

Multimodal Deep Learning

Paul Liang

Machine Learning Department
Carnegie Mellon University

**Carnegie
Mellon
University**

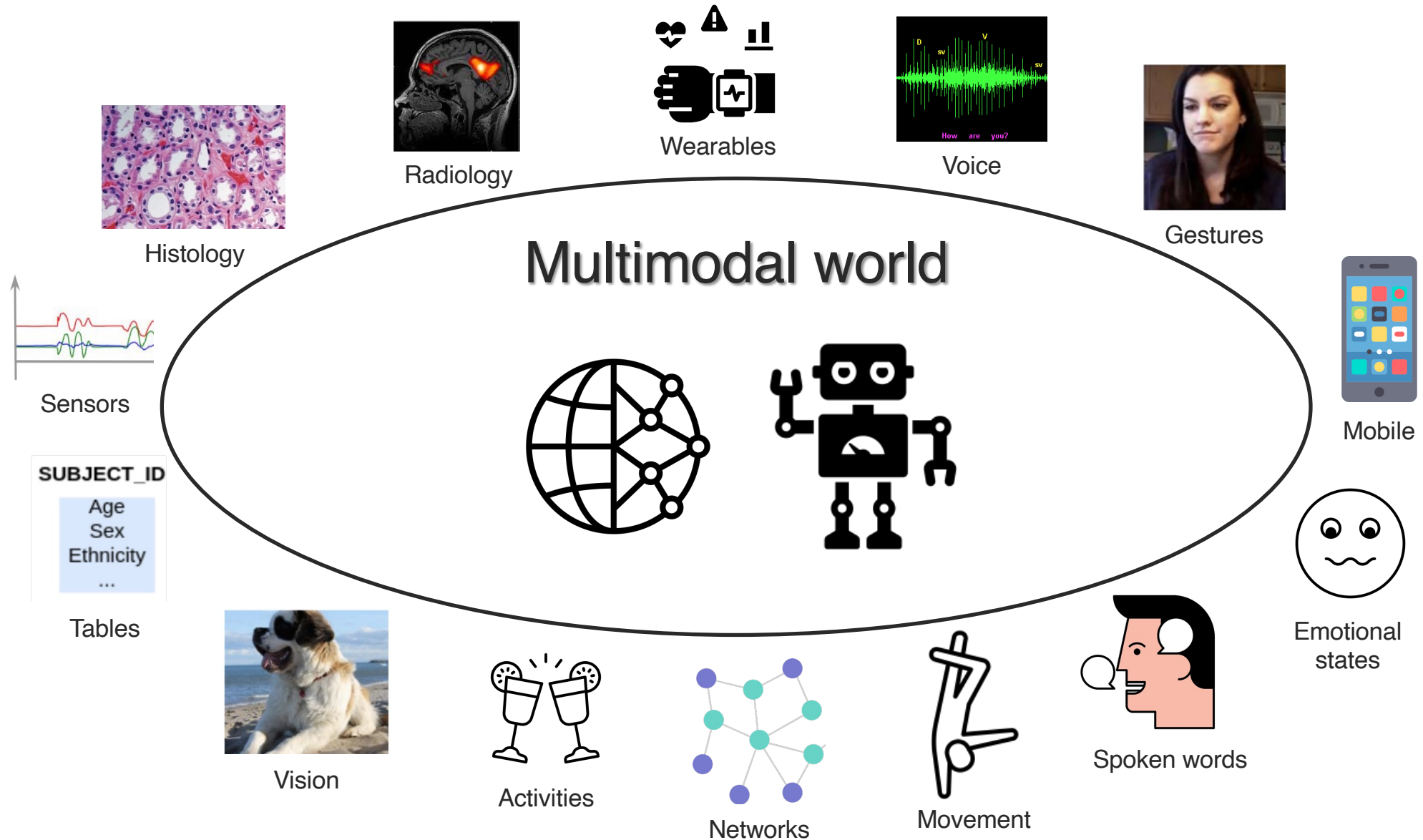
<https://www.cs.cmu.edu/~pliang/>
pliang@cs.cmu.edu

<https://github.com/pliang279>

 @pliang279



Multimodal Artificial Intelligence



Multimodal Human Communication

Understanding human language and gestures

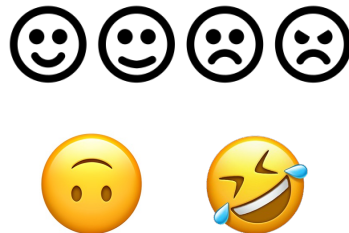
Language: *And he I don't think he got mad when hah I don't know maybe.* *Too much too fast, I mean we basically just get introduced to this character...* *All I can say is he's a pretty average guy.*

Vision: *Gaze aversion* *Uninformative* *Contradictory smile*

Acoustic: (frustrated voice) (angry voice) (disappointed voice)

Diverse annotations

- Sentiment
- Emotions
- Sarcasm
- Humor
- Personalities



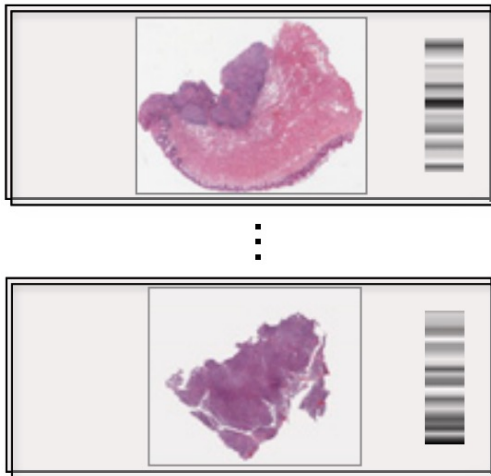
Practical applications



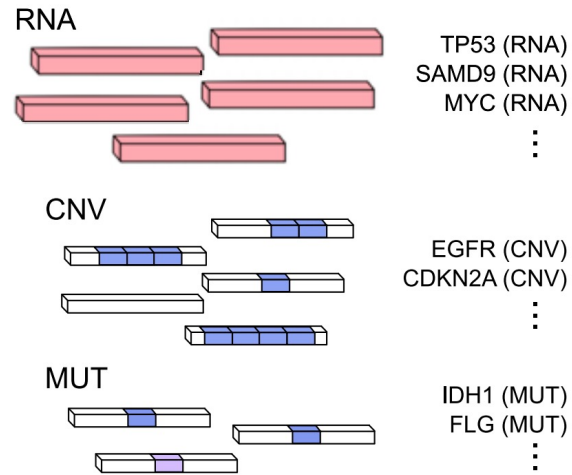
Healthcare Modalities

Medicine and healthcare

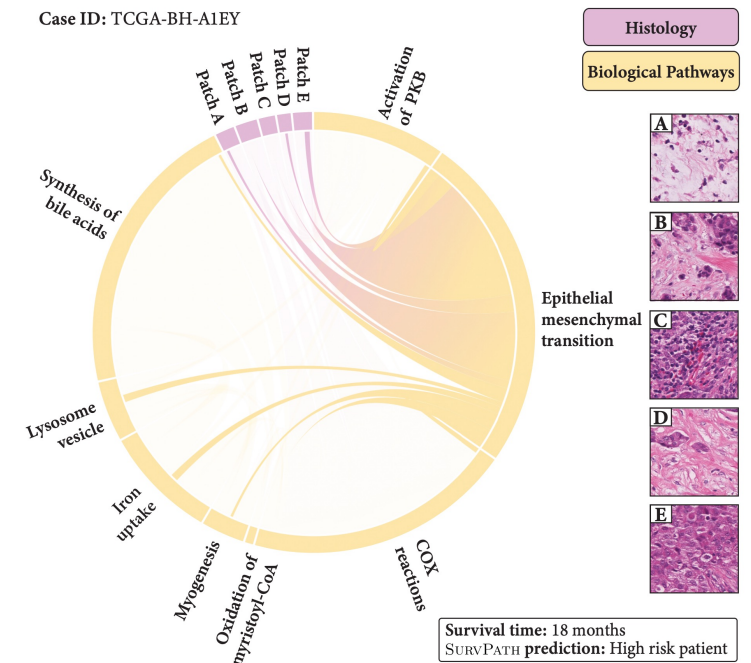
Histology images



Genomics profile



Transcriptomics

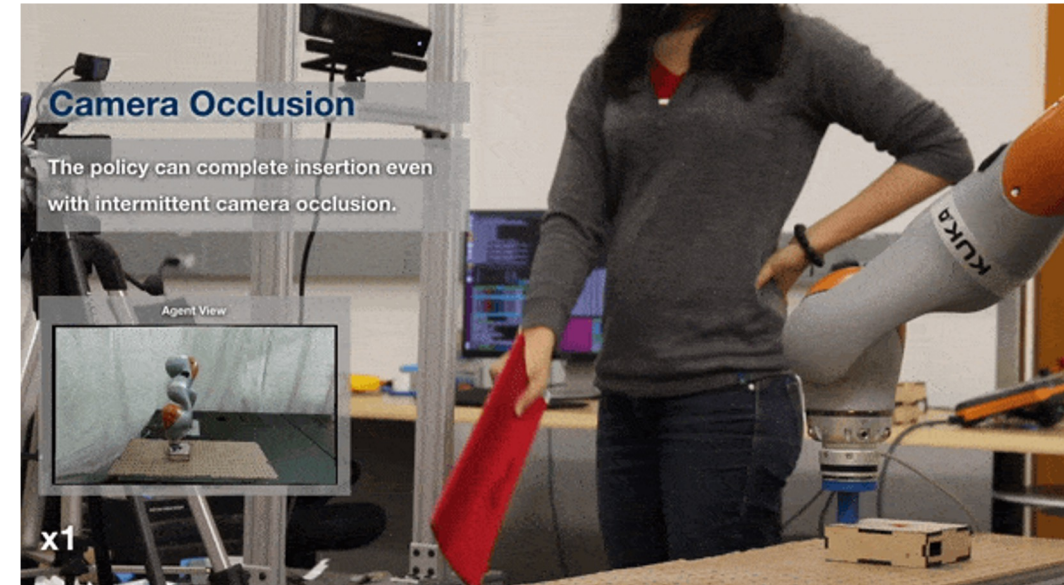
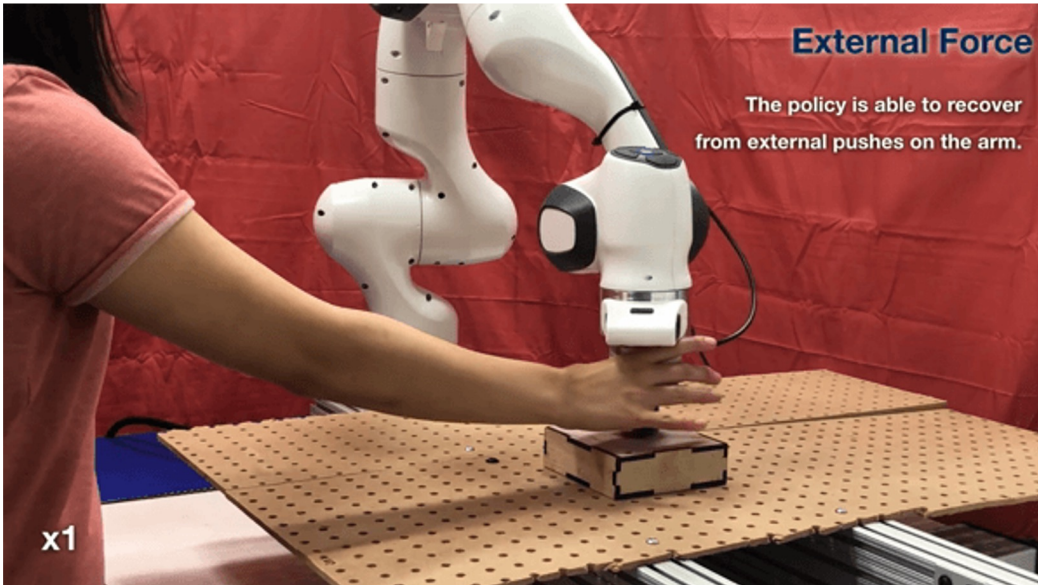


[Jaume et al., Modeling Dense Multimodal Interactions Between Biological Pathways and Histology for Survival Prediction. 2023]

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

Multisensory Robotic Intelligence

Multisensor fusion in robotics



+ robustness

Multimodal Machine Learning – Surveys, Tutorials and Courses

Foundations and Recent Trends in Multimodal Machine Learning

Paul Liang, Amir Zadeh and Louis-Philippe Morency

- ✓ 6 core challenges
- ✓ 50+ taxonomic classes
- ✓ 700+ referenced papers

<https://arxiv.org/abs/2209.03430>

Tutorials: ICML 2023, CVPR 2022, NAACL 2022

Graduate-level courses:

Multimodal Machine Learning (12th edition)

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

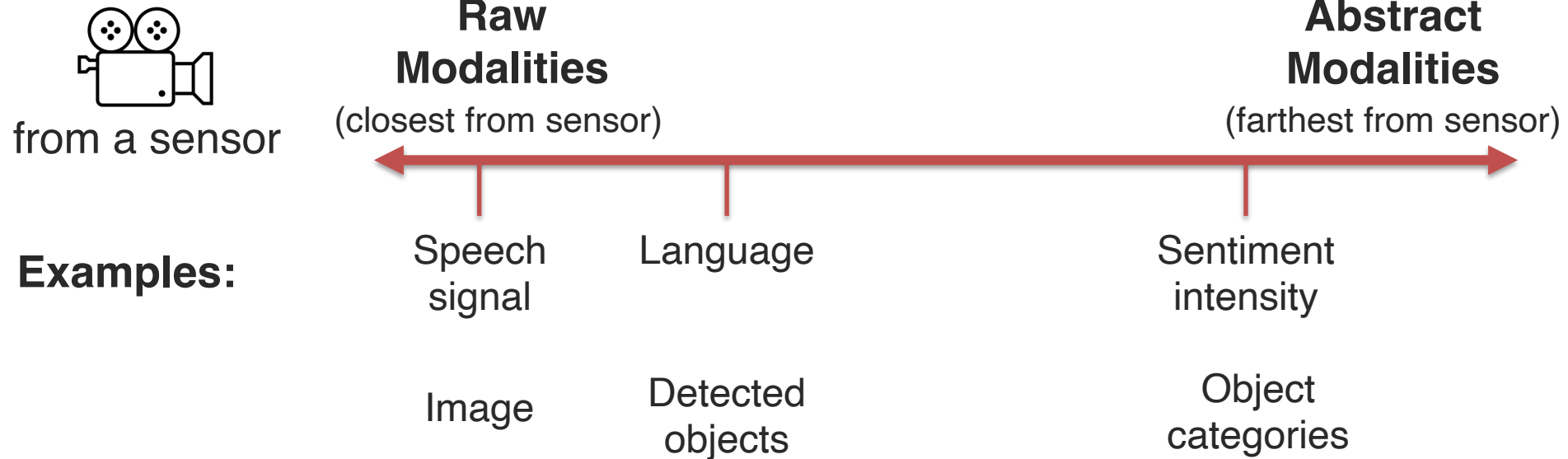
Advanced Topics in Multimodal ML

<https://cmu-multicomp-lab.github.io/adv-mmml-course/spring2023/>

What is a Modality?

Definition

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



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

What is Multimodal?

A dictionary definition...

Multimodal: with multiple modalities

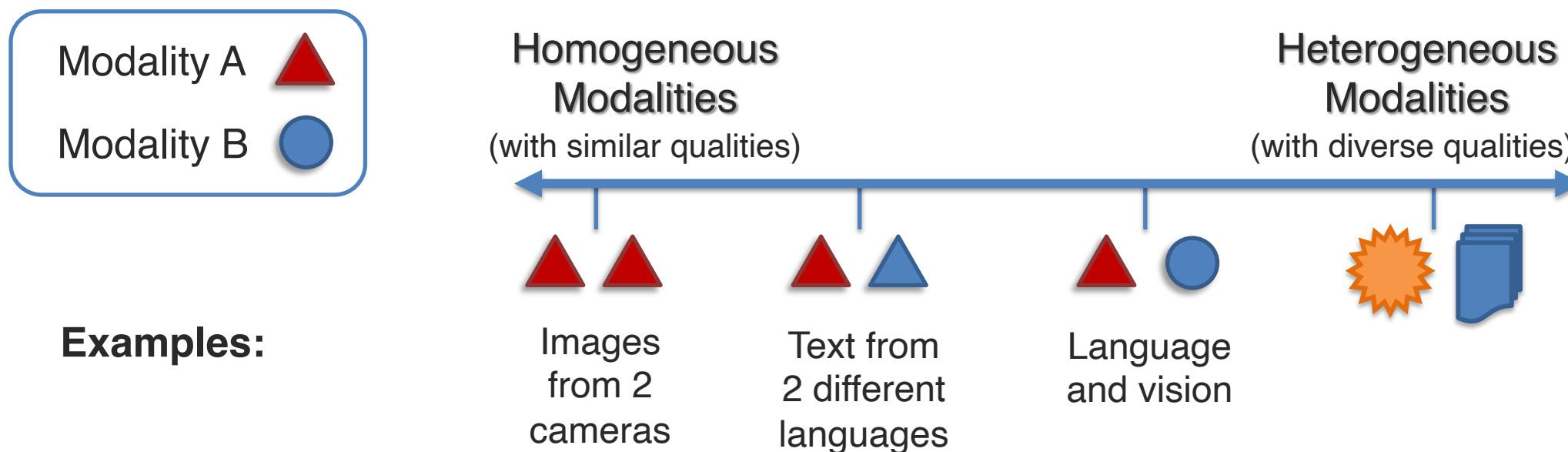
A research-oriented definition...

***Multimodal* is the scientific study of**
heterogeneous and interconnected data
Connected + Interacting

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

Heterogeneous Modalities

Heterogeneous: Diverse qualities, structures and representations.



Abstract modalities are more likely to be homogeneous

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

Dimensions of Heterogeneity

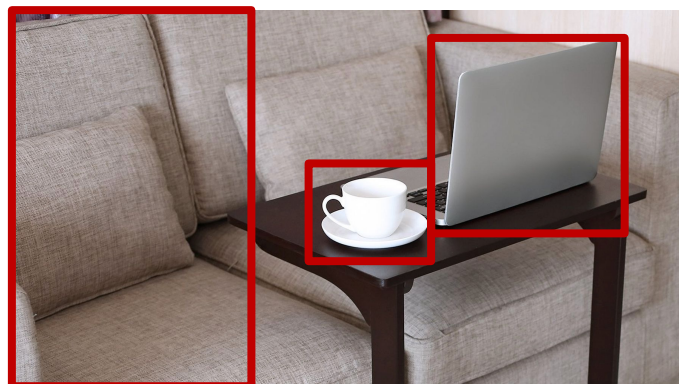
Information present in different modalities will often show diverse qualities, structures, and representations.



*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

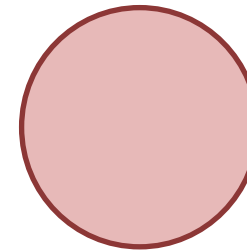
Information present in different modalities will often show diverse qualities, structures, and representations.



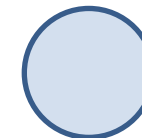
*A teacup on the right of a laptop
in a clean room.*

3 **Information:** entropy and density

$H(\blacktriangle)$

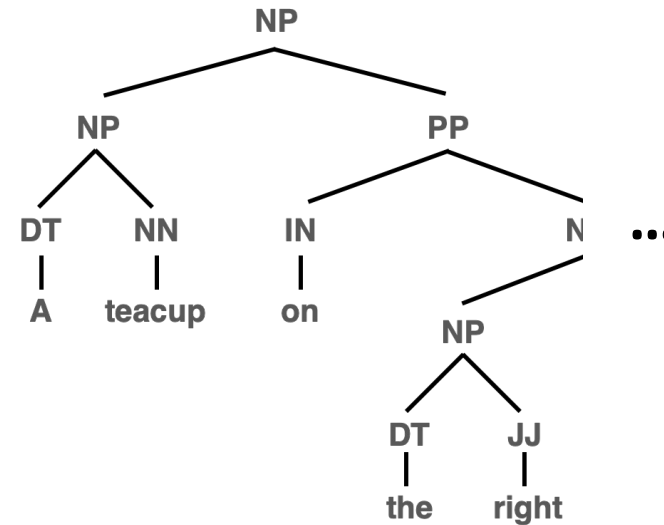


$H(\bullet)$

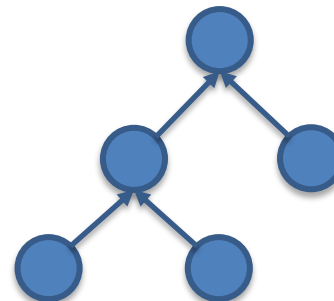
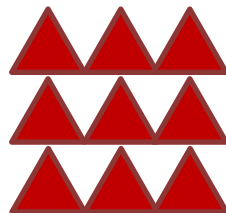


Dimensions of Heterogeneity

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



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



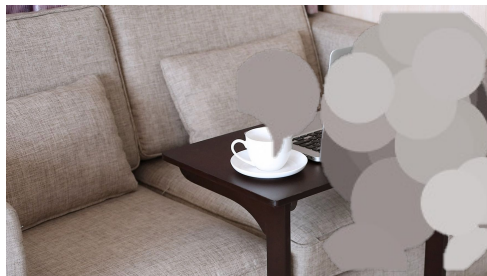
Dimensions of Heterogeneity

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

Connected Modalities

Connected: Shared information that relates modalities



Connected Modalities

Connected: Shared information that relates modalities



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



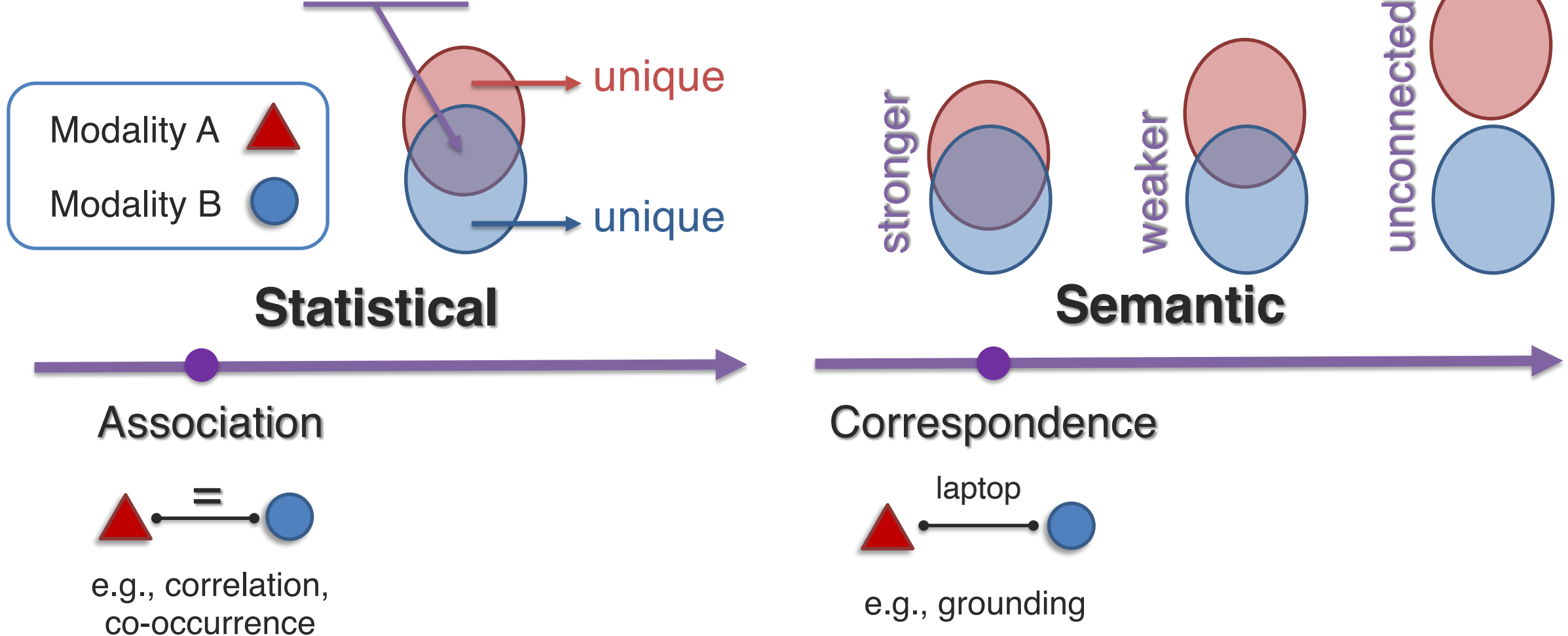
teacup



laptop

Connected Modalities

Connected: Shared information that relates modalities



Connected Modalities

Connected: Shared information that relates modalities



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



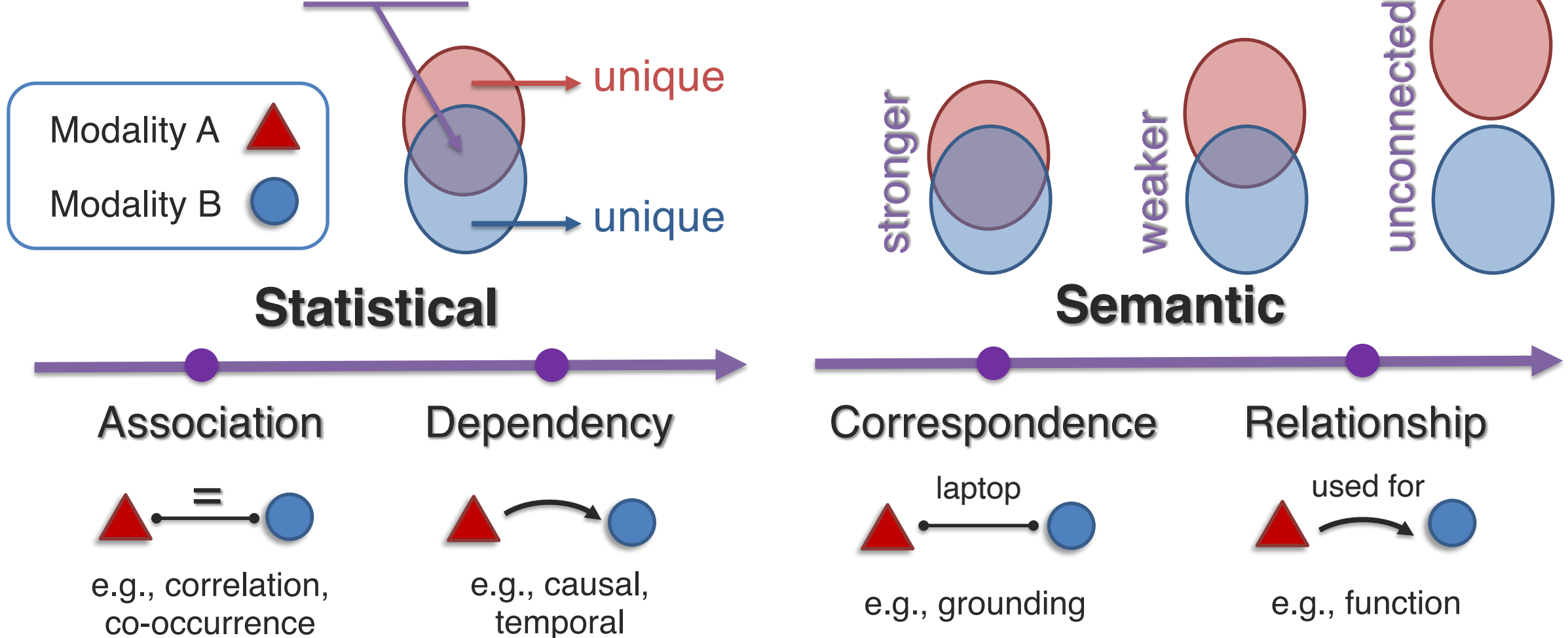
clean



room

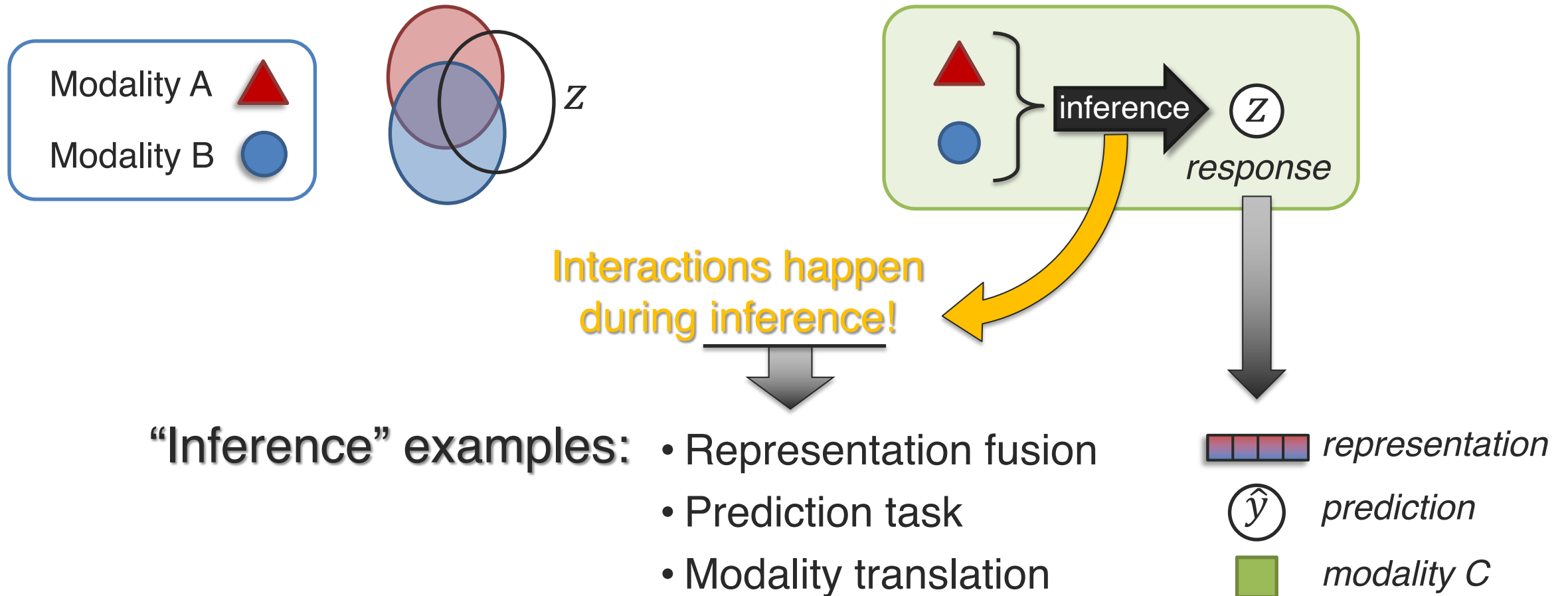
Connected Modalities

Connected: Shared information that relates modalities



Interacting Modalities

Interacting: process affecting each modality, creating new response



Interacting Modalities

Interactions: How multimodal information changes when modalities are combined for a response.

Is this indoors?



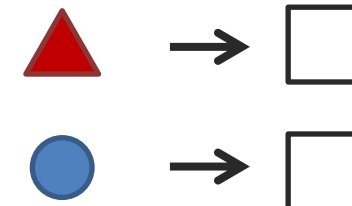
Yes!

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



Yes!

Redundant



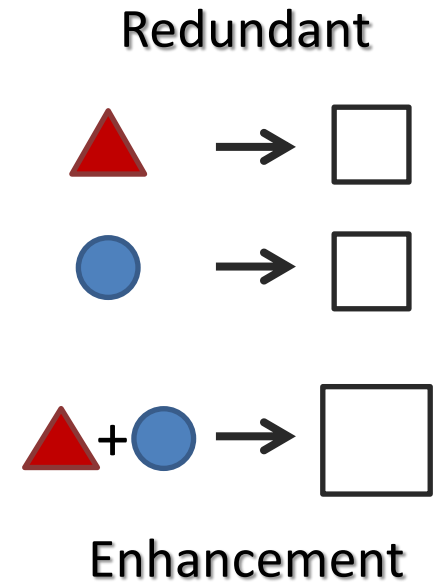
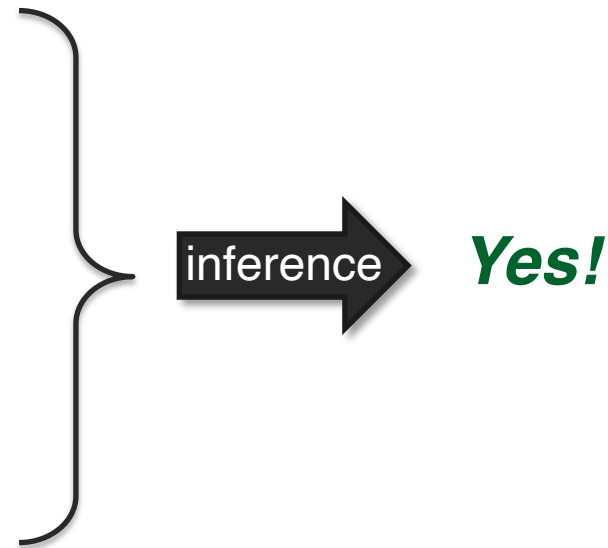
Interacting Modalities

Interactions: How multimodal information changes when modalities are combined for a response.

Is this indoors?



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



Interacting Modalities

Interactions: How multimodal information changes when modalities are combined for a response.

Is this a living room?



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

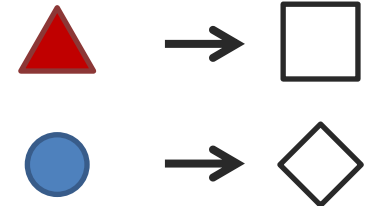


Yes!



No, probably study room.

Non-redundant



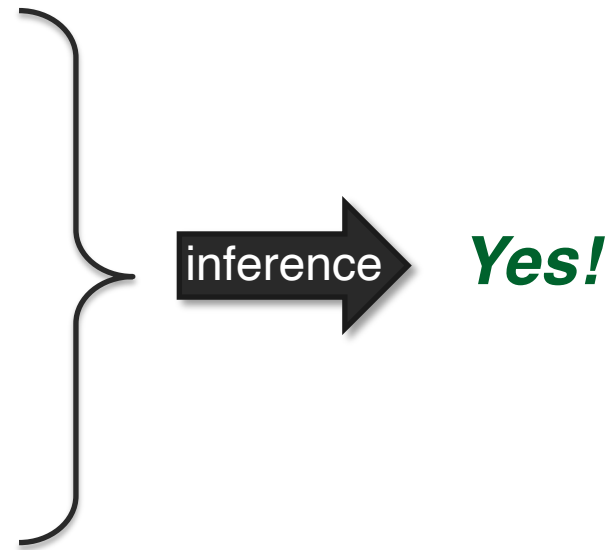
Interacting Modalities

Interactions: How multimodal information changes when modalities are combined for a response.

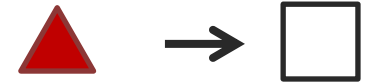
Is this a living room?



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



Non-redundant



Dominance

Interacting Modalities

Interactions: How multimodal information changes when modalities are combined for a response.

Should I work here?



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

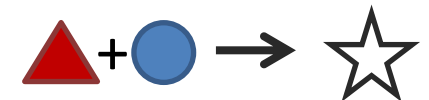


Maybe? Comfy sofa but table's too small.



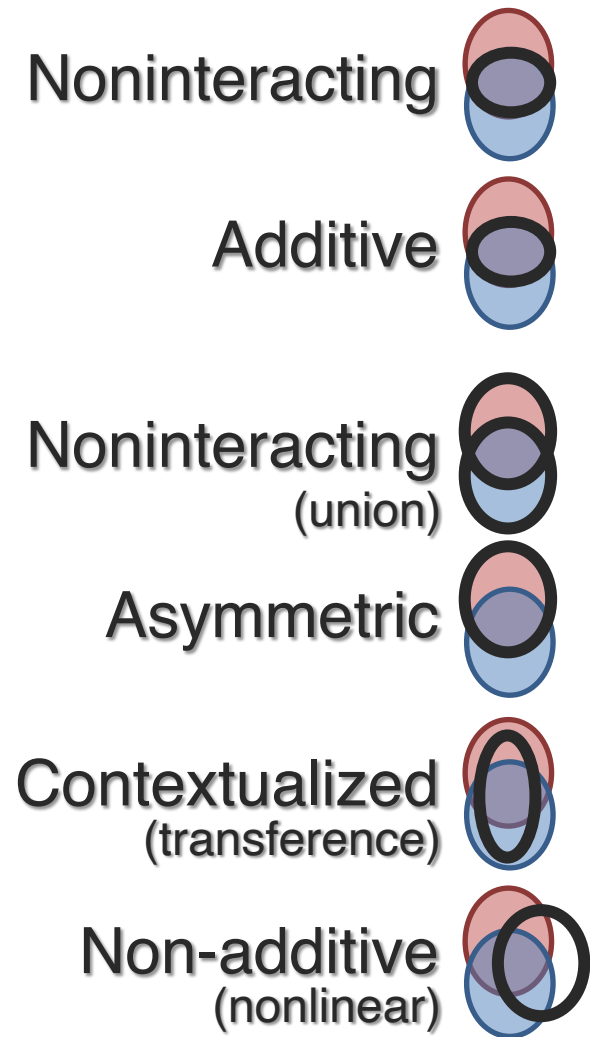
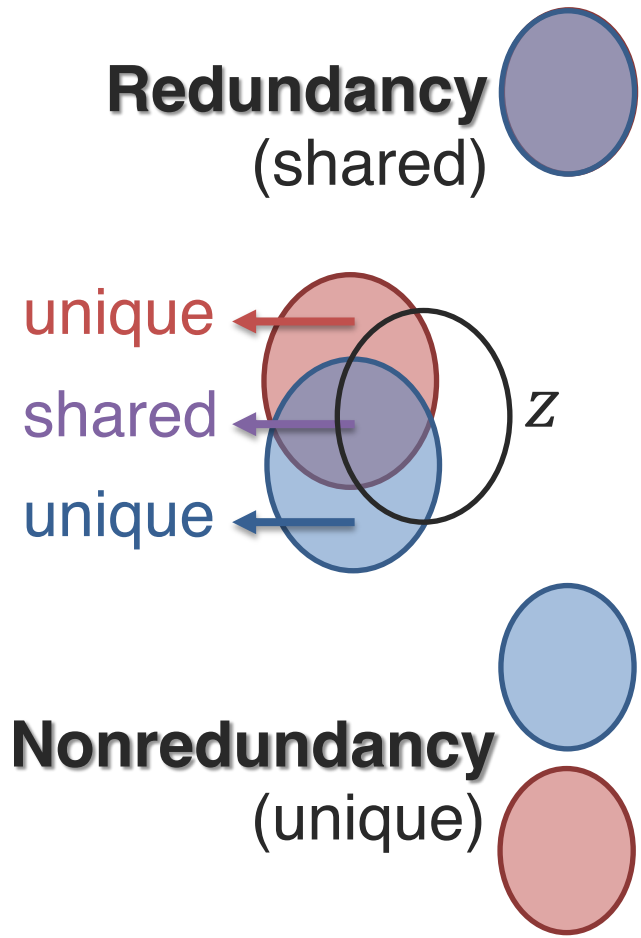
Maybe? Clean and there's tea.

Non-redundant



Emergence

Cross-modal Interaction Mechanics



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

*What is
Multimodal?*



Why is it hard?



What is next?

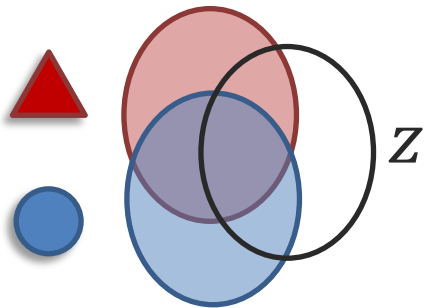
Heterogeneous



Connected

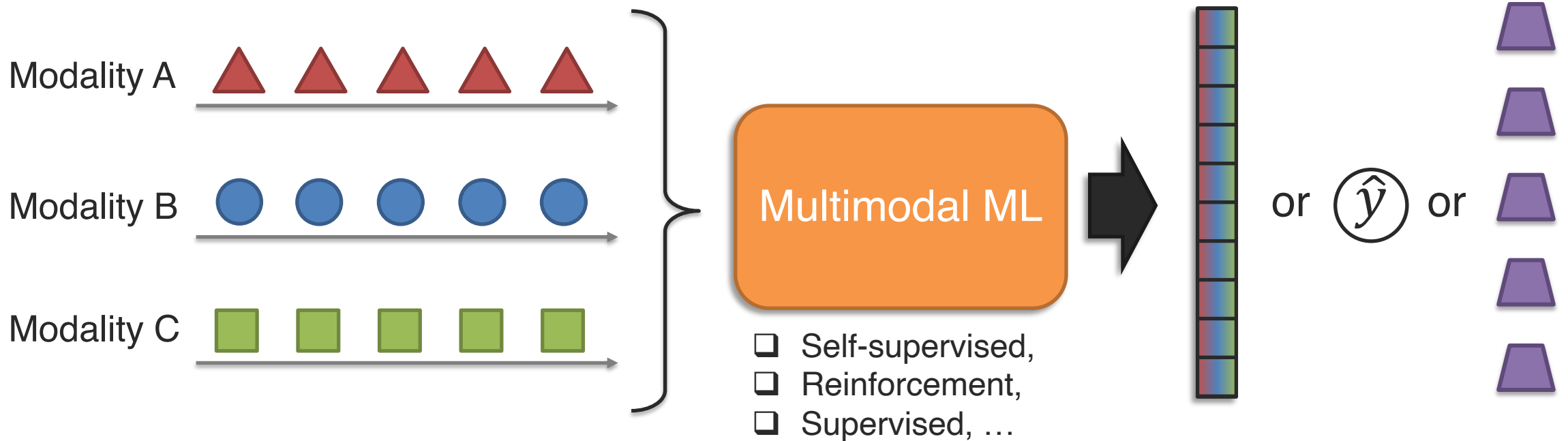


Interacting



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

Multimodal Machine Learning



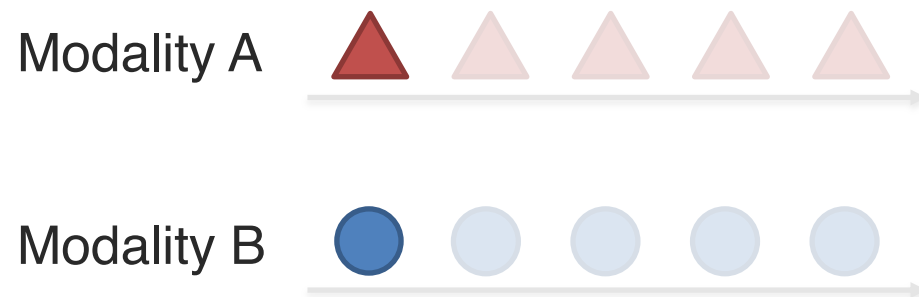
*What are the **core multimodal technical challenges**, understudied in conventional machine learning?*

Challenge 1: Representation

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

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

Individual elements:



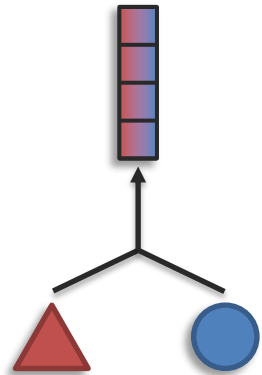
*It can be seen as a “local” representation
or
representation using holistic features*

Challenge 1: Representation

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

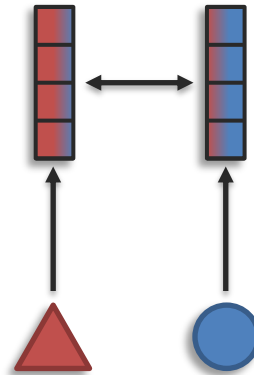
Sub-challenges:

Fusion



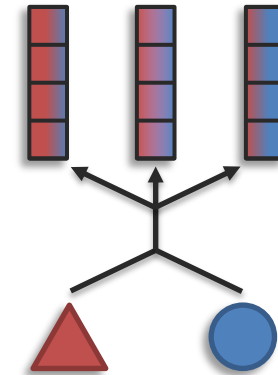
modalities $>$ # representations

Coordination



modalities = # representations

Fission



modalities $<$ # representations

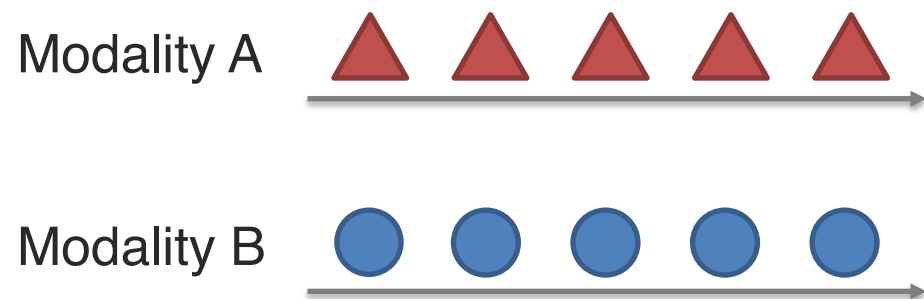
[Liang, Zadeh, Morency, Foundations and Trends in Multimodal Machine Learning. Tutorials at ICML 2023, CVPR 2022, NAACL 2022]

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

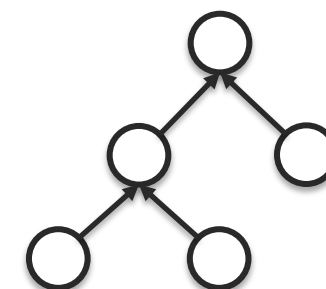
Elements with temporal structure:



Other structured examples:



Spatial



Hierarchical

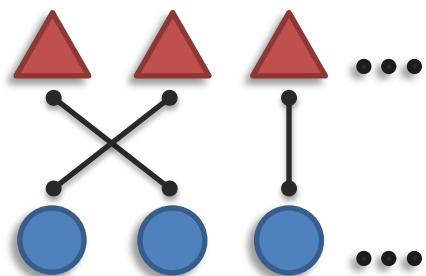
[Liang, Zadeh, Morency, Foundations and Trends in Multimodal Machine Learning. Tutorials at ICML 2023, CVPR 2022, NAACL 2022]

Challenge 2: Alignment

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

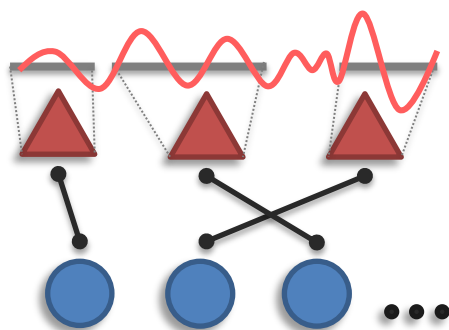
Sub-challenges:

Discrete Alignment



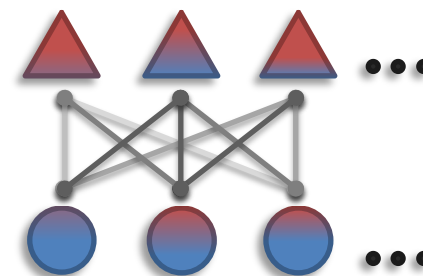
Discrete elements
and connections

Continuous Alignment



Segmentation and
continuous warping

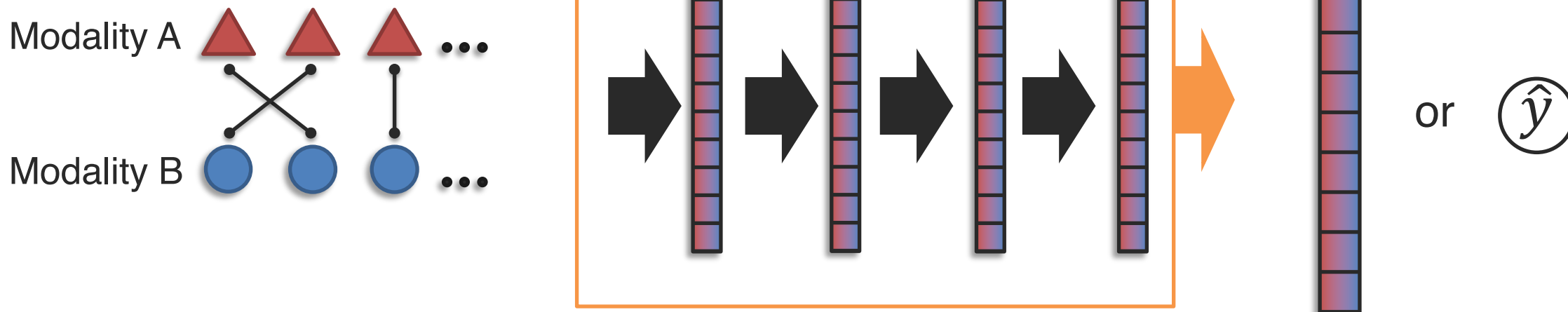
Contextualized Representation



Alignment + representation

Challenge 3: Reasoning

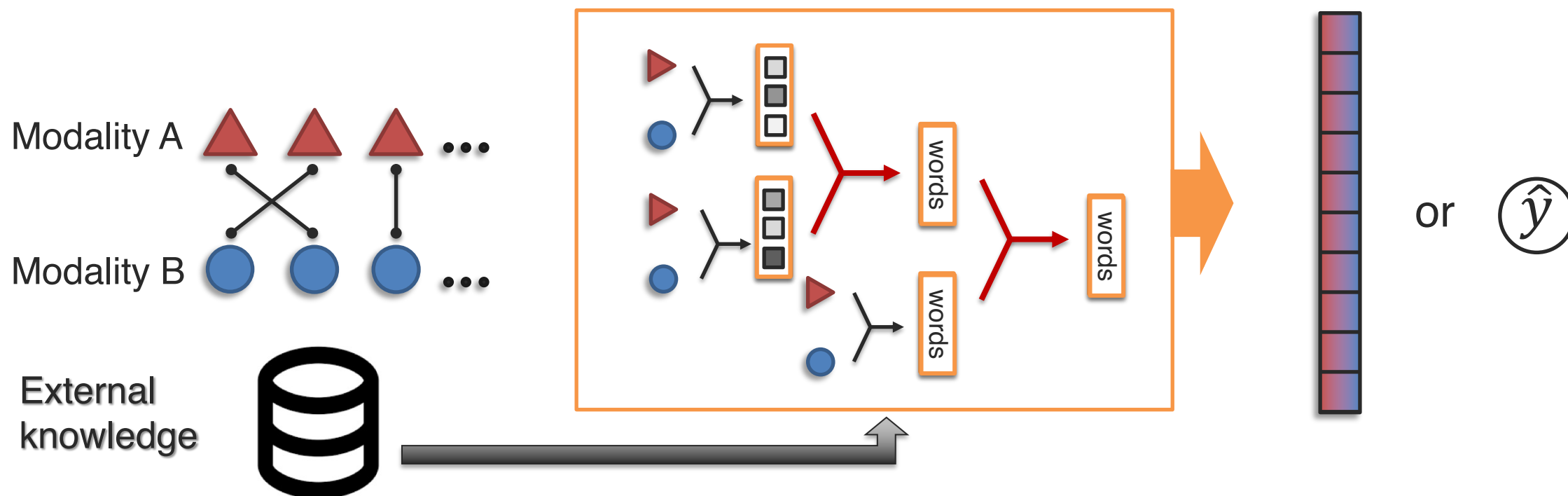
Definition: Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure



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

Challenge 3: Reasoning

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

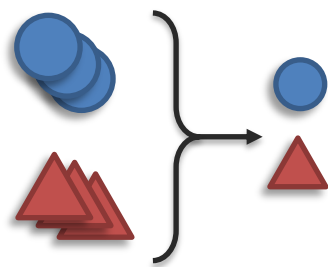


Challenge 4: Generation

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

Sub-challenges:

Summarization



Reduction



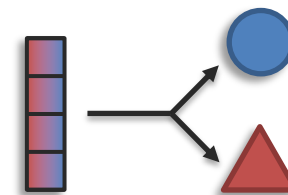
Translation



Maintenance



Creation



Expansion



Information:
(content)

Challenge 4: Generation

An astronaut riding a horse in the style of Andy Warhol.

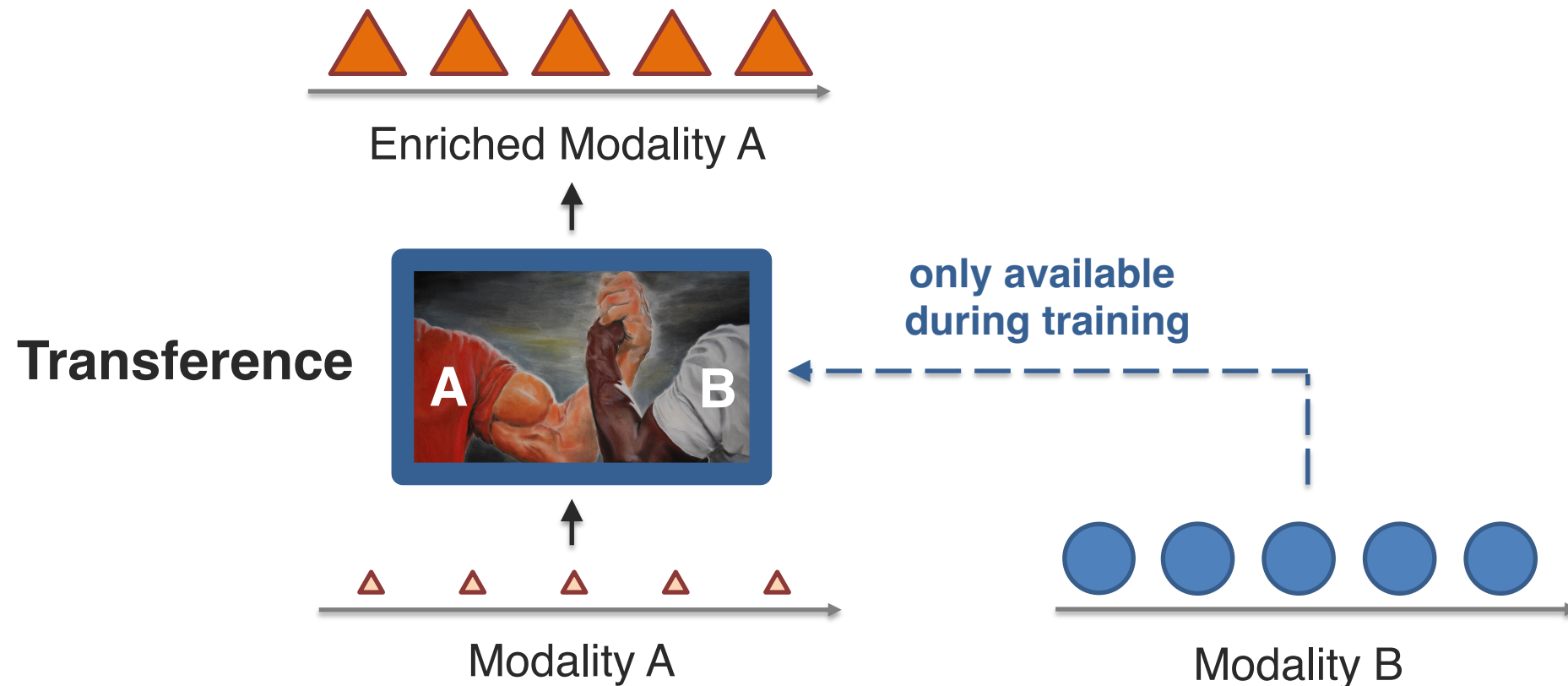


A bowl of soup that is a portal to another dimension as digital art.



Challenge 5: Transference

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

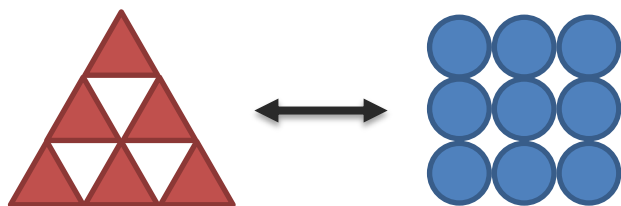


Challenge 6: Quantification

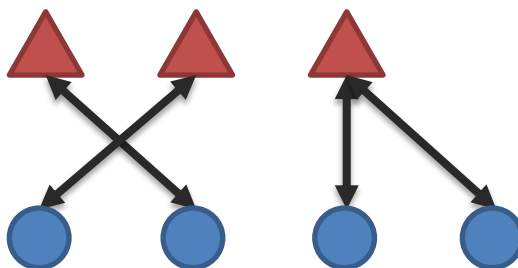
Definition: Empirical and theoretical study to better understand heterogeneity, cross-modal interactions and the multimodal learning process

Sub-challenges:

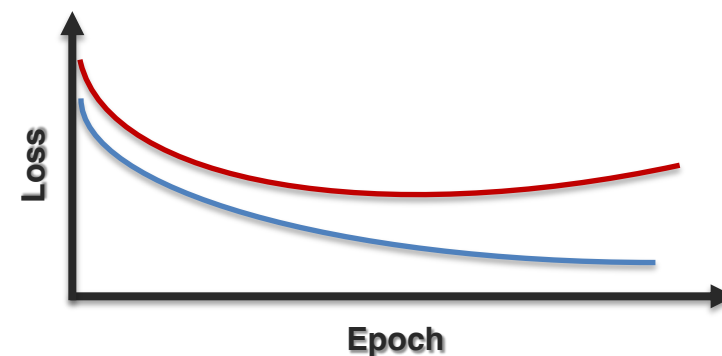
Heterogeneity



Connections & Interactions

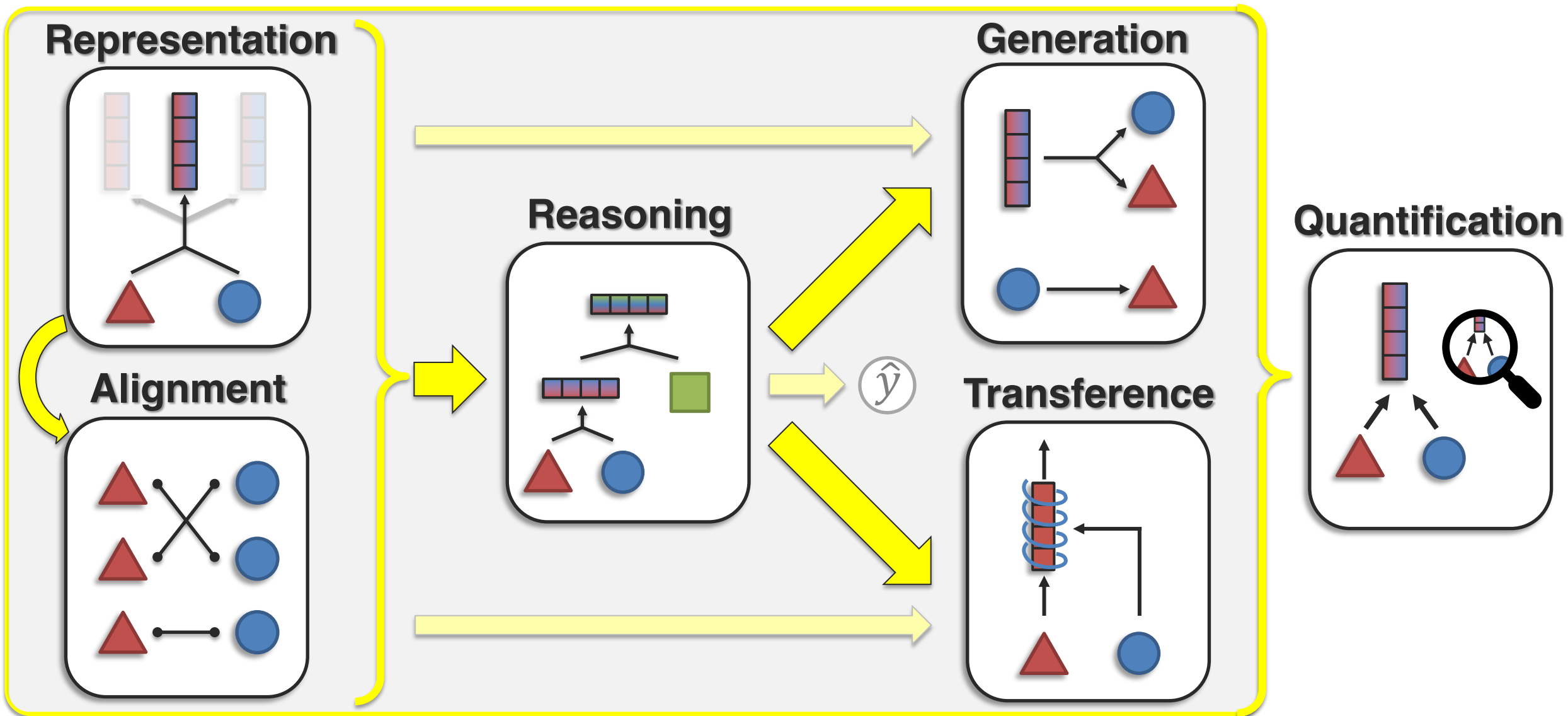


Learning

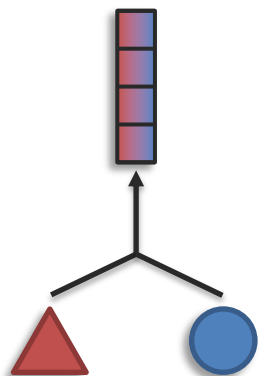


[Liang, Zadeh, Morency, Foundations and Trends in Multimodal Machine Learning. Tutorials at ICML 2023, CVPR 2022, NAACL 2022]

Core Multimodal Challenges

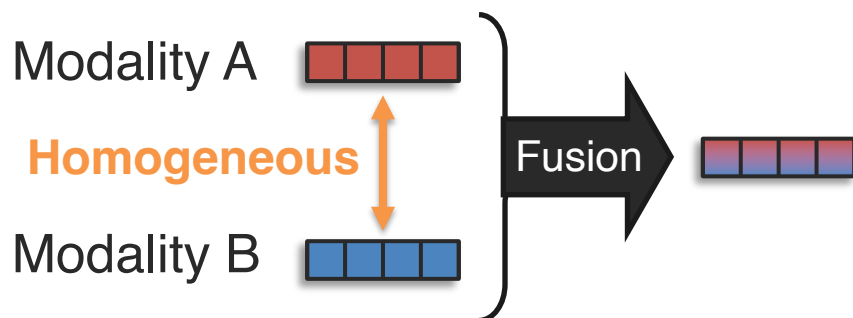


Sub-Challenge: Representation Fusion

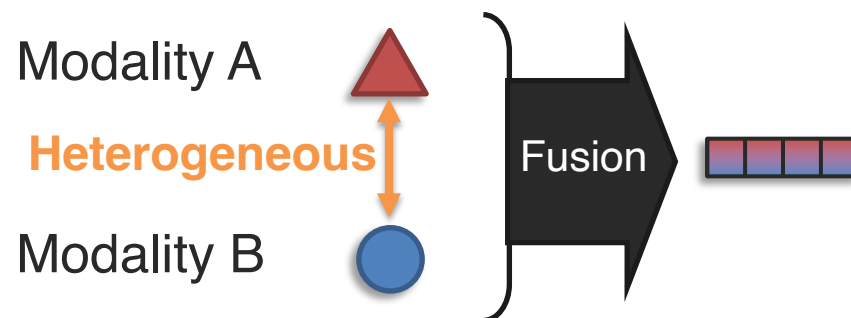


Definition: Learn a joint representation that models cross-modal interactions between individual elements of different modalities

Basic fusion:

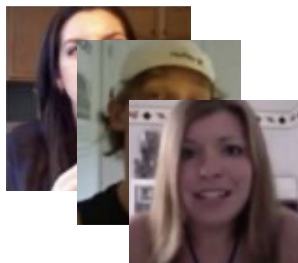


Raw-modality fusion:



From Additive to Multiplicative

300 book reviews



y : audience score

x_A : percentage of smiling

x_B : professional status
(0=non-critic, 1=critic)

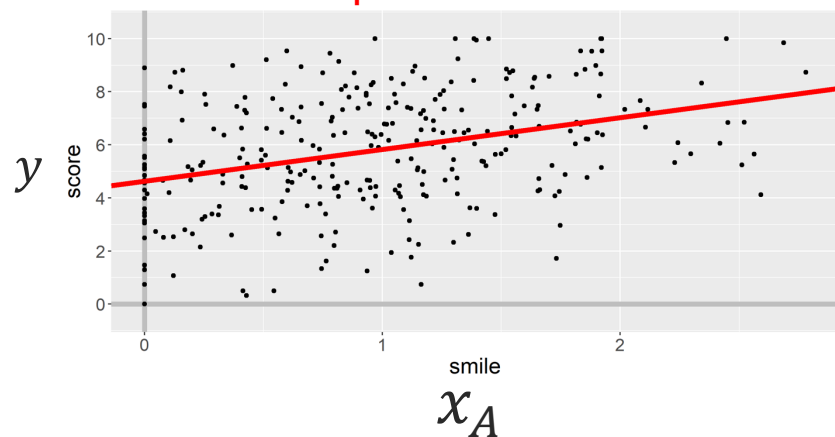
H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

$$y = w_0 + \boxed{w_1} x_A + \epsilon$$

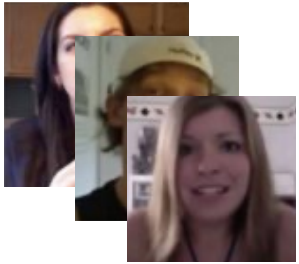
slope



	Estimate	95% CI
w_0	4.63	[4.20, 5.06]
w_1	1.20	[0.83, 1.57]

From Additive to Multiplicative

300 book reviews



y : audience score

x_A : percentage of smiling

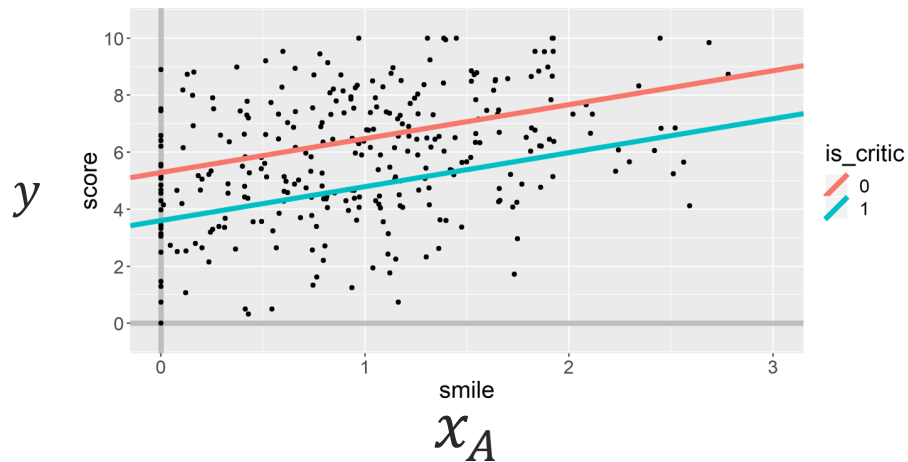
x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

$$y = w_0 + w_1 x_A + w_2 x_B + \epsilon$$



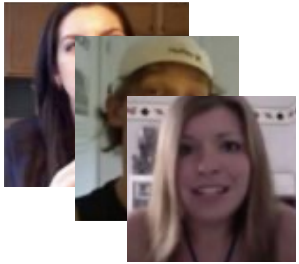
	Estimate	95% CI
w_0	5.29	[4.86, 5.73]
w_1	1.19	[0.85, 1.53]
w_2	-1.69	[-2.14, -1.24]

➔ Positive effect

➔ Negative effect

From Additive to Multiplicative

300 book reviews



y : audience score

x_A : percentage of smiling

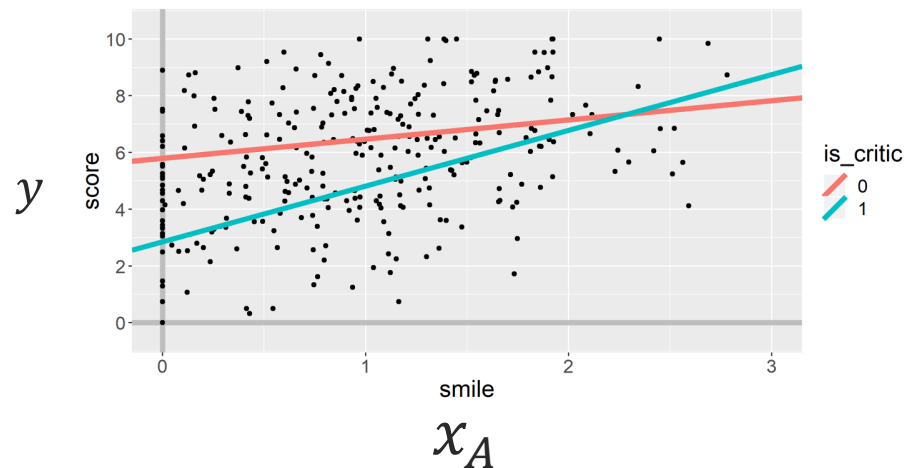
x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

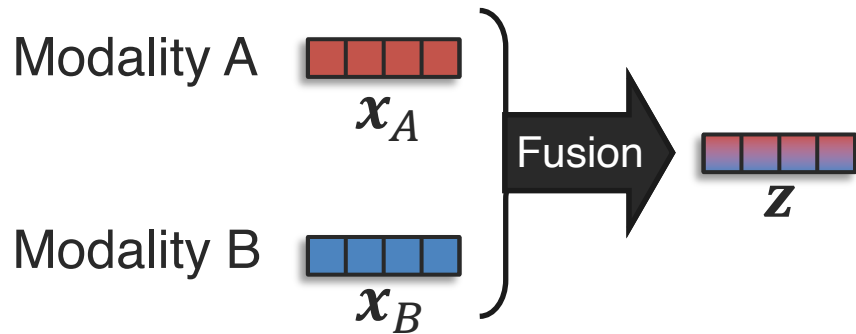
$$y = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$



	Estimate	95% CI
w_0	5.79	[5.29, 6.29]
w_1	0.68	[0.25, 1.11]
w_2	-2.94	[-3.73, -2.15]
w_3	1.29	[0.61, 1.97]

➔ **Multiplicative interaction!**

Basic Fusion – Additive Interactions

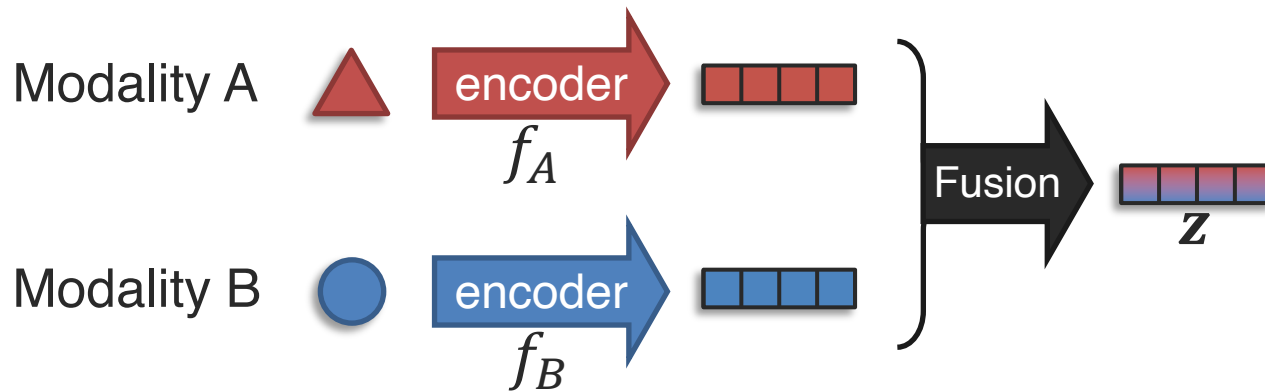


Additive fusion:

$$z = w_1 x_A + w_2 x_B$$

→ 1-layer neural network
can be seen as additive

With unimodal encoders:

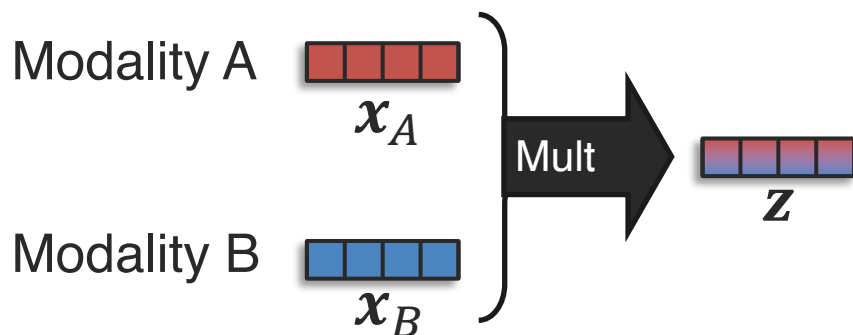


Additive fusion:

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

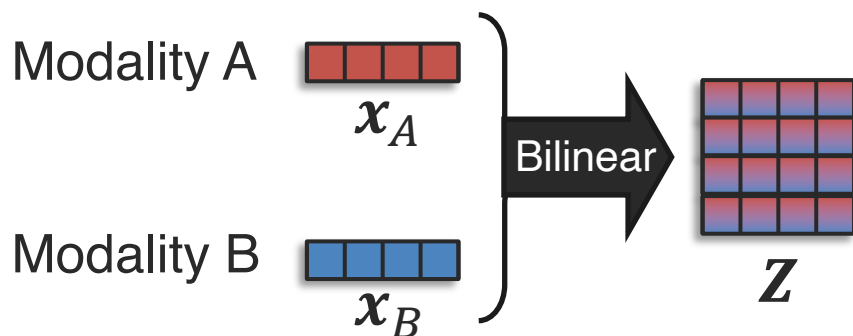
→ It could be seen as an
ensemble approach
(late fusion)

Multiplicative Interactions



Simple multiplicative fusion:

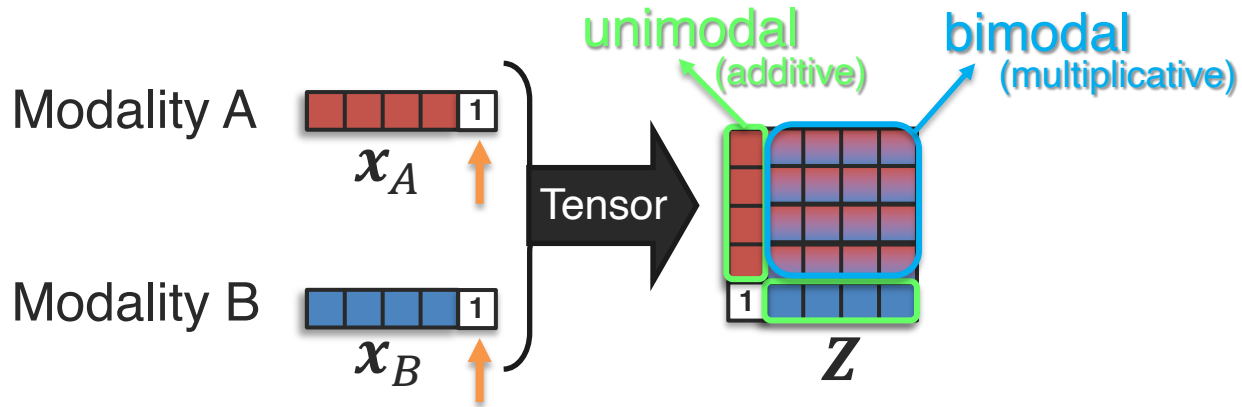
$$z = w(x_A \times x_B)$$



Bilinear Fusion:

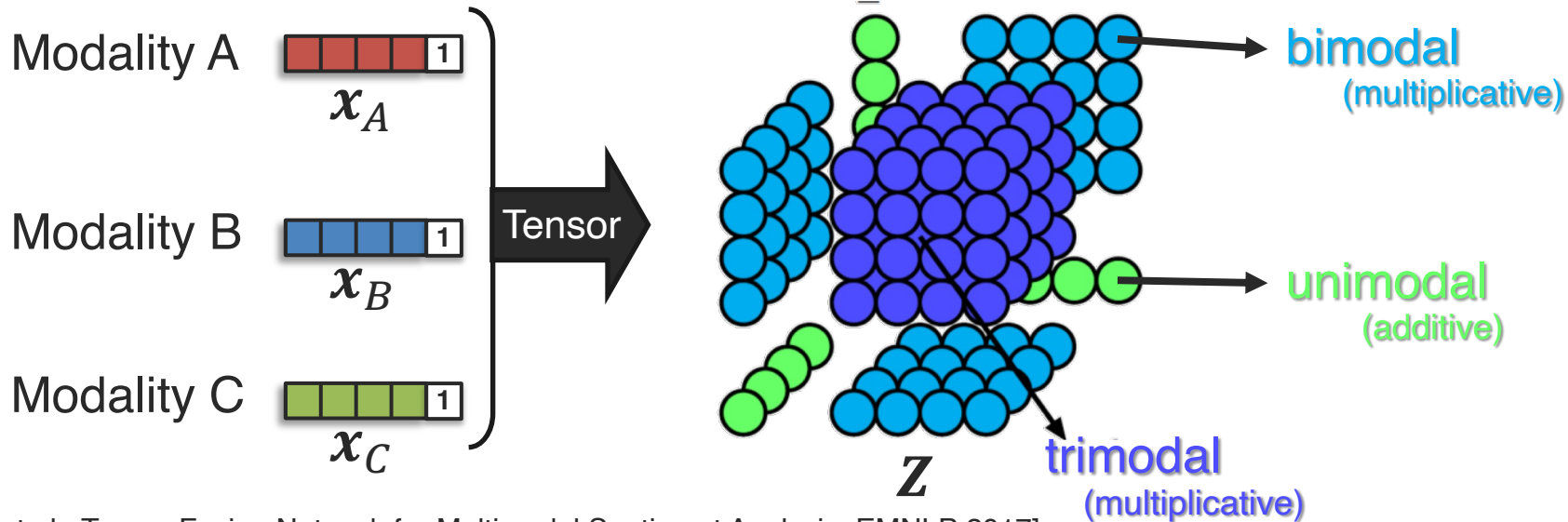
$$Z = W(x_A^T \cdot x_B)$$

Tensor Fusion



Tensor Fusion (bimodal):

$$Z = w([x_A \ 1]^T \cdot [x_B \ 1])$$

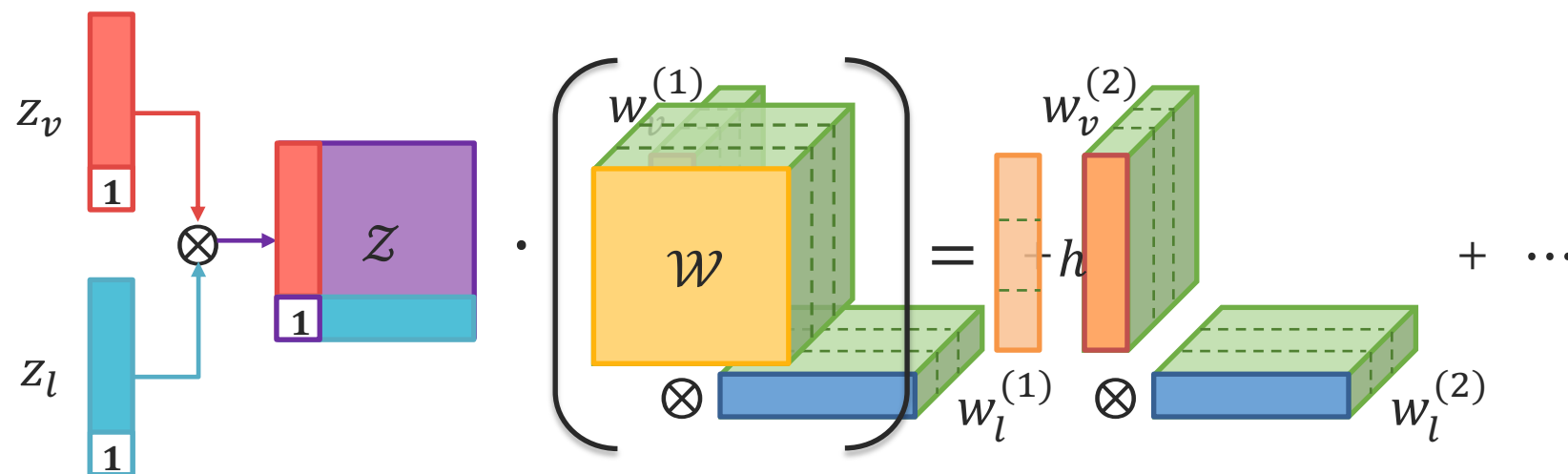


... but the weight matrix may end up quite large!

[Zadeh et al., Tensor Fusion Network for Multimodal Sentiment Analysis. EMNLP 2017]

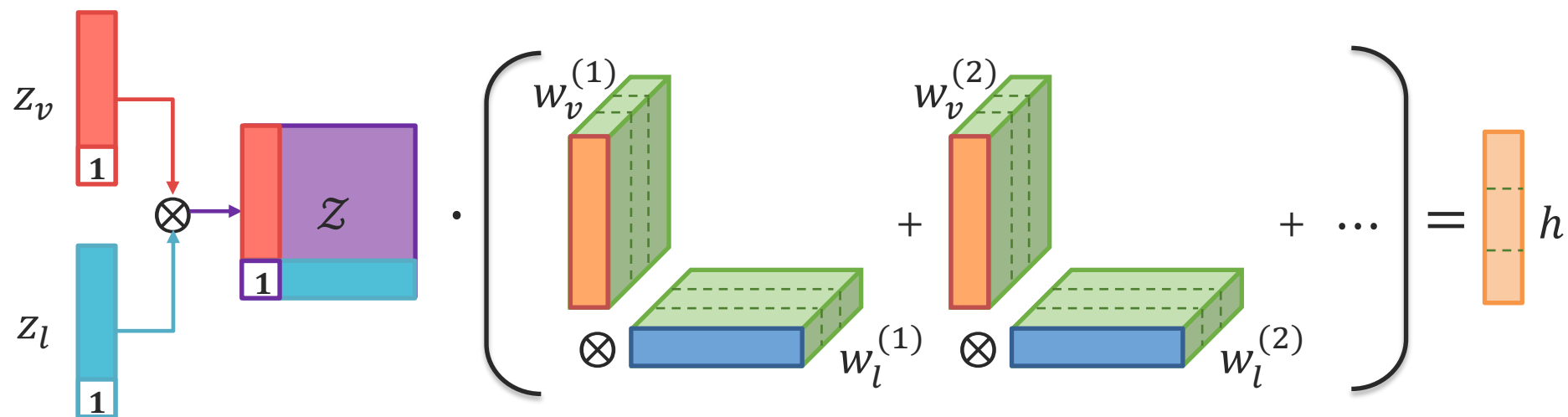
[Hou et al., Deep Multimodal Multilinear Fusion with High-order Polynomial Pooling. NeurIPS 2019]

Low-rank Tensor Fusion



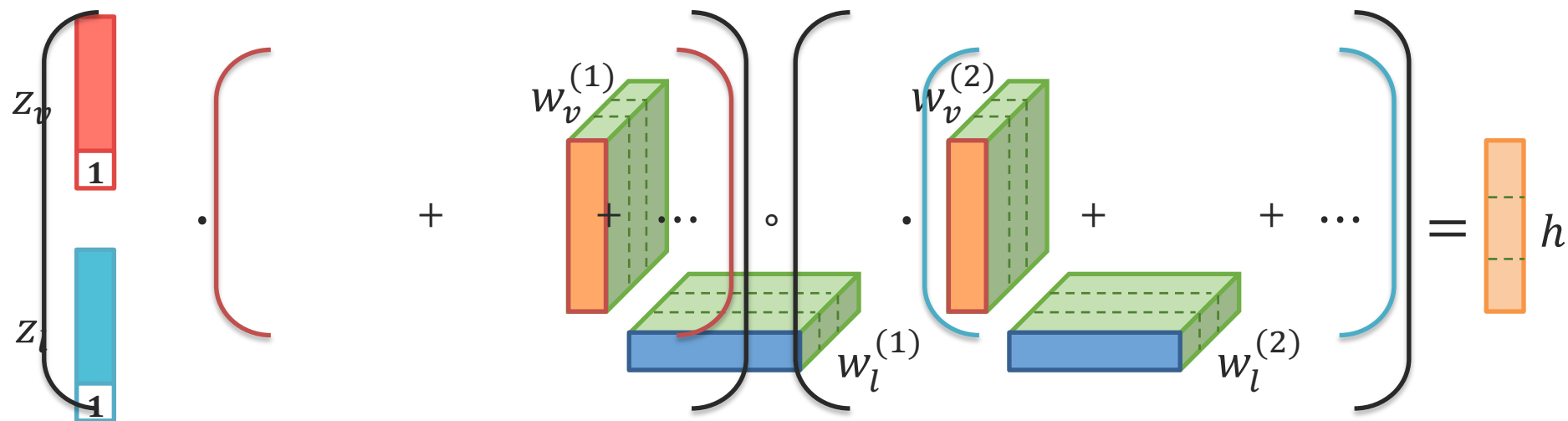
[Liu et al., Efficient Low-rank Multimodal Fusion with Modality-Specific Factors. ACL 2018]

Low-rank Tensor Fusion



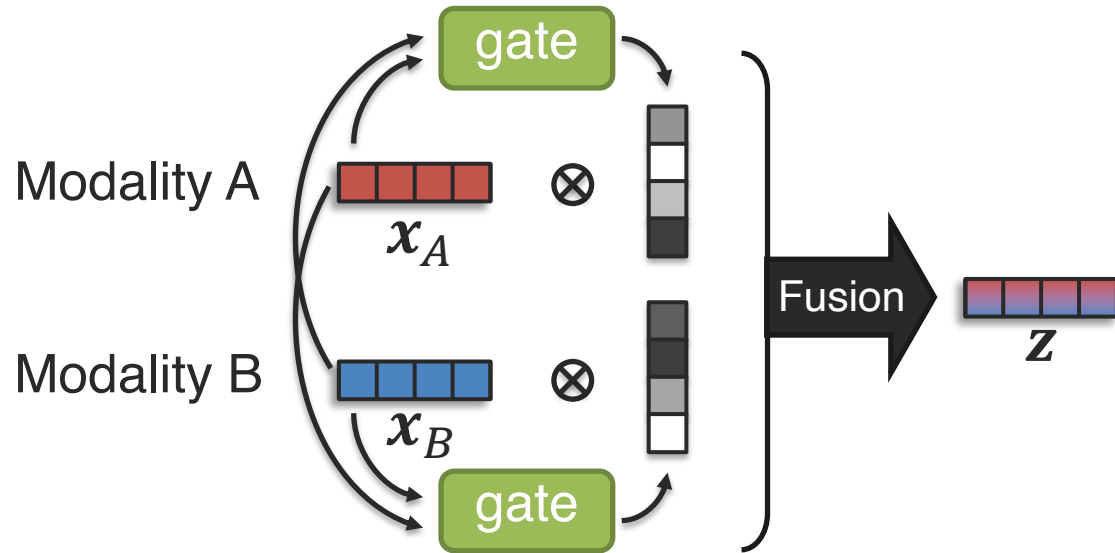
[Liu et al., Efficient Low-rank Multimodal Fusion with Modality-Specific Factors. ACL 2018]

Low-rank Tensor Fusion



[Liu et al., Efficient Low-rank Multimodal Fusion with Modality-Specific Factors. ACL 2018]

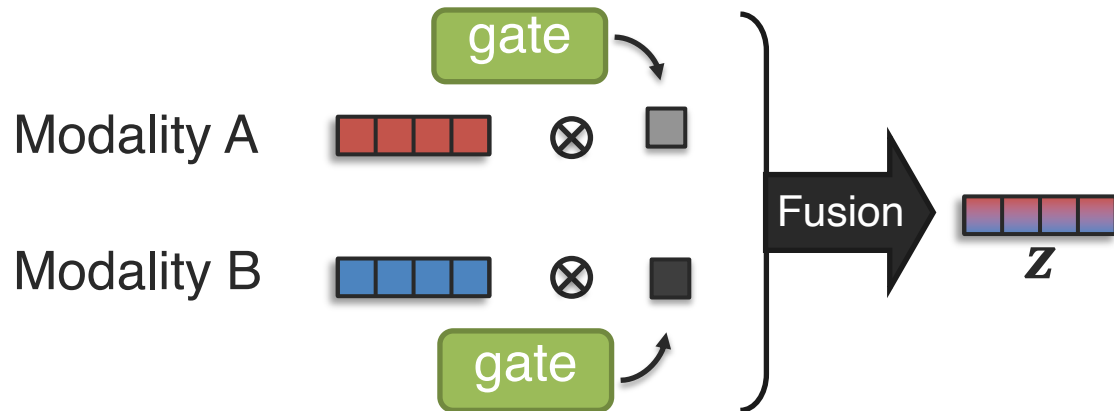
Gated Fusion



Example with additive fusion:

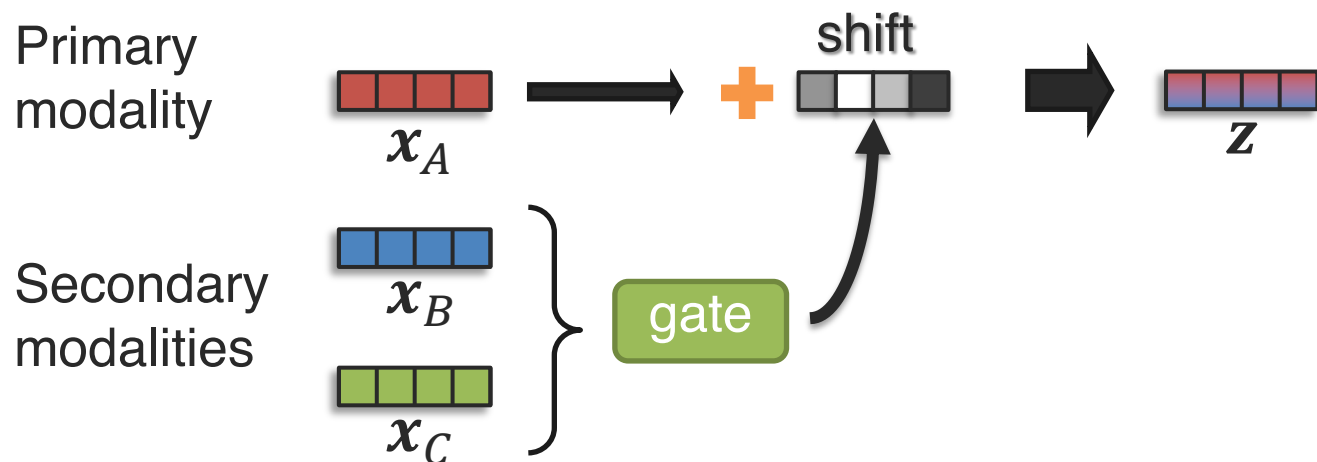
$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

$\rightarrow g_A$ and g_B can be seen as attention functions



\rightarrow Gating output can be one weight for the whole modality

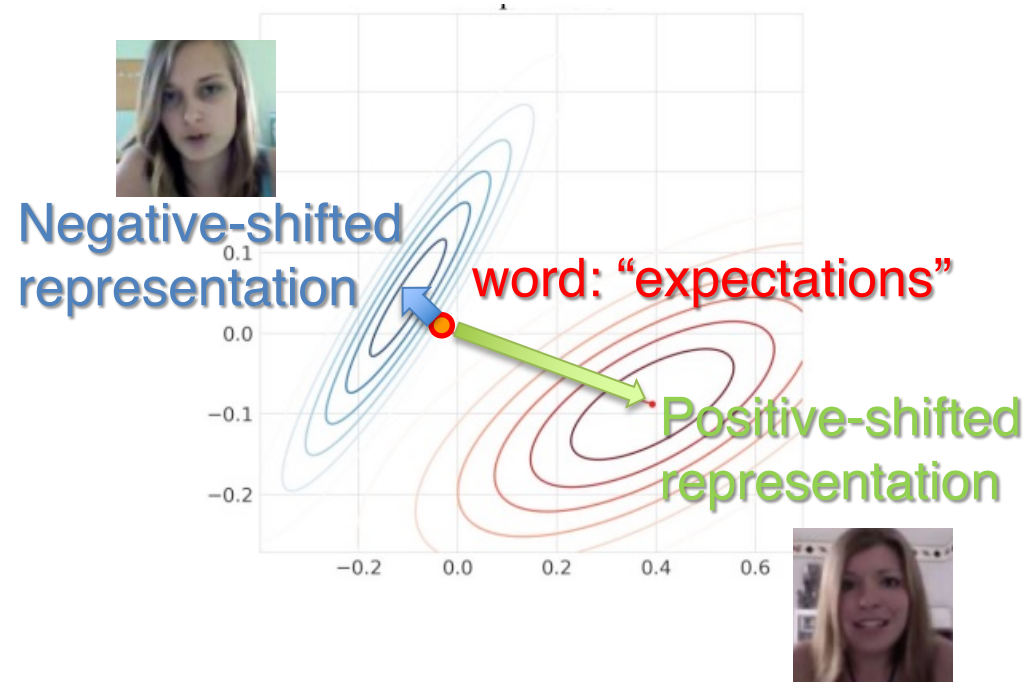
Modality-Shifting Fusion



Example with language modality:

Primary modality: language

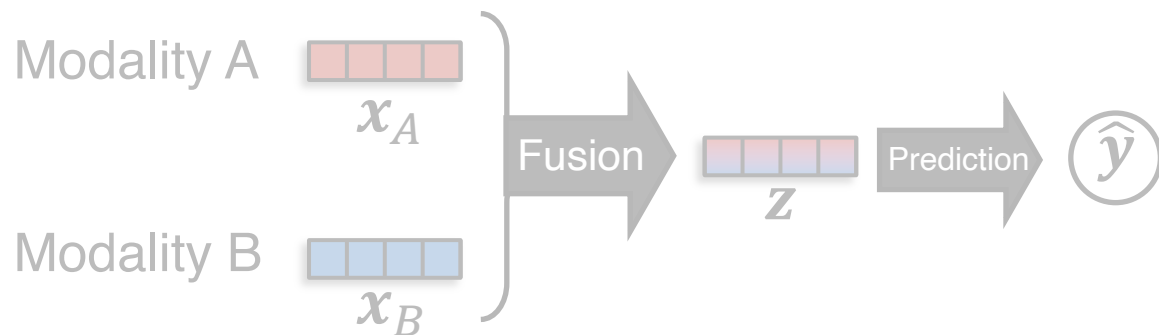
Secondary modalities: acoustic and visual



[Wang et al., Words Can Shift: Dynamically Adjusting Word Representations Using Nonverbal Behaviors, AAAI 2019]

[Rahman et al., Integrating Multimodal Information in Large Pretrained Transformers, ACL 2020]

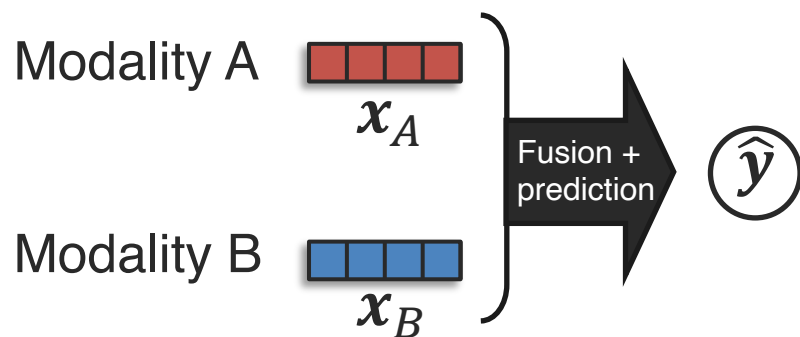
Nonlinear Fusion



Nonlinear fusion:

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

For any nonlinear model

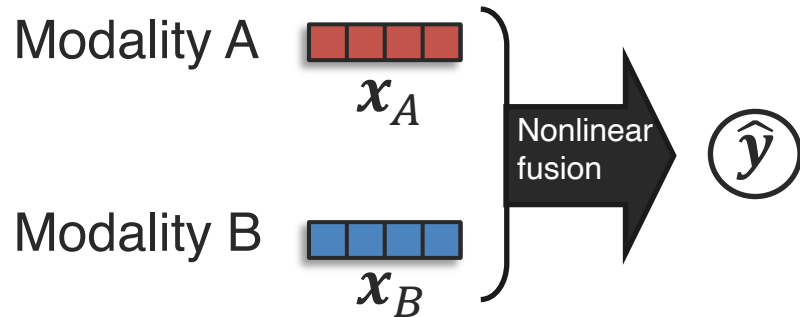


→ This could be seen as *early fusion*:

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

... 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):

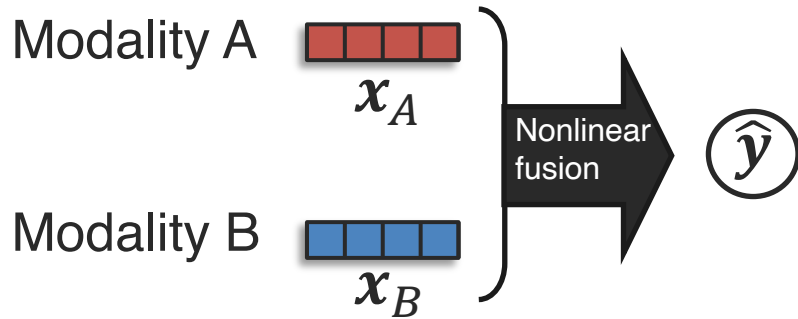
$$\tilde{f}(\mathbf{x}_A, \mathbf{x}_B) = \underbrace{\mathbb{E}_{\mathbf{x}_B} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_A(\mathbf{x}_A)} + \underbrace{\mathbb{E}_{\mathbf{x}_A} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_B(\mathbf{x}_B)}$$

Modality A + Modality B

Additive fusion
(approximation)

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

Measuring Non-Additive Interactions



Nonlinear fusion:

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

EMAP
projection

Additive fusion:

$$\hat{y}' = \hat{f}_A(x_A) + \hat{f}_B(x_B) + \hat{\mu}$$

		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
	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

Always a good baseline!

Differences are small!!!

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

Learning Non-additive Bimodal and Trimodal Interactions

Idea: prioritize simpler interactions

Multimodal Residual Optimization

Unimodal
(additive)

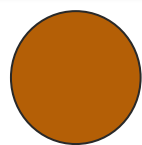
Bimodal
(non-additive)

Trimodal
(non-additive)

residual

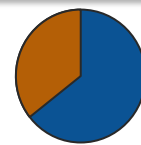
residual

$$\mathcal{L}(y, \hat{y}_{uni}) + \mathcal{L}(y - \hat{y}_{uni}, \hat{y}_{bi}) + \mathcal{L}(y - \hat{y}_{uni} - \hat{y}_{bi}, \hat{y}_{tri})$$



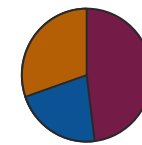
\hat{y}_{uni}

+



\hat{y}_{bi}

+



\hat{y}_{tri}

=

\hat{y}

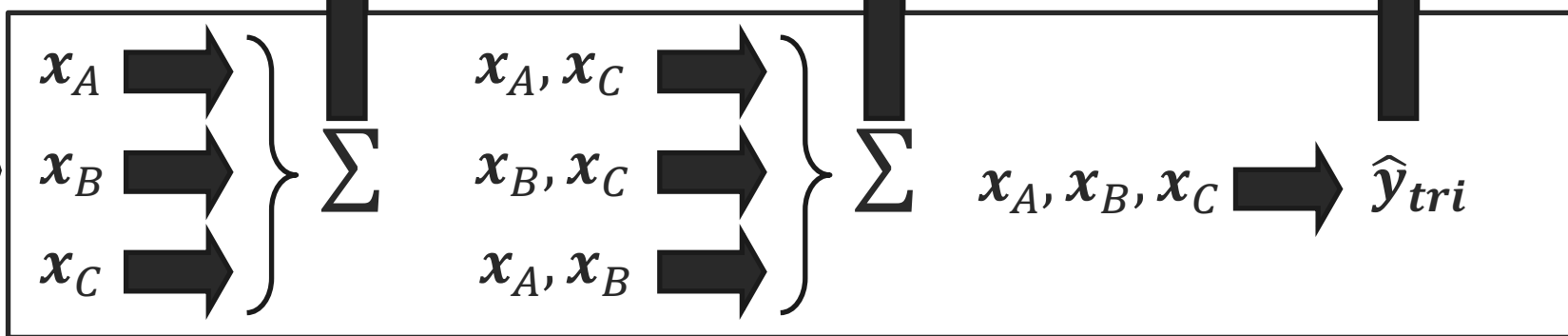
Modality A



Modality B



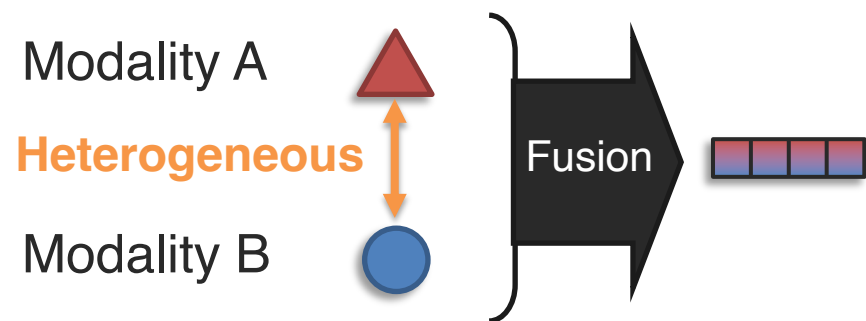
Modality C



[Wortwein et al., Beyond Additive Fusion: Learning Non-Additive Multimodal Interactions, Findings-EMNLP 2022]

Fusion with Heterogeneous Modalities

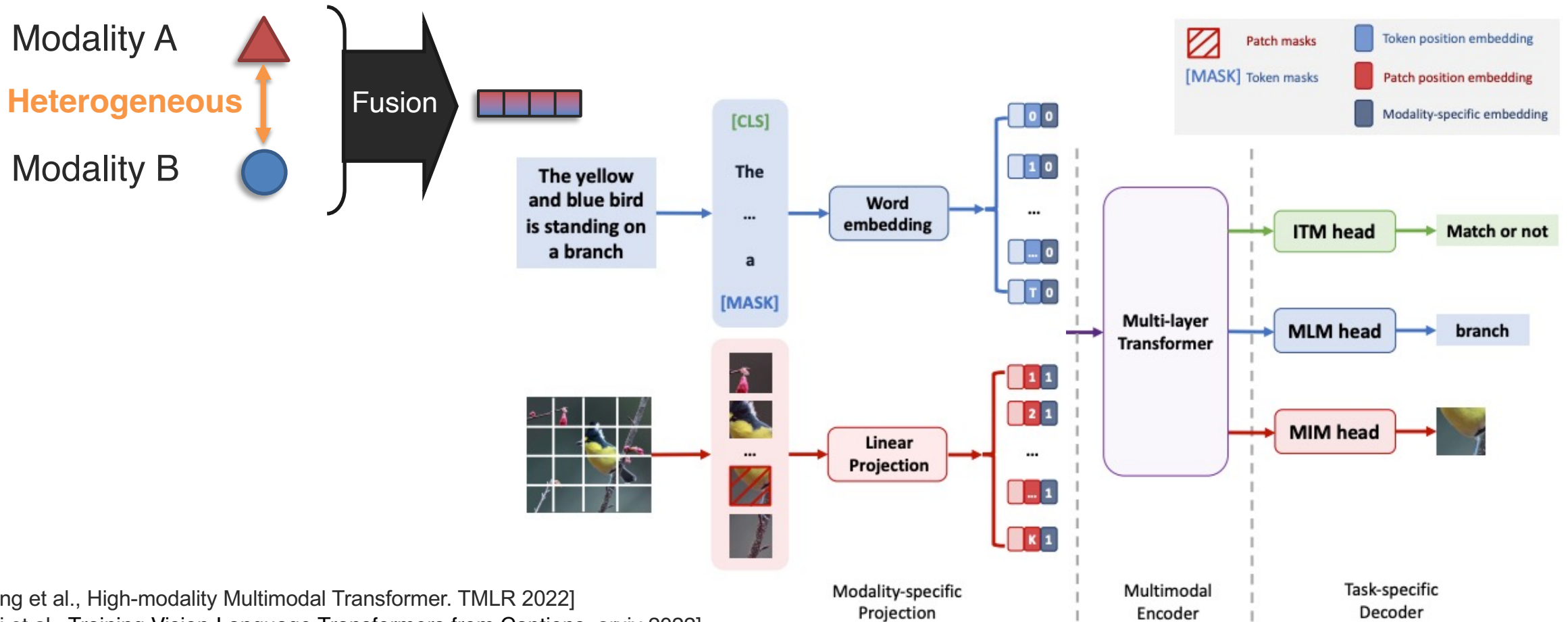
Example: From feature fusion to early fusion



[Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale. ICLR 2021]

Fusion with Heterogeneous Modalities

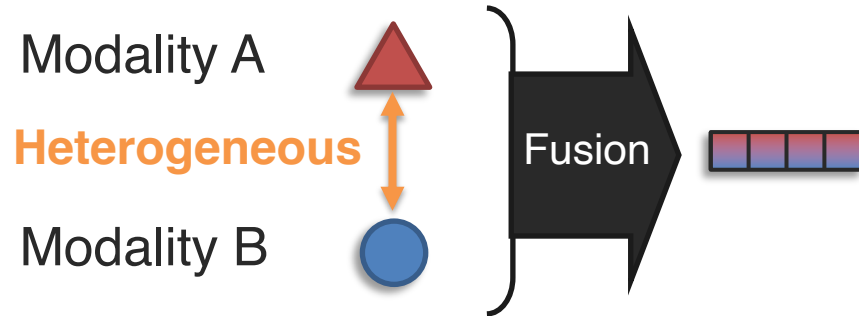
Example: From feature fusion to early fusion



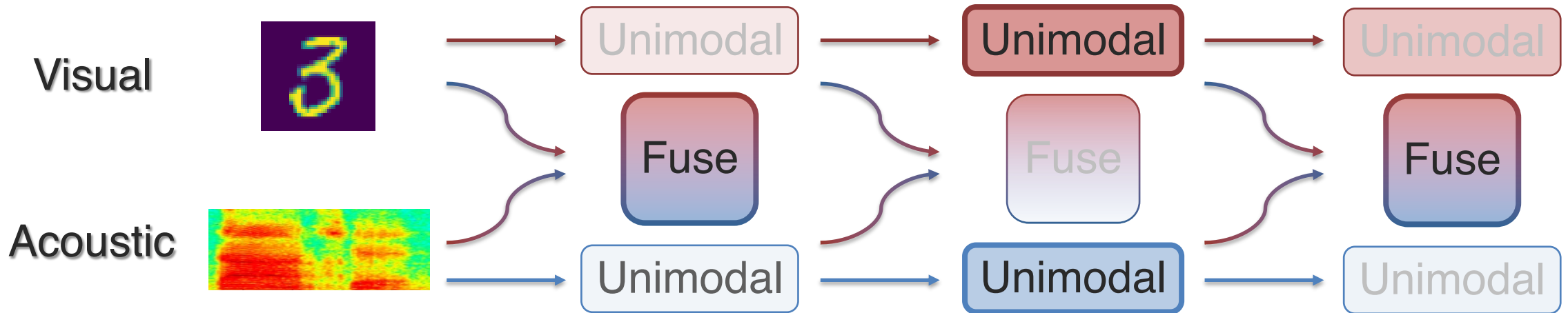
[Liang et al., High-modality Multimodal Transformer. TMLR 2022]

[Gui et al., Training Vision-Language Transformers from Captions. arxiv 2022]

Dynamic Early Fusion



Idea: Deciding when to fuse in early fusion



[Xue and Marculescu, Dynamic Multimodal Fusion, arxiv 2022]

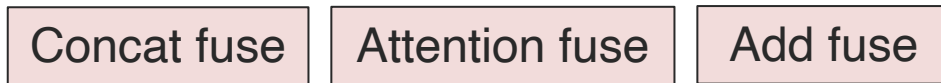
Dynamic Early Fusion

Fusion fully learned from optimization and data

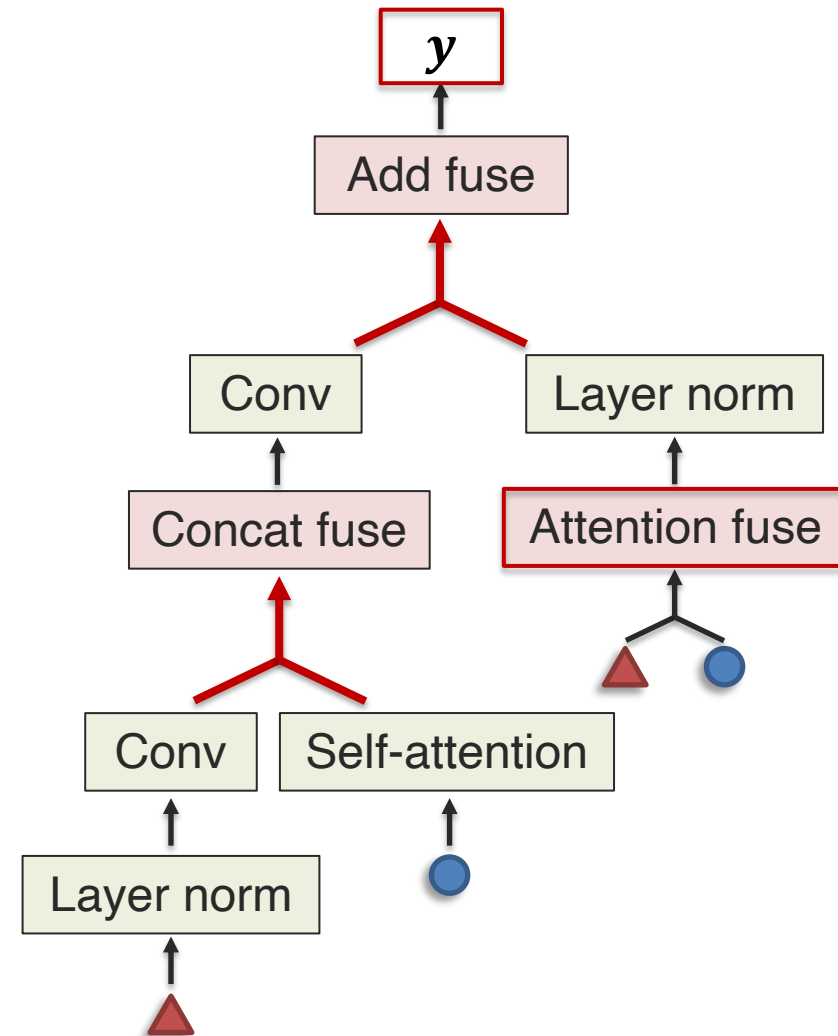
1. Define basic representation building blocks



2. Define basic fusion building blocks



3. Automatically search for composition using neural architecture search



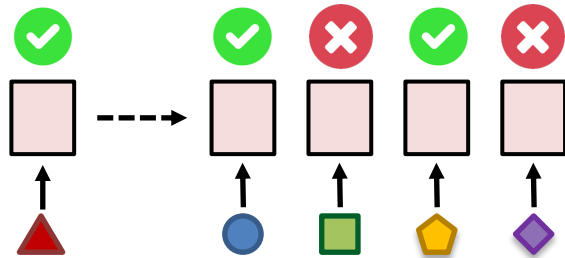
[Xu et al., MUFASA: Multimodal Fusion Architecture Search for Electronic Health Records. AAAI 2021]

[Liu et al., DARTS: Differentiable Architecture Search. ICLR 2019]

Heterogeneity-aware Fusion

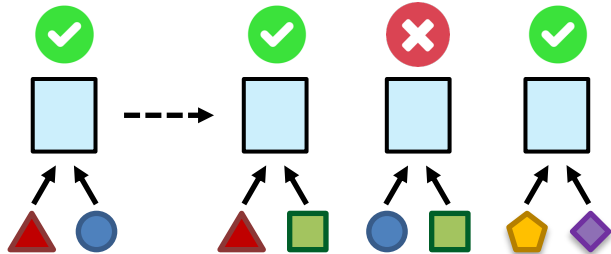
Information transfer, transfer learning perspective

1a. Estimate modality heterogeneity via transfer



(Implicitly captures heterogeneity)

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

3. Determine parameter clustering

$$U_1 = \{U_1, U_2, U_4\}$$

$$U_2 = \{U_3\}$$

$$U_3 = \{U_5\}$$

$$C_1 = \{C_{12}, C_{13}, C_{45}\}$$

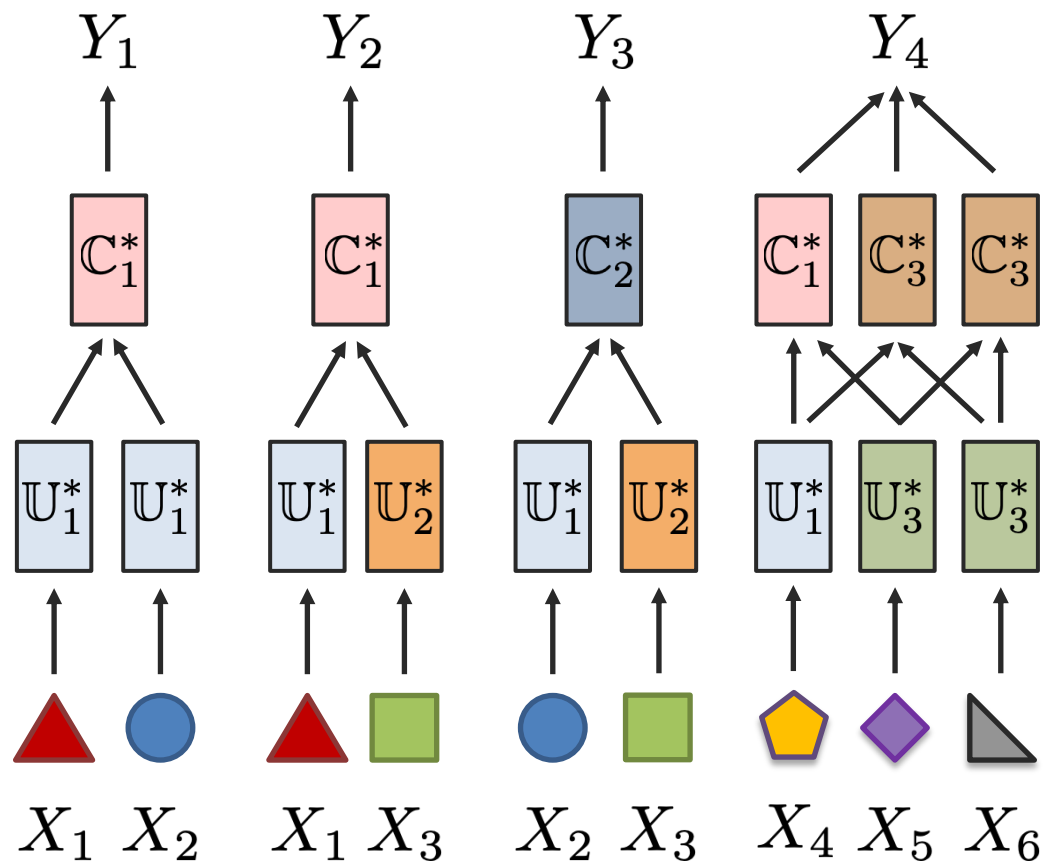
$$C_2 = \{C_{23}\}$$

[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]

Improving Optimization

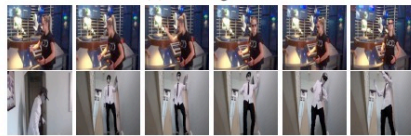
Kinetics dataset



(a) headbanging



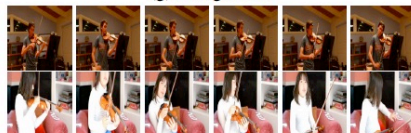
(c) shaking hands



(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?**

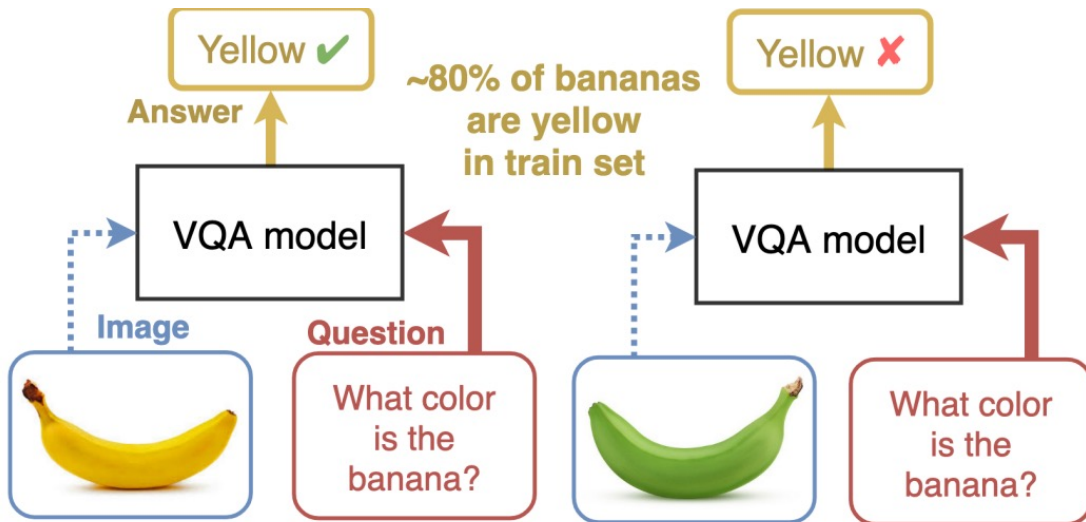
[Wang et al., What Makes Training Multi-modal Classification Networks Hard? CVPR 2020]

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

Improving Optimization

Information heterogeneity and unimodal biases

Finding: VQA models answer the question without looking at the image



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



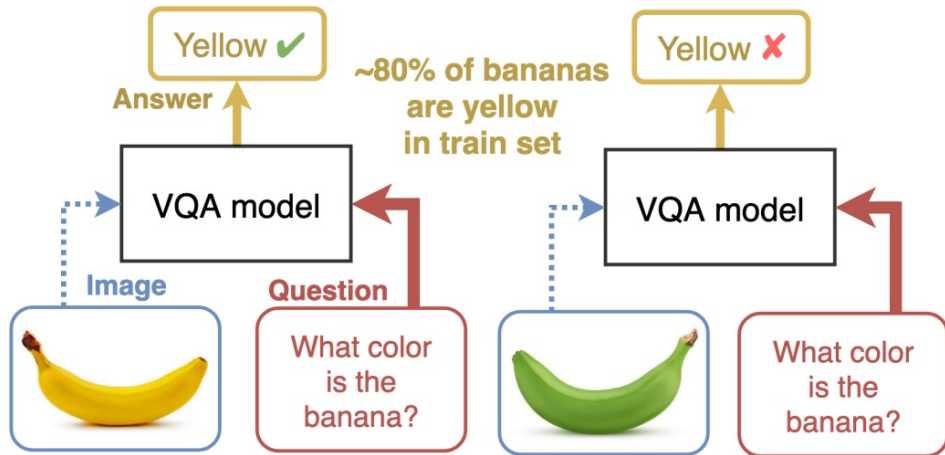
[Goyal et al., Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. CVPR 2017]

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

Improving Optimization

Information heterogeneity and modality collapse

VQA models answer the question without looking at the image

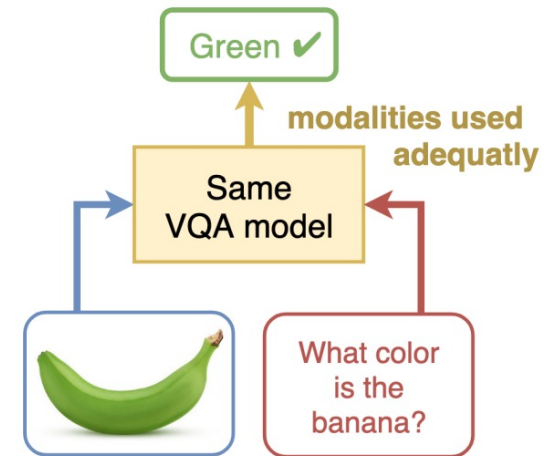


Balancing modalities

Balancing training



Not the case when trained with RUBi



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

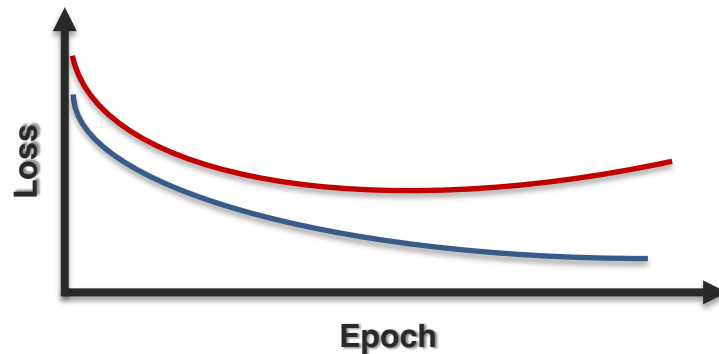
[Goyal et al., Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. CVPR 2017]

Improving Optimization

Relevance heterogeneity

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

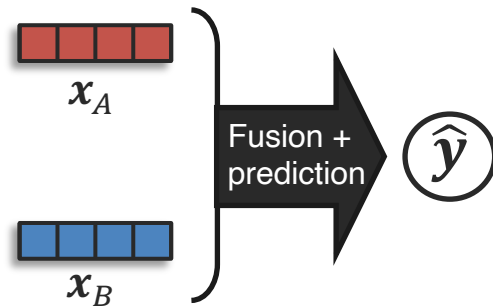
[Wang et al., What Makes Training Multi-modal Classification Networks Hard? CVPR 2020]

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

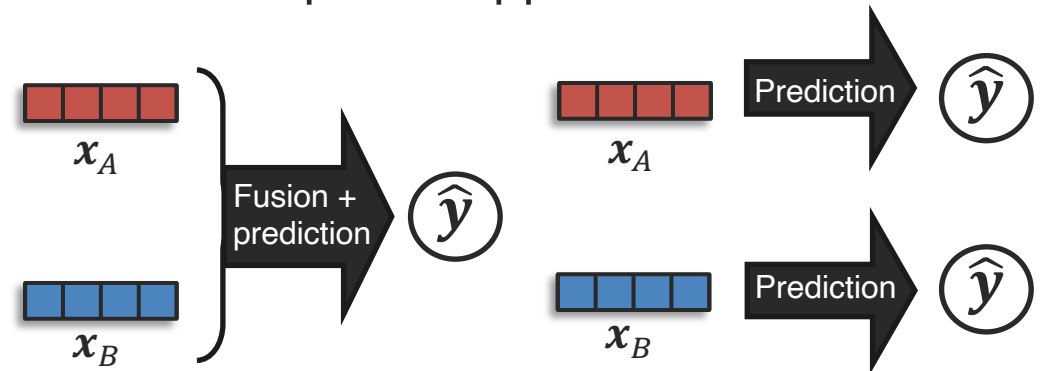
Improving Optimization

Relevance heterogeneity

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]

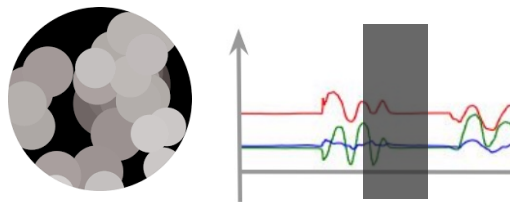
[Wu et al., Characterizing and Overcoming the Greedy Nature of Learning in Multi-modal Deep Neural Networks. ICML 2022]

Improving Robustness

Heterogeneity in noise

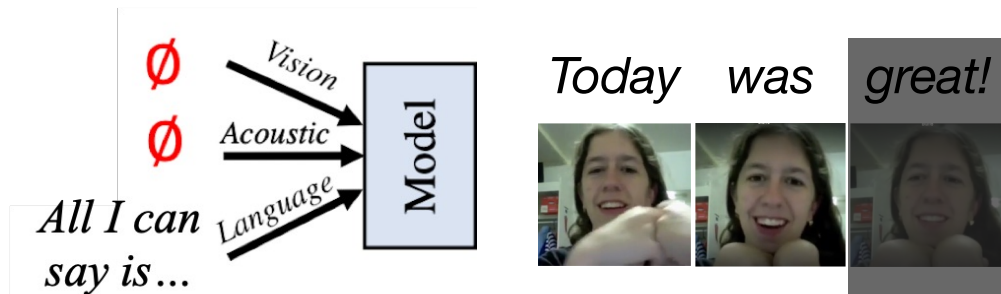
Modality-specific robustness

noise → **nosie**



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

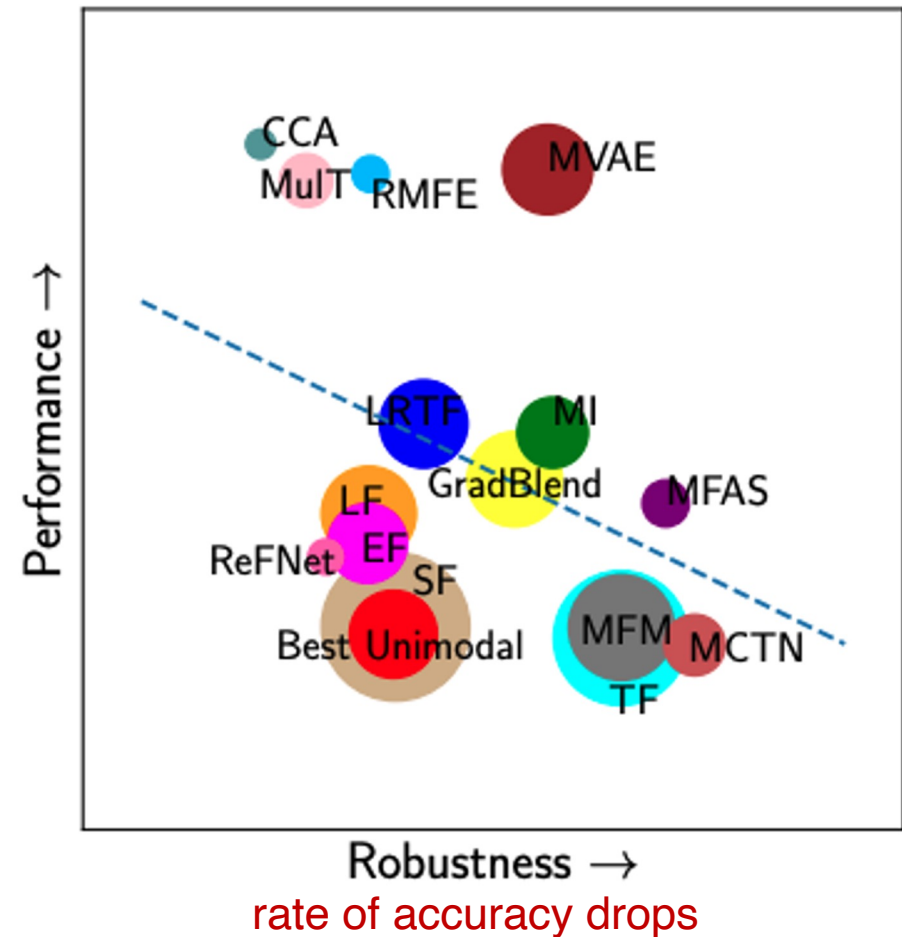
Multimodal robustness



[Zadeh et al., 2020]

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

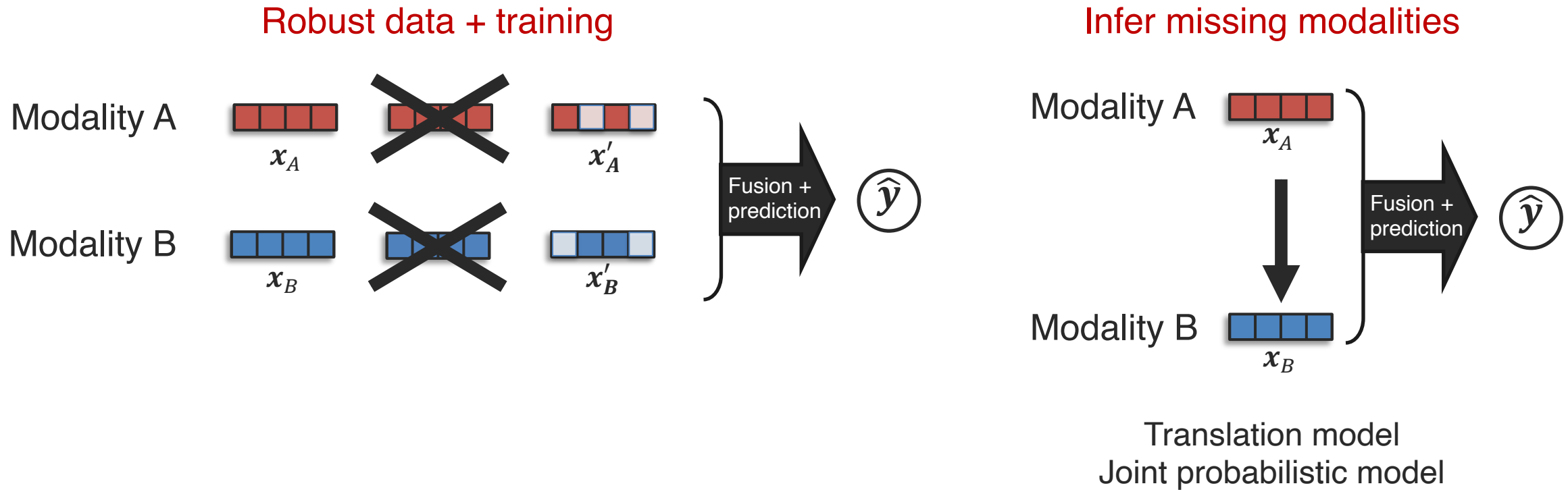
Tradeoffs between performance and robustness



Robustness →
rate of accuracy drops

Improving Robustness

Several approaches towards more robust models



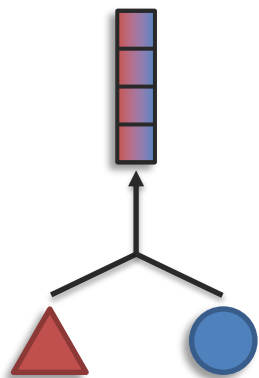
[Ngiam et al., Multimodal Deep Learning. ICML 2011]

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

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

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

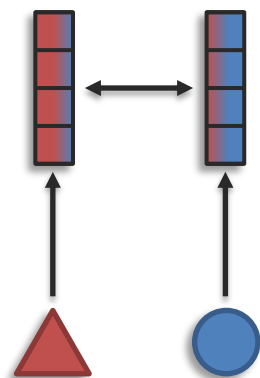
Sub-Challenge 1a: Representation Fusion



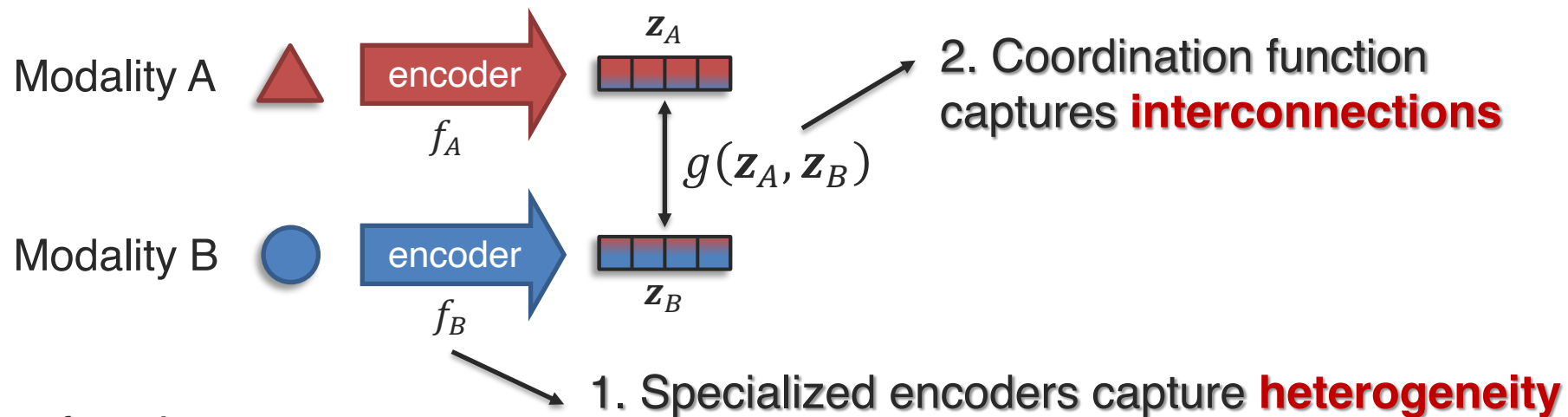
Definition: Learn a joint representation that models cross-modal interactions between individual elements of different modalities



Sub-Challenge: Representation Coordination



Definition: Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions

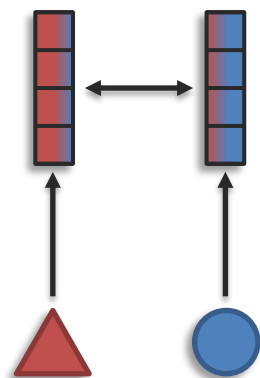


Learning with coordination function:

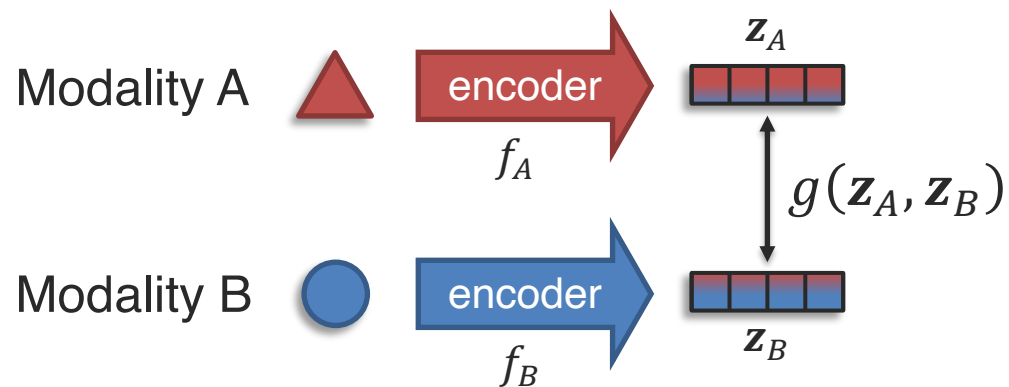
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

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

Sub-Challenge: Representation Coordination



Definition: Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions



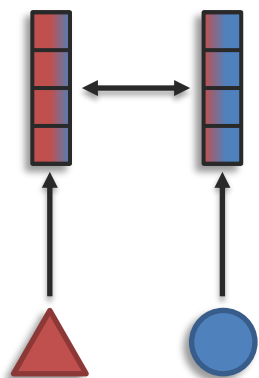
Learning with coordination function:

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

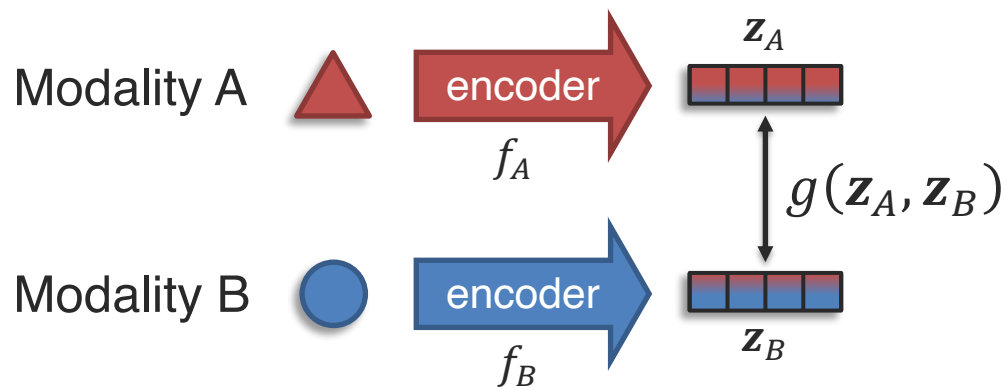
with model parameters θ_g , θ_{f_A} and θ_{f_B}

① Cosine similarity:
$$g(\mathbf{z}_A, \mathbf{z}_B) = \frac{\mathbf{z}_A \cdot \mathbf{z}_B}{\|\mathbf{z}_A\| \|\mathbf{z}_B\|}$$

Sub-Challenge: Representation Coordination



Definition: Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions



Learning with coordination function:

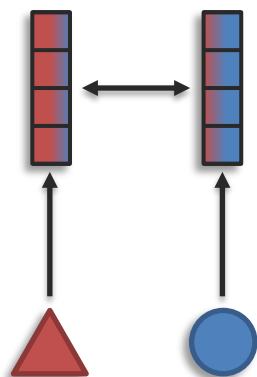
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

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

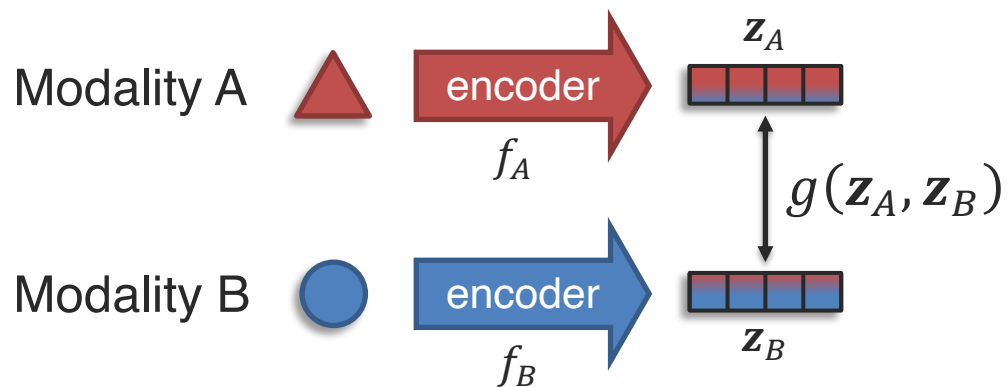
② Kernel similarity functions:

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

Sub-Challenge: Representation Coordination



Definition: Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions



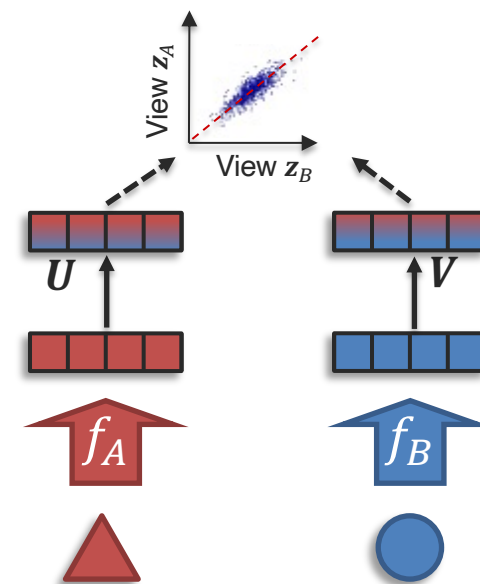
Learning with coordination function:

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

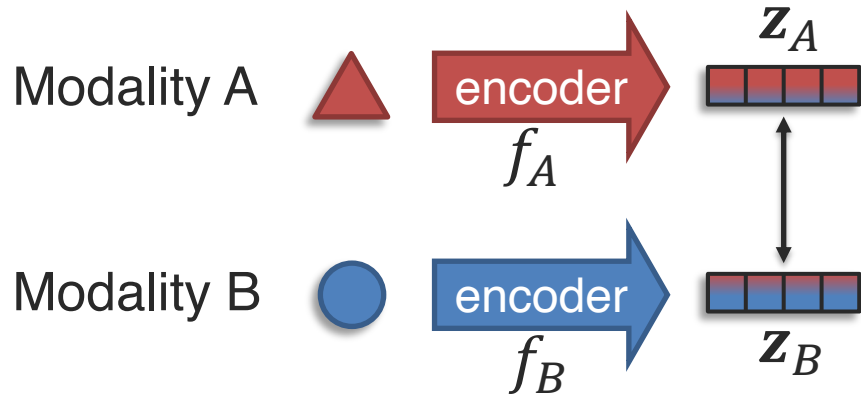
with model parameters θ_g , θ_{f_A} and θ_{f_B}

③ Canonical Correlation Analysis:

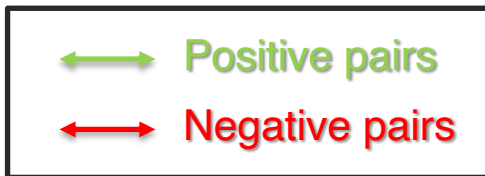
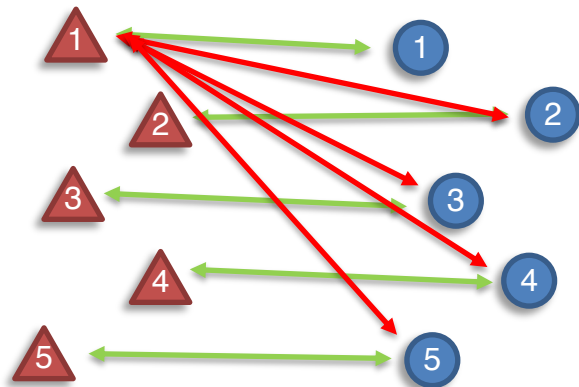
$$\operatorname{argmax}_{V, U, f_A, f_B} \operatorname{corr}(z_A, z_B)$$



Coordination with Contrastive Learning



Paired data: $\{\triangle, \circ\}$
(e.g., images and text descriptions)



Contrastive loss:

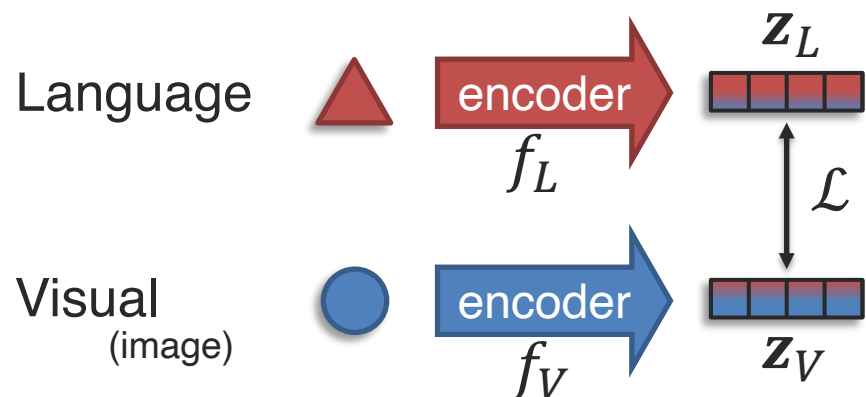
 brings **positive pairs** closer and pushes **negative pairs** apart

Simple contrastive loss:

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

Similarity functions are often cosine similarity

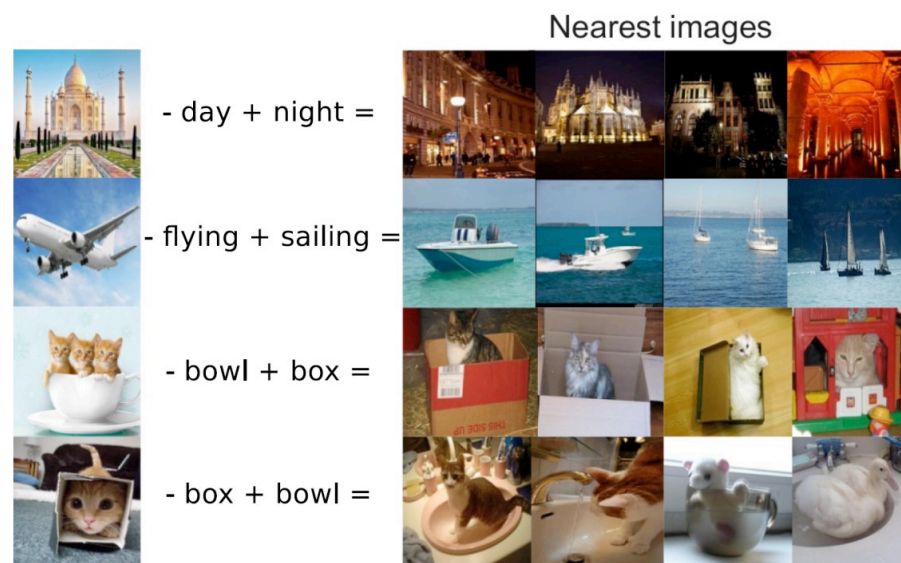
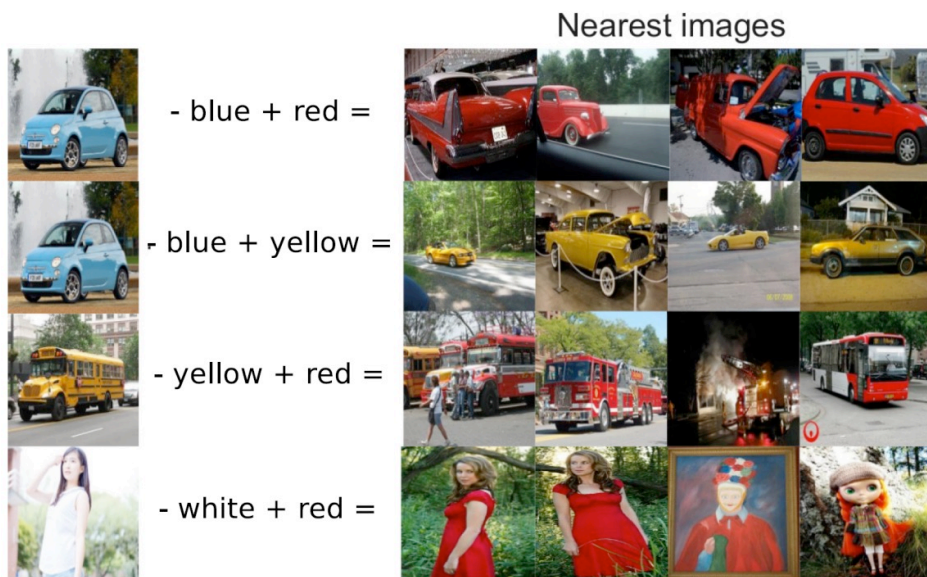
Example – Visual-Semantic Embeddings



Two contrastive loss terms:

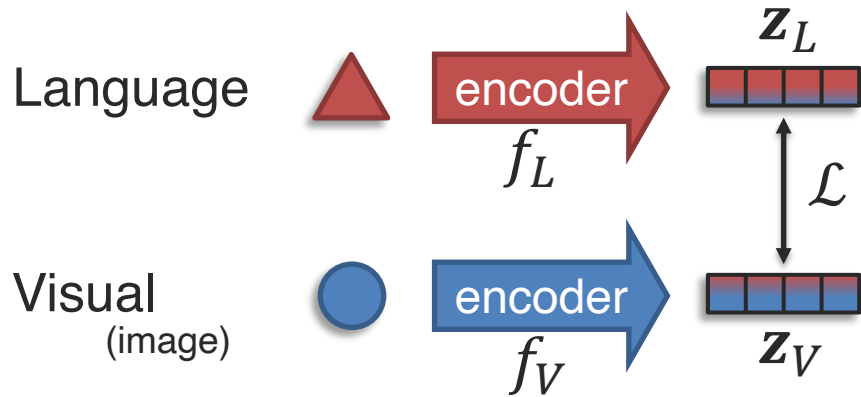
$$\max\{0, \alpha + \text{sim}(z_L, z_V^+) - \text{sim}(z_L, z_V^-)\}$$

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



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

Example – CLIP (Contrastive Language–Image Pre-training)



Popular contrastive loss: InfoNCE

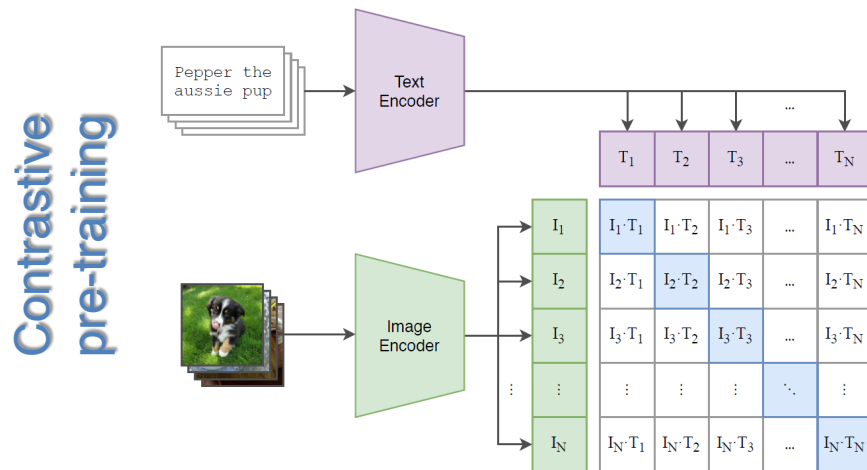
$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(z_A^i, z_B^i)}{\sum_{j=1}^N \text{sim}(z_A^i, z_B^j)}$$

Similarity function can be cosine similarity

positive pairs

negative pairs and positive pairs

Positive and negative pairs:

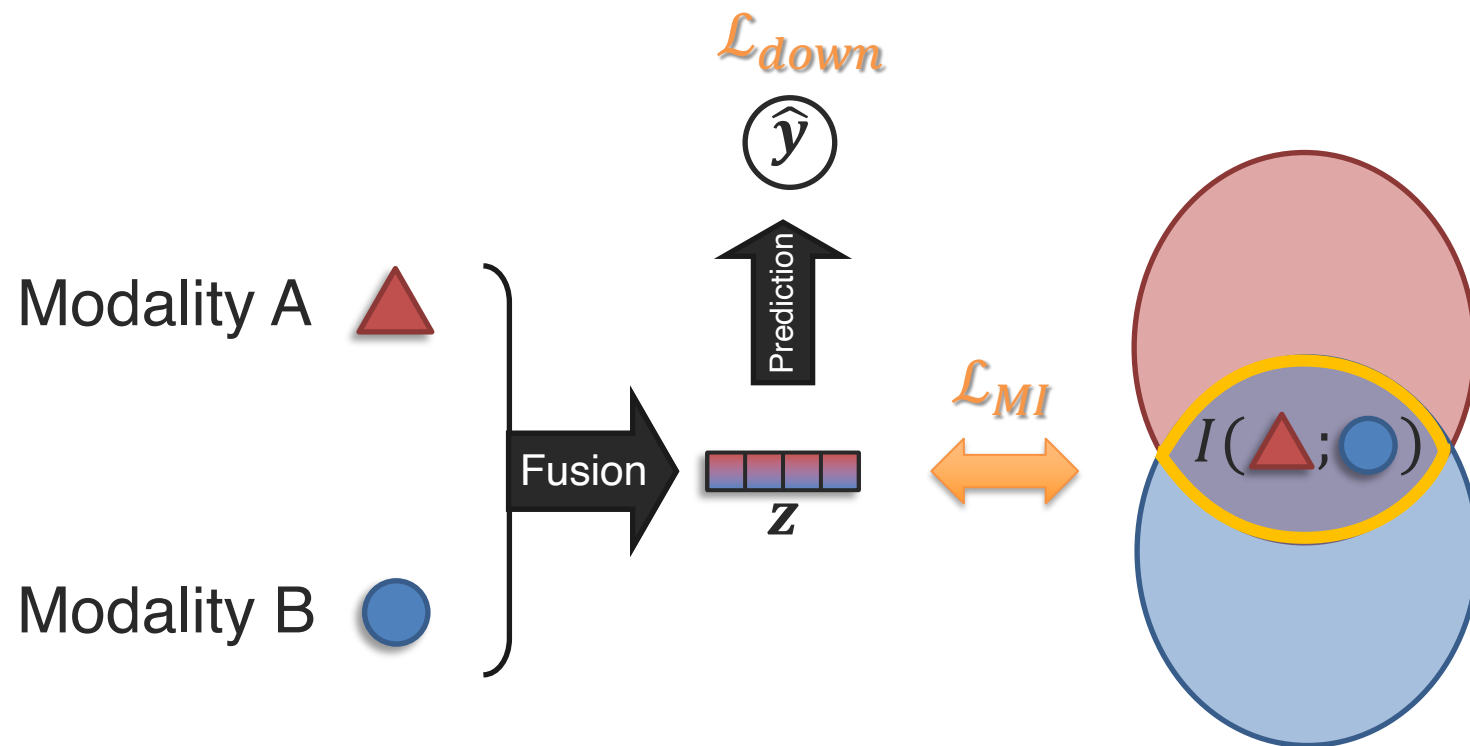


CLIP encoders (f_L and f_V) are great for language-vision tasks

z_L and z_V are coordinated but not identical representation spaces

[Radford et al., Learning Transferable Visual Models From Natural Language Supervision, ICML 2021]

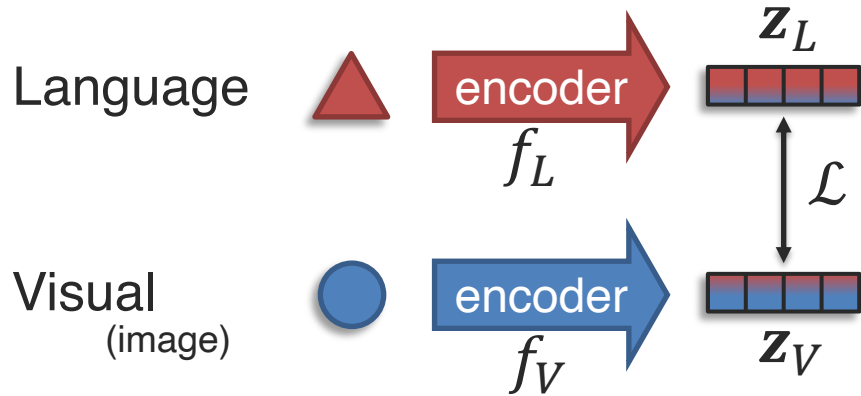
Multimodal Fusion with Mutual Information



Assumption?

Information present in both modalities is most important for the downstream task

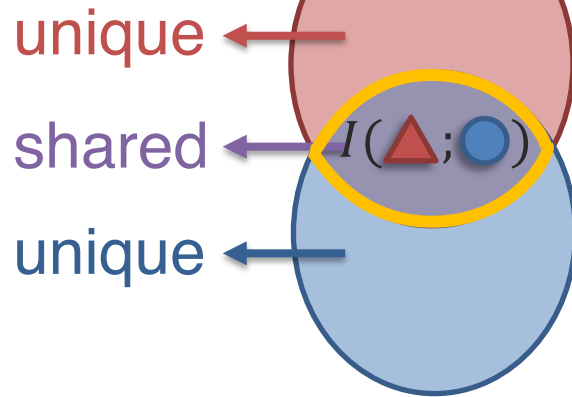
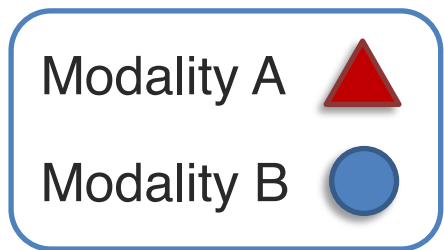
Contrastive Learning and Connected Modalities



Popular contrastive loss: InfoNCE

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

Connected modalities:

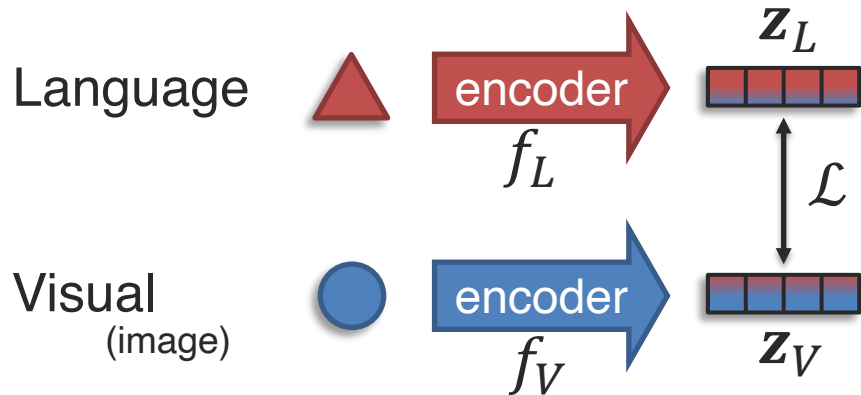


CLIP focuses on shared connections

Mutual information $I(X; Y)$

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

Contrastive Learning and Mutual Information



InfoNCE:

$$\mathcal{L} = -\mathbb{E} \left[\log \frac{f(\mathbf{x}_A^i, \mathbf{x}_B^i)}{\sum_{j=1}^N f(\mathbf{x}_A^i, \mathbf{x}_B^j)} \right]$$

critic function

Critic function f is trained to be a binary classifier distinguishing $\mathbf{x}_A, \mathbf{x}_B \sim p(\mathbf{x}_A, \mathbf{x}_B)$ vs $\mathbf{x}_A, \mathbf{x}_B \sim p(\mathbf{x}_A)p(\mathbf{x}_B)$

InfoNCE/CL:

- 'Captures' mutual information
- Optimizes a lower bound on mutual information

At optimal loss, $f^*(\mathbf{x}_A, \mathbf{x}_B) = \frac{p(\mathbf{x}_A, \mathbf{x}_B)}{p(\mathbf{x}_A)p(\mathbf{x}_B)}$.

Plugging f^* back into \mathcal{L} gives:

$$\mathcal{L}^* \geq \mathbb{E} \left[\log \frac{p(\mathbf{x}_A)p(\mathbf{x}_B)}{p(\mathbf{x}_A, \mathbf{x}_B)} N \right] = -I(X_A, X_B) + \log N$$

In other words:

$$I(X_A, X_B) \geq \log N - \mathcal{L}^*$$

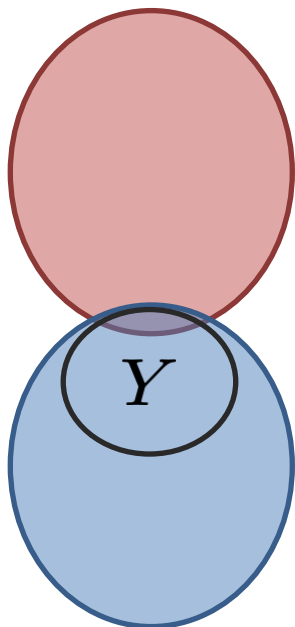
Multiview Redundancy and Contrastive Learning



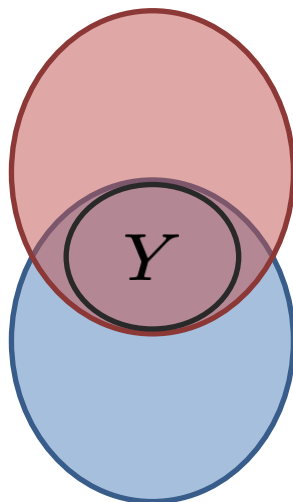
How much information should be shared?

Multi-view redundancy: $I(X_1; X_2) = I(X_1; Y)$

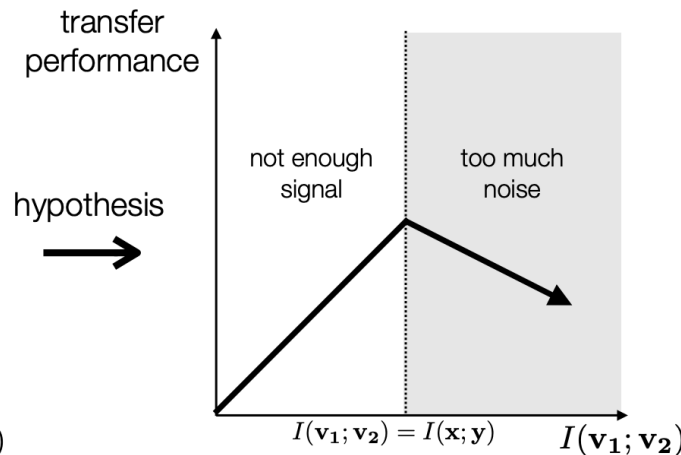
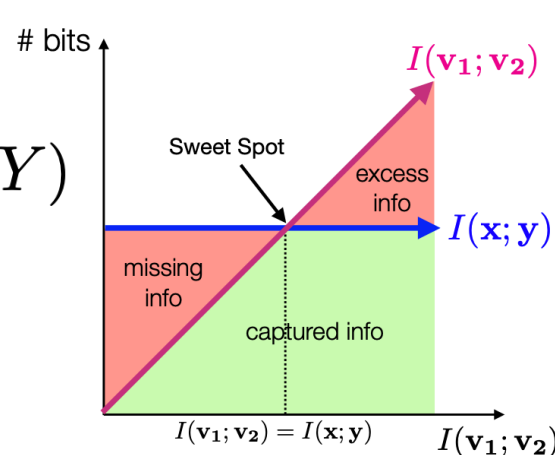
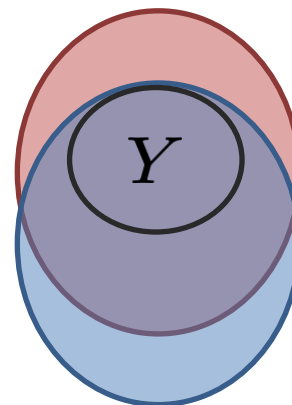
Not enough signal



Just right



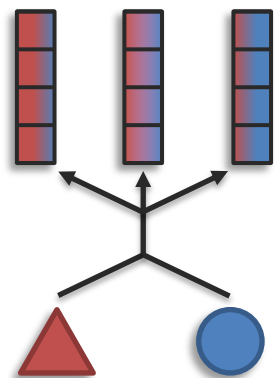
Too much noise



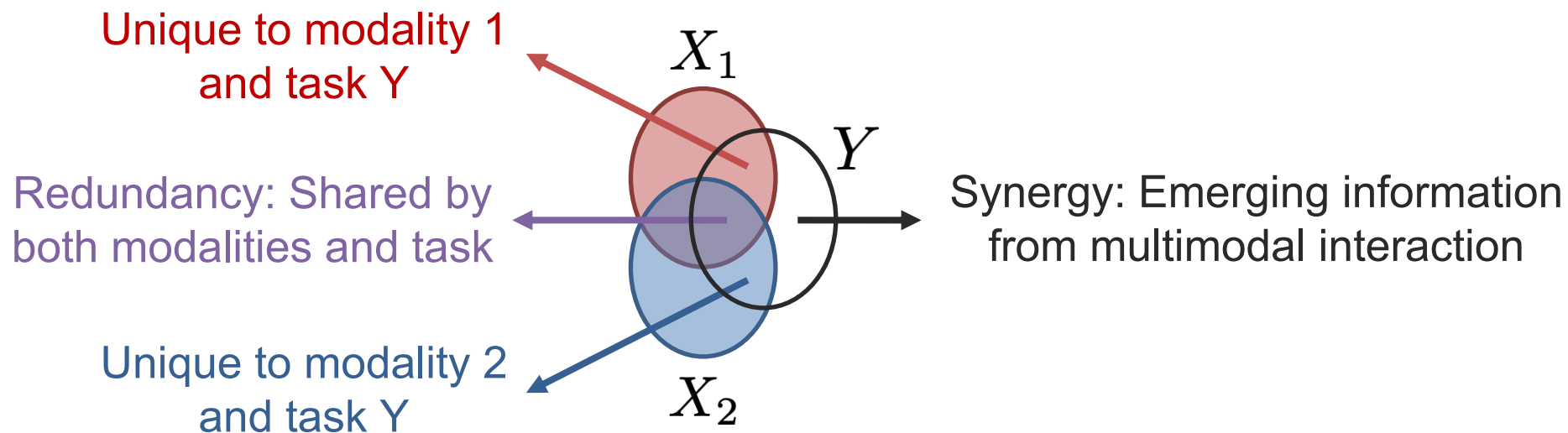
Multi-view redundancy may not hold for multimodal problems!

[Tian et al., What makes for Good Views for Contrastive Learning? NeurIPS 2020]
 [Tosh et al., Contrastive Learning, Multi-view Redundancy, and Linear models. ALT 2021]

Sub-Challenge 1c: Representation Fission



Definition: Learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering



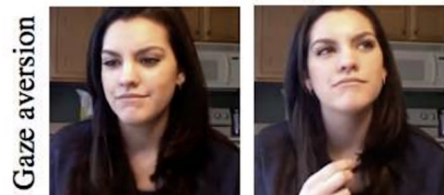
Quantifying Interactions

These interactions can be estimated efficiently:

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

I don't know maybe.

Vision:



(frustrated voice)

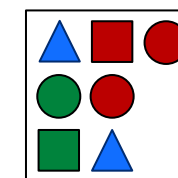
Acoustic:

Sheldon :

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

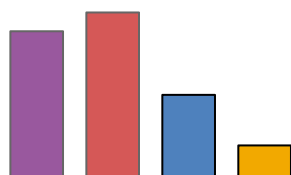


- **Text** : suggests a compliment.
- **Audio** : neutral tone.
- **Video** : straight face.



Is there a red shape above a circle?

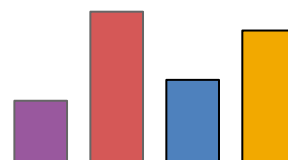
Sentiment



$R U_\ell U_{av} S$

Language/Agreement

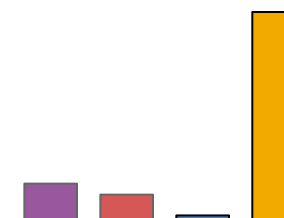
Sarcasm



$R U_\ell U_{av} S$

Multimodal Transformer

VQA



$R U_\ell U_i S$

Multiplicative/Transformer

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

Can also be used to choose most appropriate models – can they be used to better train/design new models?

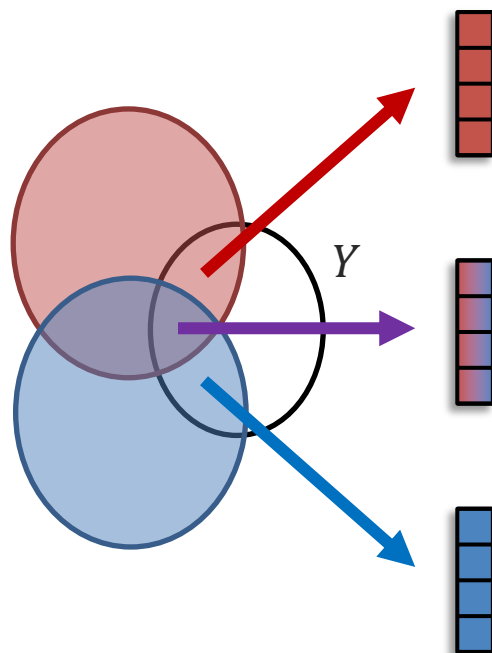
[Liang et al., Quantifying & Modeling Feature Interactions: An Information Decomposition Framework, NeurIPS 2023]

Factorized Contrastive Learning

Modeling task-relevant unique information



Can you please pass the cow?



- 2) Maximize task-relevant **unique** information

$$I(\mathbf{Z}; Y | \bullet)$$

- 1) Maximize task-relevant **shared** information

$$I(\mathbf{Z}; \bullet; Y) \text{ and } I(\mathbf{Z}; \blacktriangle; Y)$$

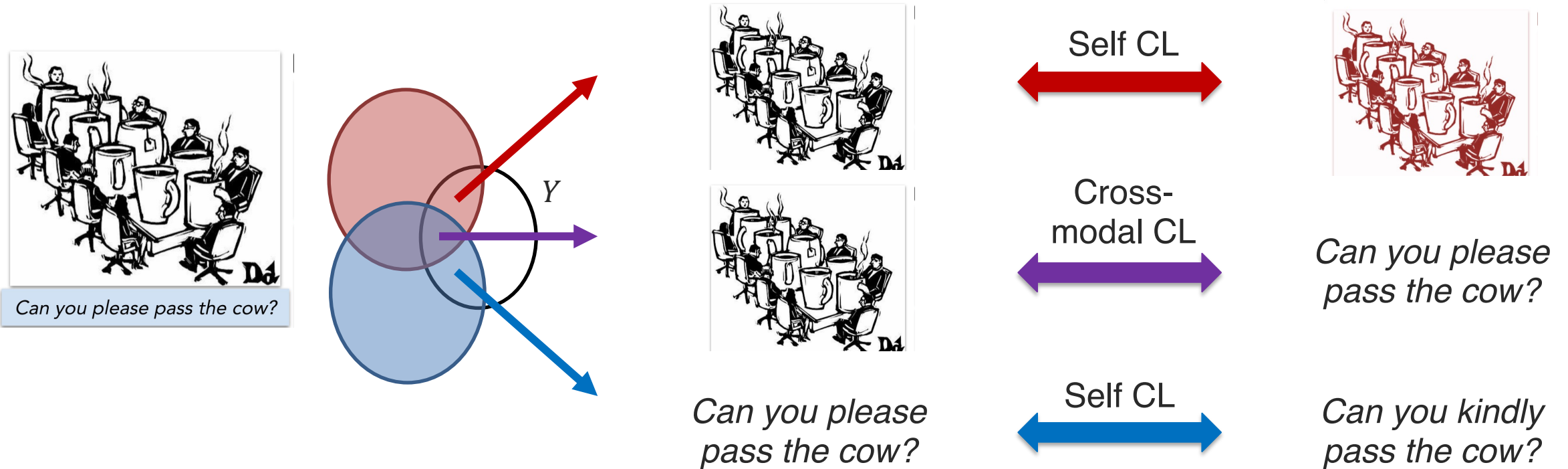
- 3) Maximize task-relevant **unique** information

$$I(\mathbf{Z}; Y | \blacktriangle)$$

[Liang et al., Factorized Contrastive Learning: Going Beyond Multi-view Redundancy, NeurIPS 2023]

Factorized Contrastive Learning

Modeling task-relevant unique information



Approximate task-relevance Y using multi-view data augmentations
New scalable lower and upper bounds on mutual information

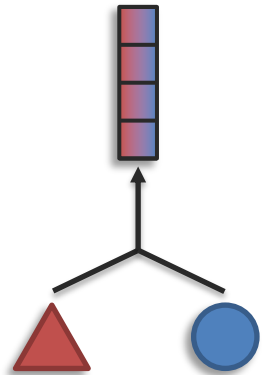
[Liang et al., Factorized Contrastive Learning: Going Beyond Multi-view Redundancy, NeurIPS 2023]

Challenge 1: Representation

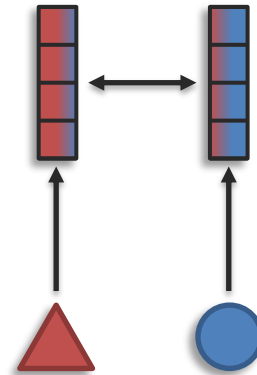
Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

Sub-challenges:

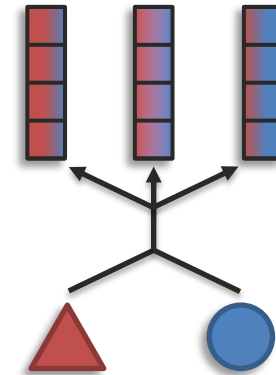
Fusion



Coordination



Fission



What is Multimodal?



Why is it hard?



What is next?

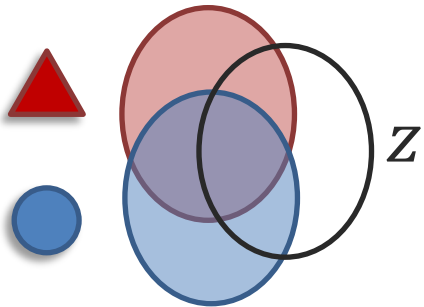
Heterogeneous



Connected



Interacting



Representation

Alignment

Reasoning

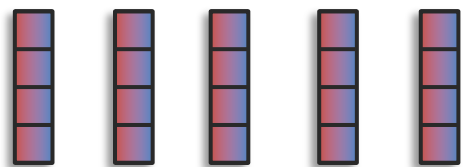
Generation

Transference

Quantification

Future Direction: Heterogeneity

Homogeneity



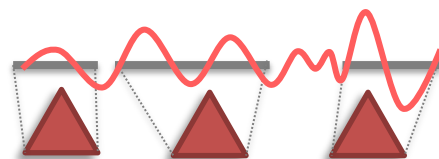
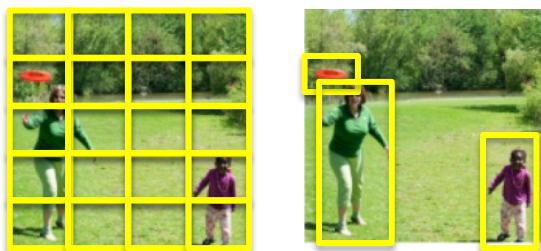
vs

Heterogeneity



Examples:

Arbitrary tokenization



Beyond differentiable interactions

Causal, logical, brain-inspired
Theoretical studies

MultiBench

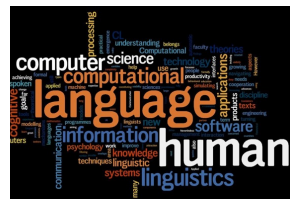
<https://github.com/pliang279/MultiBench>

Future Direction: High-modality

Few modalities



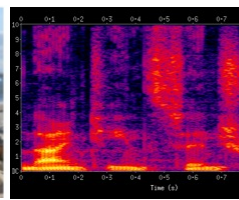
High-modality



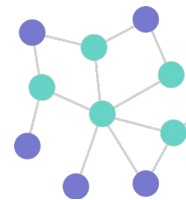
Language



Vision



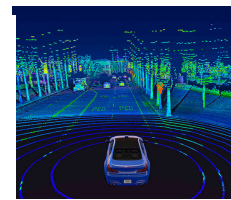
Audio



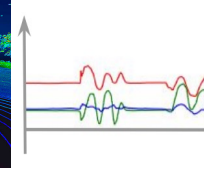
Graphs



Control



LIDAR



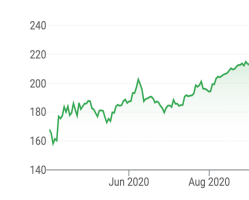
Sensors



Set

SUBJECT_ID
Age
Sex
Ethnicity
...

Table



Financial



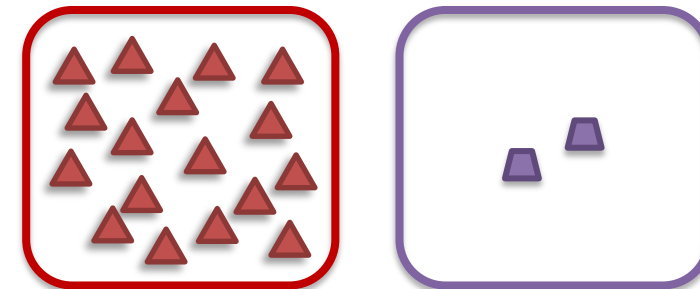
Medical

Examples:

Non-parallel learning

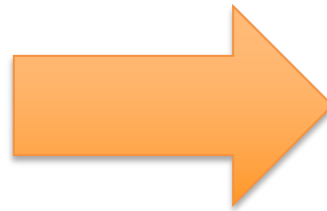
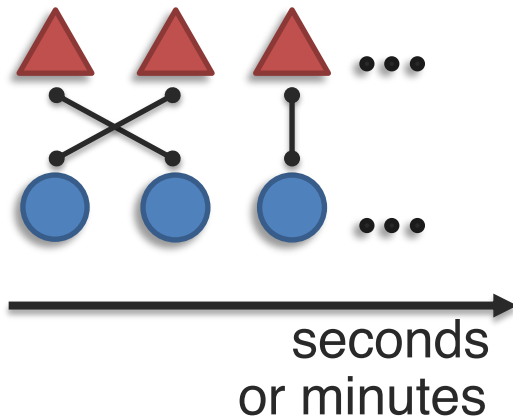


Limited resources



Future Direction: Long-term

Short-term



Long-term



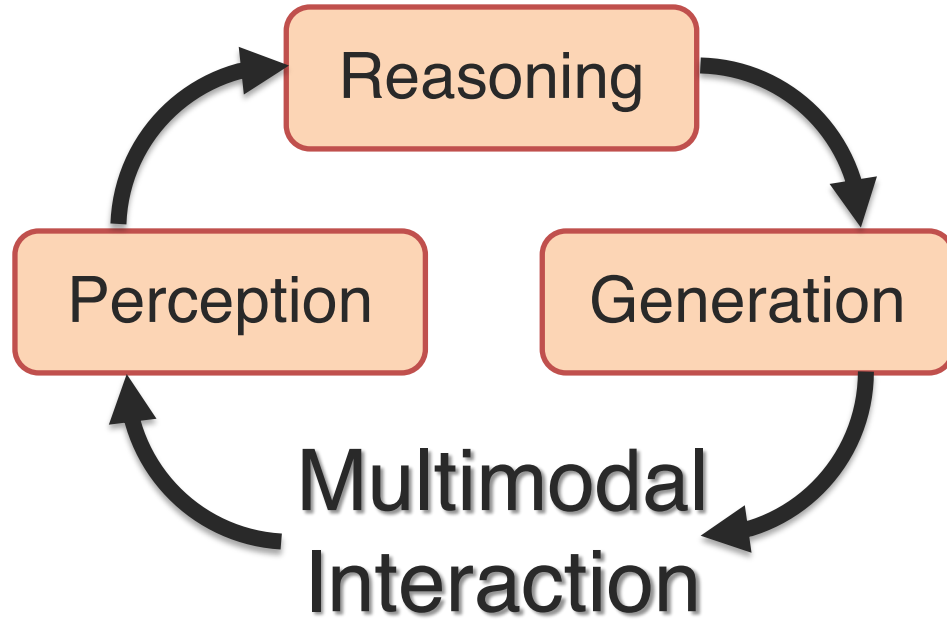
Examples:

Compositionality

Memory

Personalization

Future Direction: Interaction



Social Intelligence



Examples:

Multi-Party

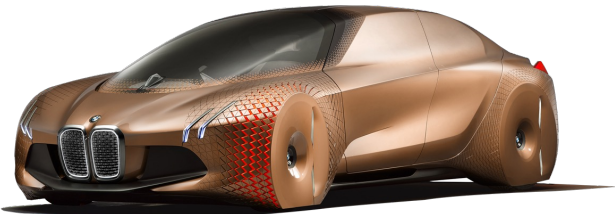
Causality

Ethical

Future Direction: Real-world



Healthcare
Decision Support



Intelligent Interfaces and
Vehicles



Online Learning
and Education

Examples:

Robustness

Fairness

Generalization

Interpretation

What is Multimodal?

Heterogeneous



Connected



Interacting



Why is it hard?

Representation

Alignment

Reasoning

Generation

Transference

Quantification



What is next?

High-modality

Heterogeneity

Long-term

Interaction

Real-world

<https://cmu-multicomp-lab.github.io/mmml-course/fall2023/>

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

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

pliang@cs.cmu.edu

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