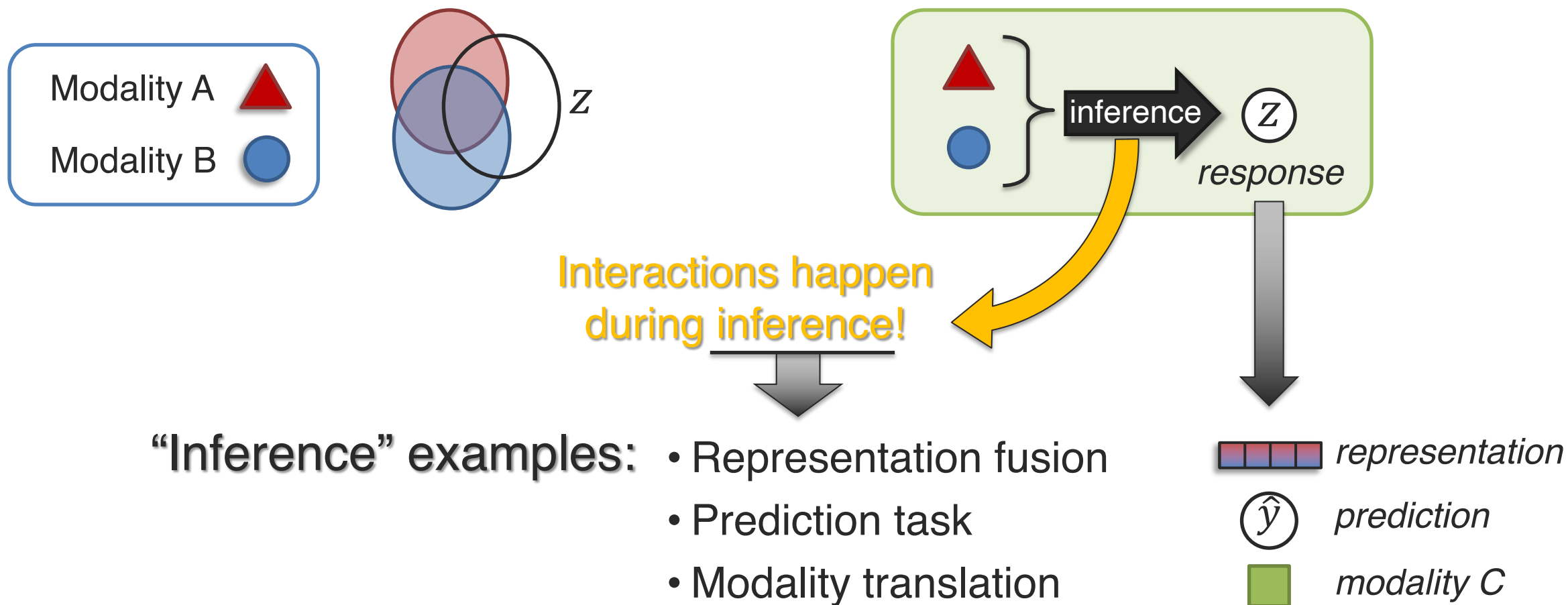**Correspondence**

laptop

e.g., grounding

**Relationship**

used for

e.g., function

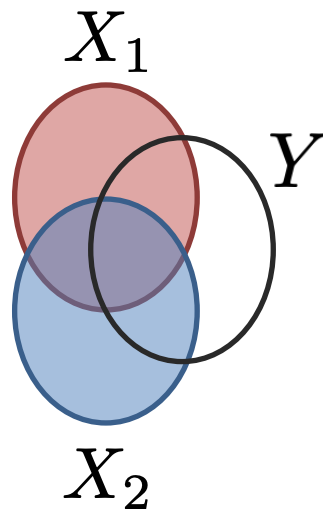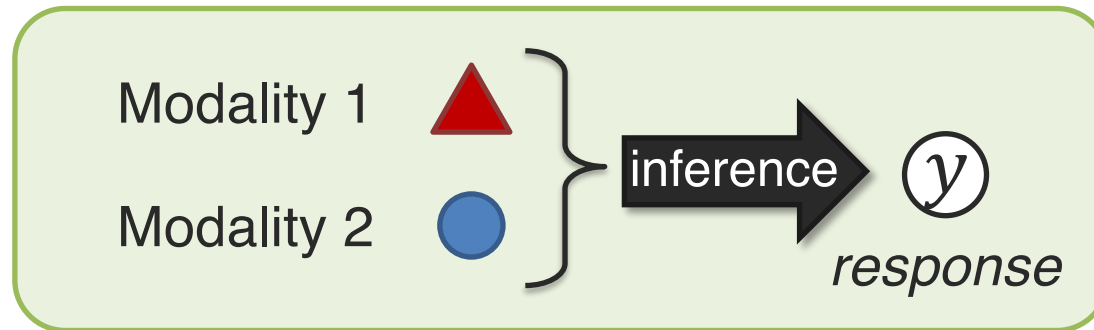# Sub-Challenge 6b: Cross-modal Interactions

**Interacting:** process affecting each modality, creating new response

Modality A ▲

Modality B ●

$z$

▲
●
} **inference** ➡ $z$
*response*

**Interactions happen during inference!**

"Inference" examples:
- Representation fusion
- Prediction task
- Modality translation

*representation*

$\hat{y}$ *prediction*

*modality C*

# Part 1: Multimodal Interactions

**Interactions:** Understanding *commonalities* between modalities and how they *combine* to provide information for a task.



[Liang, Zadeh, and Morency. Foundations and Trends on Multimodal Machine Learning. ICML 2023, CVPR 2022, NAACL 2022 Tutorials]

# Multimodal Interactions

**Interactions:** Understanding *commonalities* between modalities and how they *combine* to provide information for a task.



Redundancy: Shared by both modalities and task

[Liang, Zadeh, and Morency. Foundations and Trends on Multimodal Machine Learning. ICML 2023, CVPR 2022, NAACL 2022 Tutorials]
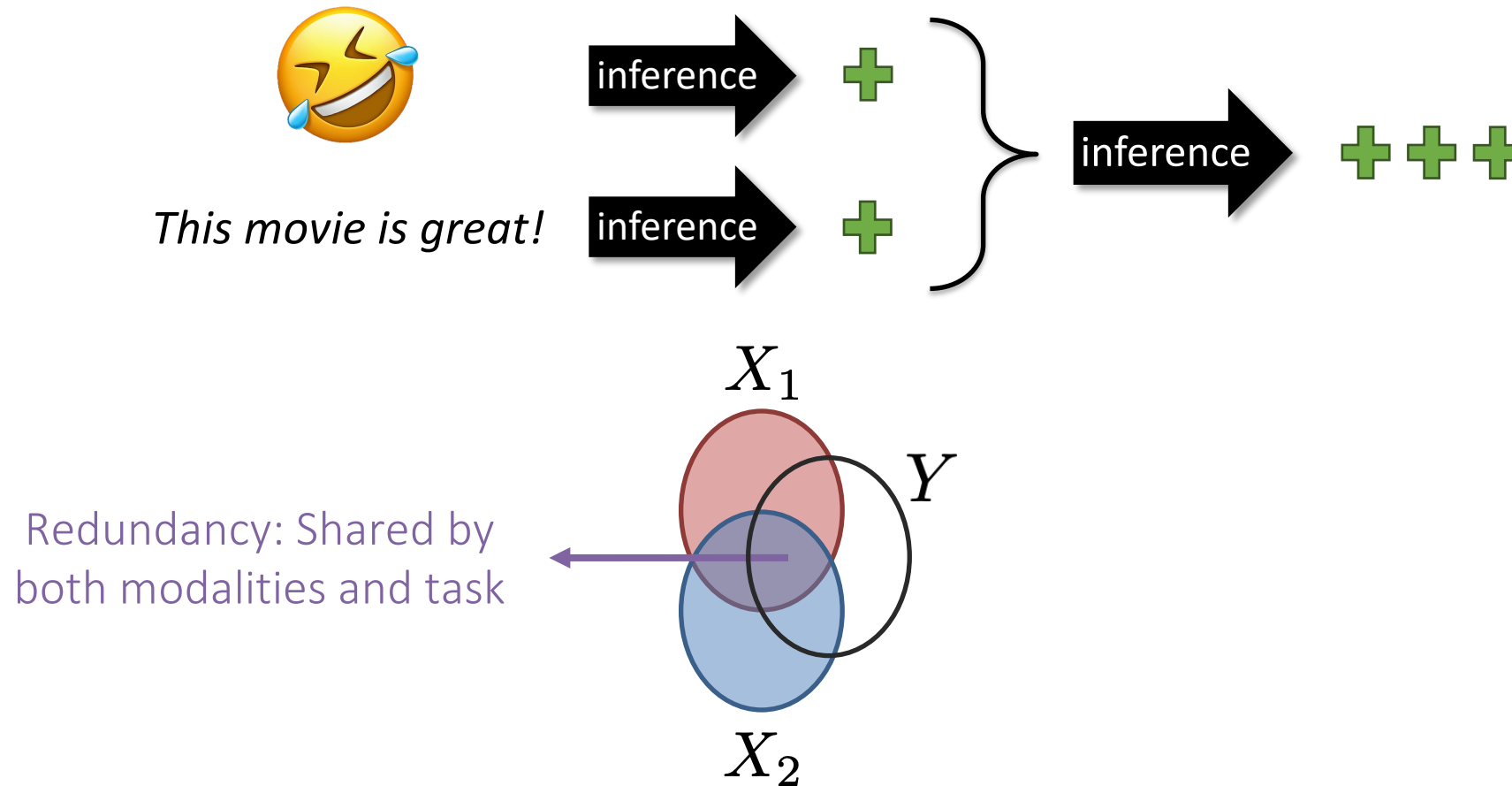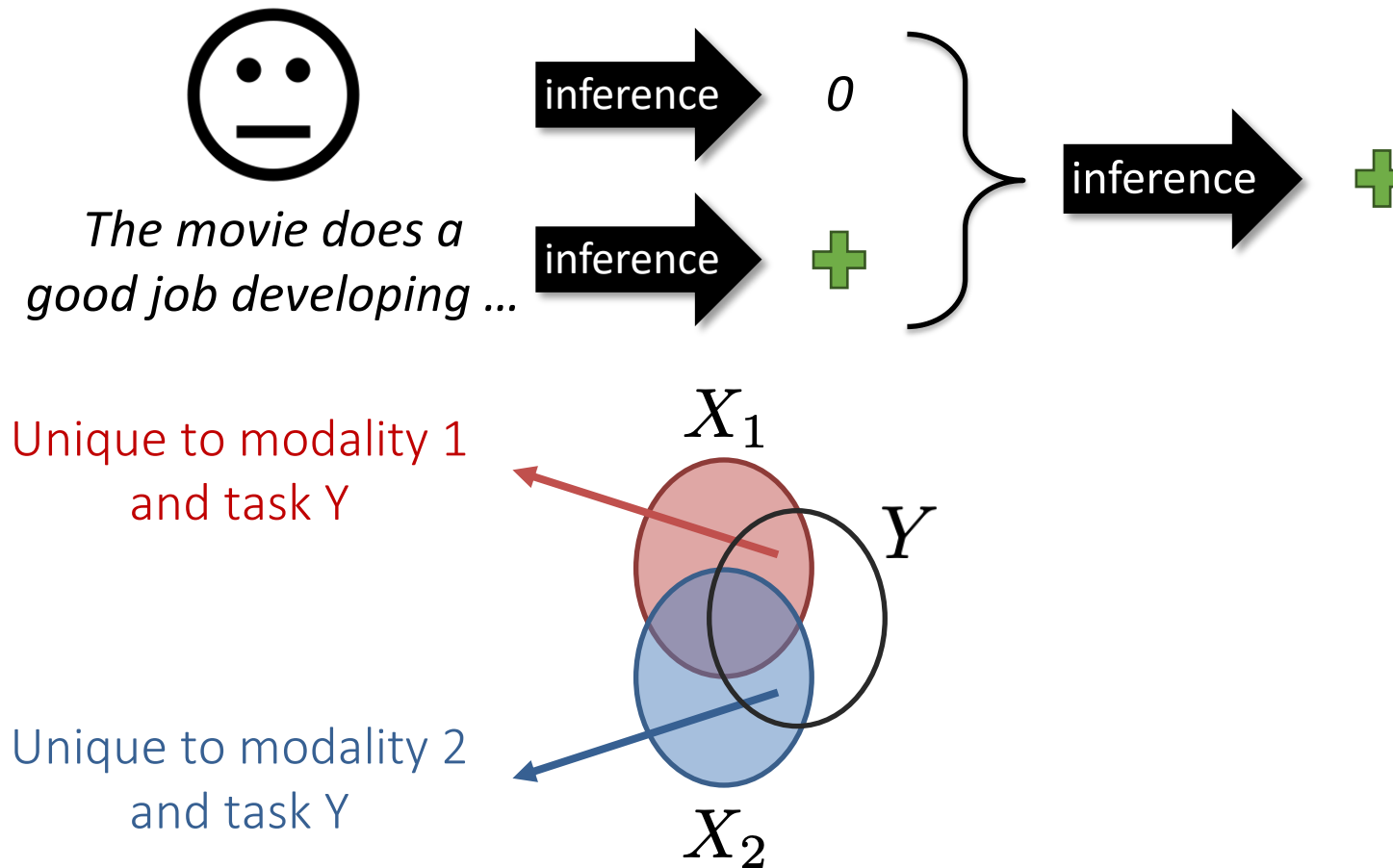
# Multimodal Interactions

**Interactions:** Understanding *commonalities* between modalities and how they *combine* to provide information for a task.



Unique to modality 1 and task Y

Unique to modality 2 and task Y

$X_1$

$X_2$

$Y$

inference → 0

inference → +

inference → +

*The movie does a good job developing ...*

[Liang, Zadeh, and Morency. Foundations and Trends on Multimodal Machine Learning. ICML 2023, CVPR 2022, NAACL 2022 Tutorials]

# Multimodal Interactions

**Interactions:** Understanding *commonalities* between modalities and how they *combine* to provide information for a task.
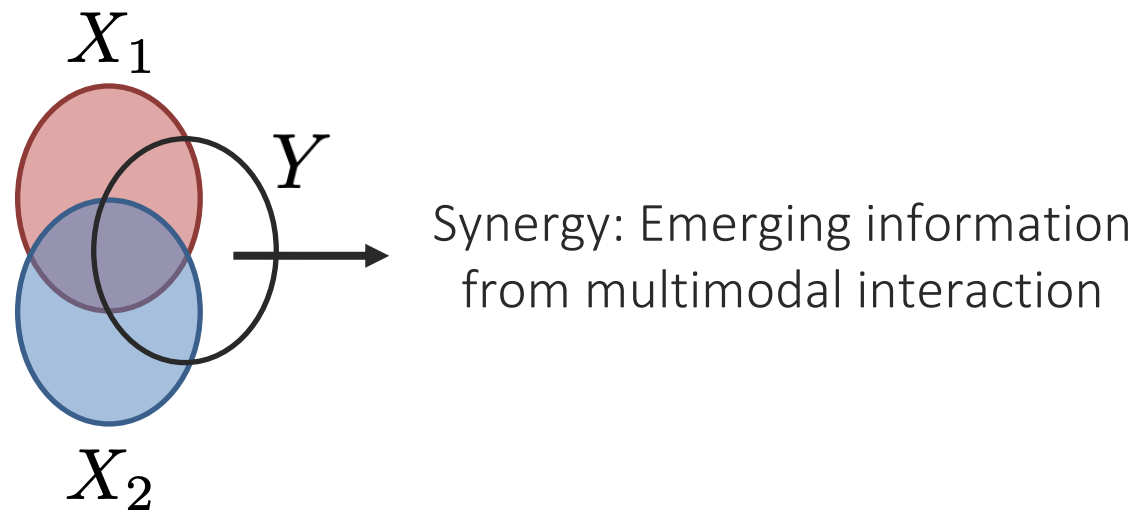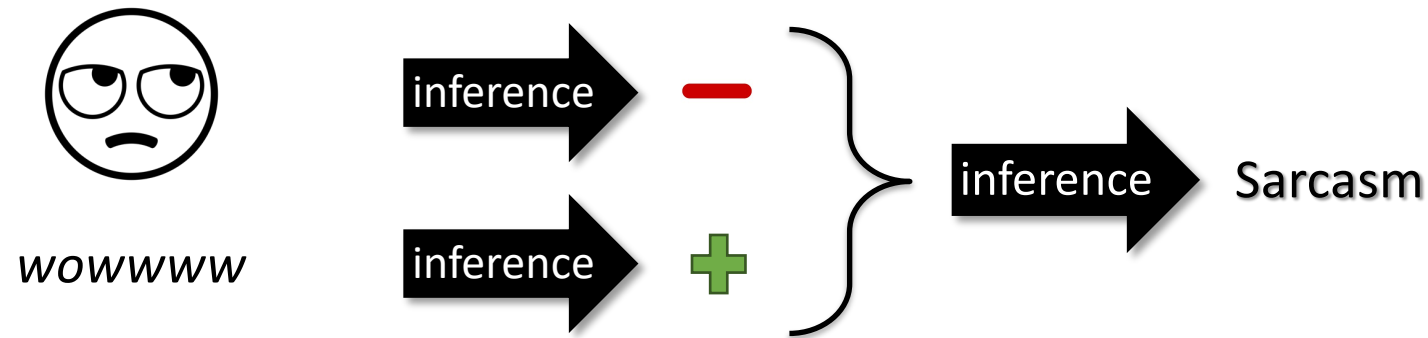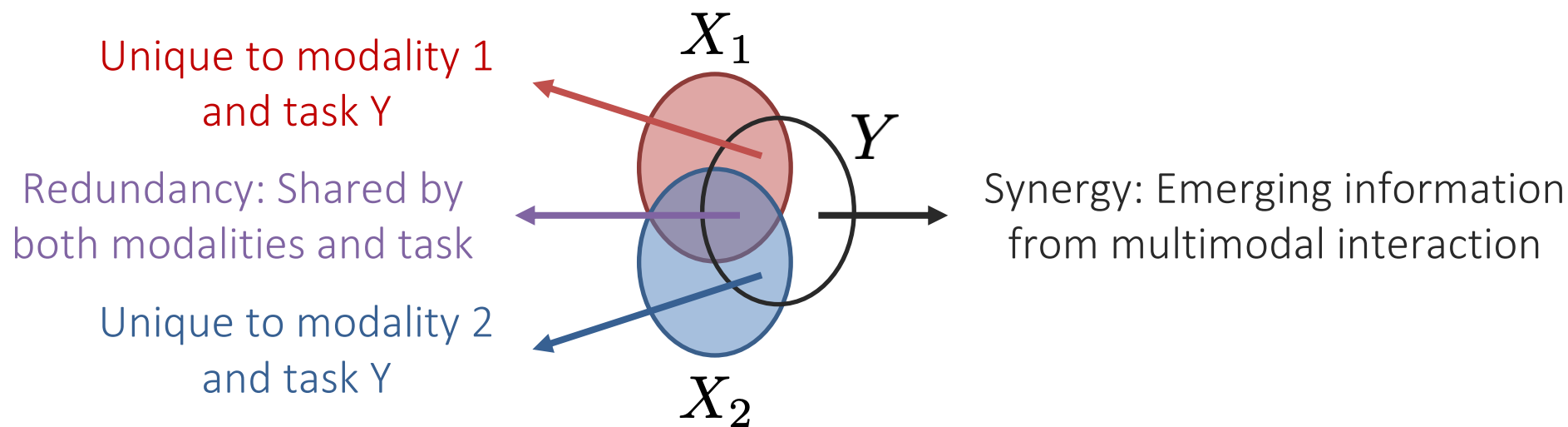


*wowwww*

inference　➖

inference　➕

inference　**Sarcasm**

$X_1$

$Y$

$X_2$

Synergy: Emerging information from multimodal interaction

[Liang, Zadeh, and Morency. Foundations and Trends on Multimodal Machine Learning. ICML 2023, CVPR 2022, NAACL 2022 Tutorials]

# Quantifying Multimodal Interactions

**Fundamental questions in multimodal learning**

*What interactions are in my data?*

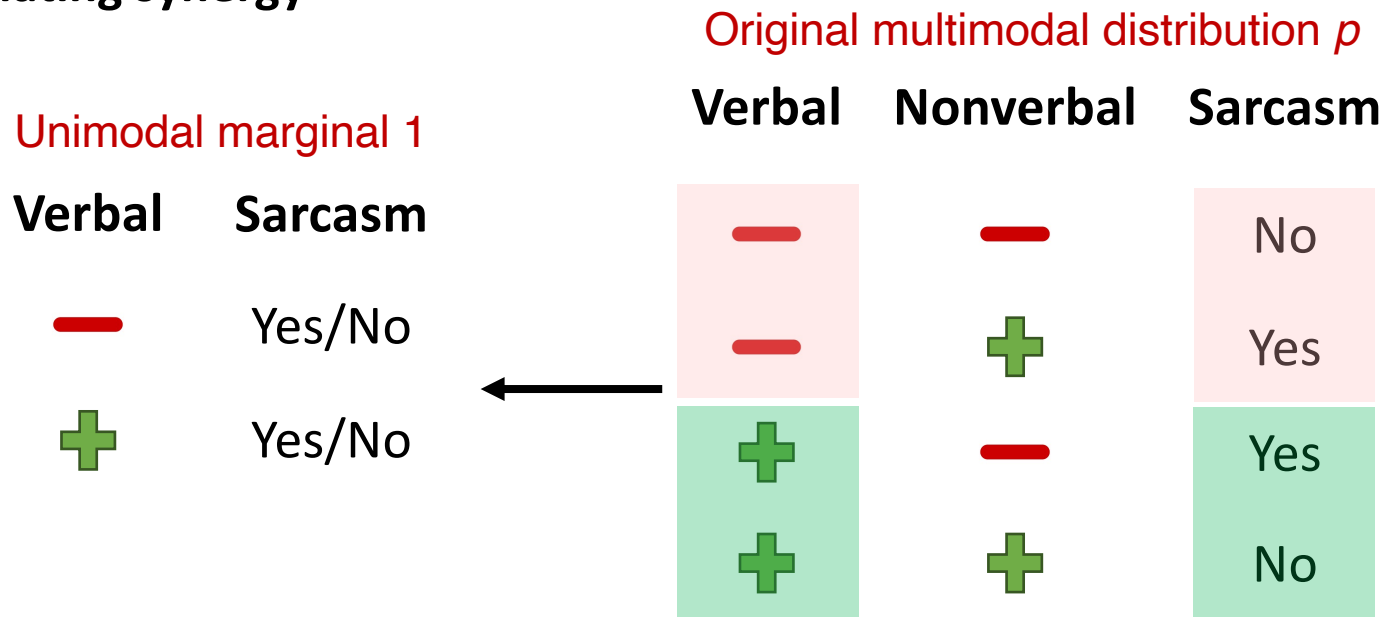*What interactions do models learn?*

*What models are suitable for my data?*



Unique to modality 1 and task Y

Redundancy: Shared by both modalities and task

Unique to modality 2 and task Y

$X_1$

$Y$

$X_2$

Synergy: Emerging information from multimodal interaction

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Mathematical Framework for Multimodal Interactions

**Estimating synergy**

Original multimodal distribution $p$

Unimodal marginal 1

**Verbal    Nonverbal    Sarcasm**

**Verbal       Sarcasm**

— Yes/No

➕ Yes/No

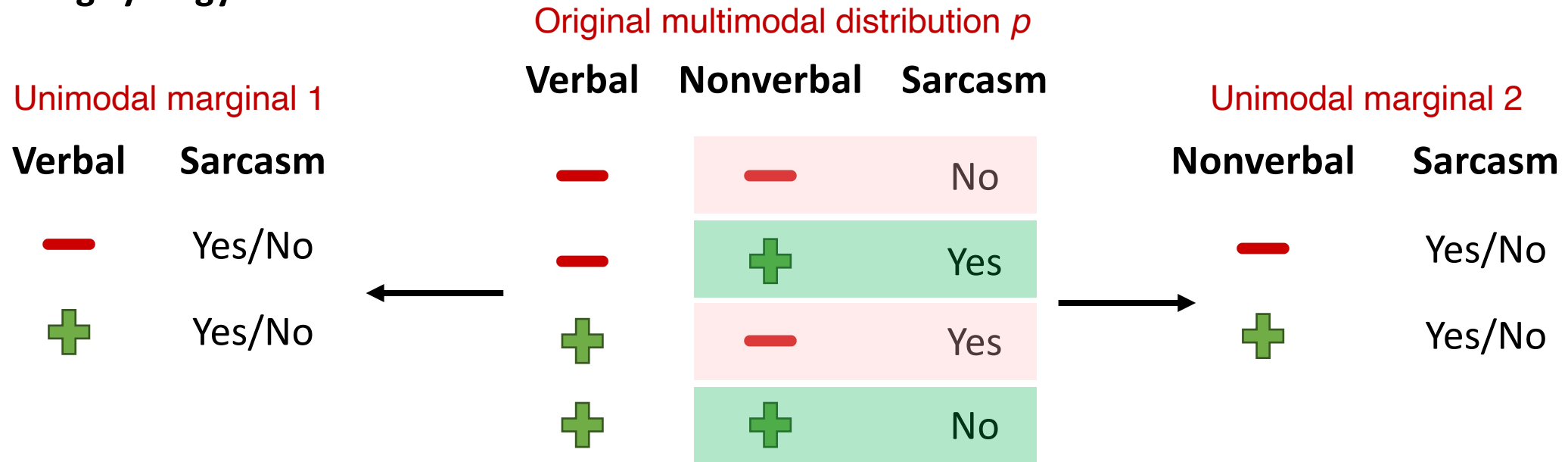| Verbal | Nonverbal | Sarcasm |
|---|---|---|
| — | — | No |
| — | ➕ | Yes |
| ➕ | — | Yes |
| ➕ | ➕ | No |

Synergy = Original multimodal information about the task
   − multimodal information given by the ***worst*** distribution combining the same modalities

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Mathematical Framework for Multimodal Interactions

**Estimating synergy**

Original multimodal distribution *p*

**Unimodal marginal 1**

| **Verbal** | **Nonverbal** | **Sarcasm** |
|---|---|---|

**Unimodal marginal 2**

| **Verbal** | **Sarcasm** |
|---|---|

| **Nonverbal** | **Sarcasm** |
|---|---|

— Yes/No

— — No

— Yes/No

+ Yes/No

— + Yes

+ Yes/No

+ — Yes

+ + No

Synergy = Original multimodal information about the task
— multimodal information given by the ***worst*** distribution combining the same modalities

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Mathematical Framework for Multimodal Interactions

**Many ways of combining these 2 unimodal marginals into a multimodal distribution!**

| Unimodal marginal 1 | | Original multimodal distribution $p$ | | | Unimodal marginal 2 | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| **Verbal** | **Sarcasm** | **Verbal** | **Nonverbal** | **Sarcasm** | **Nonverbal** | **Sarcasm** |
| ▬ | Yes/No | ▬ | ▬ | No | ▬ | Yes/No |
| ✚ | Yes/No | ▬ | ✚ | Yes | ✚ | Yes/No |
| | | ✚ | ▬ | Yes | | |
| | | ✚ | ✚ | No | | |

Multimodal is **very** informative about task    $I_p(\{X_2, X_2\}; Y) = 1$

Synergy = Original multimodal information about the task
– multimodal information given by the ***worst*** distribution combining the same modalities

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Mathematical Framework for Multimodal Interactions

**Many ways of combining these 2 unimodal marginals into a multimodal distribution!**

Another distribution $q$
combining modalities

Unimodal marginal 1                                                                         Unimodal marginal 2

| **Verbal** | **Sarcasm** | **Verbal** | **Nonverbal** | **Sarcasm** | **Nonverbal** | **Sarcasm** |
|---|---|---|---|---|---|---|
| — | Yes/No | — | — | Yes/No | — | Yes/No |
| + | Yes/No | — | + | Yes/No | + | Yes/No |
| | | + | — | Yes/No | | |
| | | + | + | Yes/No | | |

Multimodal is **less**
informative about task     $I_q(\{X_2, X_2\}; Y) = 0$

Synergy = Original multimodal information about the task
– multimodal information given by the **_worst_** distribution combining modalities = 1 − 0 = **1**

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Mathematical Framework for Interactions

**More formally as partial information decomposition:** [Bertschinger et al., 2014]



$q$ must be a ***coupling*** of the unimodal marginals:

$$\Delta_p = \{q(x_1, x_2, y) : q(x_1, y) = p(x_1, y), q(x_2, y) = p(x_2, y)\}$$

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

Task-relevant
multimodal info

Task-relevant multimodal
info without synergy:

$$S_{q^*} = I_{q^*}(X_1, X_2; Y) - \min_{q \in \Delta_p} I_q(X_1, X_2; Y) = 0$$

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Estimating Partial Information Decomposition

Equivalent formulation as max-entropy optimization:

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

$$\Delta_p = \{q(x_1, x_2, y) : q(x_1, y) = p(x_1, y), q(x_2, y) = p(x_2, y)\}$$

If $X_1, X_2, Y$ have small discrete support: exact solution via convex programming.

**Clustering**

**Convex programming with linear constraints**



$x_1$

Text

$x_2$

Image

$|\mathcal{X}_2|$

$|\mathcal{X}_1|$

$|\mathcal{Y}|$

$$Q^* = \arg \max_Q H_Q(Y|X_1, X_2)$$

CVXPY

$R, U_1, U_2, S$

s.t. $\sum_{x_2} Q = p(x_1, y), \sum_{x_1} Q = p(x_2, y),$

$Q \geq 0, \sum_{x_1, x_2, y} Q = 1.$

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Estimating Partial Information Decomposition

Equivalent formulation as
max-entropy optimization:

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

$$\Delta_p = \{q(x_1, x_2, y) : q(x_1, y) = p(x_1, y), q(x_2, y) = p(x_2, y)\}$$

If $X_1, X_2, Y$ high-dimensional & continuous: an approximate neural network estimator.



[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Quantifying Multimodal Datasets

**1. Dataset quantification:**

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



Language: *And he I don't think he got mad when hah I don't know maybe.*
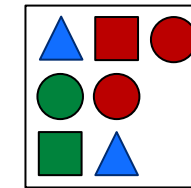
Vision:

Acoustic:

Gaze aversion

(frustrated voice)

Sentiment

**Sheldon :**

Its just a *privilege* to watch your mind at work.

· **Text** : suggests a compliment.
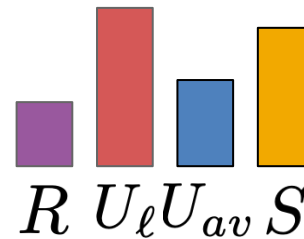· **Audio** : neutral tone.
· **Video** : straight face.
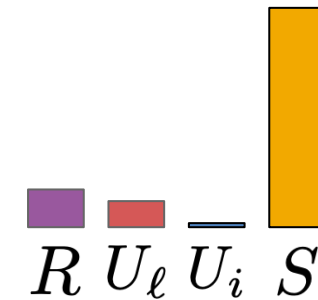
Sarcasm

*Is there a red shape above a circle?*

VQA

$R\ U_\ell U_{av}\ S$

$R\ U_\ell U_{av}\ S$

$R\ U_\ell\ U_i\ S$

Also matches human judgment of interactions, and other sanity checks on synthetic datasets

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Model Selection

**1. Dataset quantification:**

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

**Interaction polytope**

$(1, 0, 0, 0)$

$(0, 1, 0, 0)$

$(0.1, 0.7, 0.2, 0.3)$

$(0, 0, 1, 0)$

$(0, 0, 0, 1)$

Can be done with synthetic data

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Model Selection

**2. Model quantification:**

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

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

**Interaction polytope**



$(1, 0, 0, 0)$

$(0, 1, 0, 0)$

$(0.1, 0.7, 0.2, 0.3)$

$(0, 0, 1, 0)$

$(0, 0, 0, 1)$

Model families trained
on synthetic data

- Unimodal models
- Ensemble
- Multiplicative interactions
- and many more…

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# 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**



**Interaction polytope**

$(1, 0, 0, 0)$

$(0, 1, 0, 0)$

$(0.1, 0.7, 0.2, 0.3)$

$(0, 0, 1, 0)$

$(0, 0, 0, 1)$

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Model Selection

**3. Model selection:**

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



**Language:** *And he I don't think he got mad when hah I don't know maybe.*

**Vision:**

**Acoustic:** (frustrated voice)

Sentiment

$R \; U_\ell U_{av} \; S$

Language/Agreement

**Sheldon:**
Its just a *privilege* to watch your mind at work.

- **Text**: suggests a compliment.
- **Audio**: neutral tone.
- **Video**: straight face.

Sarcasm

$R \; U_\ell U_{av} \; S$

Multimodal Transformer

Is there a red shape above a circle?

VQA

$R \; U_\ell \; U_i \; S$

Multiplicative/Transformer

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Application 1: Mental Health

**Daily mood prediction as a stepping-stone towards real-time assessment of suicide ideation.**

Text + app + keystroke interactions

| Slower implies positive | Faster implies positive |
|---|---|
| just | why, thank, haha |
| next, was, into, people | making, work, idk |
| stuff, cute, phone, want, talk, see | they, send, dont, man, going |
| don't, talk | think, you, all, love |

+ words like 'love', 'thanks', 'haha' become more positive when typed faster

- words like 'don't, 'just' become more negative when typed faster

$R \; U_\ell U_{ak} \; S$

[Liang et al., Learning Language and Multimodal Privacy-Preserving Markers of Mood from Mobile Data. ACL 2021]

# Application 2: Computational Pathology



Histology images

Genomics profile

RNA
TP53 (RNA)
SAMD9 (RNA)
MYC (RNA)

CNV
EGFR (CNV)
CDKN2A (CNV)

MUT
IDH1 (MUT)
FLG (MUT)

$R\ U_h\ U_g\ S$

Glioma: Genomics unimodal

$R\ U_h\ U_g\ S$

Pancreas: Histology + genomics interaction

Understanding the models and adoption in practice by doctors

[Liang et al., Quantifying & Modeling Multimodal Interactions: An Information Decomposition Framework. NeurIPS 2023]

# Implications of Studying Multimodal Interactions

$X_1$

Uniqueness

Redundancy

Synergy

Uniqueness

$Y$

$X_2$

Optimizing these interactions as training objectives:

*CLIP/contrastive learning*

$X_1$

$Y$

$X_2$

Visualizing the interactions learned in individual neurons:

Lip

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

Predicting multimodal performance to decide modality utility:

$p(x_1, y)$   $f_1 : \blacktriangle \longrightarrow y_1$

$p(x_2, y)$   $f_2 : \bullet \longrightarrow y_2$

$p(x_1, x_2, y)$   ***80% accuracy***

[Liang et al., FactorCL. NeurIPS 2023]     [Liang et al., MultiViz. ICLR 2023]     [Liang et al., Semi-supervised. arXiv 2023]

# 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 $f(x_A, x_B) = f_A(x_A) + f_B(x_B)$.

Modality A $\quad x_A$

Modality B $\quad x_B$

Fusion + prediction $\quad \widehat{y}$

$$\tilde{f}(x_A, x_B) = \underbrace{\underset{x_B}{\mathbb{E}}[f(x_A, x_B)]}_{f_A(x_A)} + \underbrace{\underset{x_A}{\mathbb{E}}[f(x_A, x_B)]}_{f_B(x_B)} - \underbrace{\underset{x_A, x_B}{\mathbb{E}}[f(x_A, x_B)]}_{\mu}$$

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.

$\mu$ 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 $f(x_A, x_B) = f_A(x_A) + f_B(x_B)$.

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

Natural second-order extension of gradient-based approaches!

[Liang et al., MultiViz: Visualizing and Understanding Multimodal Models. ICLR 2023]

# Quantifying Cross-modal Interactions

**Identifying individual cross-modal interactions**

| CLEVR | VQA 2.0 | Flickr-30k | CMU-MOSEI |
|:-----:|:-------:|:----------:|:---------:|



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

*How many **birds**?*

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

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

Correspondence

Relationships

[Liang et al., MultiViz: Visualizing and Understanding Multimodal Models. ICLR 2023]

# Quantifying Cross-modal Interactions

**Classification of cross-modal interactions**

Modality A $\quad$ $x_A$

Fusion + prediction $\rightarrow$ $\widehat{y}$

**Unimodal importance**

(e.g., GradCAM, LIME, SHAP)

Modality B $\quad$ $x_B$

$I_A$

$I_B$

**Analysis**

$\|I_A\| \gg \|I_B\|$ Dominance

$I_A \cdot I_B > 0$ Complementary

$I_A \cdot I_B < 0$ Conflict

(A) sentiment — L, V, A — dominant, complement, conflict

Language is often **dominant** in multimodal sentiment analysis

[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.

[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
in a single instance.

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

# 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 Quantification

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

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



Quantification output

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

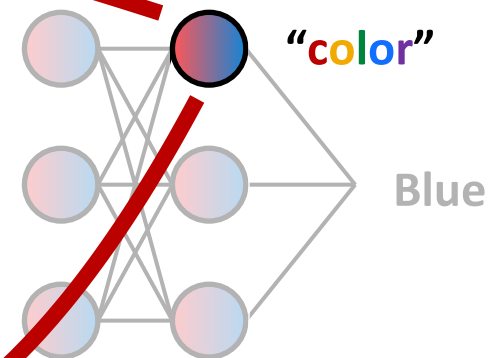# Directly Evaluating Quantification

**Direct evaluation: Create datasets for each tested quality, but limited to synthetic data**

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



Redundancy      Unique 1      Unique 2      Synergy

# Directly Evaluating Quantification

**Direct evaluation: Create datasets for each tested quality, but limited to synthetic data**

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

<span style="color:red">Can be done with synthetic data</span>

**Interaction polytope**



$(1,0,0,0)$     $(0,1,0,0)$

$(0.1, 0.7, 0.2, 0.3)$

$(0,0,1,0)$     $(0,0,0,1)$

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**Indirect evaluation**

*Find some downstream quality that practitioners find useful and can be easily evaluated.*



[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**Indirect evaluation: Model simulation**



1. Model simulation
*Can humans reproduce model predictions with high accuracy and agreement?*

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**Indirect evaluation: Model simulation**



55.0%   65.0%   61.7%   71.7%   **81.7%**

U   U + C   U + C + Local R   U + C + Local R + Global R   U + C + Local R + Global R + P

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

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**Indirect evaluation: Model error analysis and debugging**



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

*Fix bugs*

*Find bugs*

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**Indirect evaluation: Model error analysis and debugging**



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

Local analysis

3. Multimodal representations

"**color**"

Blue

*What color is the Salisbury Rd sign?*   *What color is the building?*   *What color are the checkers on the wall?*

Global analysis

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

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**Indirect evaluation: Model error analysis and debugging**

*"Models pick up cross-modal interactions but fail in identifying color!"* → *Add targeted examples involving color.*

+1.4%    +0.2%    **+30.5%**

Random    Uncertainty    MultiViz

*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: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Indirectly Evaluating Quantification

**More ways for indirect evaluation:**
- Model selection: given fixed budget, try randomly or try models in order based on what quantification tells me.
- Data/modality selection: given fixed budged, collect random data or collect based on what quantification tells me.
- If quantification gives theoretical result, check how well the theory matches experiments.

[Liang et al., MultiViz: Towards Visualizing and Understanding Multimodal Models. ICLR 2023]

# Challenges: Quantifying Multimodal Interactions

**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** (shared)

unique
shared
unique

$z$

**Nonredundancy** (unique)

Noninteracting

Additive

Noninteracting (union)

Asymmetric

Contextualized (transference)

Non-additive (nonlinear)

signal    response

a+b → ☐    Equivalence

a+b → ☐    Enhancement

a+b → ☐ and ◯    Independence

a+b → ☐    Dominance

a+b → ☐ (or ☐)    Modulation

a+b → △    Emergence

**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

zebras ⟶ **prediction**

**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?

[Agarwal et al., Towards Causal VQA: Revealing & Reducing Spurious Correlations by Invariant & Covariant Semantic Editing. CVPR 2020]

# 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?**

[Wang et al., What Makes Training Multi-modal Classification Networks Hard? CVPR 2020]

# 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

[Wang et al., What Makes Training Multi-modal Classification Networks Hard? CVPR 2020]

# 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

[Wang et al., What Makes Training Multi-modal Classification Networks Hard? CVPR 2020]

# Challenges

**Open challenges:**
- Learning, generalization, and optimization in high-dimensional settings (p >> 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**

**Heterogeneity**

**Interactions**

**Learning**

**Open challenges**

# 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 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.

Modality A ▲

Modality B ●

**Homogeneous Modalities**
(with similar qualities)

**Heterogeneous Modalities**
(with diverse qualities)

**Examples:**

Images from 2 cameras

Text from 2 different languages

Language and vision

**Abstract modalities are more likely to be homogeneous**

**Connected Modalities**

Connected: Shared information that relates modalities

Modality A ▲
Modality B ●

unique
unique

stronger
weaker
unconnected

**Statistical**

Association — Dependency

e.g., correlation, co-occurrence

e.g., causal, temporal

**Semantic**

Correspondence — Relationship

laptop

e.g., grounding

used for

e.g., function

Interacting: process affecting each modality, creating new response



Modality A ▲

Modality B ●

$z$

inference → $z$

*response*

Interactions happen during inference!

"Inference" examples:
- Representation fusion
- Prediction task
- Modality translation

▭▭▭▭ *representation*

$\hat{y}$ *prediction*

▭ *modality C*

# Interacting Modalities



**Redundancy** (shared)

unique

shared

*z*

unique

**Nonredundancy** (unique)

Noninteracting

Additive

Noninteracting (union)

Asymmetric

Contextualized (transference)

Non-additive (nonlinear)

signal     response

a+b → □   Equivalence

a+b → □   Enhancement

a+b → □ and ○   Independence

a+b → □   Dominance

a+b → □ (or □)   Modulation

a+b → △   Emergence

**Heterogeneous**
+
**Connected**
+
**Interacting**

$z$

**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:**

Modality A

Modality B

*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 | Coordination | Fission |
|--------|--------------|---------|



# modalities **>** # representations          # modalities **=** # representations          # 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:**

Modality A

Modality B

**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**

**Continuous Alignment**

**Contextualized Representation**



Discrete elements and connections

Segmentation and continuous warping

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**

**Translation**

**Creation**

**Information:** (content)

Reduction

Maintenance

Expansion

# Challenge 5: Transference

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



Enriched Modality A

**Transference**

A    B

**only available during training**

Modality A

Modality B

# 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

**Future Direction: High-modality**

**MultiBench**
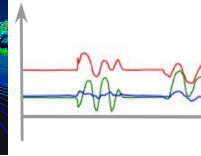https://github.com/pliang279/MultiBench

Few modalities → High-modality

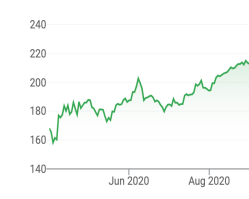Language   Vision   Audio   Graphs   Control   LIDAR   Sensors   Set   Table   Financial   Medical

**Examples:**

**Non-parallel learning**

**Limited resources**

# Future Direction: Long-term

## Short-term



seconds
or minutes

## Long-term



**Examples:**

**Compositionality**          **Memory**          **Personalization**

**Future Direction: Interaction**

Social Intelligence



Reasoning

Perception

Generation

Multimodal
Interaction

**Examples:**

**Multi-Party**          **Causality**          **Ethical**
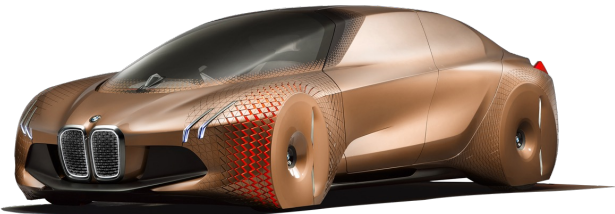
**Future Direction: Real-world**



Healthcare
Decision Support



Intelligent Interfaces and
Vehicles



Online Learning
and Education

**Examples:**

**Robustness**            **Fairness**            **Generalization**
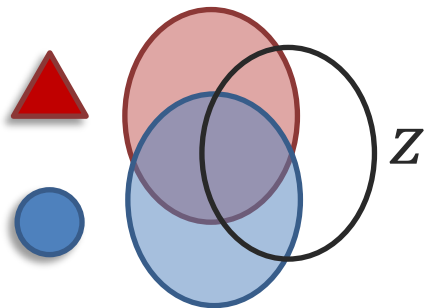
# What is Multimodal? → Why is it hard? → What is next?

| What is Multimodal? | Why is it hard? | What is next? |
|---|---|---|
| **Heterogeneous** ✚ **Connected** ✚ **Interacting** | **Representation** | **Heterogeneity** |
| | **Alignment** | **High-modality** |
| | **Reasoning** | **Long-term** |
| | **Generation** | **Interaction** |
| | **Transference** | |
| | **Quantification** | **Real-world** |

$z$