



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal Machine Learning

Lecture 3.2: Multimodal Coordination and Fission

Paul Liang

** Co-lecturer: Louis-Philippe Morency. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yonatan Bisk.*

Objectives of today's class

- Representation coordination
 - Coordination functions
 - Kernel similarity functions
 - Canonical correlation analysis
 - Contrastive learning
 - Information, entropy and mutual information
- Representation fission
 - Factorized multimodal representations
 - Clustering and fine-grained fission

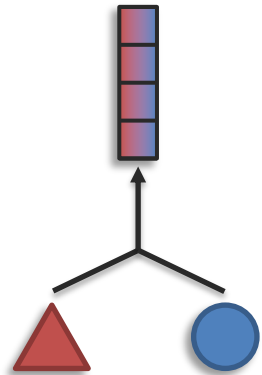
Multimodal Representation

Challenge 1: Representation

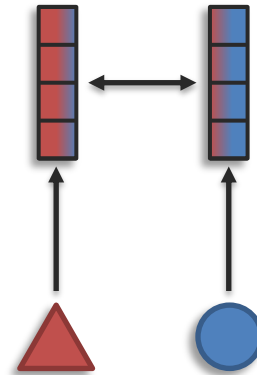
Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

Sub-challenges:

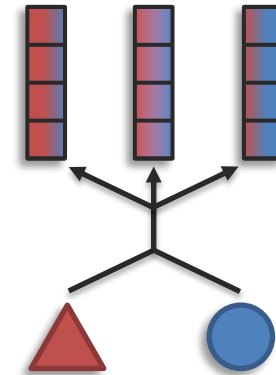
Fusion



Coordination

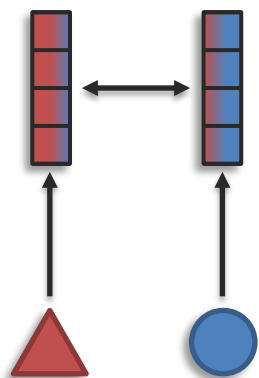


Fission



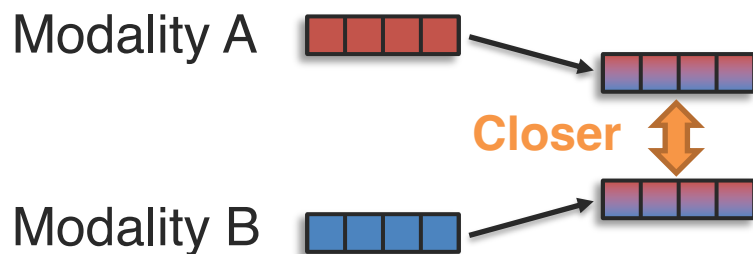
Representation Coordination

Sub-Challenge 1b: Representation Coordination

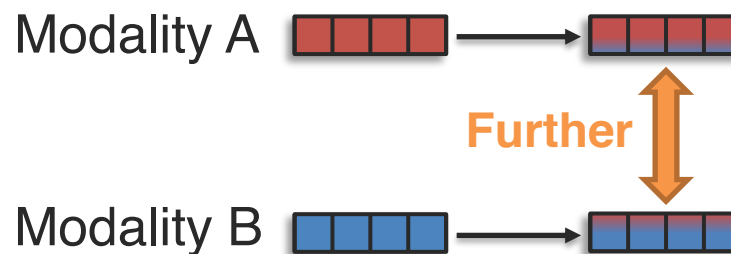


Definition: Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions

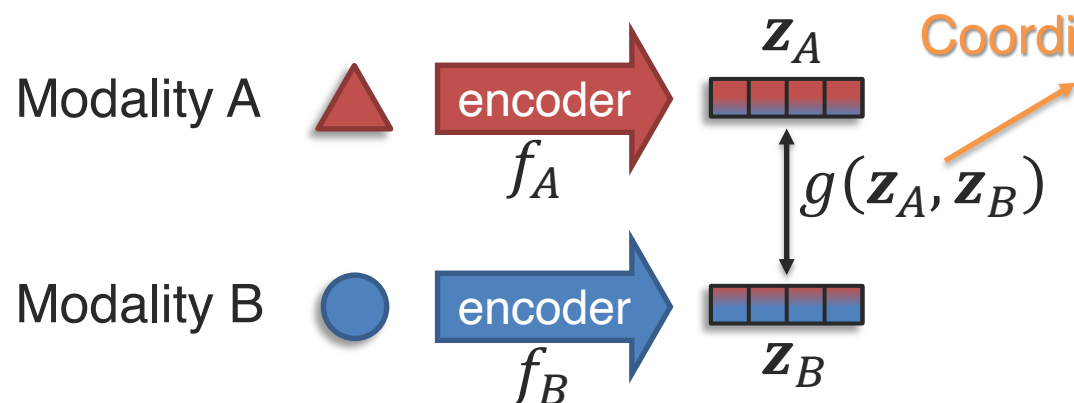
Strong Coordination:



Partial Coordination:



Coordination Function



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}

➡ Requires paired data

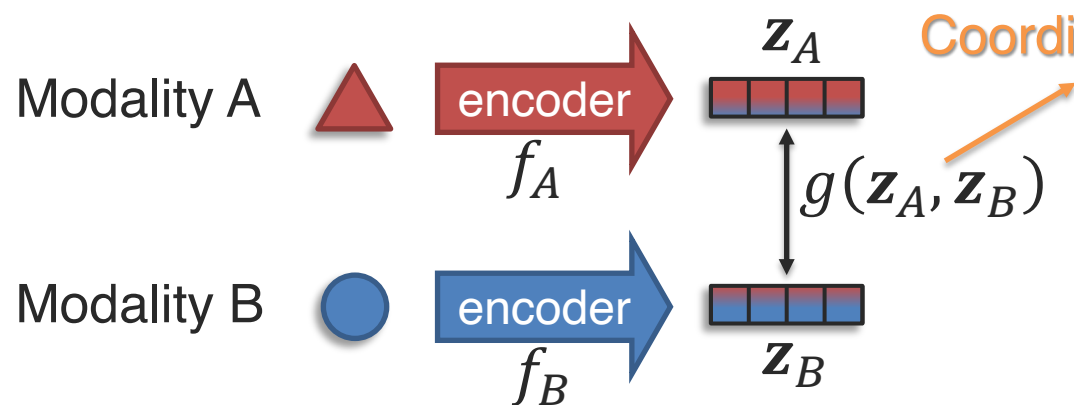
Examples of coordination function:

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

Strong coordination!

➡ For normalized inputs (e.g., $\mathbf{z}_A - \bar{\mathbf{z}}_A$), equivalent to *Pearson correlation coefficient*

Coordination Function



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}

Examples of coordination function:

② Kernel similarity functions:

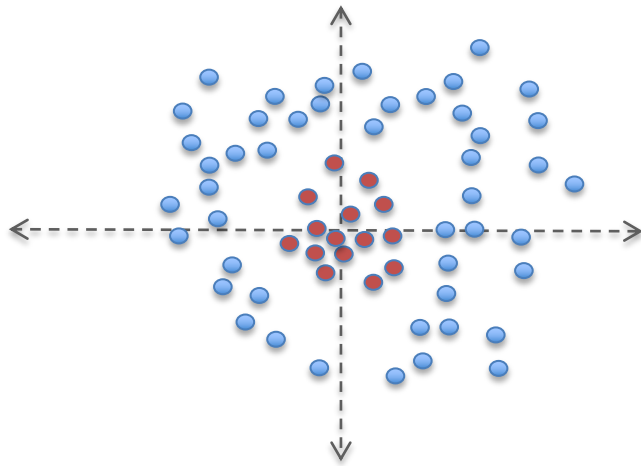
$$g(\mathbf{z}_A, \mathbf{z}_B) = k(\mathbf{z}_A, \mathbf{z}_B) \begin{cases} \bullet \text{ Linear} \\ \bullet \text{ Polynomial} \\ \bullet \text{ Exponential} \\ \bullet \text{ RBF} \end{cases}$$

➔ All these examples bring relatively strong coordination between modalities

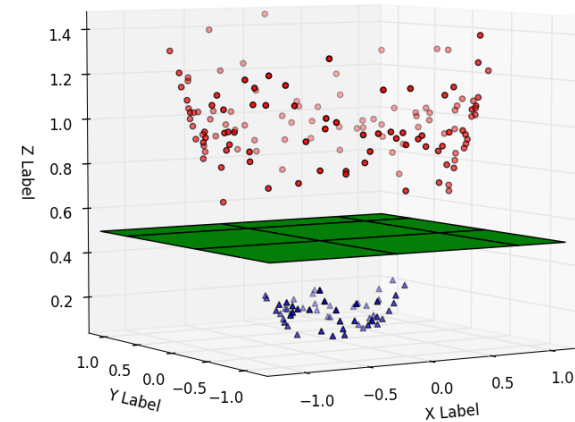
Kernel Function

A kernel function: Acts as a similarity metric between data points

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad \rightarrow \phi(\mathbf{x}) \text{ can be high-dimensional space!}$$



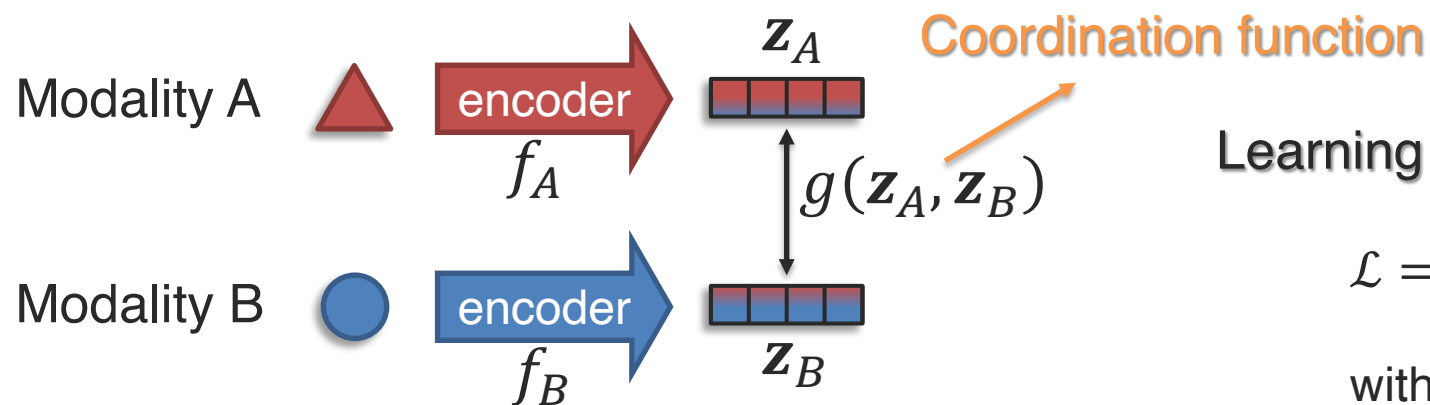
Not linearly separable in x space



Same data, but now linearly separable in $\phi(\mathbf{x})$ space

Radial Basis Function (RBF) Kernel : $K(\mathbf{x}_i, \mathbf{x}_j) = \exp -\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2$

Coordination Function



Learning with coordination function:

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

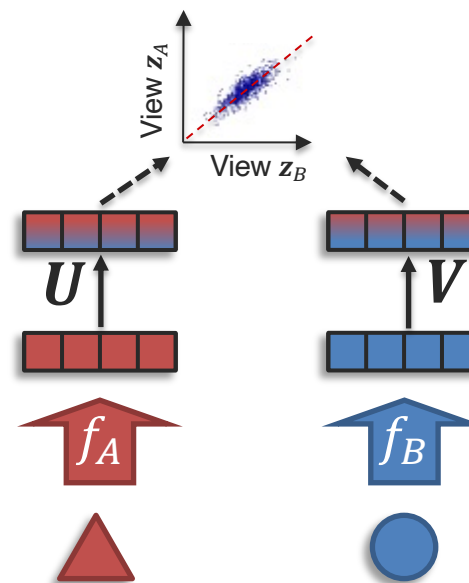
with model parameters θ_g , θ_{f_A} and θ_{f_B}

Examples of coordination function:

③ Canonical Correlation Analysis (CCA):

$$\operatorname{argmax}_{V, U, f_A, f_B} \operatorname{corr}(\mathbf{z}_A, \mathbf{z}_B)$$

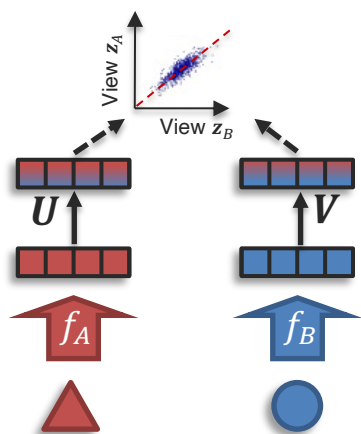
 CCA includes multiple projections, all orthogonal with each others



Correlated Projection

- 1 Learn two linear projections, one for each view, that are maximally correlated:

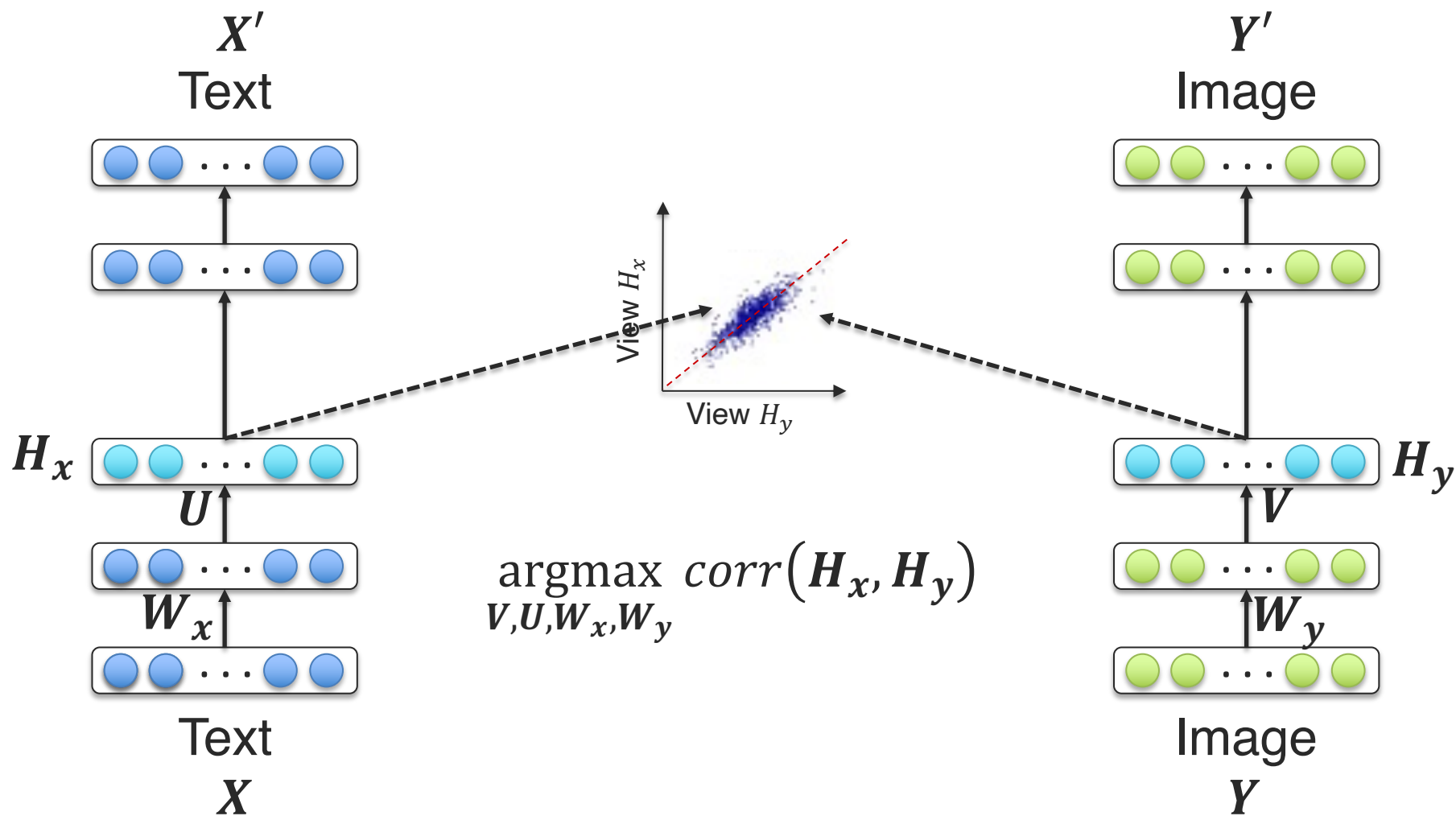
$$(\mathbf{u}^*, \mathbf{v}^*) = \operatorname{argmax}_{\mathbf{u}, \mathbf{v}} \operatorname{corr}(\mathbf{u}^T \mathbf{X}, \mathbf{v}^T \mathbf{Y})$$



Two views X, Y where same instances have the same color

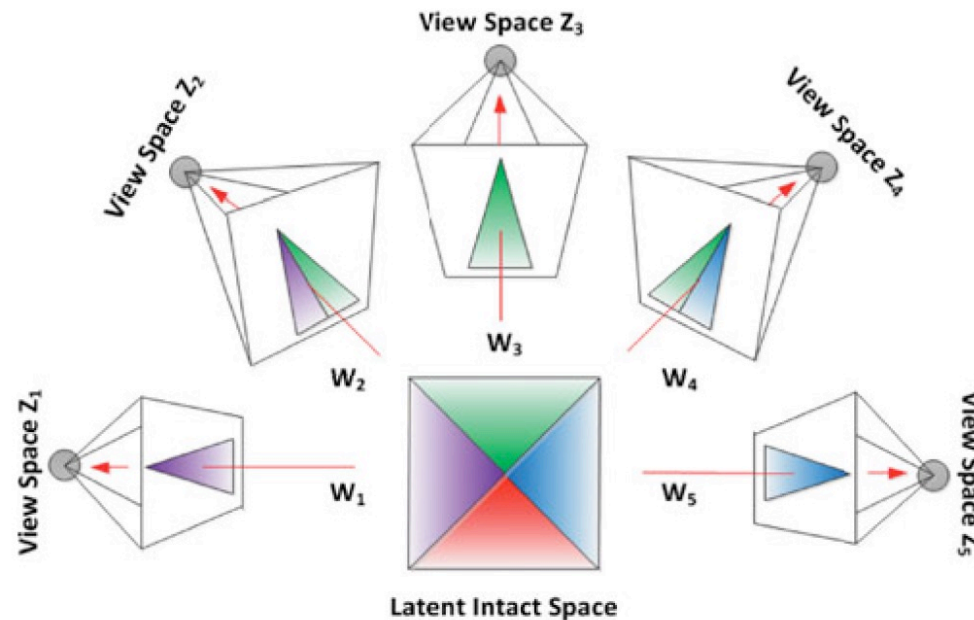
➡ Remember that X and Y consist of paired data

Deep Canonically Correlated Autoencoders (DCCA)



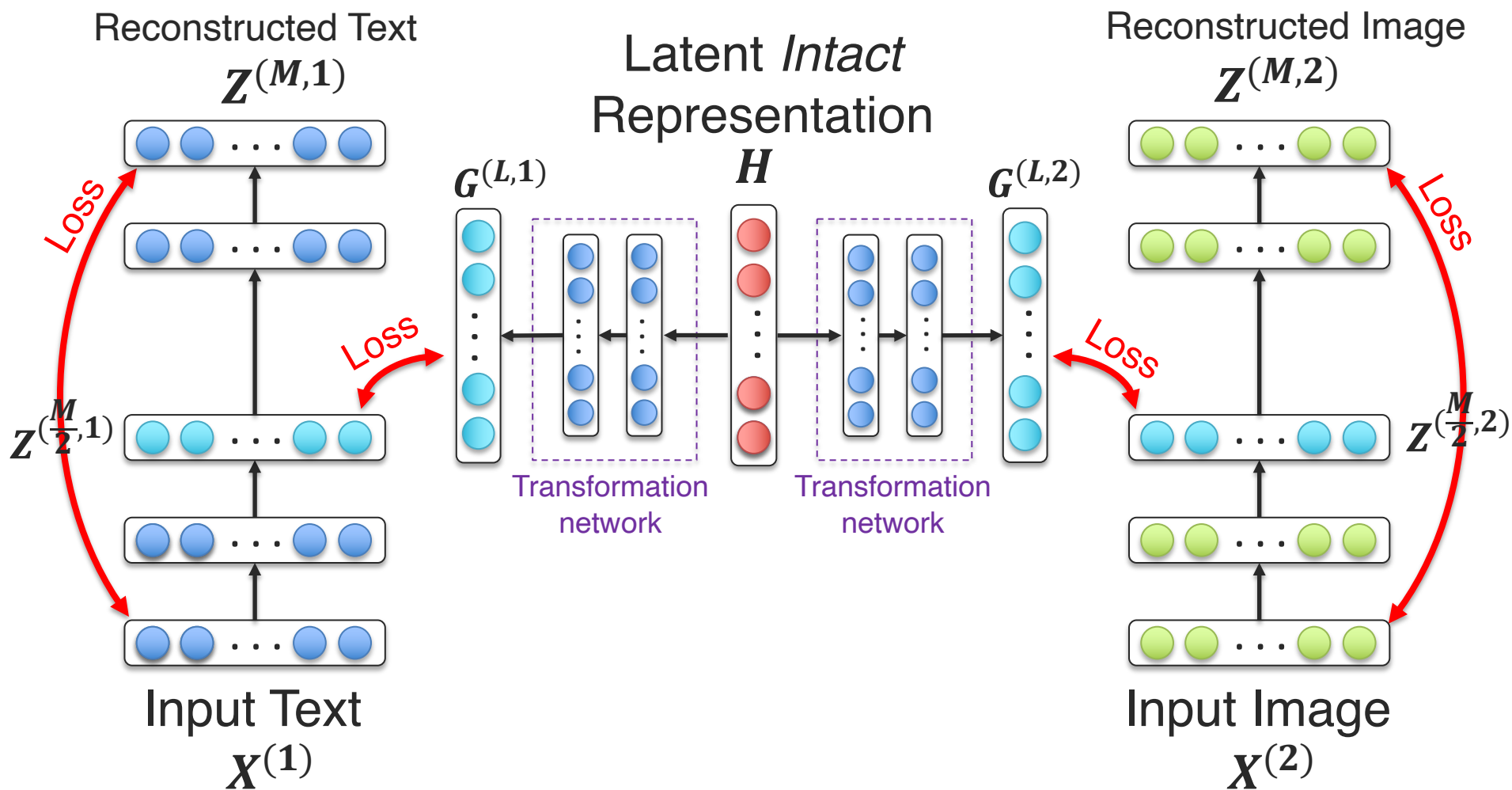
Multi-view Latent “Intact” Space

Given multiple views z_i from the same “object”:

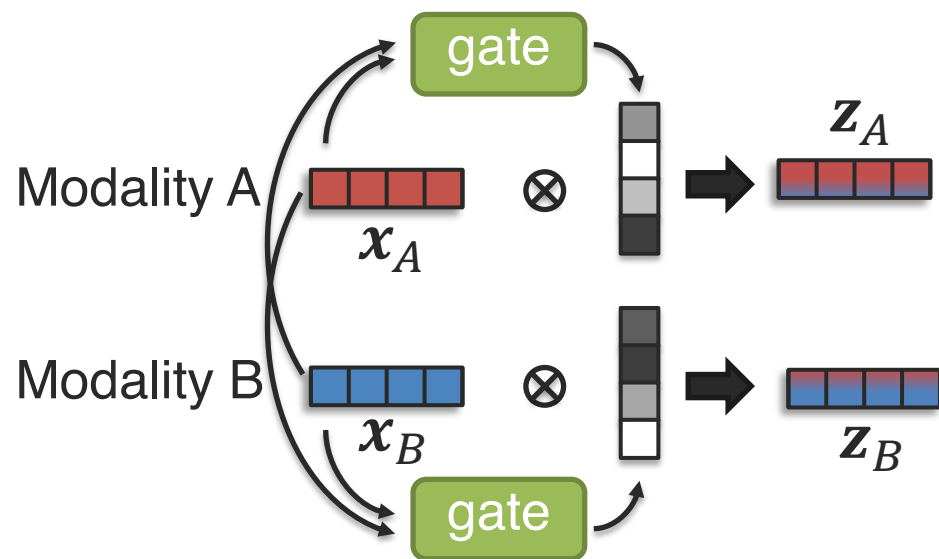


- 1) There is an “intact” representation which is *complete* and *not damaged*
- 2) The views z_i are partial (and possibly degenerated) representations of the intact representation

Auto-Encoder in Auto-Encoder Network



Gated Coordination



Gated coordination:

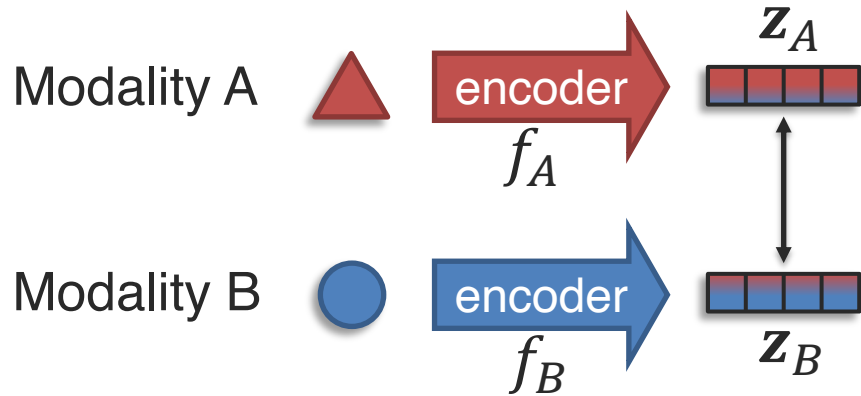
$$z_A = g_A(x_A, x_B) \cdot x_A$$

$$z_B = g_B(x_A, x_B) \cdot x_B$$

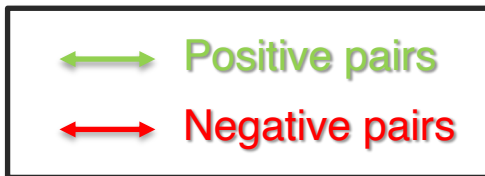
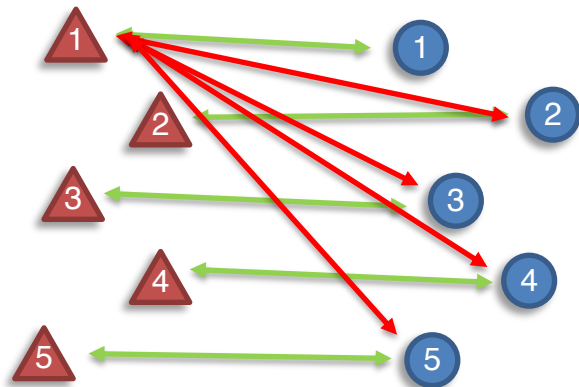
➡ Related to attention modules in transformers

More about it next week!

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

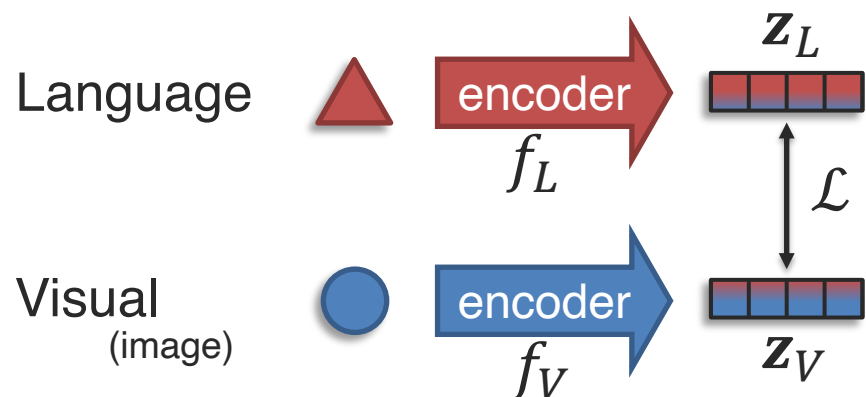
Simple contrastive loss:

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

Similarity functions are often cosine similarity

→ Similar to hinge loss

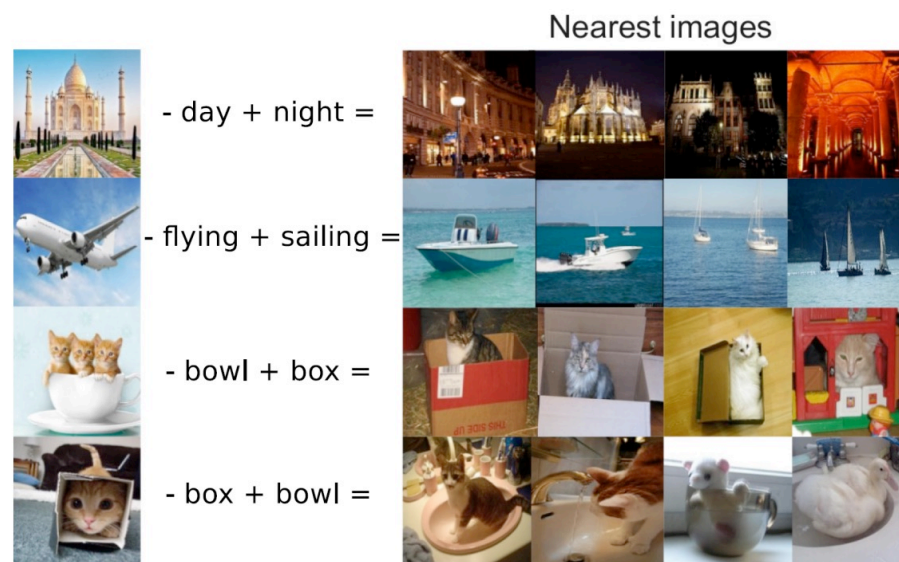
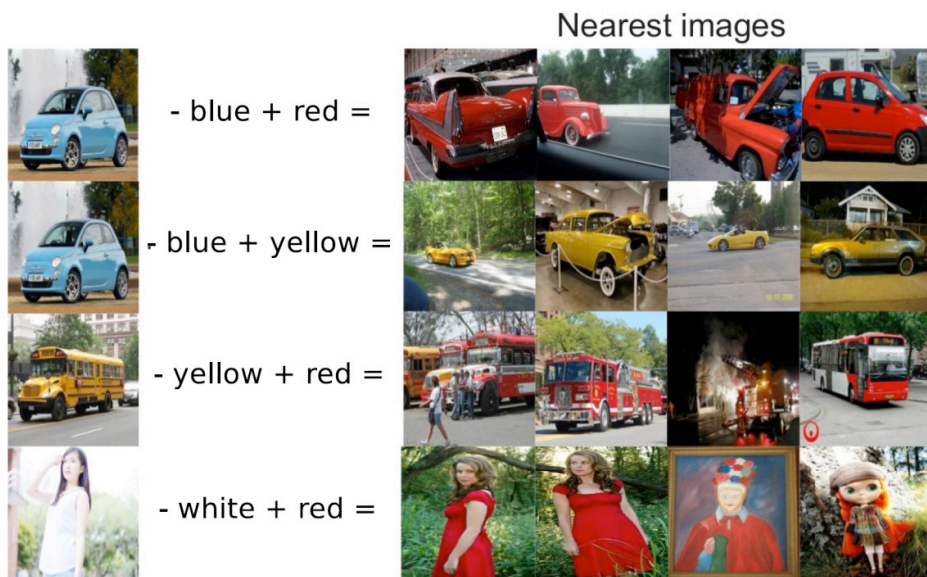
Example – Visual-Semantic Embeddings



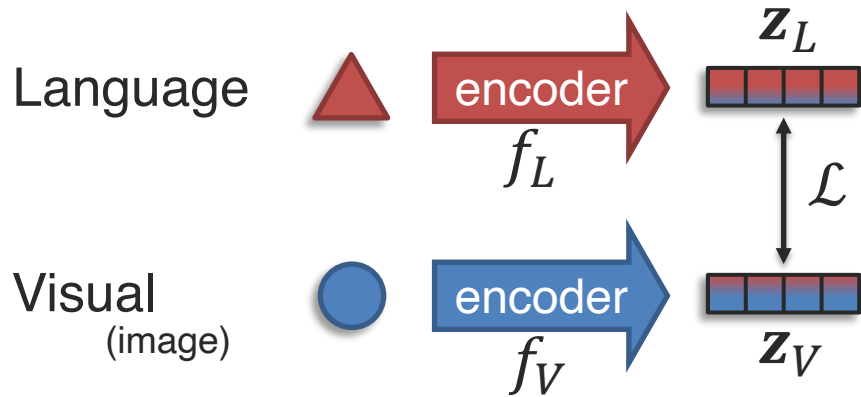
Two contrastive loss terms:

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

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



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



Popular contrastive loss: InfoNCE

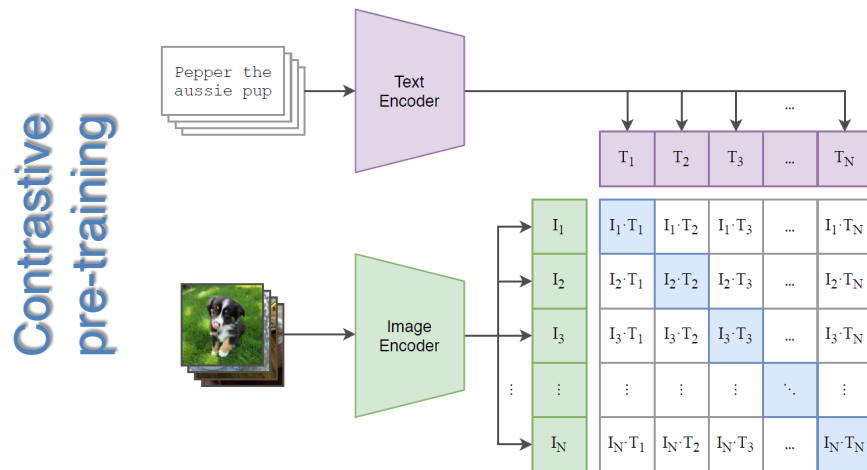
$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(z_A^i, z_B^i)}{\sum_{j=1}^N \text{sim}(z_A^i, z_B^j)}$$

Similarity function can be cosine similarity

positive pairs

negative pairs and positive pairs

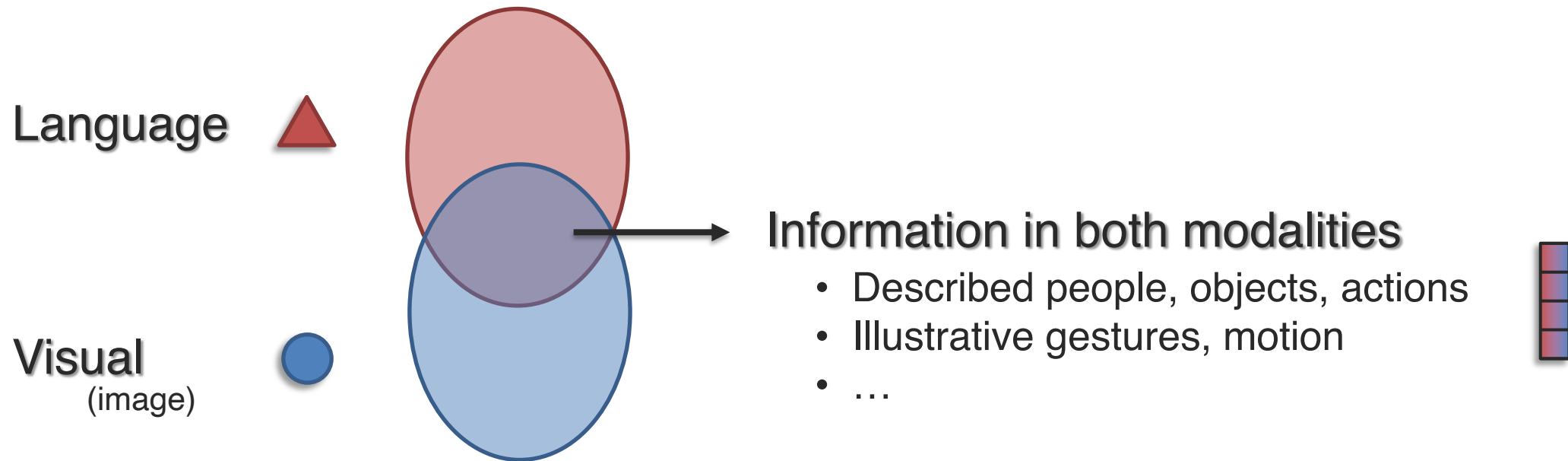
Positive and negative pairs:



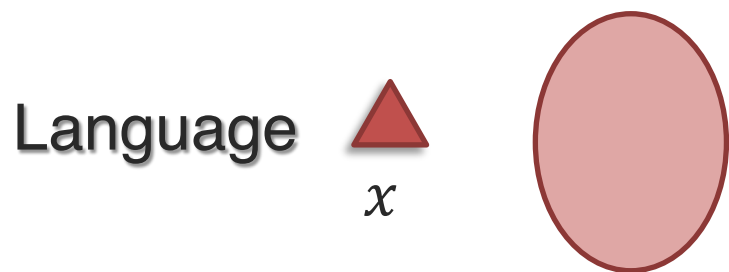
CLIP encoders (f_L and f_V) are great for language-vision tasks

z_L and z_V are coordinated but not identical representation spaces

Multimodal Coordination – Information Theory



Information and Entropy – Information Theory



How much information in the modality?

Information Theory (Shannon, 1948)

Main intuition: “Information value” of a communicated message x depends on how random its content is

x : “1,1,1,1,1,1,1,1,1,1,1,1,1”

➔ Not very random... So, low information

x : “0,1,0,1,0,0,1,1,1,0,0,1”

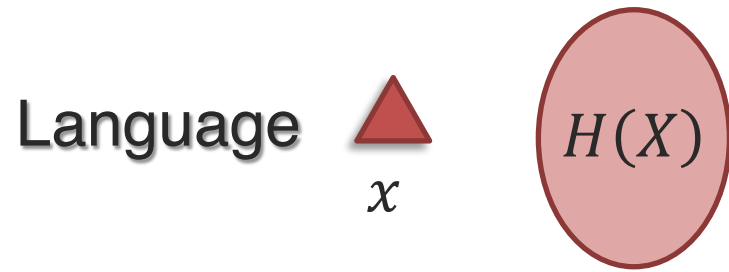
➔ More random... So, higher information

Information content $I(x)$

$$I(x) \sim \frac{1}{p(x)}$$

$$I(x) = \log\left(\frac{1}{p(x)}\right) = -\log(p(x))$$

Information and Entropy – Information Theory



How much information in the modality?

Information Theory (Shannon, 1948)

Information content $I(X) = -\log(p(X))$

➔ For discrete alphabet \mathcal{X} , then X is discrete random variable

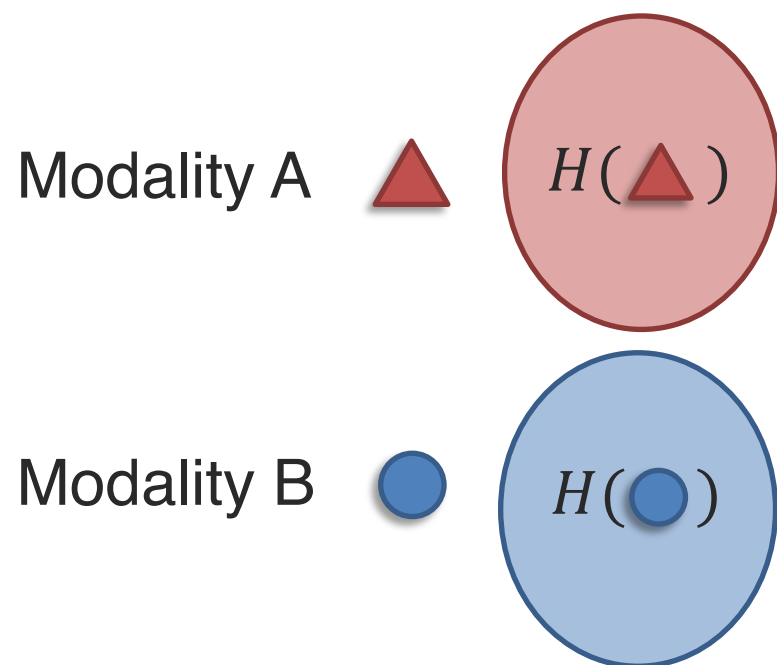
Entropy: weighted average of all possible outcomes from \mathcal{X}

$$H(X) = \mathbb{E}[I(X)] = \mathbb{E}[-\log(p(X))] = - \sum_{x \in \mathcal{X}} p(X) \log(p(X))$$

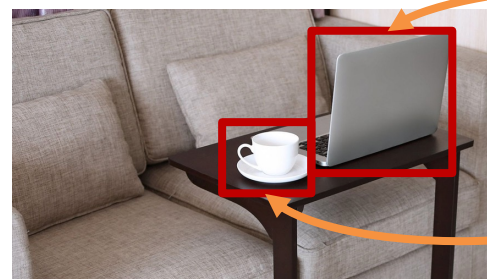
➔ Entropy can also be defined for continuous random variables

Entropy with Two Modalities

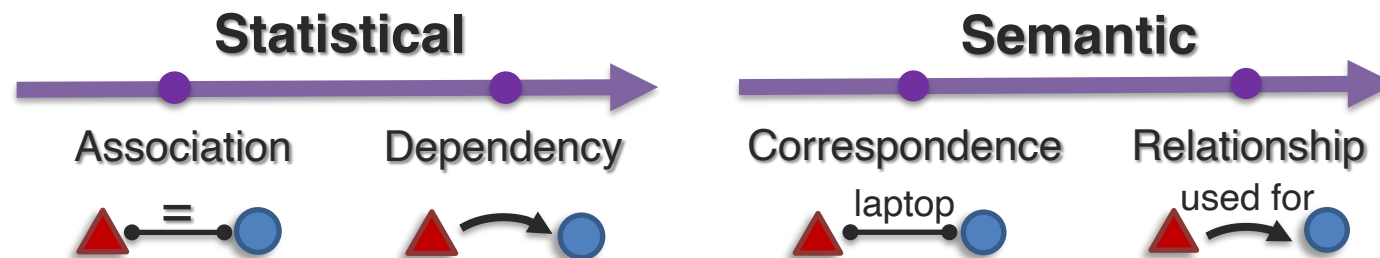
If no overlapping information



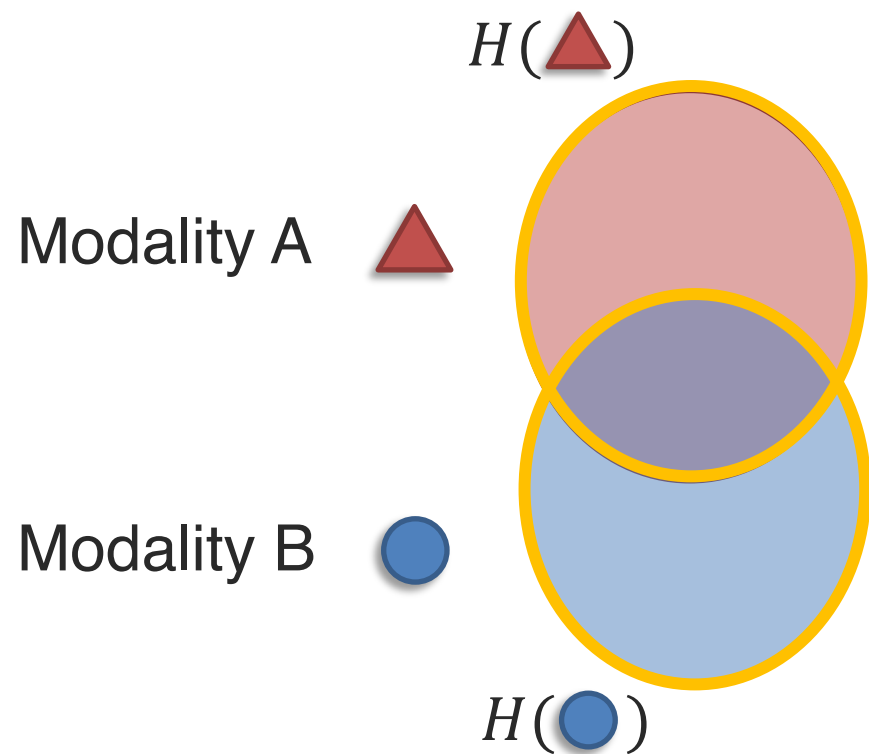
➔ But in most real-world scenarios, modalities are *inter-connected*



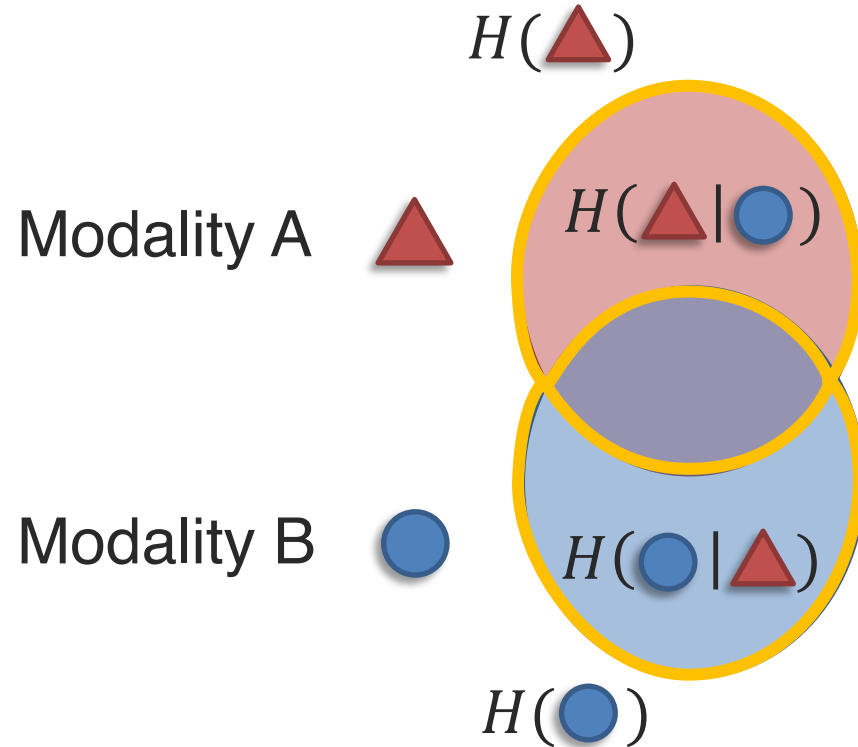
A **teacup** on the right of a **laptop** in a clean room.



Entropy with Two Modalities



Entropy with Two Modalities



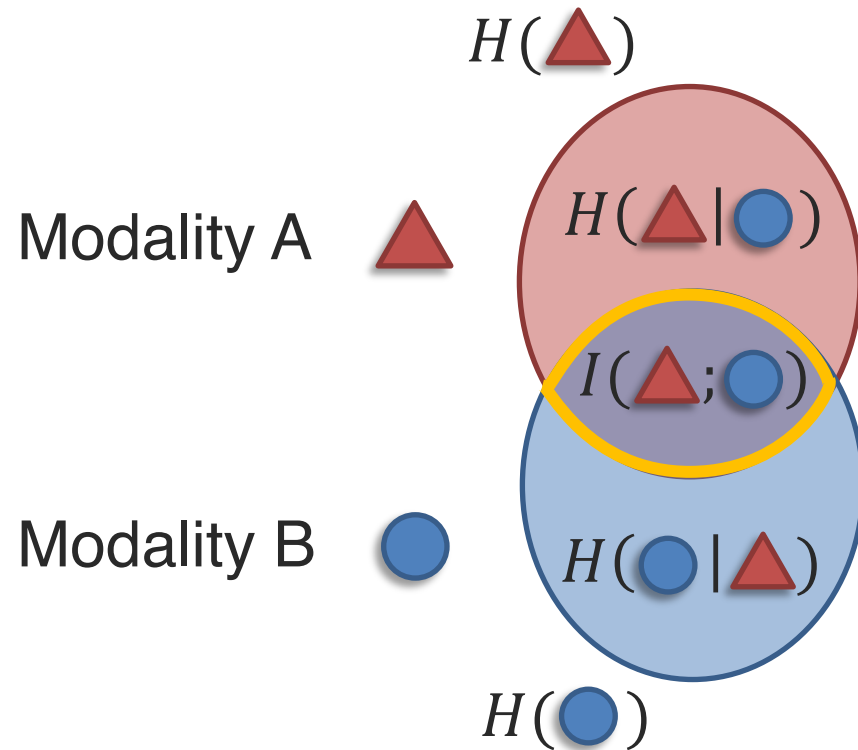
Conditional entropy $H(Y|X)$

$$H(Y|X) = -\mathbb{E}_{X,Y}[\log p(y|x)]$$

$$= -\mathbb{E}_{X,Y} \left[\log \frac{p(x,y)}{p(x)} \right]$$

If X and Y independent, $H(Y|X) = H(Y)$.
If X fully determines Y , then $H(Y|X) = 0$.

Entropy with Two Modalities



Mutual information $I(X; Y)$

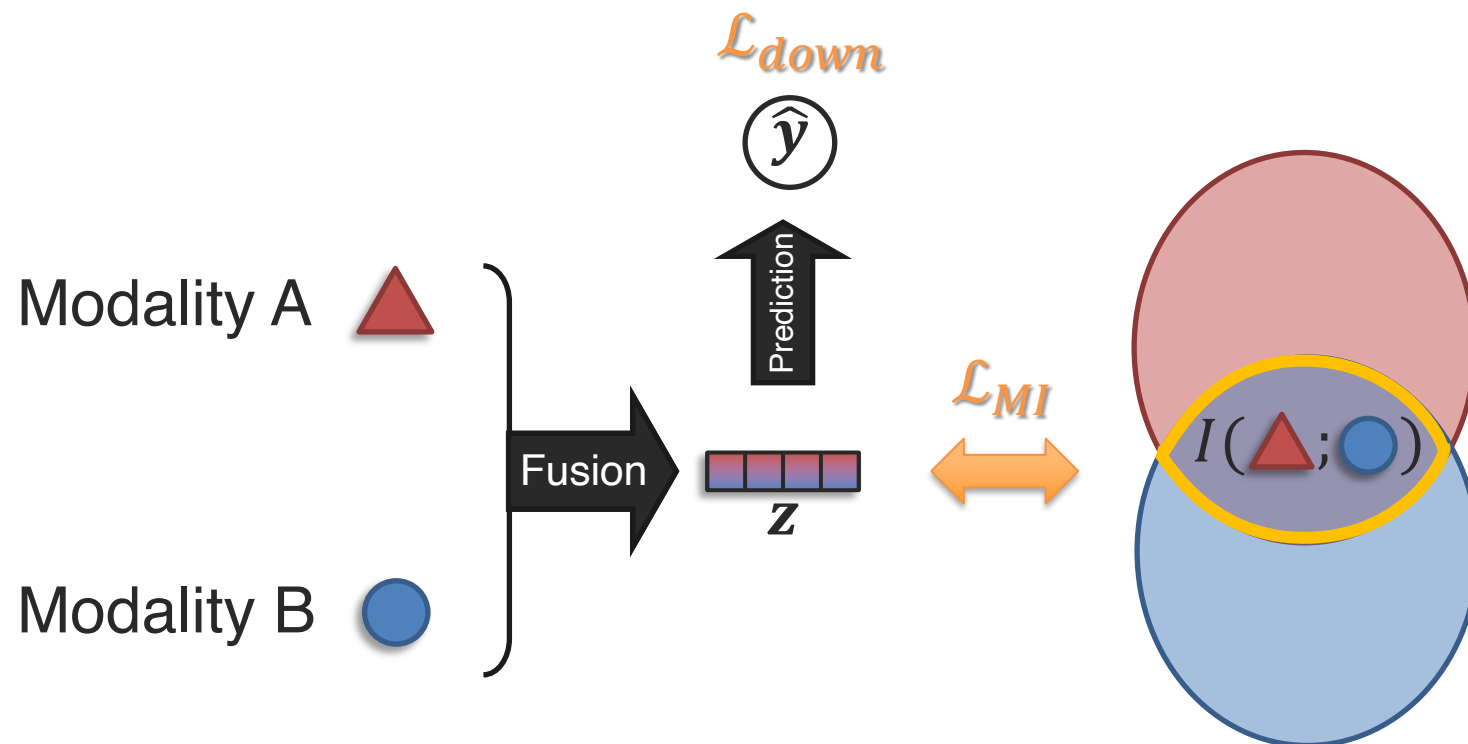
$$I(X; Y) = H(X) - H(X|Y)$$

$$= \mathbb{E}_{X,Y} \left[\log \frac{1}{P_X(x)} + \log \frac{P_{XY}(x, y)}{P_Y(y)} \right]$$

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

using KL-divergence $\leftarrow I(X; Y) = D_{KL}(P_{XY}(x, y) \parallel P_X(x)P_Y(y))$

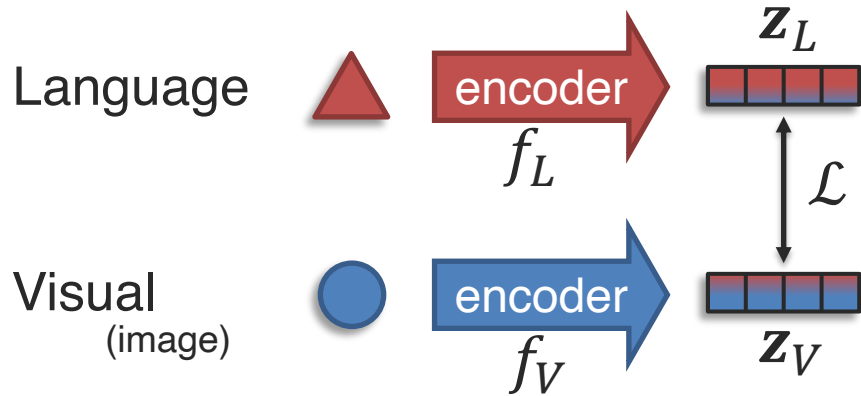
Multimodal Fusion with Mutual Information



Assumption?

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

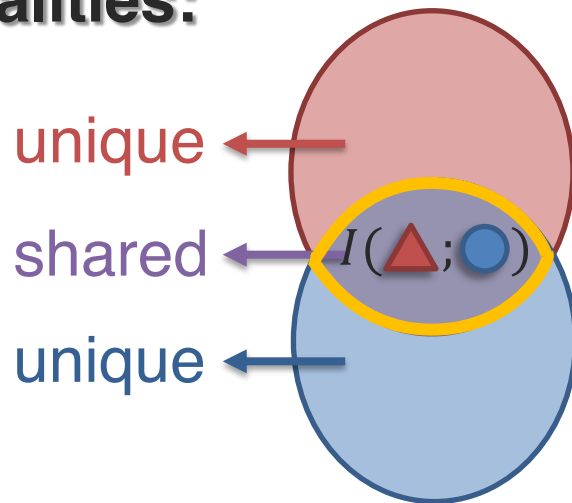
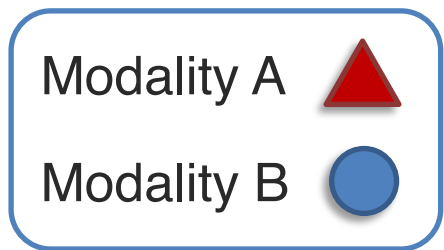
Contrastive Learning and Connected Modalities



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:

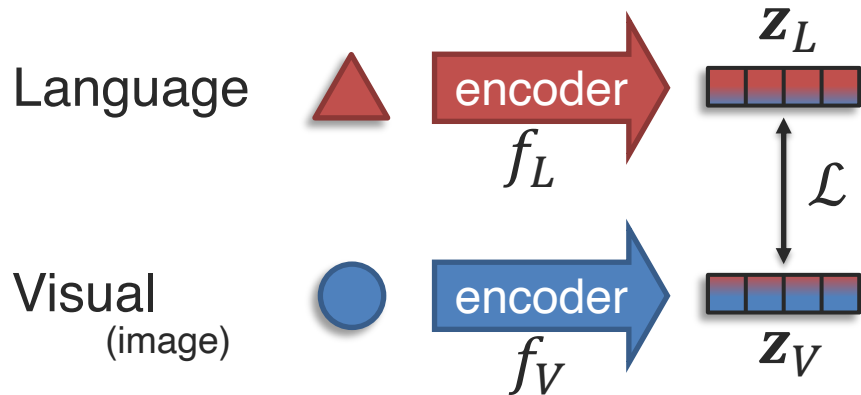


CLIP focuses on shared connections

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:

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

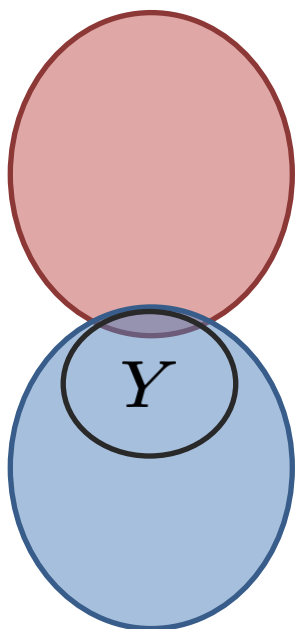
Multiview Redundancy and Contrastive Learning



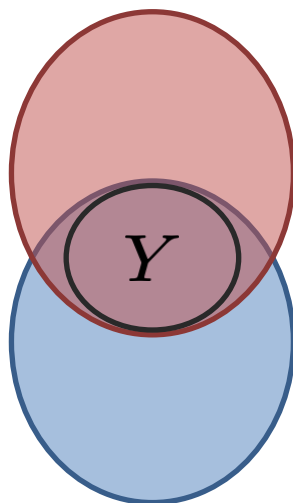
How much information should be shared?

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

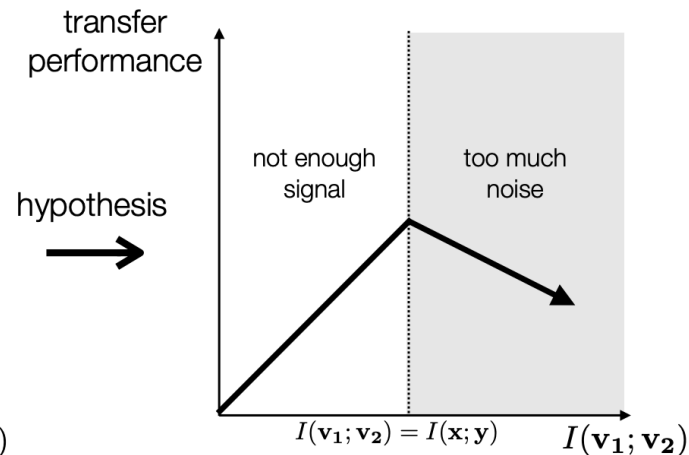
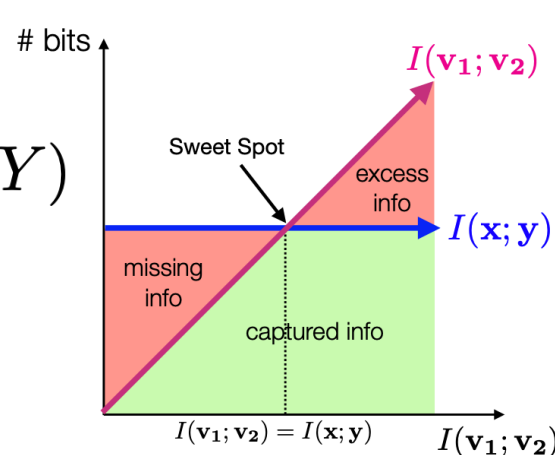
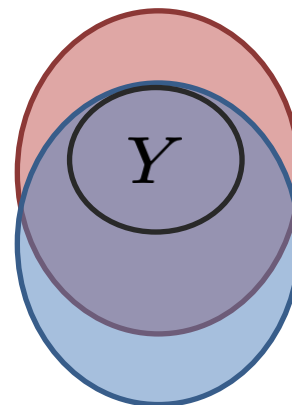
Not enough signal



Just right



Too much noise



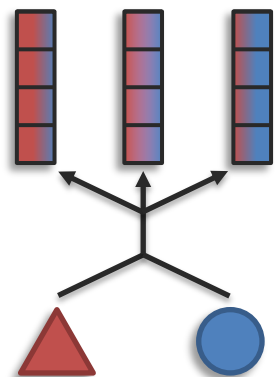
Multi-view redundancy may not hold for multimodal problems!

[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]

Representation

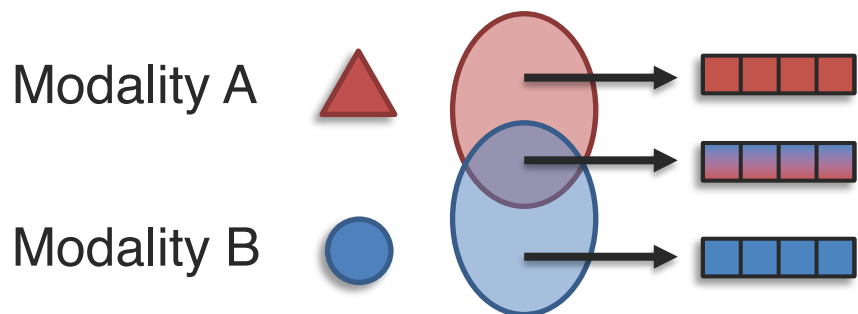
Fission

Sub-Challenge 1c: Representation Fission

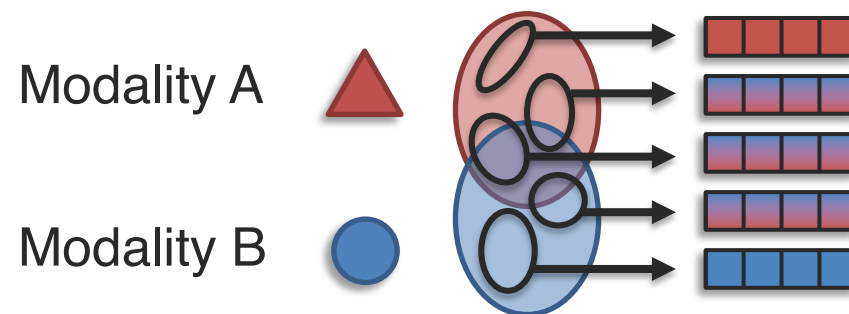


Definition: Learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering

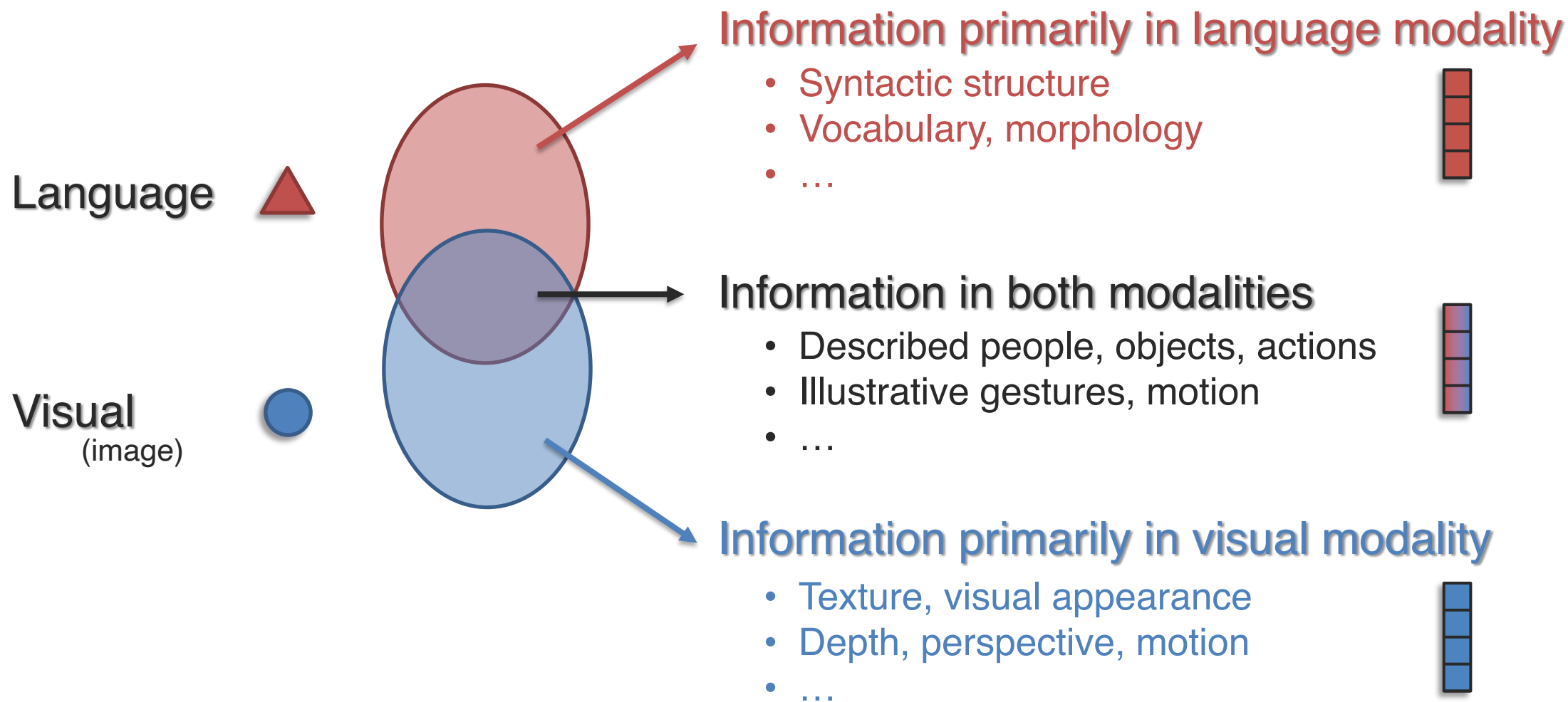
Modality-level fission:



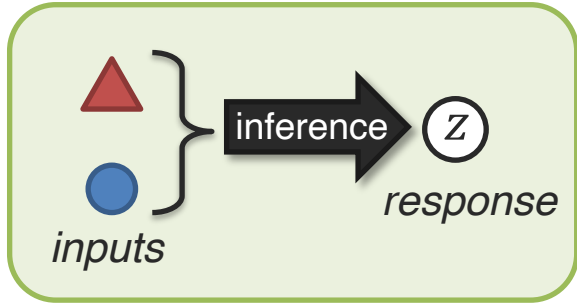
Fine-grained fission:



Modality-Level Fission



Recall Taxonomy of Interactions



Multimodal Communication



Redundancy

signal response

a → □

b → □

signal response

a+b → □

a+b → □

Equivalence

Enhancement

Nonredundancy

a → □

b → ○

a+b → □ and ○

a+b → □

a+b → □ (or □)

a+b → △

Independence

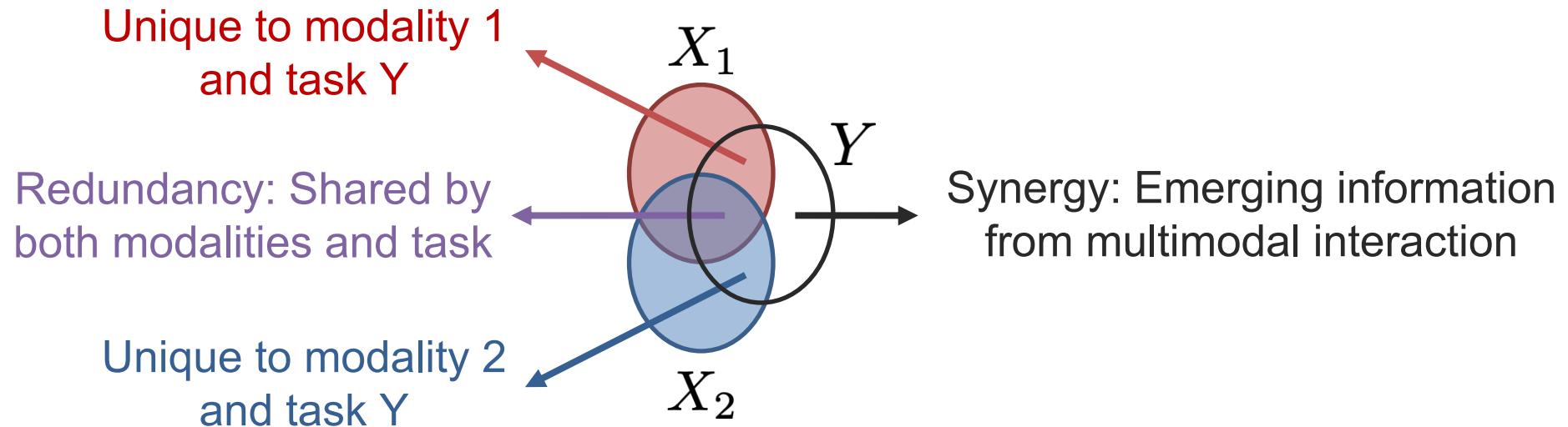
Dominance

Modulation

Emergence

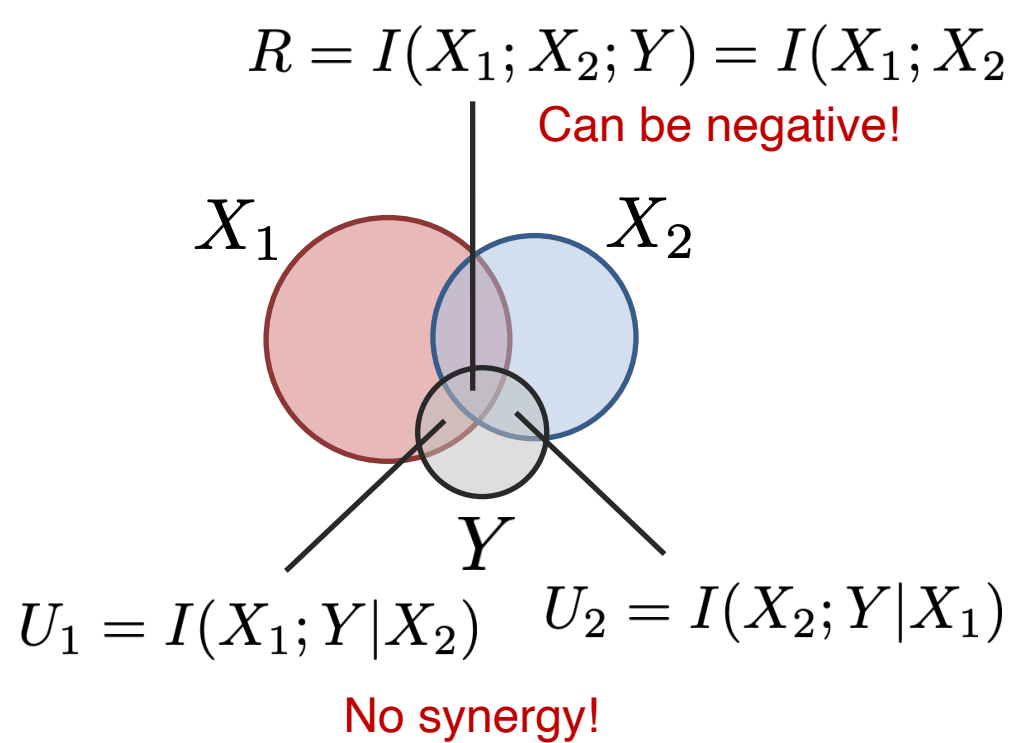
Partan and Marler (2005). *Issues in the classification of multimodal communication signals*. *American Naturalist*, 166(2)

Representation Fission via Information Theory

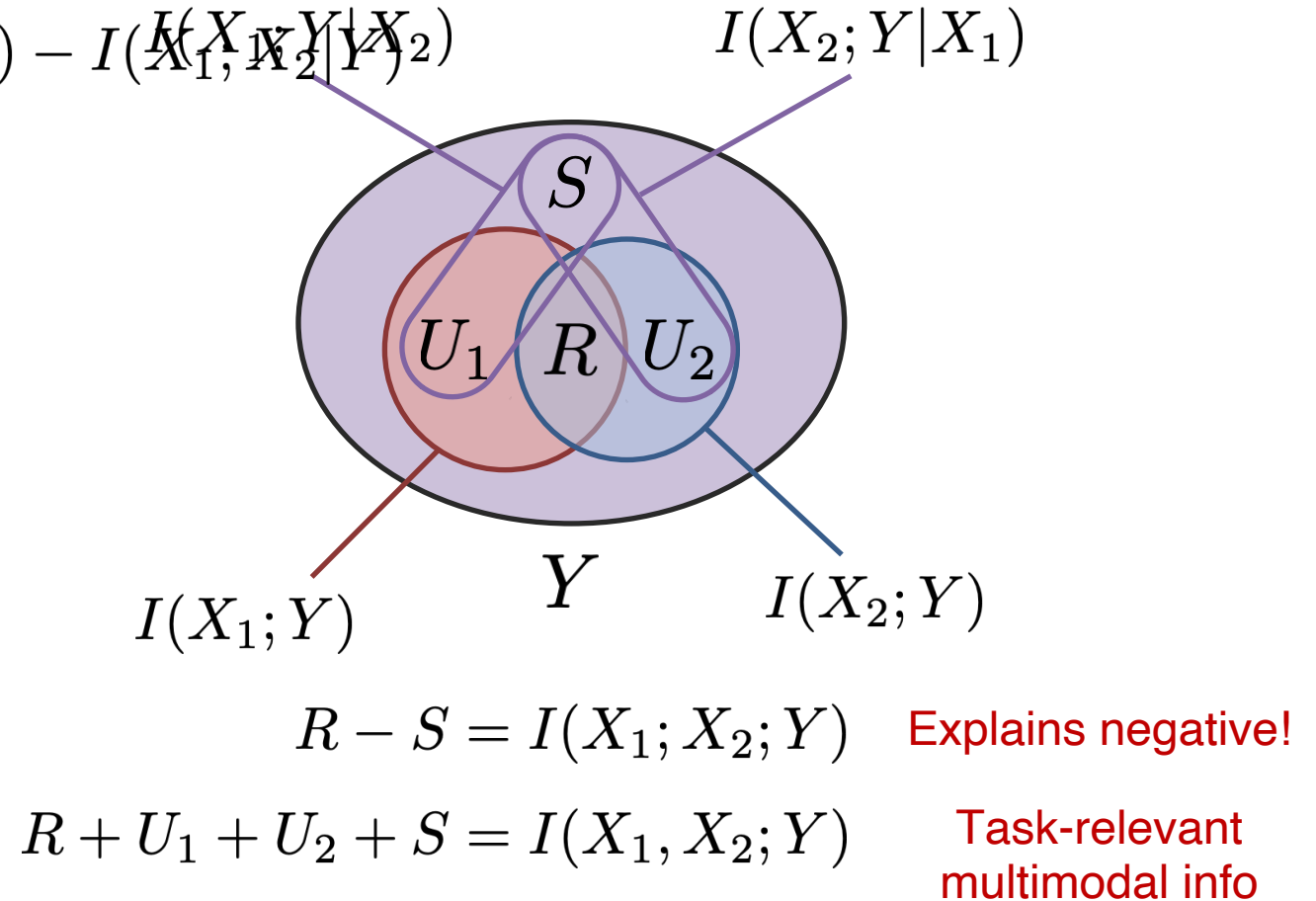


Partial Information Decomposition

Classical Information Theory



Partial Information Decomposition

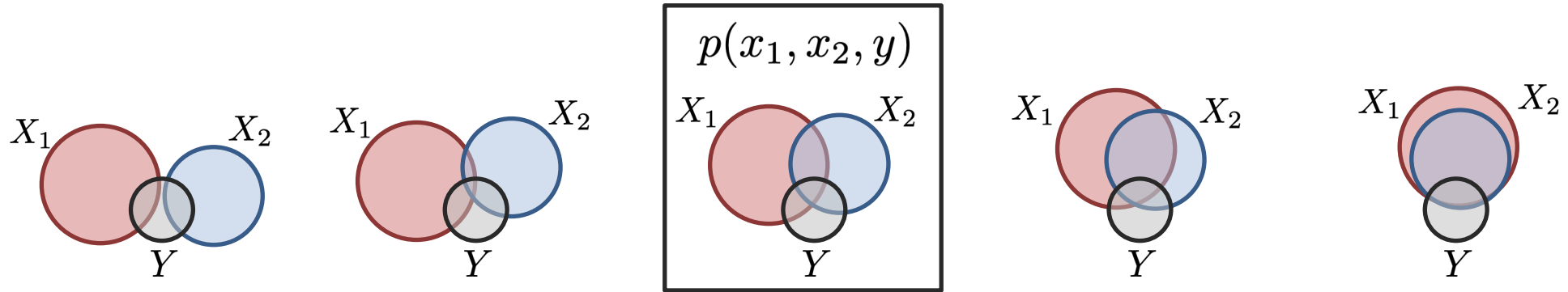


Partial Information Decomposition

One type of information decomposition

Unimodal marginal-matching distributions:

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



$$S = \underbrace{I_p(X_1, X_2; Y)}_{\text{Task-relevant multimodal info}} - \underbrace{\min_{q \in \Delta_p} I_q(X_1, X_2; Y)}_{\text{Task-relevant multimodal info without synergy}}$$

Task-relevant multimodal info Task-relevant multimodal info without synergy:

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

Partial Information Decomposition

One type of information decomposition

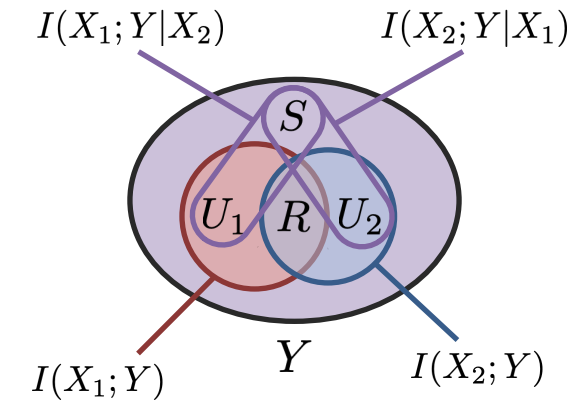
Unimodal marginal-matching distributions:

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

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

+ consistency equations relating interactions with information theory:

Only need unimodal marginals to infer redundancy and uniqueness:



$$R = \max_{q \in \Delta_p} I_q(X_1; X_2; Y) \quad U_1 = \min_{q \in \Delta_p} I_q(X_1; Y|X_2) \quad U_2 = \min_{q \in \Delta_p} I_q(X_2; Y|X_1)$$

Can be solved efficiently as a convex optimization problem

Scales to high-dimensional continuous modalities via neural networks

Quantifying Interactions

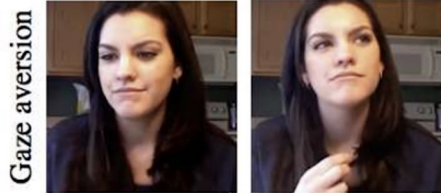


These interactions can be efficiently estimated – gives a path towards understanding interactions

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

I don't know maybe.

Vision:



(frustrated voice)

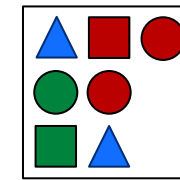
Acoustic:

Sheldon :

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

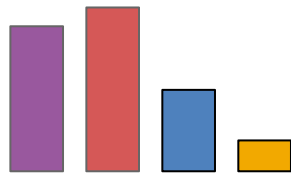


- **Text** : suggests a compliment.
- **Audio** : neutral tone.
- **Video** : straight face.



Is there a red shape above a circle?

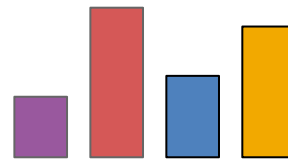
Sentiment



$R U_\ell U_{av} S$

Language/Agreement

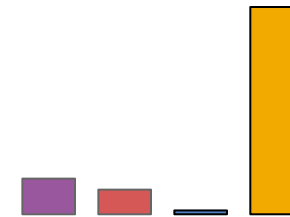
Sarcasm



$R U_\ell U_{av} S$

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. arXiv 2023]

Quantifying Interactions



Lower and upper bounds for interactions in a semi-supervised setting: $p(x_1, y), p(x_2, y), p(x_1, x_2)$

Idea 2: min-entropy couplings

Efficient approximation algorithms
[Cicalese et al., 2002, Compton 2022]

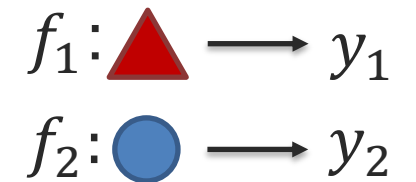
Upper bound:
$$\bar{S} = c_2 - \min_{r \in \Delta_p} H_r(X_1, X_2, Y) - \min_{q \in \Delta_p} I_q(X_1, X_2; Y)$$

$$\bar{S} = \max_{r \in \Delta_p} I_r(X_1, X_2; Y) - \min_{q \in \Delta_p} I_q(X_1, X_2; Y)$$

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

Lower bound:
$$\underline{S} = \underbrace{\alpha(f_1, f_2)}_{\text{Task-relevant multimodal info}} \cdot c_1 - \underbrace{\max(U_1, U_2)}_{\text{Task-relevant multimodal info without synergy}}$$

Idea 1: disagreement



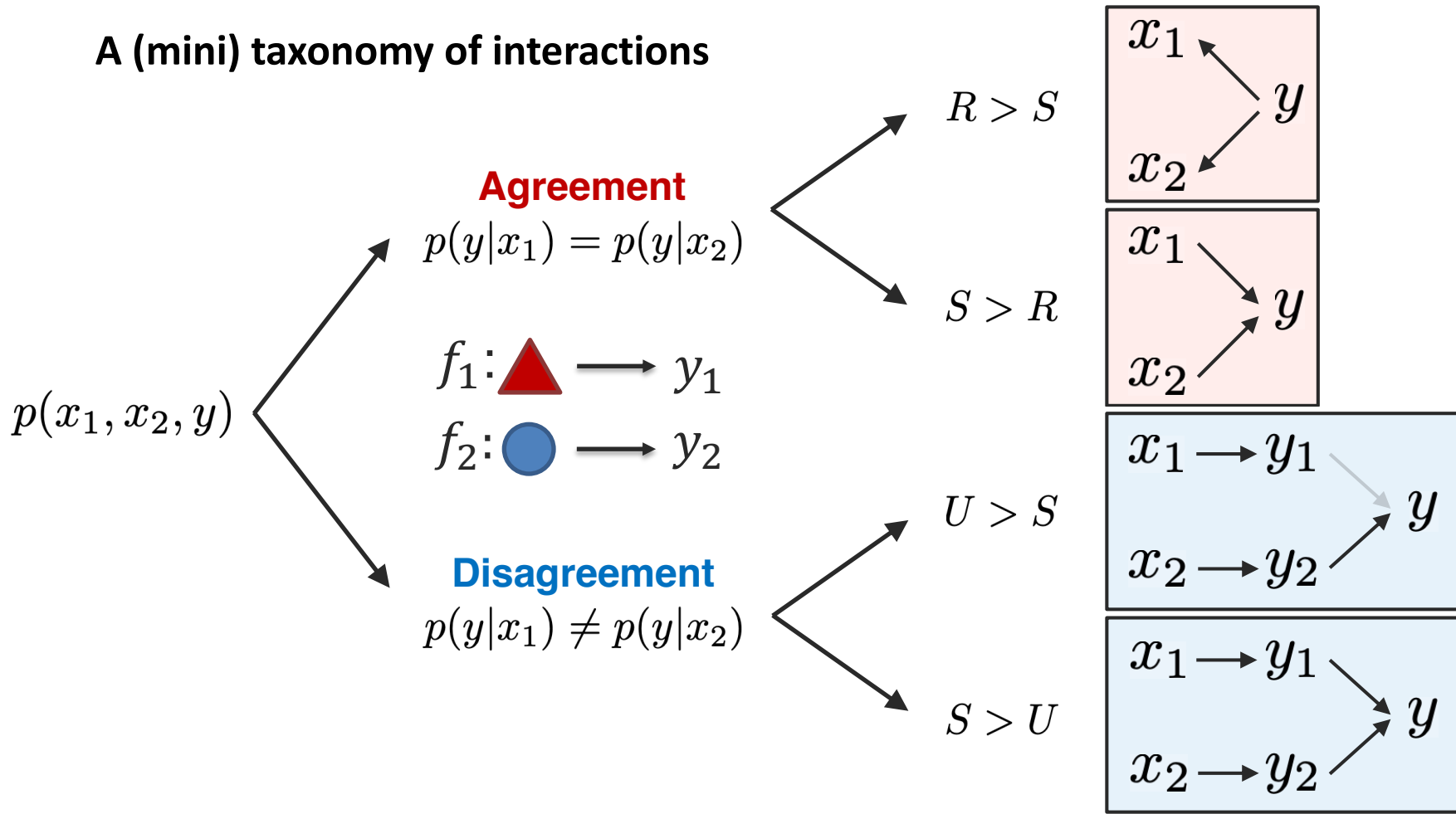
Gives theoretical results on estimating interactions and model performance for semi-supervised multimodal learning

[Liang et al., Multimodal Learning Without Labeled Multimodal Data: Guarantees and Applications, arXiv 2023]

On Agreement, Disagreement, and Synergy



A (mini) taxonomy of interactions



$I(X_1; X_2) > I(X_1; X_2|Y)$
Agreement redundancy
 Contrastive learning

$I(X_1; X_2) < I(X_1; X_2|Y)$
Agreement synergy
 Future work?

Disagreement uniqueness
 Feature selection

Disagreement synergy
 Future work?

[Blum and Mitchell. Combining Labeled and Unlabeled Data with Co-training. COLT 1998

[Peng et al., Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. TPAMI 2005]

[Liang et al., Multimodal Learning Without Labeled Multimodal Data: Guarantees and Applications, arXiv 2023]

Factorized Learning of Shared + Unique Information

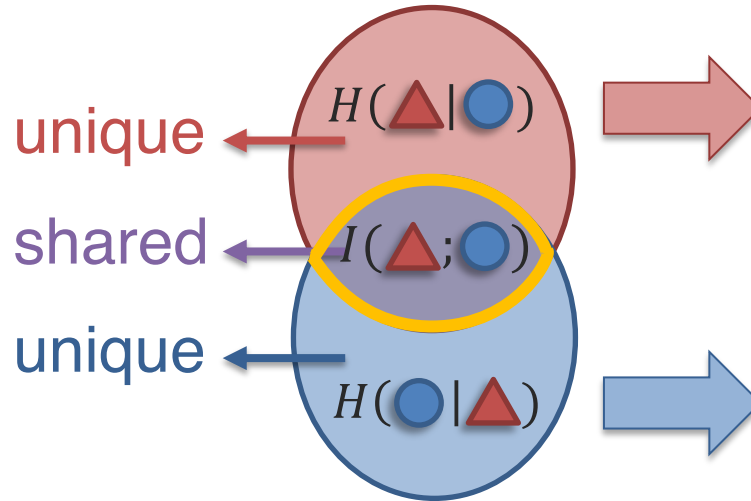
Modeling unique information



Can you please pass the cow?

Modality A 

Modality B 



1 Maximize the mutual information

$$I(\mathbf{z}; \bullet) \quad \text{and} \quad I(\mathbf{z}; \blacktriangle)$$

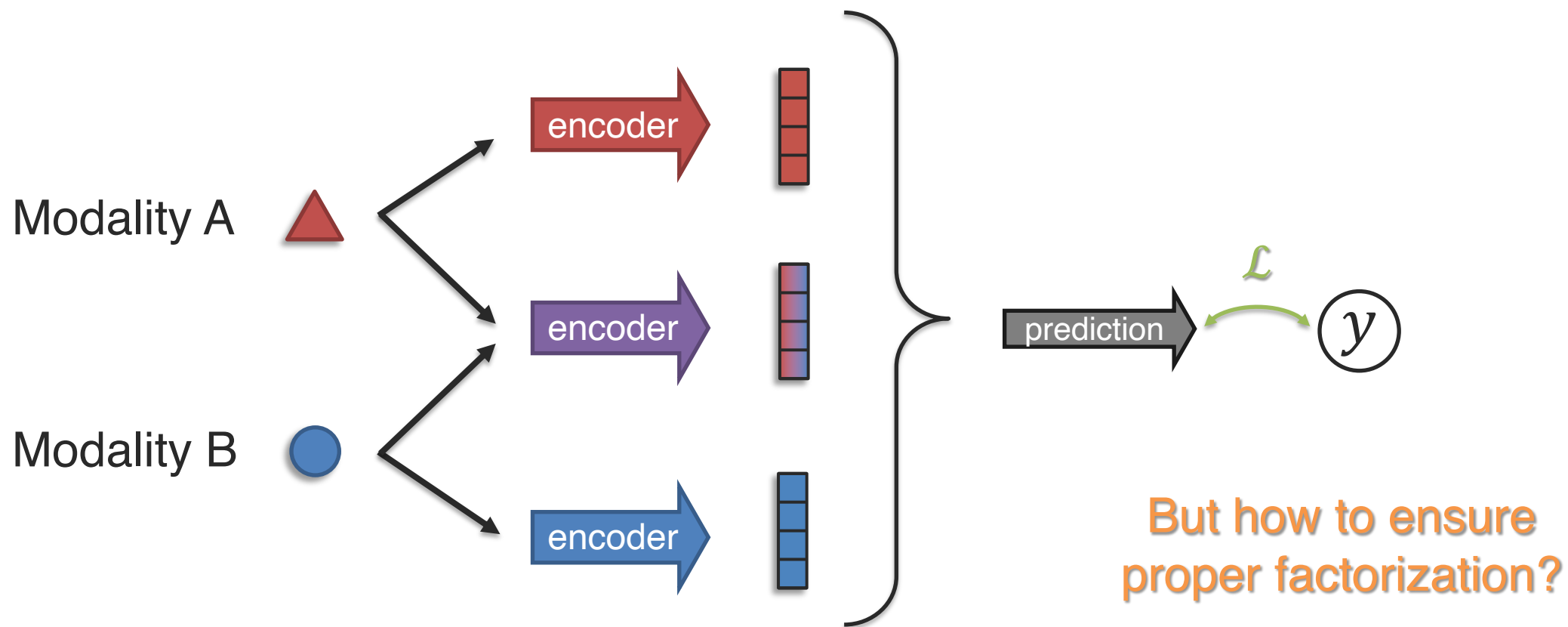
2 Minimize the conditional entropy

$$H(\mathbf{z}|\bullet) \quad \text{and} \quad H(\mathbf{z}|\blacktriangle)$$

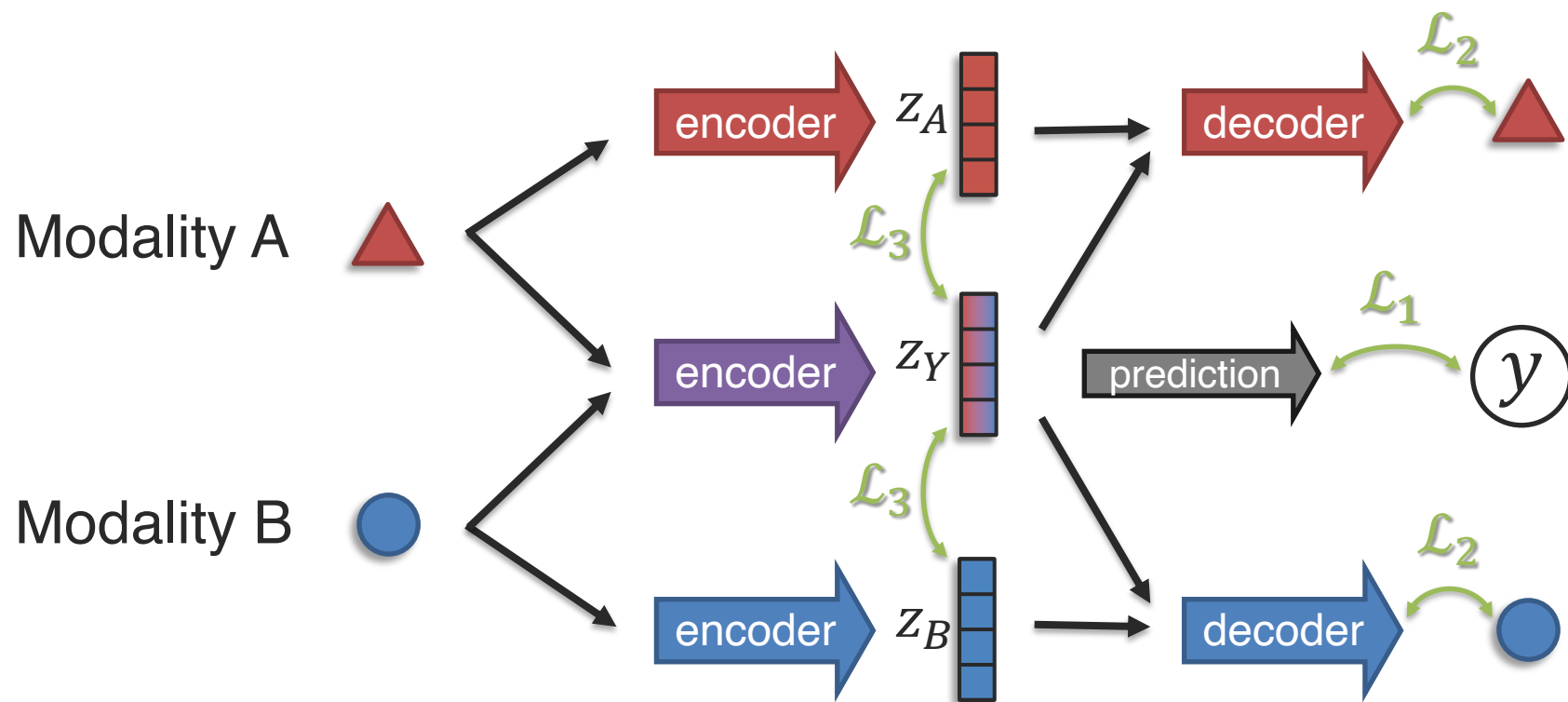
[Tsai et al., Learning Factorized Multimodal Representations. ICLR 2021]

[Wang et al., Rethinking Minimal Sufficient Representation in Contrastive Learning, CVPR 2022]

Factorized Multimodal Representations



A Generative-Discriminative Approach



$$\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_3$$

\mathcal{L}_1 : discriminative

\mathcal{L}_2 : generative

\mathcal{L}_3 : no overlap

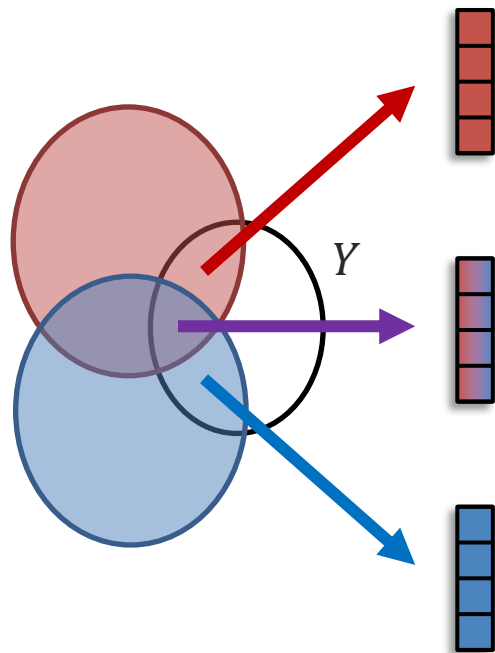
Independent priors
for z_A , z_B and z_Y

Learning Task-relevant Unique Information

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

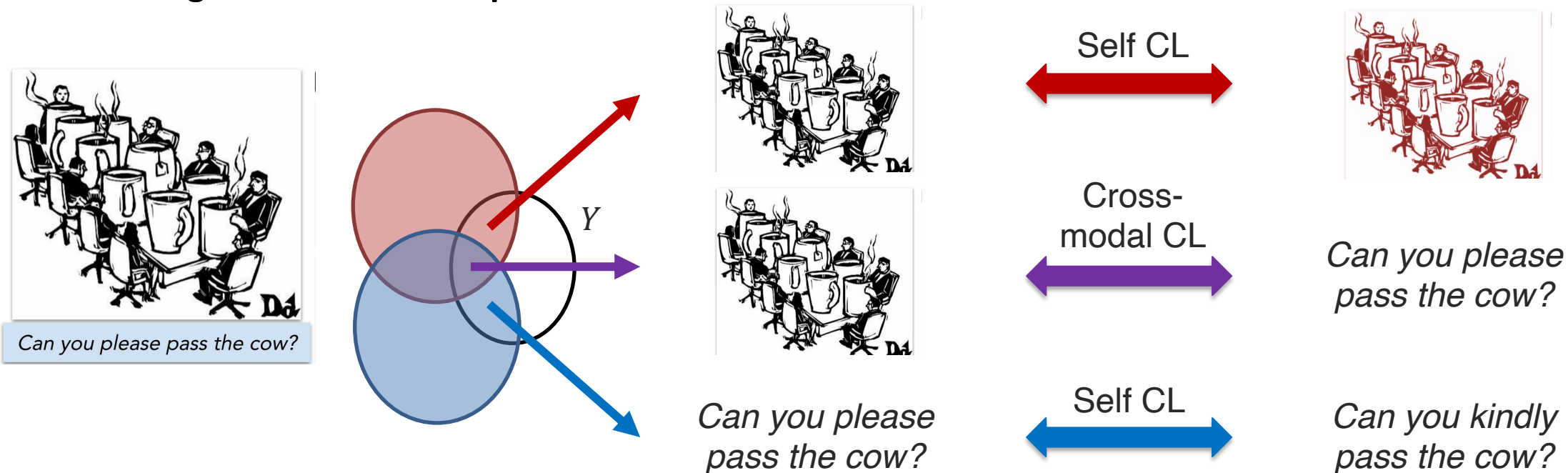
$$I(\mathbf{Z}; \bullet; Y) \quad \text{and} \quad I(\mathbf{Z}; \blacktriangle; Y)$$

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

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

Learning Task-relevant Unique Information

Modeling task-relevant unique information

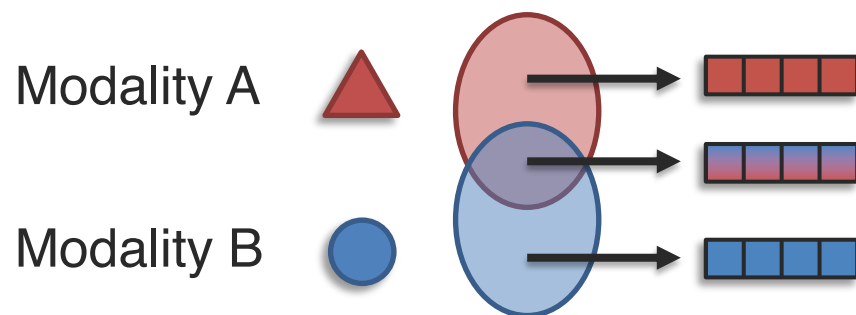


Approximate task-relevance Y using multi-view data augmentations
New scalable lower and upper bounds on mutual information

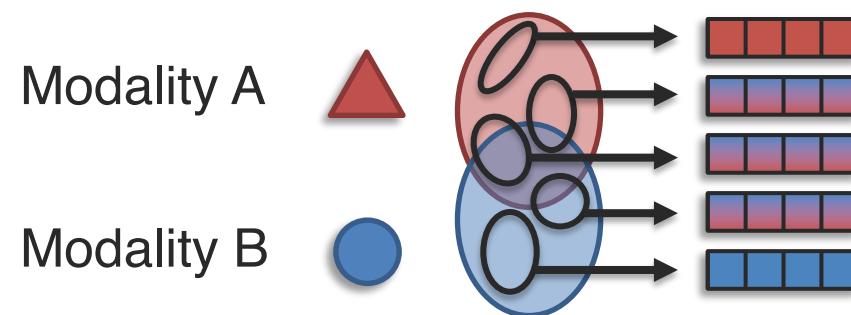
Fine-Grained Fission

How to automatically discover these internal clusters, factors?

Modality-level fission:

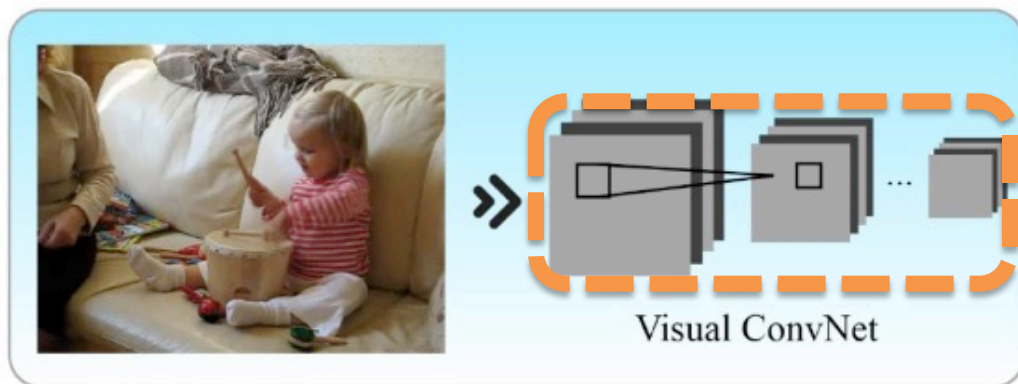


Fine-grained fission:

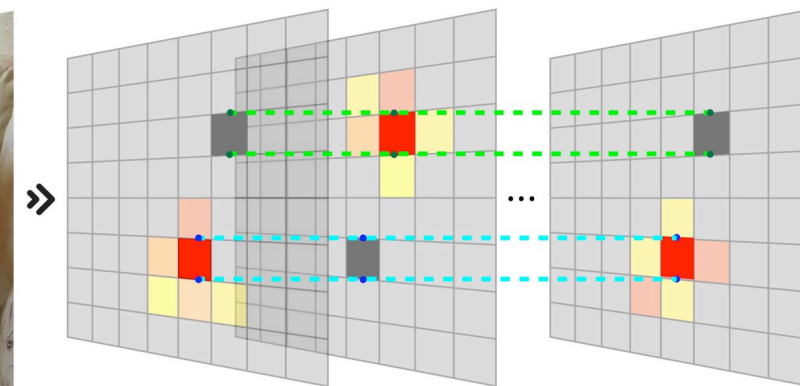
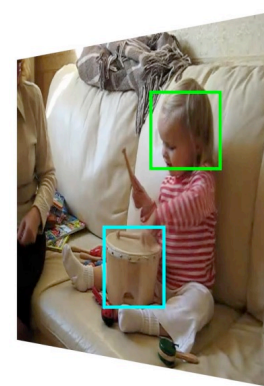
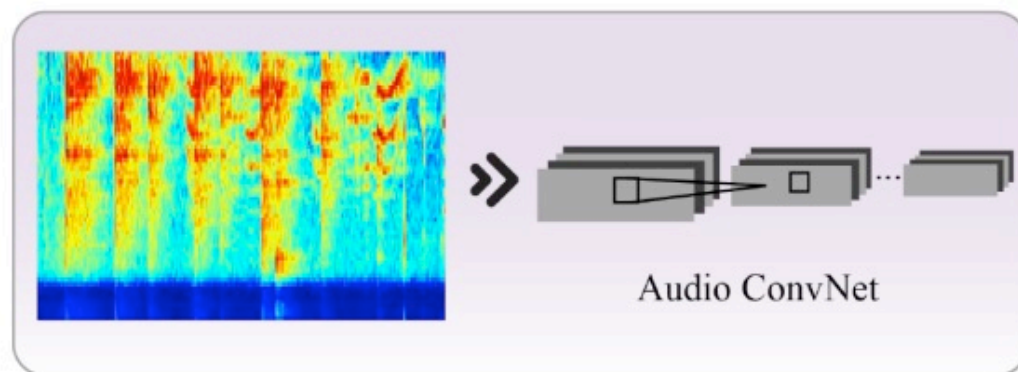


Fine-Grained Fission – A Clustering Approach

Unimodal Encoders



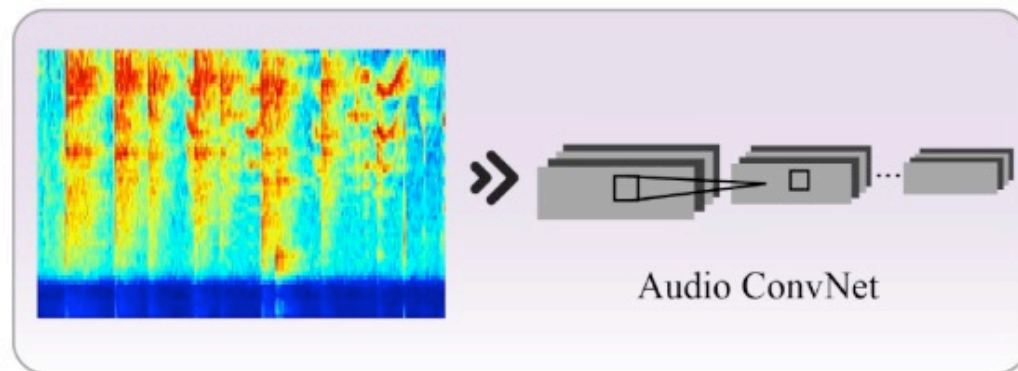
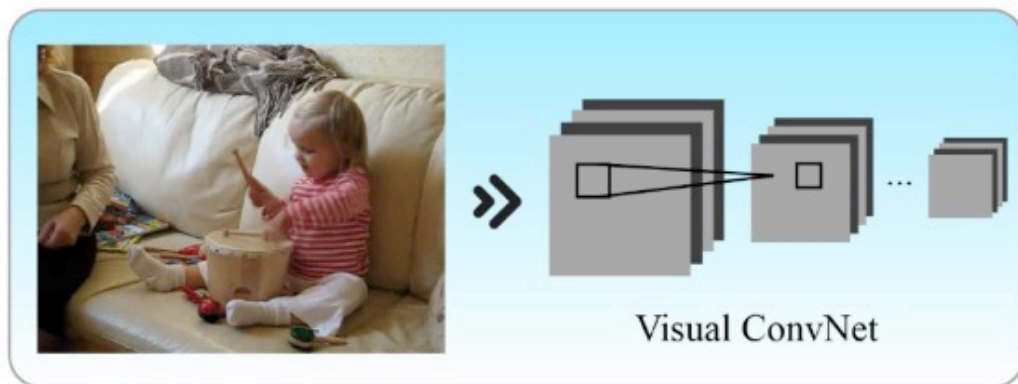
Localized activations for different objects



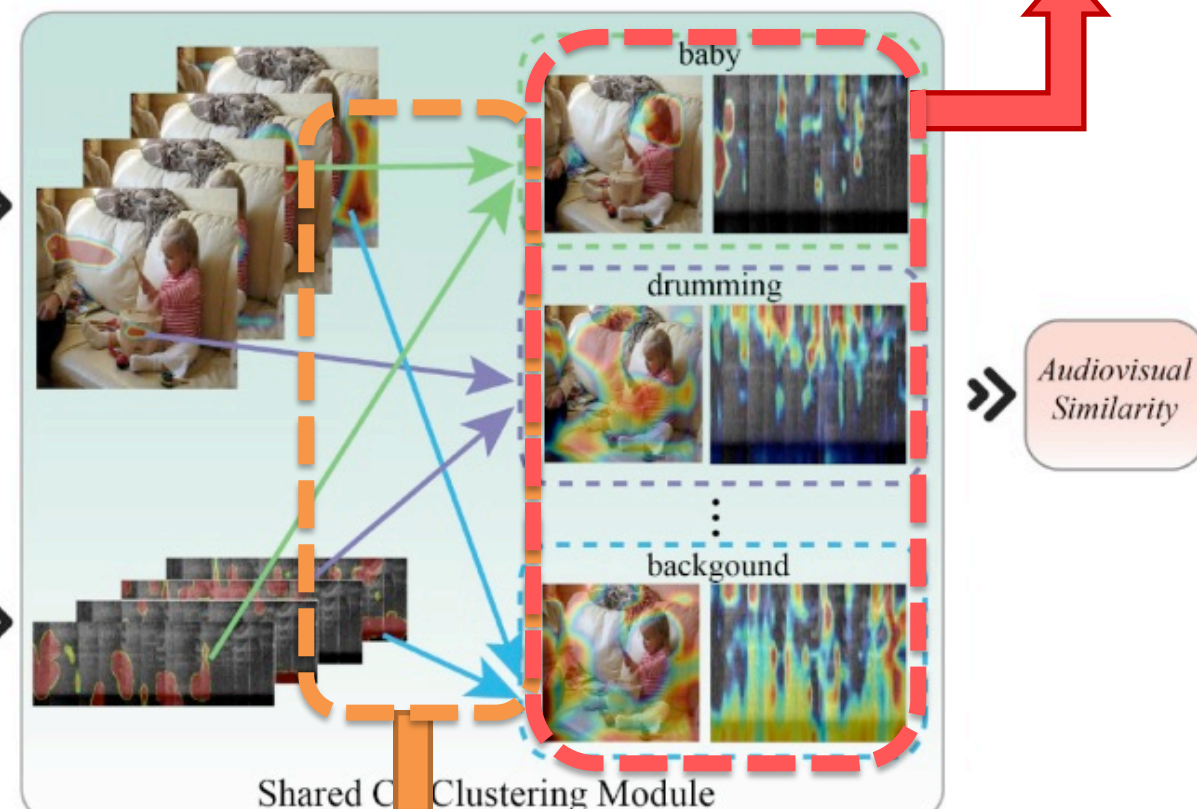
Fine-Grained Fission – A Clustering Approach

Discovers multiple audio-visual correspondences

Unimodal Encoders



Multimodal Fission

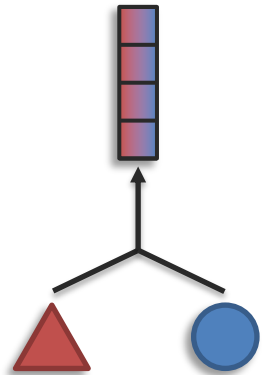


Challenge 1: Representation

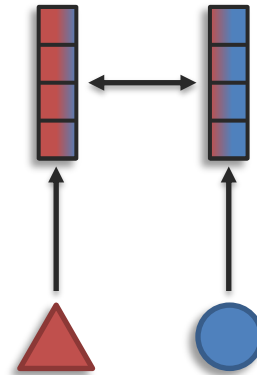
Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

Sub-challenges:

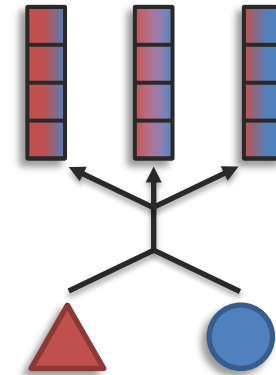
Fusion



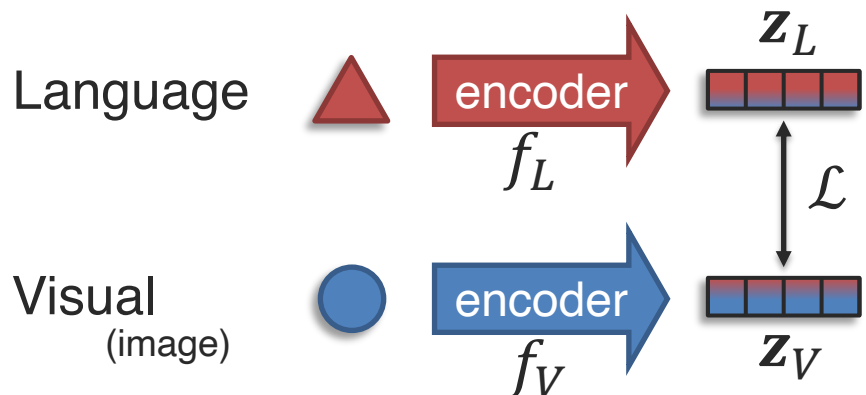
Coordination



Fission



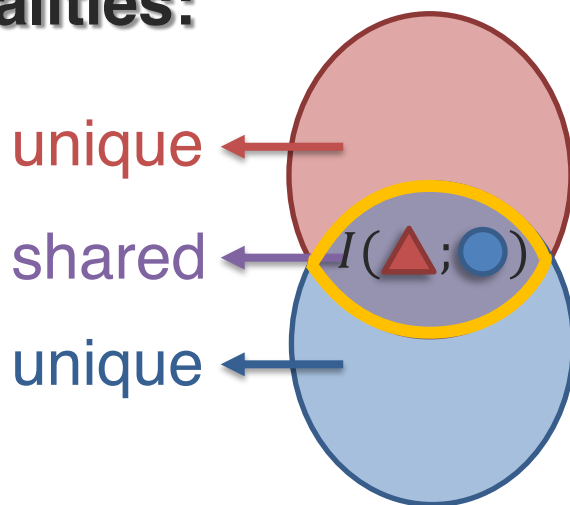
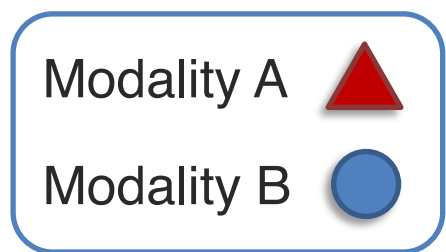
Recap: Contrastive Learning and Connected Modalities



Popular contrastive loss: InfoNCE

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

Connected modalities:

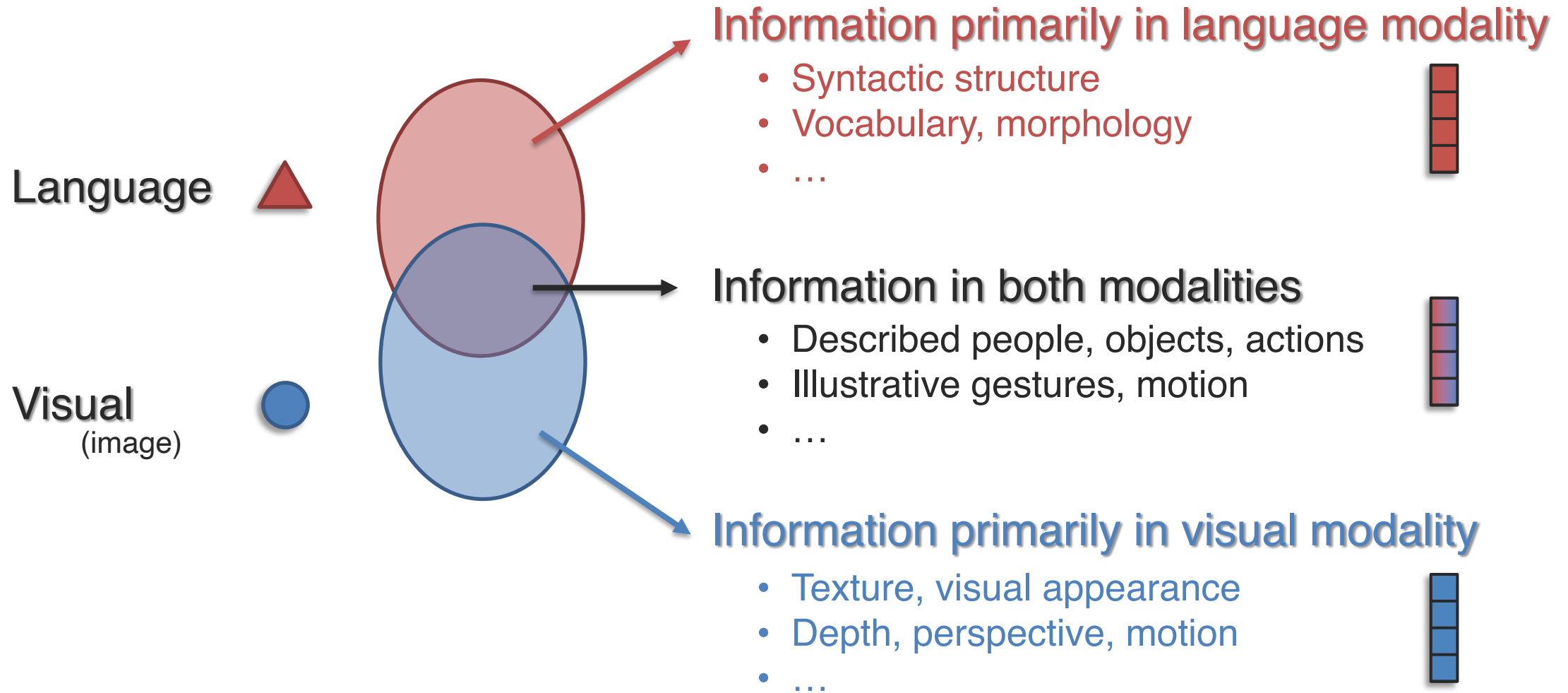


CLIP focuses on shared connections

Mutual information $I(X; Y)$

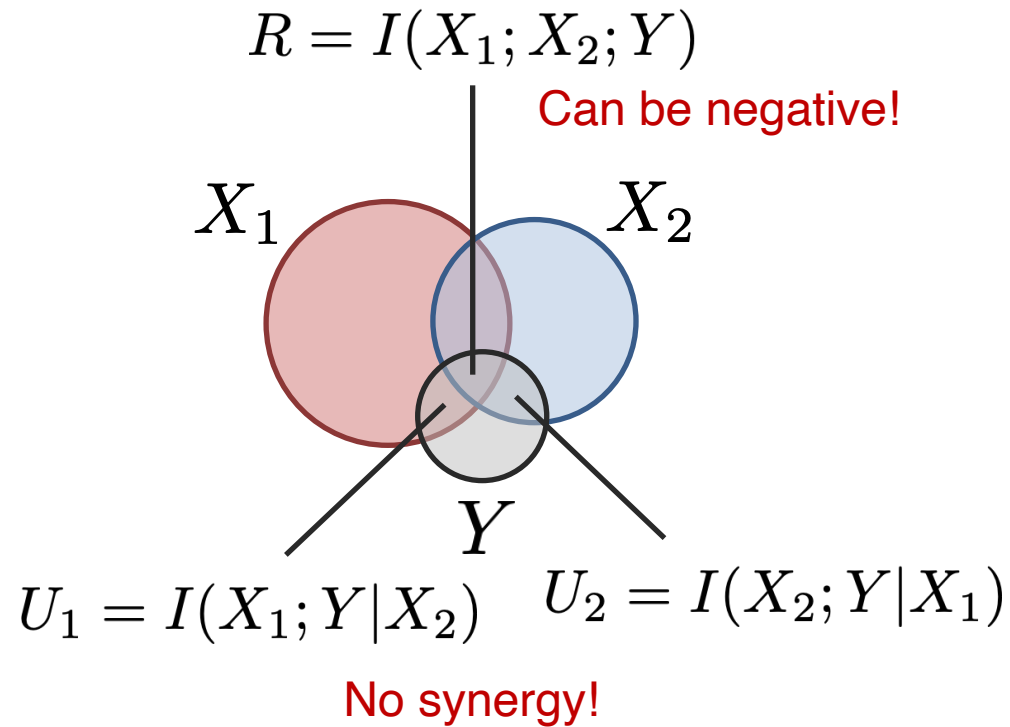
$$\mathbb{E}_{X,Y} \left[\log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} \right]$$

Recap: Modality-Level Fission



Recap: Partial Information Decomposition

Classical Information Theory



Partial Information Decomposition

