



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Multimodal Machine Learning

## Lecture 4.1: Multimodal Representations (Part 2)

Louis-Philippe Morency

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yanatan Bisk.*

# Administrative Stuff

# First Project Assignment

---

Due date: Sunday 9/25 at 8m

Four main sections:

- Introduction
- Related work
- Experimental setup
- Research ideas

Follows ICML paper format

➔ The two main sections are related work and research ideas

➔ # teammates = # research ideas

➔ Page limit depends on team size:

- 3 students : 4 pages + references
- 4 students : 4.5 pages + references
- 5 students : 5 pages + references
- 6 students : 5.5 pages + references

## Team Meetings with Instructor

---

- No lecture on Tuesday 9/27
- 15-mins meeting with instructor
  - Optional, but highly suggested
  - Not all teammates are required to attend
  - Prepare 2 slides to summarize your research ideas
- Meetings on Tuesday 9/27 and Wednesday 9/28
- Signup form:  
<https://calendly.com/morency/student-meetings>



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Multimodal Machine Learning

## Lecture 4.1: Multimodal Representations (Part 2)

Louis-Philippe Morency

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yo*

# Objectives of today's class

---

- Representation fusion
  - Multimodal auto-encoder
  - Fusion from raw modalities
- Representation coordination
  - Coordination functions
    - Kernel similarity functions
    - Canonical correlation analysis
  - Contrastive learning
- Representation fission
  - Factorized multimodal representations
  - Information, entropy and mutual information
  - Clustering and fine-grained fission

# Multimodal Representation

---

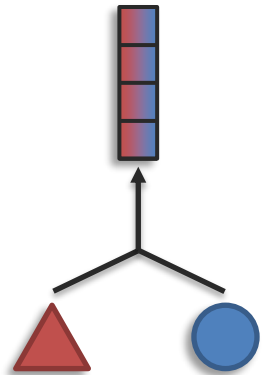
# Challenge 1: Representation

---

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

## Sub-challenges:

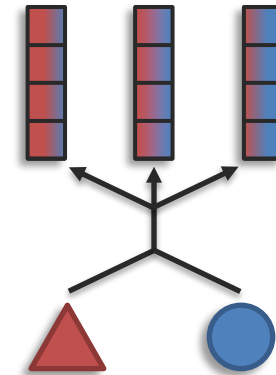
### Fusion



### Coordination

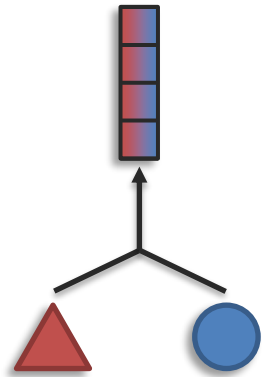


### Fission





# Sub-Challenge 1a: Representation Fusion

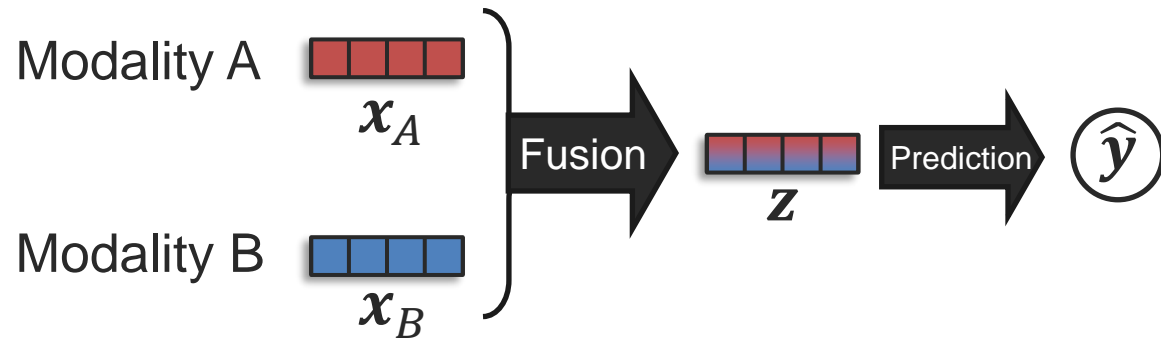


**Definition:** Learn a joint representation that models cross-modal interactions between individual elements of different modalities



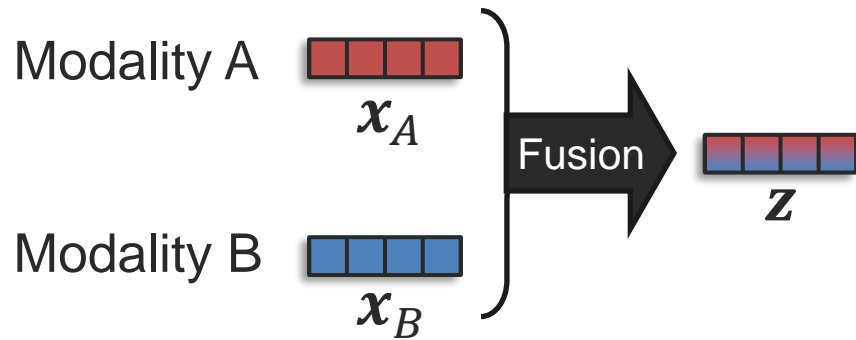
# Learning Fusion Representations

---



How to learn fusion models?

# Learning Fusion Representations

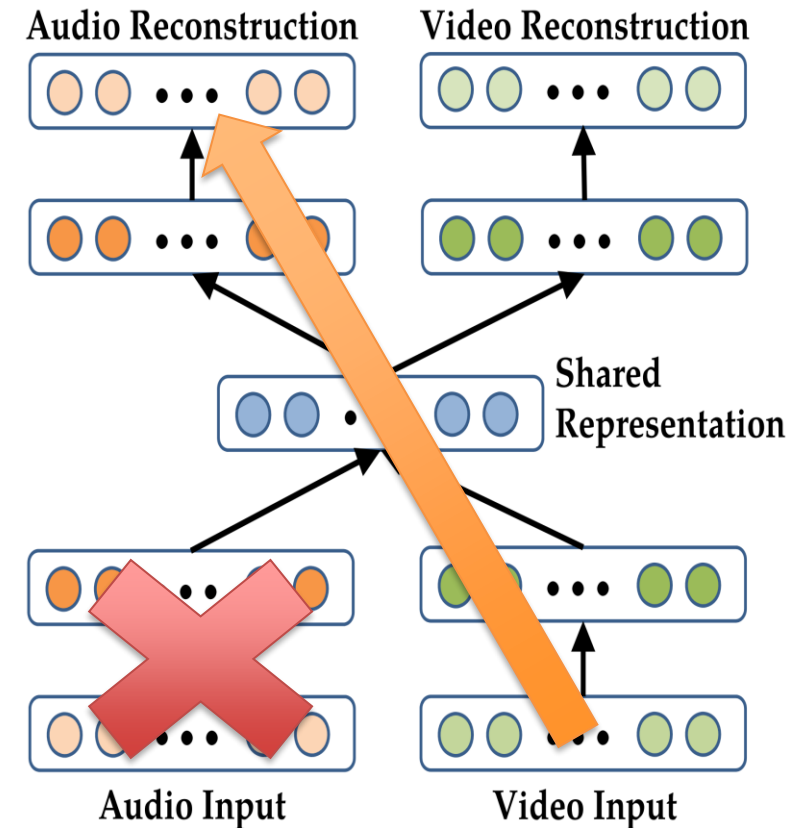


How to learn fusion models?

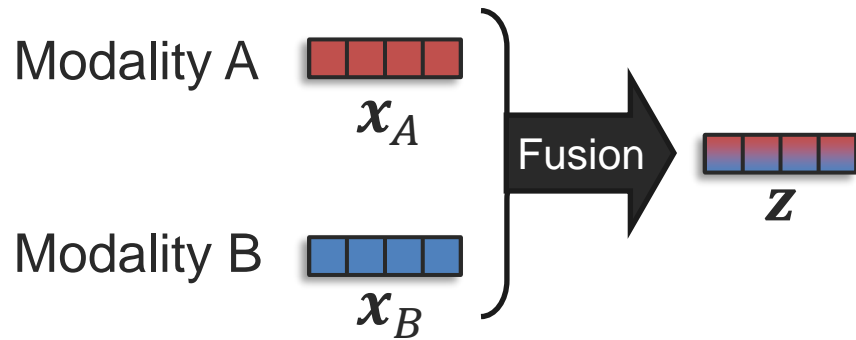
What will be the loss function?

Can it hallucinate the other modality?

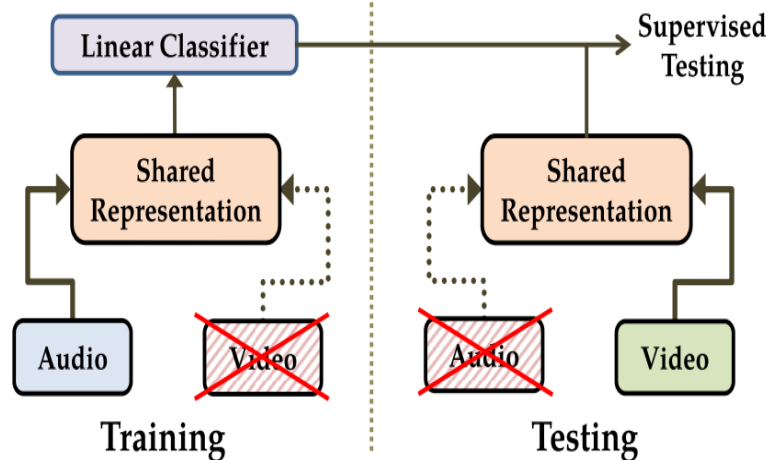
## Multimodal Autoencoder



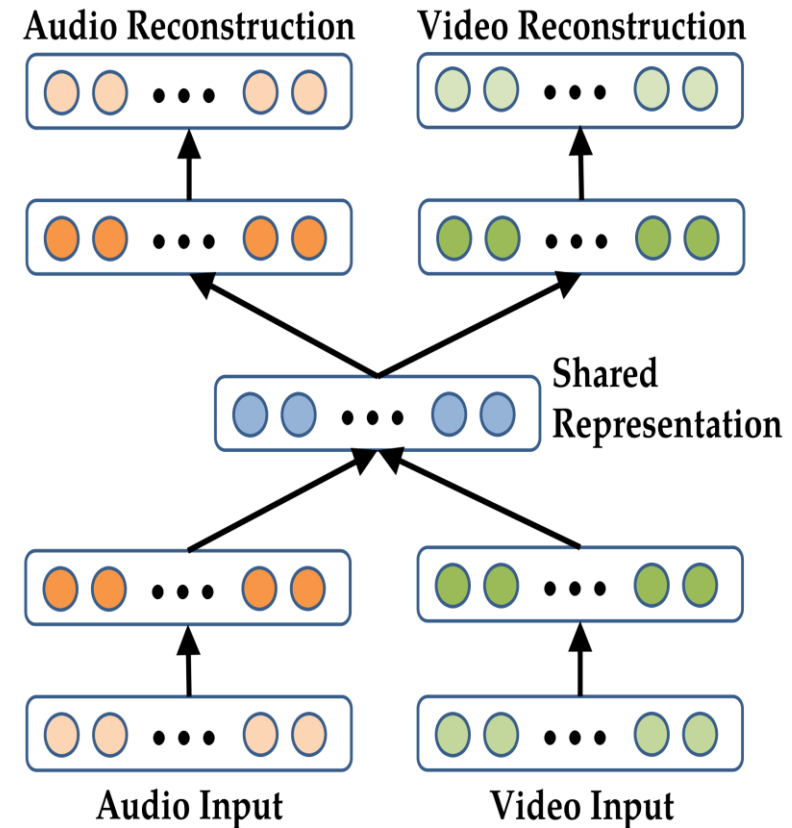
# Learning Fusion Representations



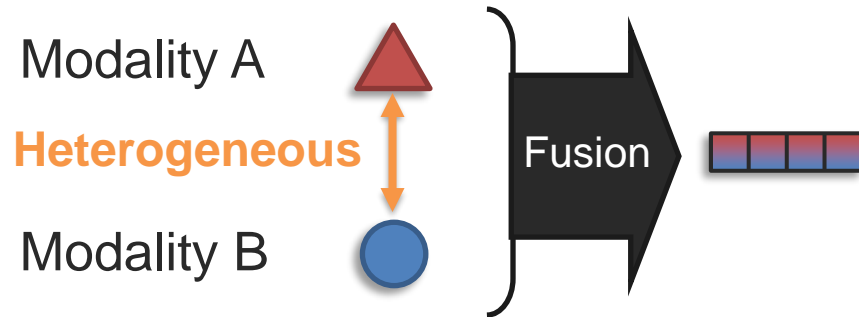
Interesting experiment: “Hearing to see”  
(zero-shot cross-modal adaptation)



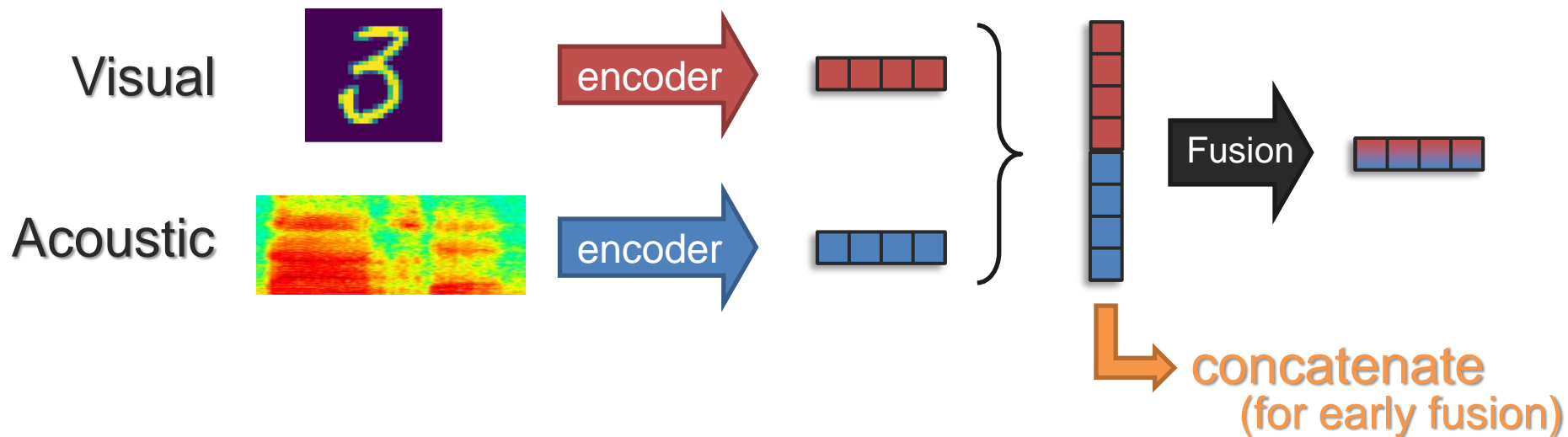
## Multimodal Autoencoder



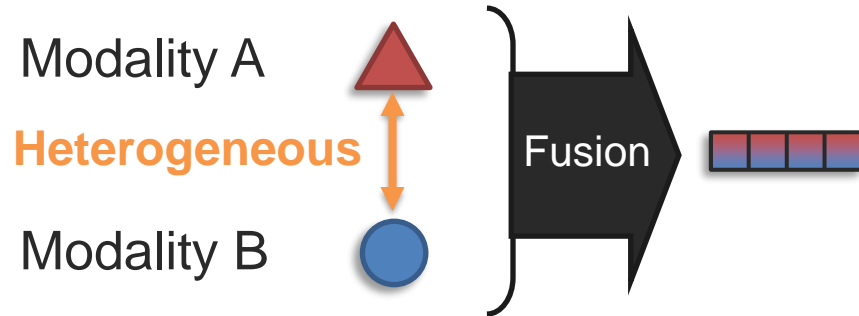
# Fusion with Raw Modalities



## Example: From Early Fusion...

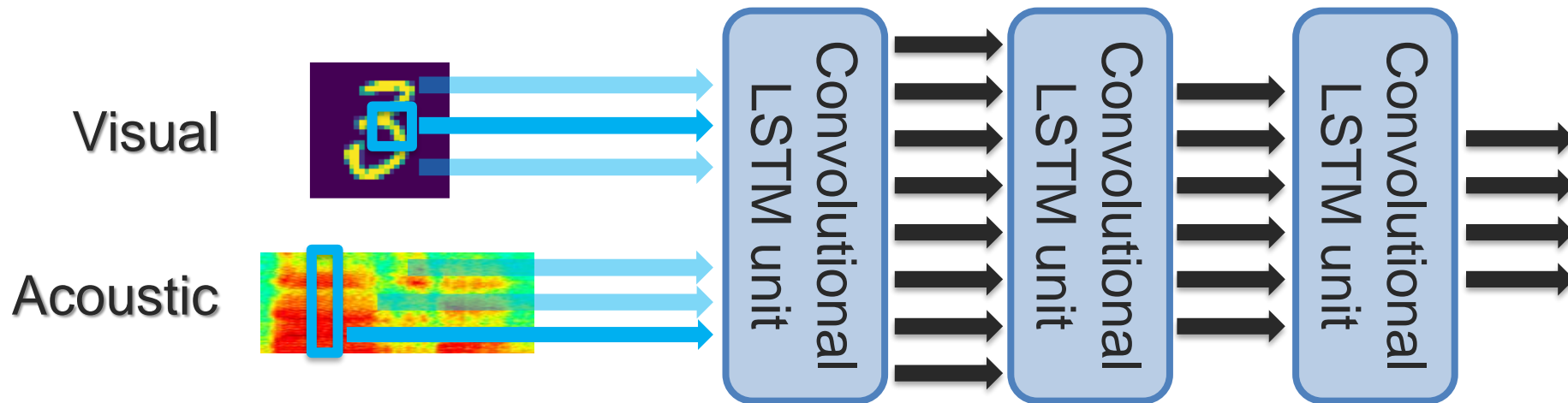


# Fusion with Raw Modalities



**Open Challenge!**

**Example: From Early Fusion... to Very Early Fusion (inspired by human brain)**

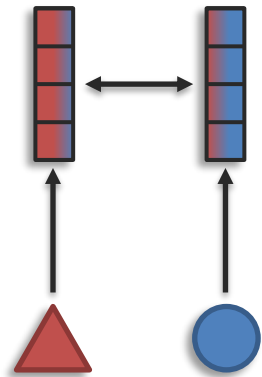


Barnum, et al. "On the Benefits of Early Fusion in Multimodal Representation Learning." *arxiv* 2022

# Representation Coordination

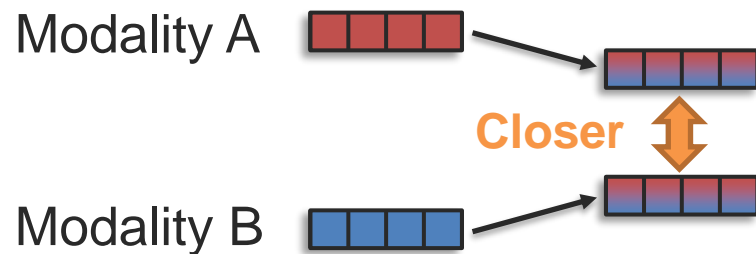
---

# Sub-Challenge 1b: Representation Coordination

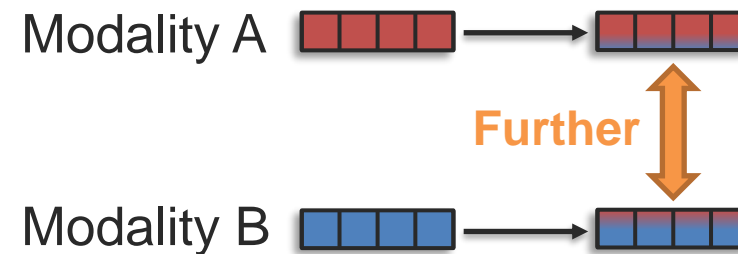


**Definition:** Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions

Strong Coordination:

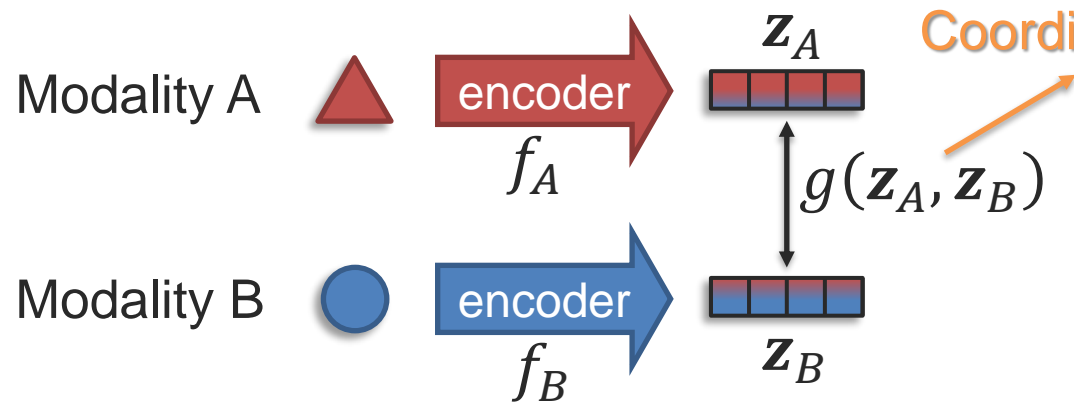


Partial Coordination:





# Coordination Function



Learning with coordination function:

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

with model parameters  $\theta_g$ ,  $\theta_{f_A}$  and  $\theta_{f_B}$

➡ Requires paired data

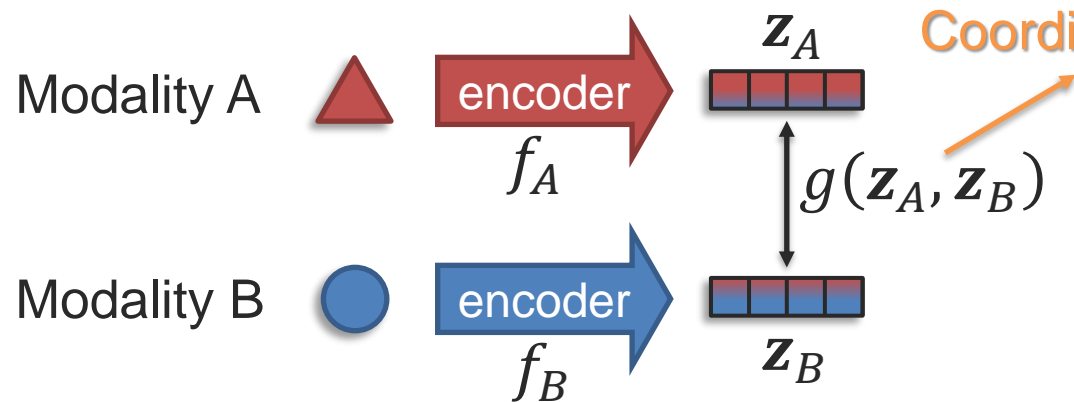
## Examples of coordination function:

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  $\theta_g$ ,  $\theta_{f_A}$  and  $\theta_{f_B}$

## Examples of coordination function:

② Kernel similarity functions:

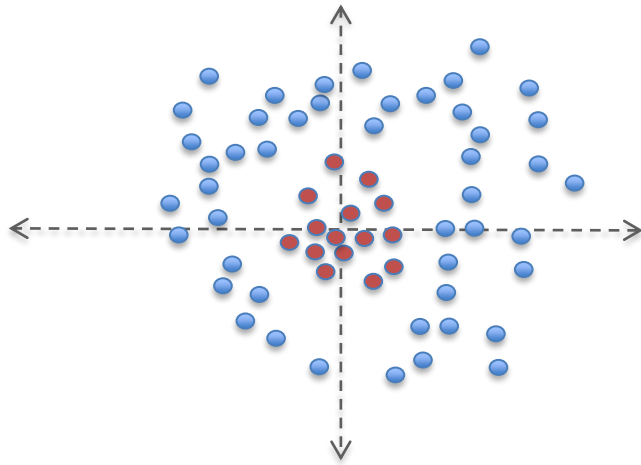
$$g(\mathbf{z}_A, \mathbf{z}_B) = k(\mathbf{z}_A, \mathbf{z}_B) \left\{ \begin{array}{l} \bullet \text{ Linear} \\ \bullet \text{ Polynomial} \\ \bullet \text{ Exponential} \\ \bullet \text{ RBF} \end{array} \right.$$

➡ All these examples bring relatively strong coordination between modalities

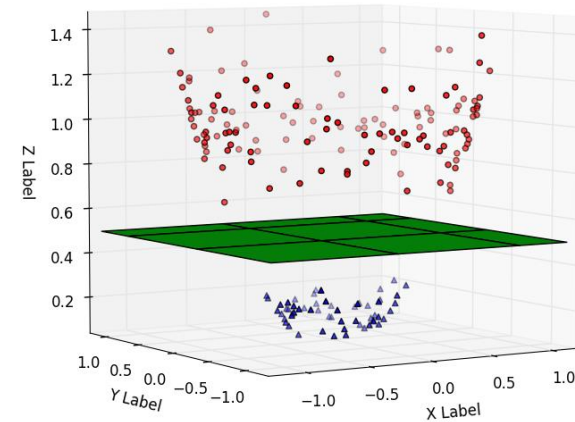
# Kernel Function

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

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



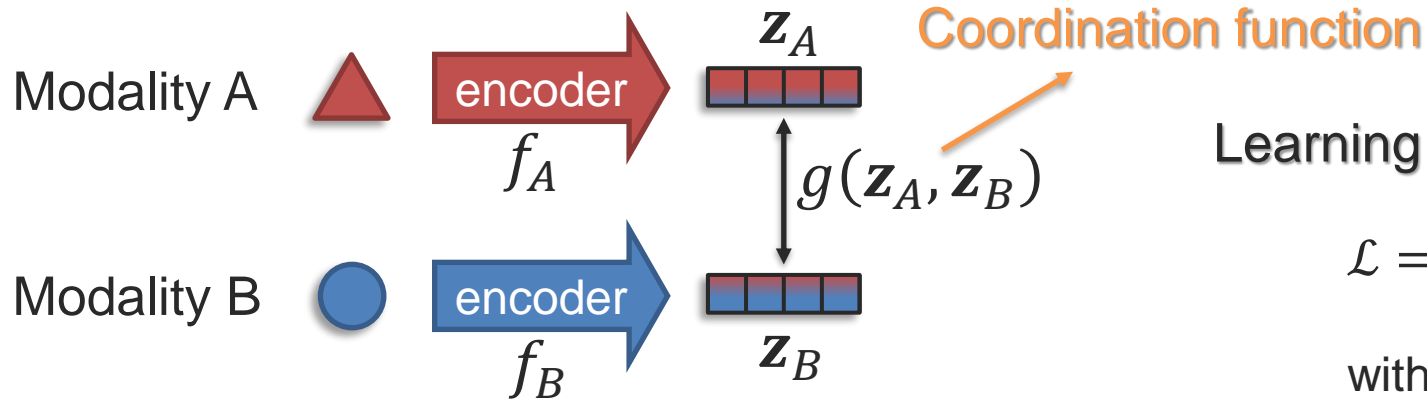
Not linearly separable in  $x$  space



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

$\rightarrow$  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:

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

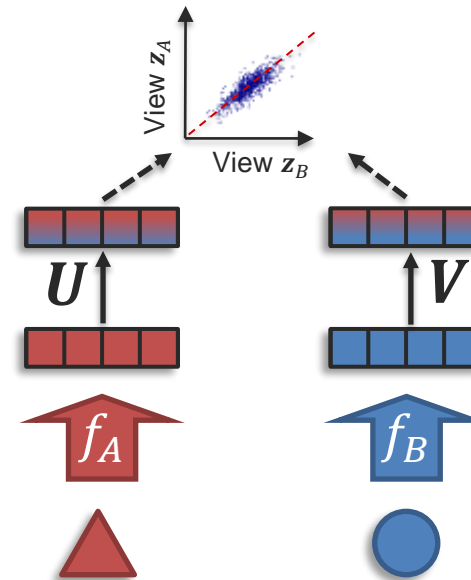
with model parameters  $\theta_g$ ,  $\theta_{f_A}$  and  $\theta_{f_B}$

## Examples of coordination function:

③ Canonical Correlation Analysis (CCA):

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

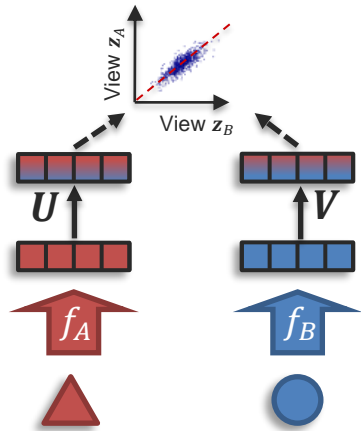
→ CCA includes multiple projections, all orthogonal with each others



# 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

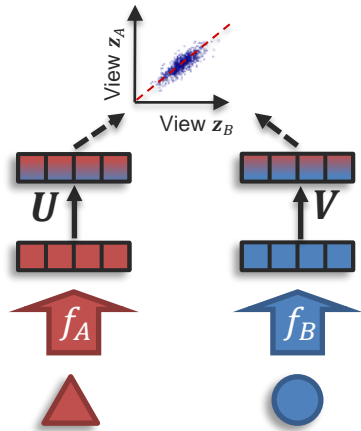
➡ Remember that  $X$  and  $Y$  consist of paired data

# Canonical Correlation Analysis

- ② We want these multiple projection pairs to be orthogonal (“canonical”) to each other:

$$\mathbf{u}_{(i)}^T \Sigma_{XY} \mathbf{v}_{(j)} = \mathbf{u}_{(j)}^T \Sigma_{XY} \mathbf{v}_{(i)} = \mathbf{0} \quad \text{for } i \neq j$$

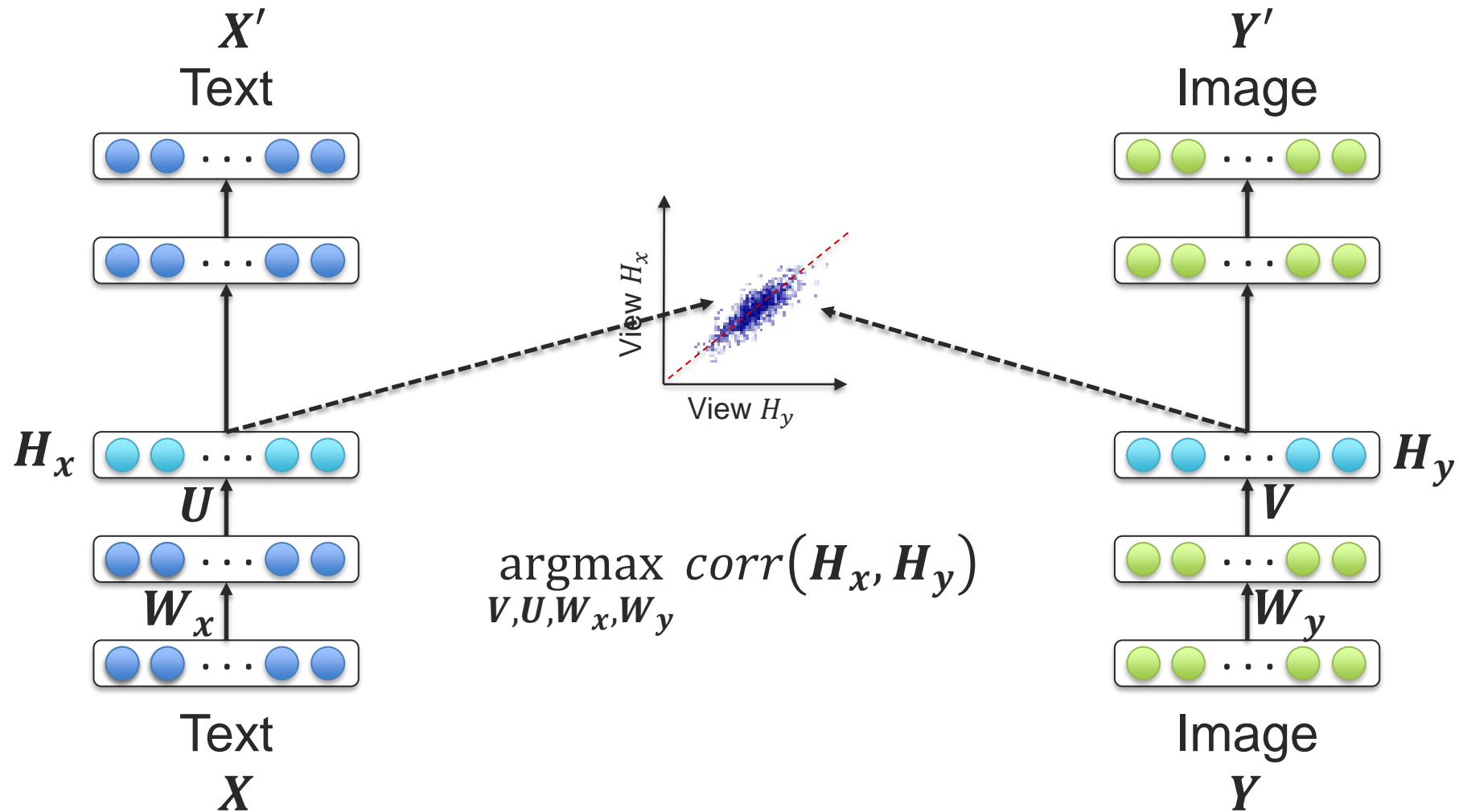
$$|U \Sigma_{XY} V| = \text{tr}(U \Sigma_{XY} V) \quad \text{where } U = [\mathbf{u}_{(1)}, \mathbf{u}_{(2)}, \dots, \mathbf{u}_{(k)}]$$
$$\text{and } V = [\mathbf{v}_{(1)}, \mathbf{v}_{(2)}, \dots, \mathbf{v}_{(k)}]$$



- ③ Since this objective function is invariant to scaling, we can constraint the projections to have unit variance:

$$U^T \Sigma_{XX} U = I \quad V^T \Sigma_{YY} V = I$$

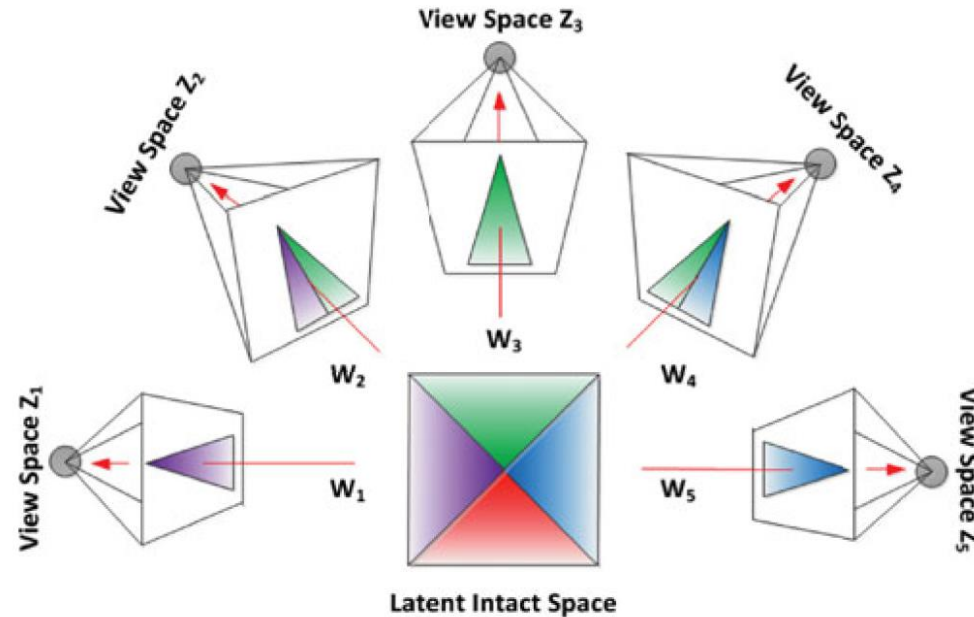
# Deep Canonically Correlated Autoencoders (DCCA)



# Multi-view Latent “Intact” Space

---

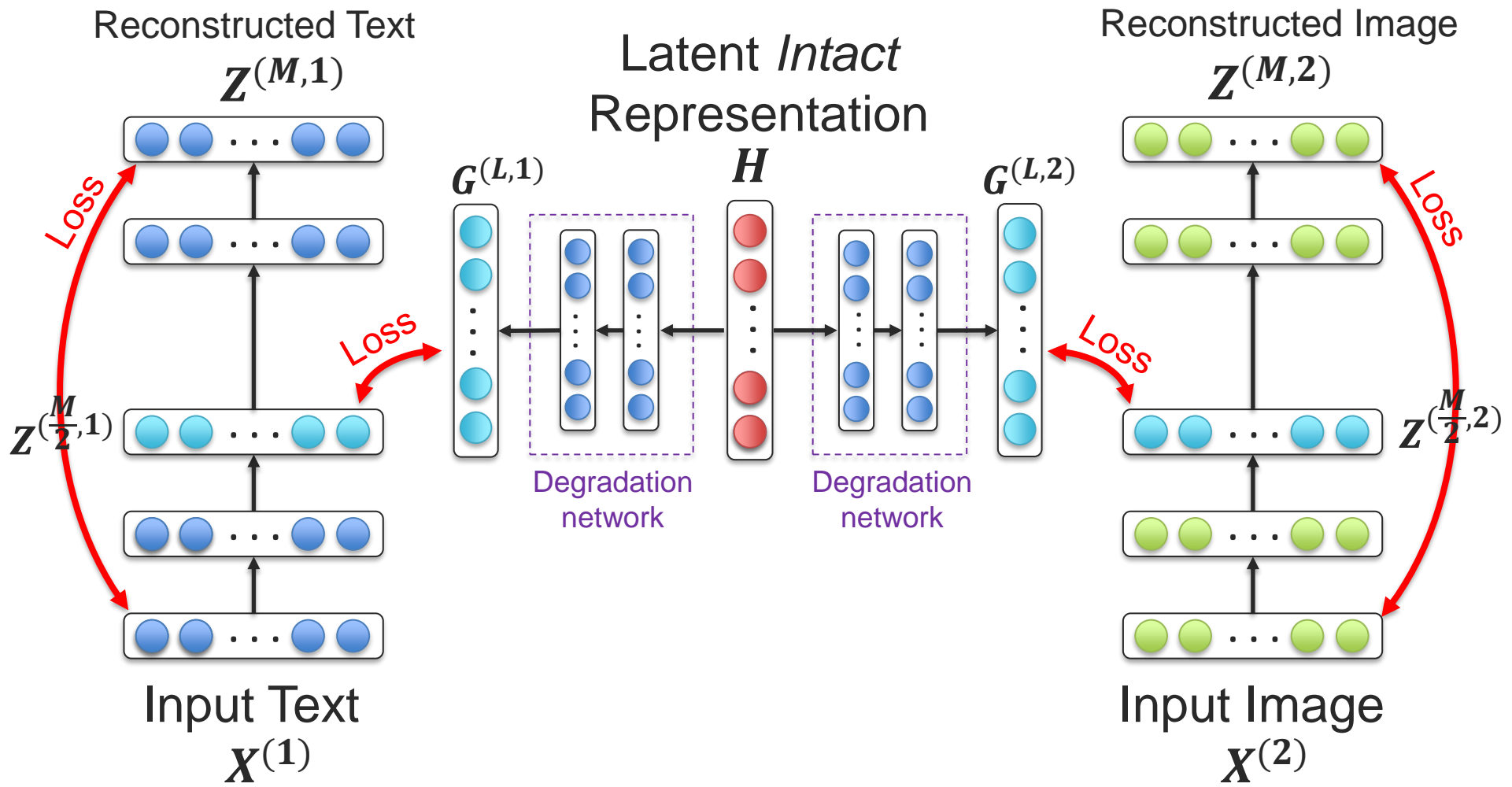
Given multiple views  $z_i$  from the same “object”:



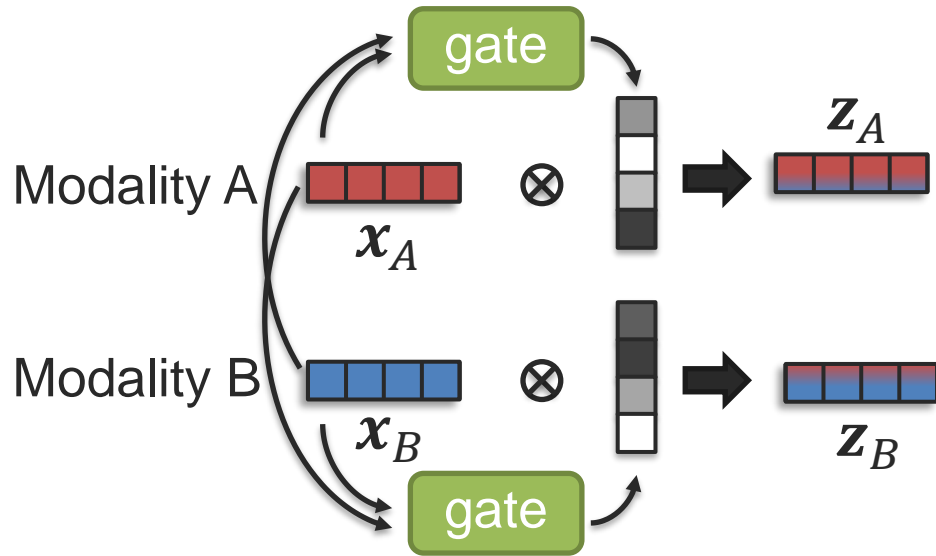
- 1) There is an “intact” representation which is *complete* and *not damaged*
- 2) The views  $z_i$  are partial (and possibly degenerated) representations of the intact representation



# Auto-Encoder in Auto-Encoder Network



# Gated Coordination



Gated coordination:

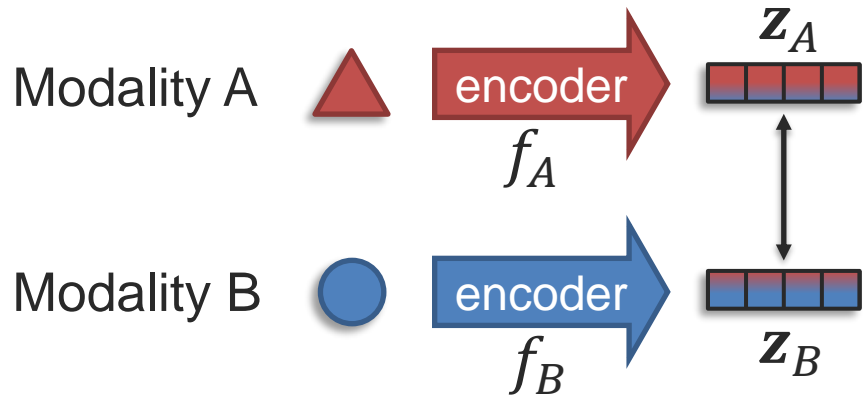
$$z_A = g_A(x_A, x_B) \cdot x_A$$

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

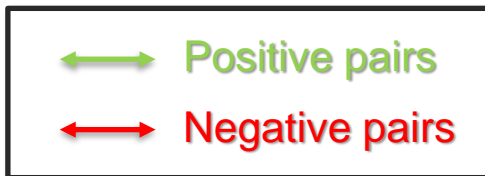
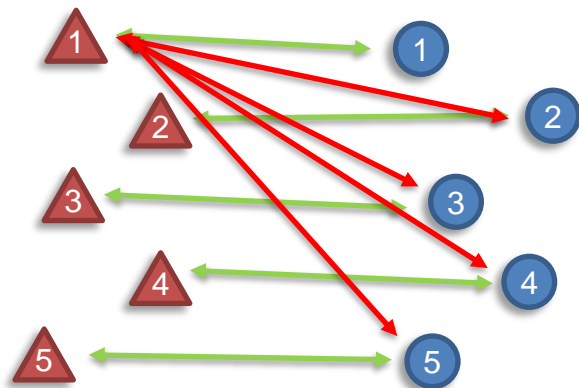
➡ Related to attention modules in transformers

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

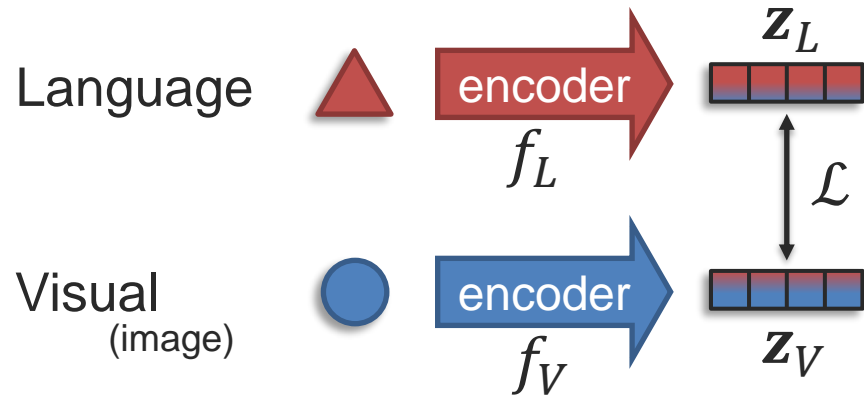
Simple contrastive loss:

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

Similarity functions are often cosine similarity

→ Similar to hinge loss

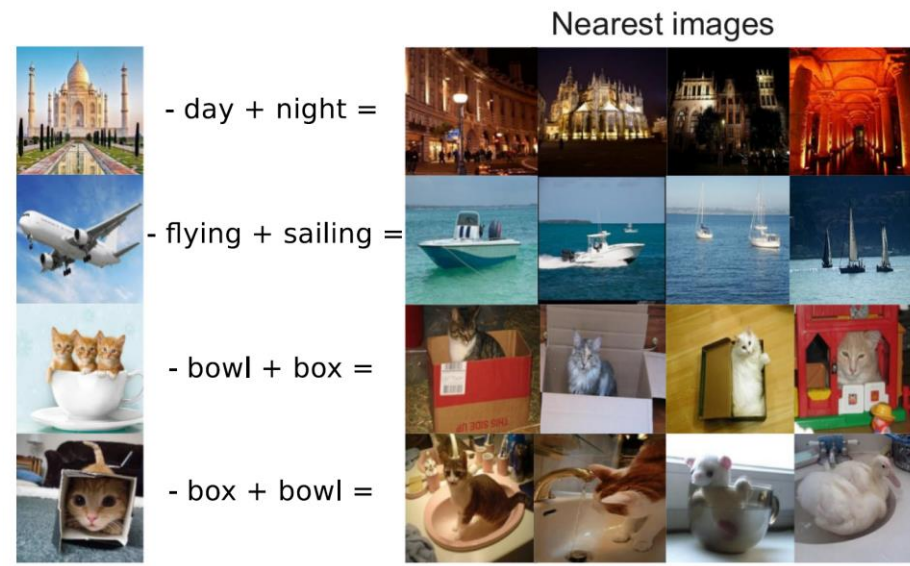
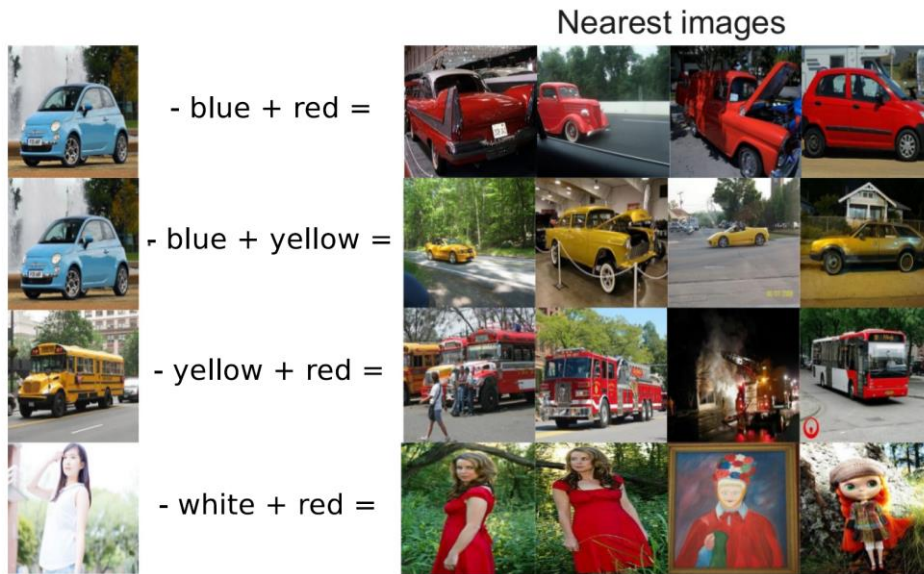
# Example – Visual-Semantic Embeddings



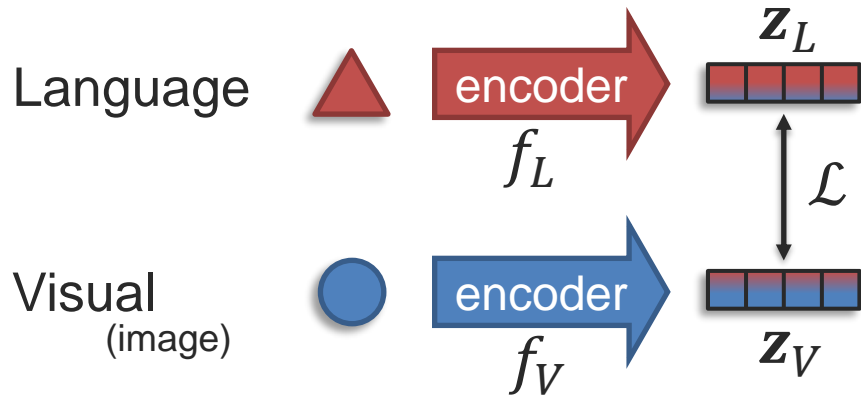
Two contrastive loss terms:

$$\max\{0, \alpha + \text{sim}(\mathbf{z}_L, \mathbf{z}_V^+) - \text{sim}(\mathbf{z}_L, \mathbf{z}_V^-)\}$$

$$+ \max\{0, \alpha + \text{sim}(\mathbf{z}_V, \mathbf{z}_L^+) - \text{sim}(\mathbf{z}_V, \mathbf{z}_L^-)\}$$



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



Popular contrastive loss: InfoNCE

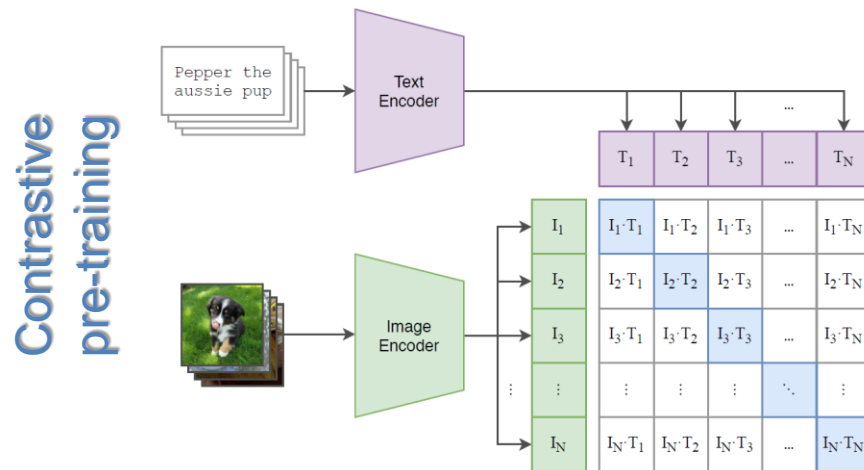
$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(z_A^i, z_B^i)}{\sum_{j=1}^N \text{sim}(z_A^i, z_B^j)}$$

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

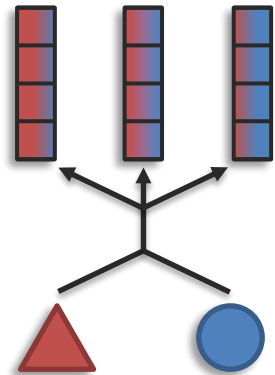
$z_L$  and  $z_V$  are coordinated but not identical representation spaces

# Representation

# Fission

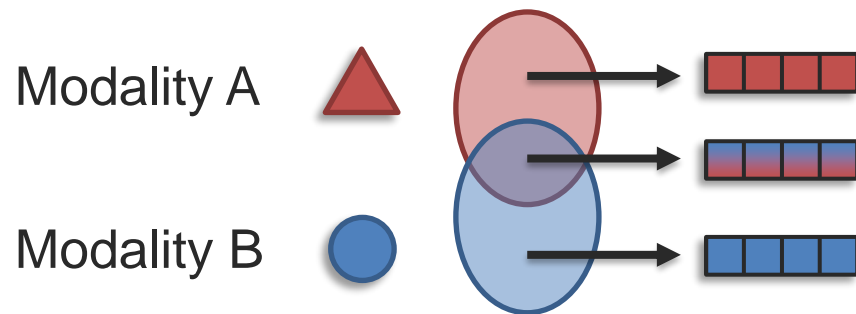
---

# Sub-Challenge 1c: Representation Fission

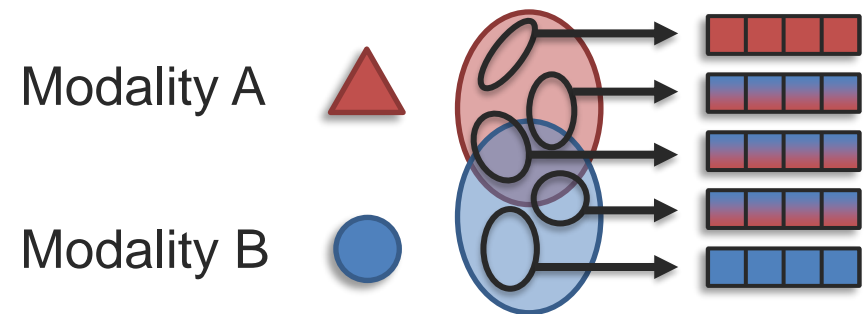


**Definition:** learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering

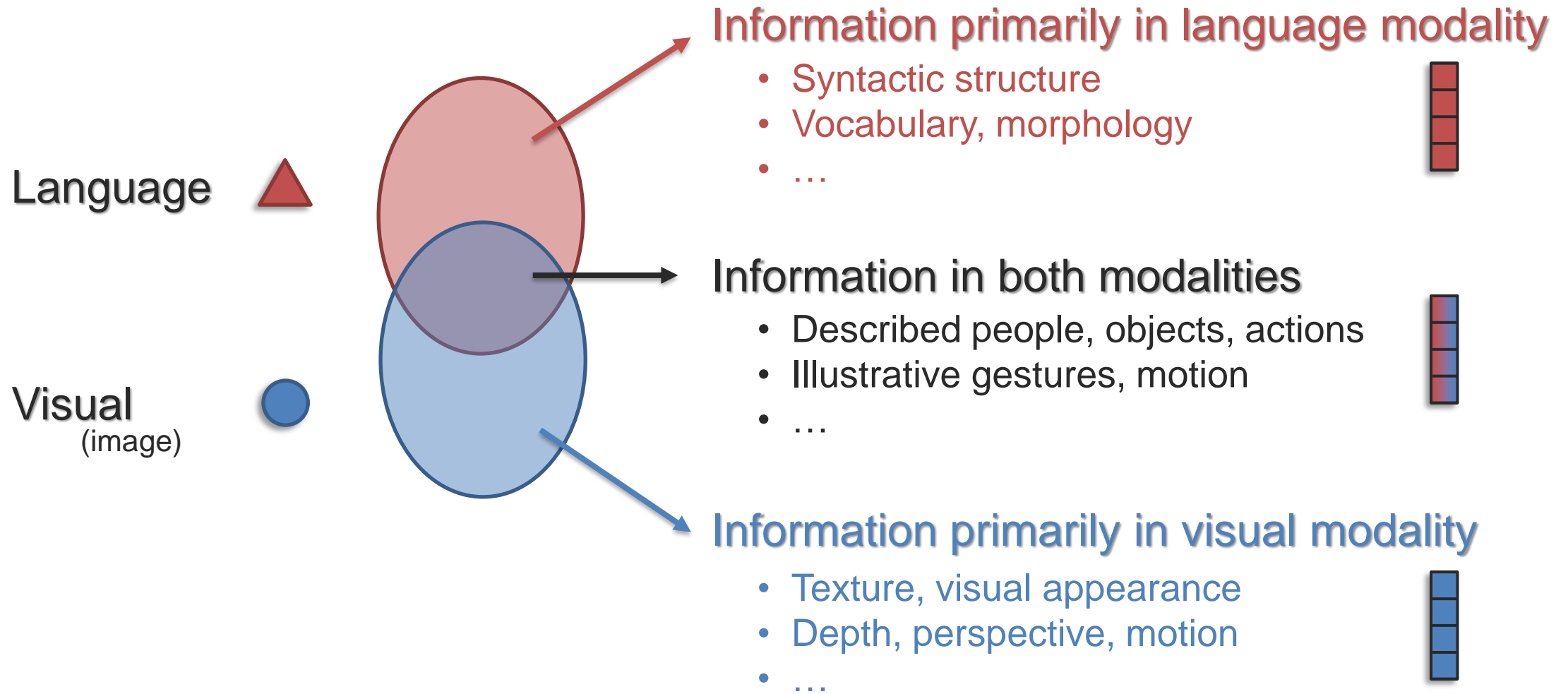
Modality-level fission:



Fine-grained fission:



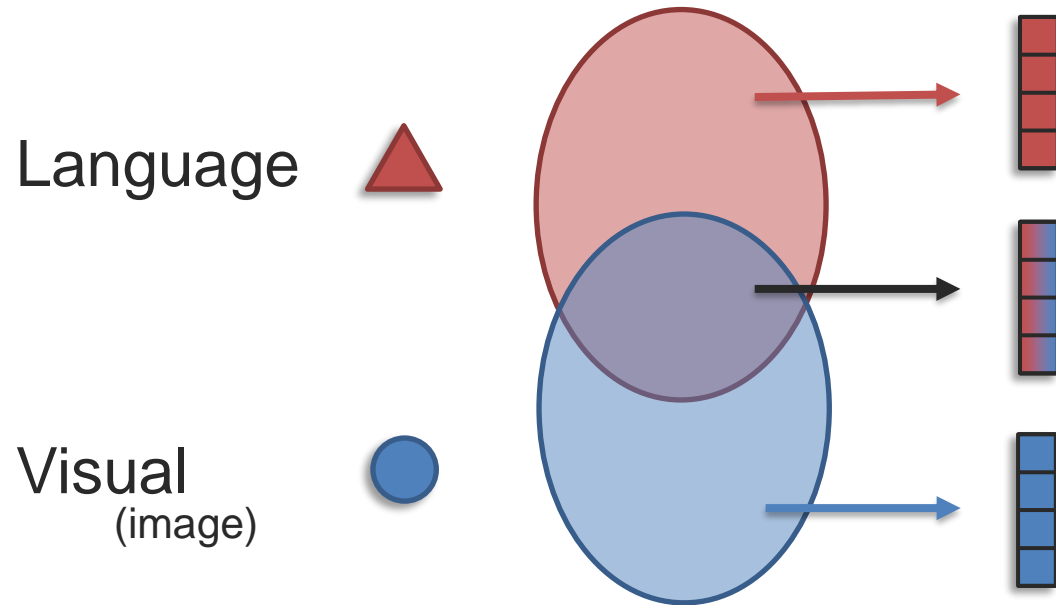
# Modality-Level Fission





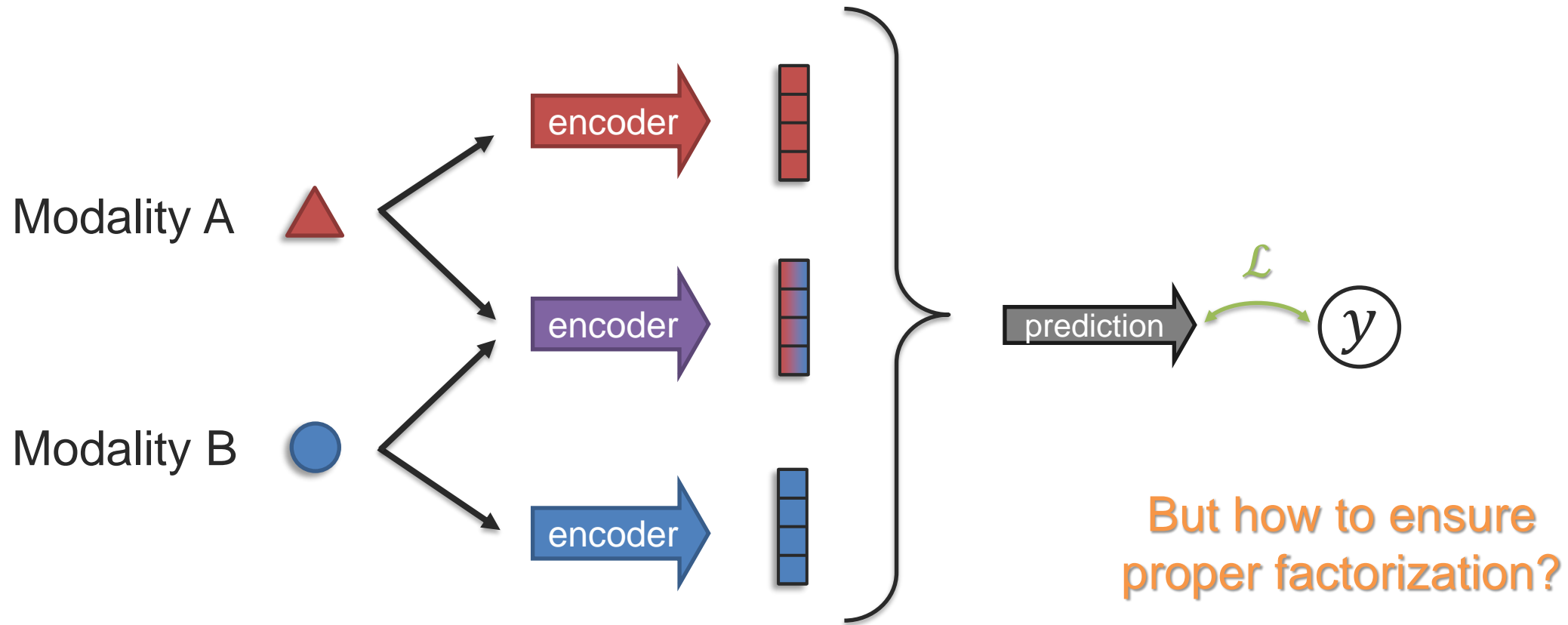
# Modality-Level Fission

---



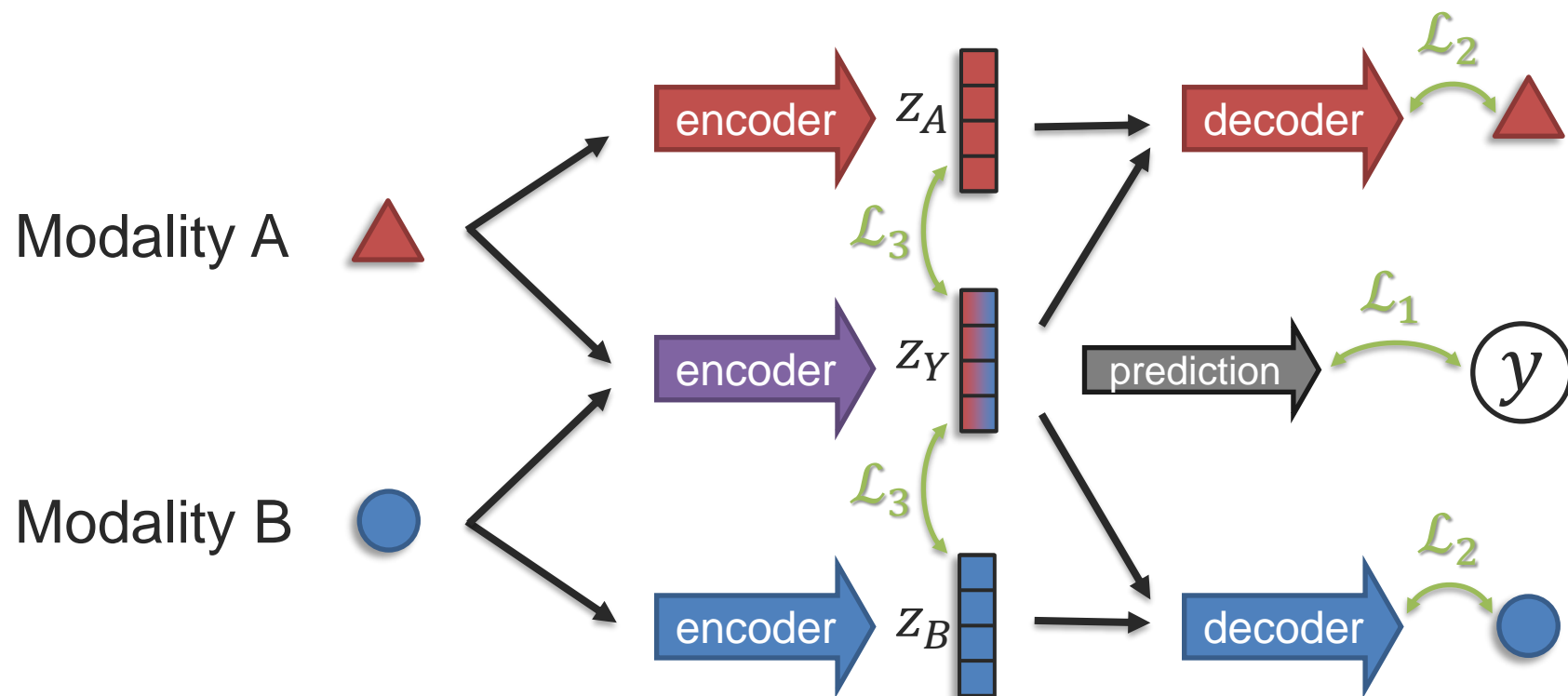
How to learn factorized multimodal representations?

# A Discriminative Approach – Factorized Multimodal Representations



But how to ensure proper factorization?

# A Generative-Discriminative Approach



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

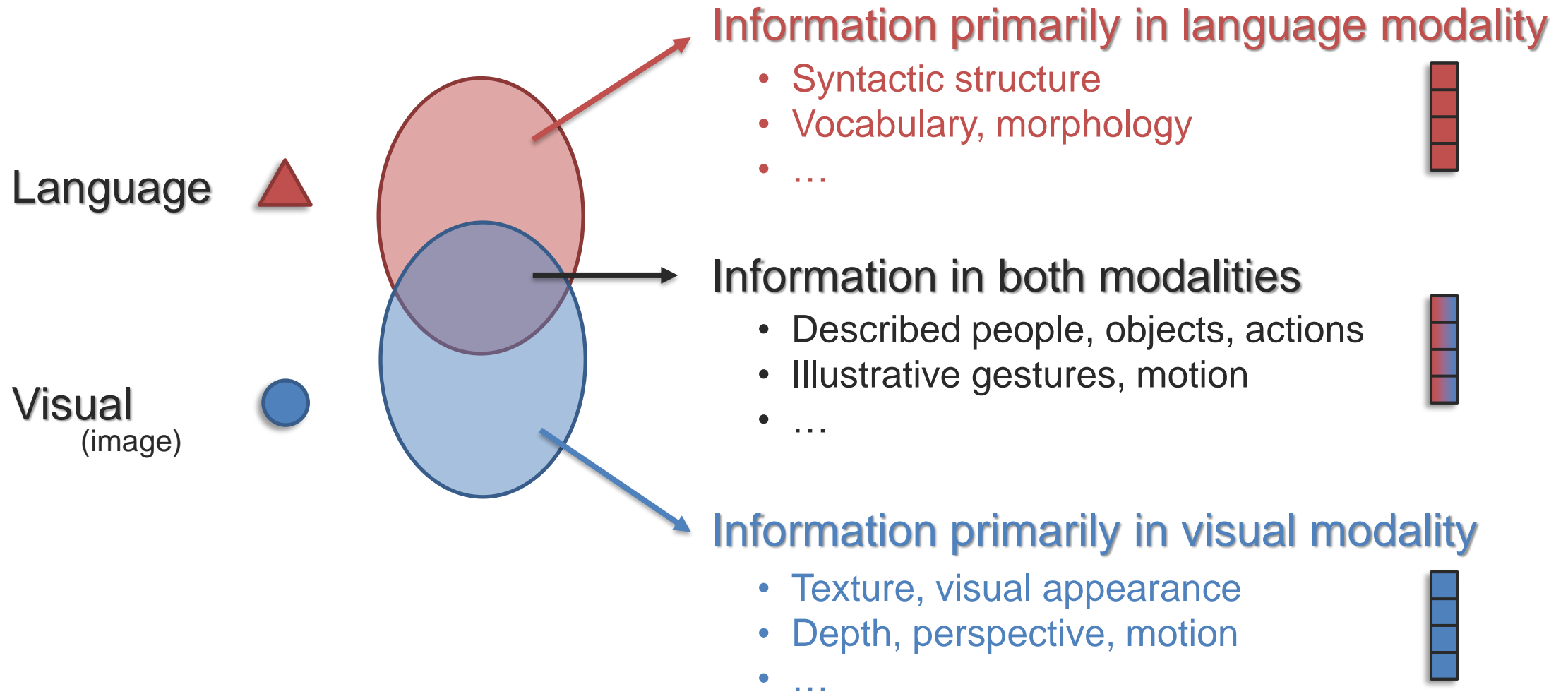
$\mathcal{L}_1$ : discriminative

$\mathcal{L}_2$ : generative

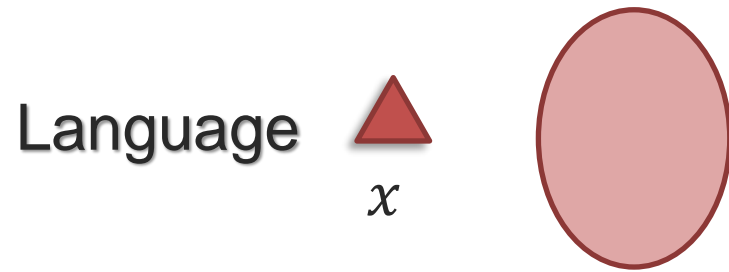
$\mathcal{L}_3$ : no overlap

Separate priors  
for  $z_A$ ,  $z_B$  and  $z_Y$

# Modality-Level Fission – 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 surprising its content is

$x$ : “12, 34, 45, 62 was not a winning combination”

➔ Not surprising... So, low information

$x$ : “11, 28, 38, 58 was a winning combination”

➔ Low chances... So, higher information

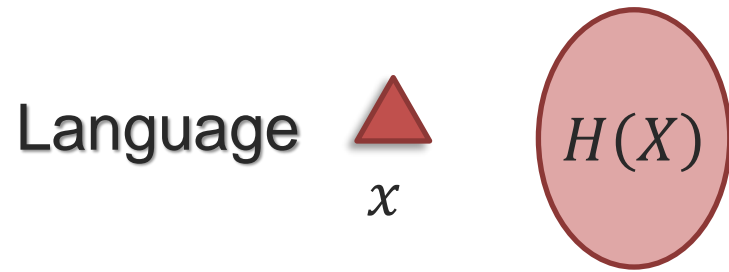
Information content  $I(x)$

$$I(x) \sim \frac{1}{p(x)} \quad \text{➔ But how to scale?}$$

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

# 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

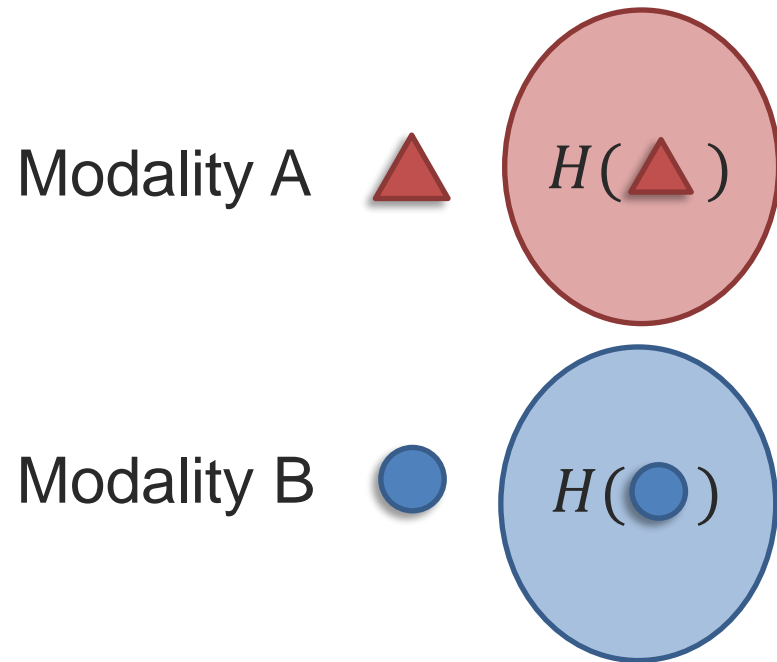
Entropy: weighted average of all possible outcomes from  $\mathcal{X}$

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

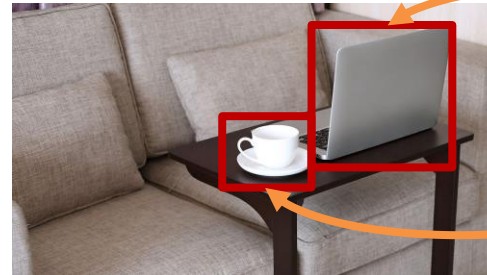
➔ Entropy can also be defined for continuous random variables

# Entropy with Two Modalities

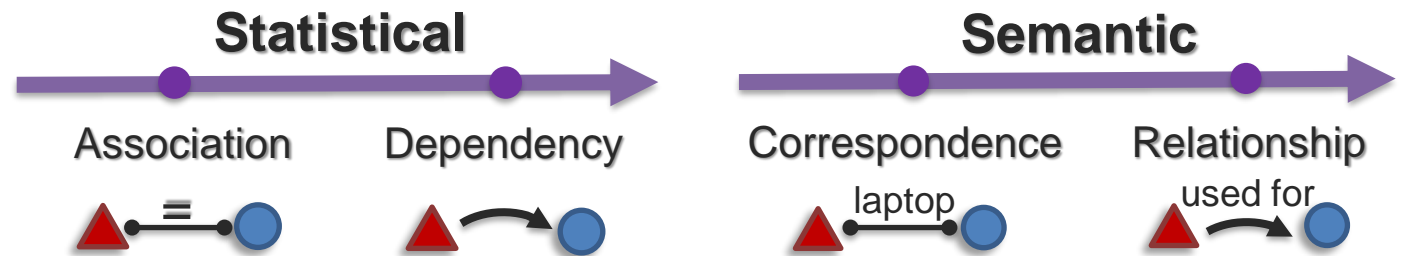
If no overlapping information



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

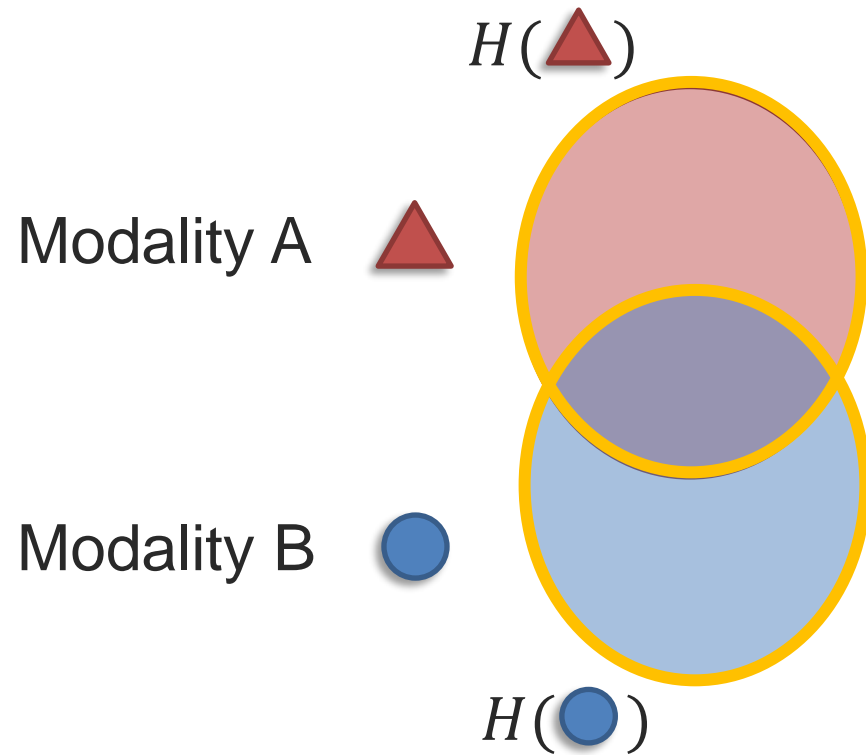


A **teacup** on the right of a **laptop** in a clean room.



# Entropy with Two Modalities

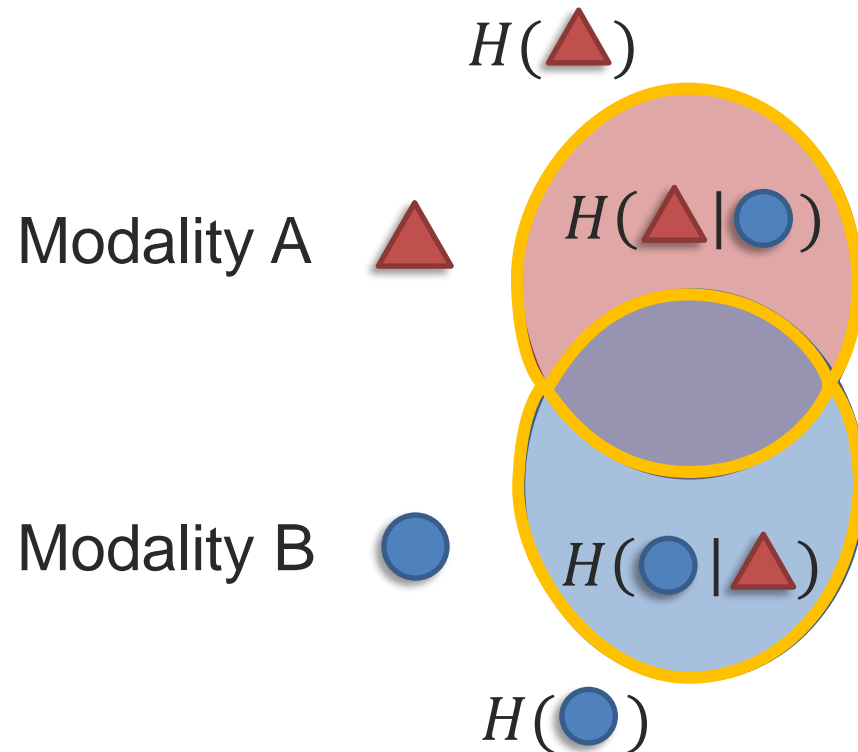
---





# Entropy with Two Modalities

---

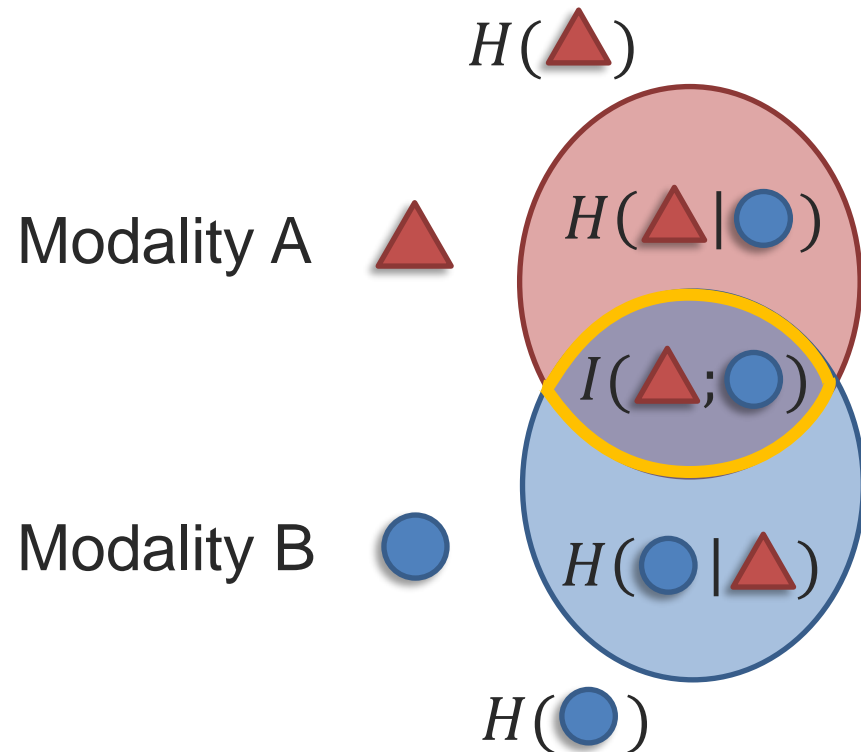


Conditional entropy  $H(Y|X)$

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

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

# Entropy with Two Modalities



Mutual information  $I(X; Y)$

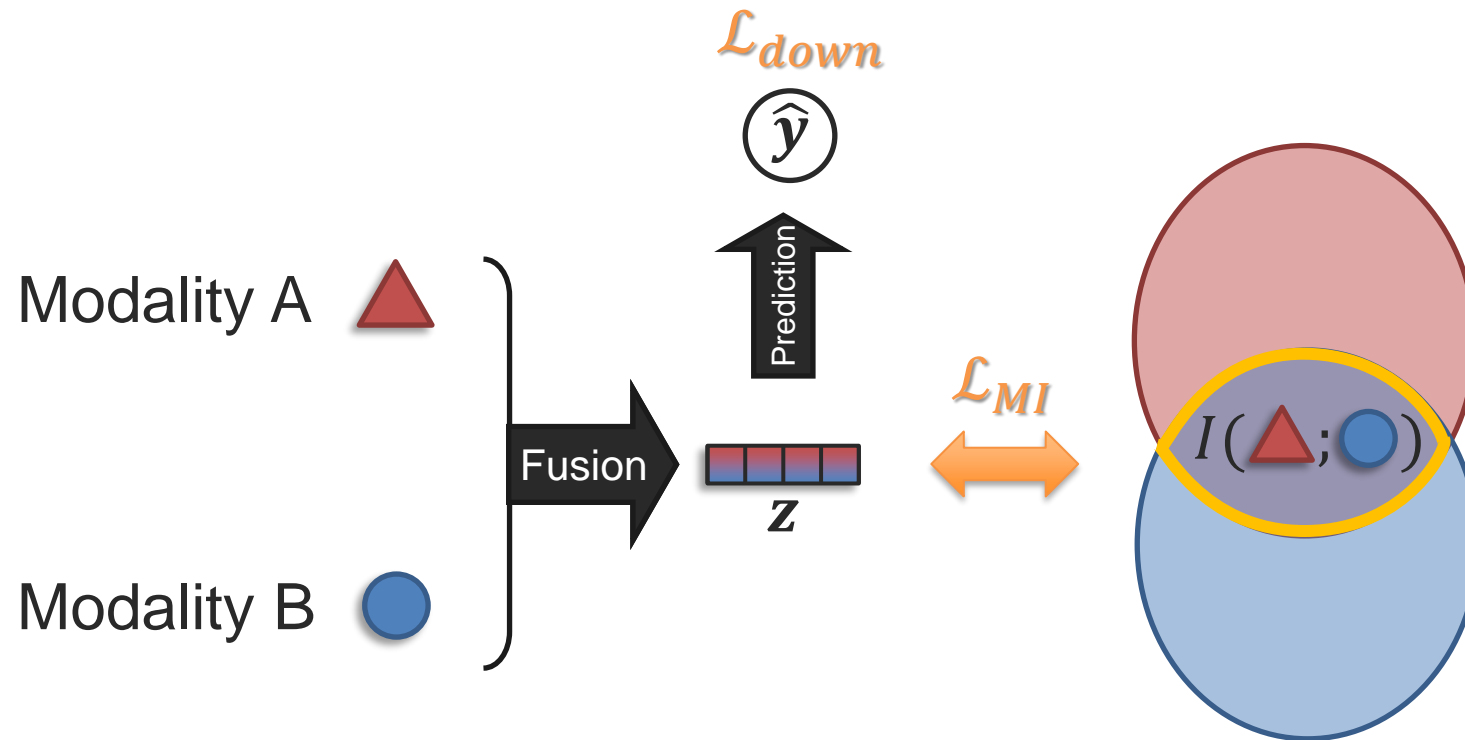
$$I(X; Y) = H(X) - H(X|Y)$$

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

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

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

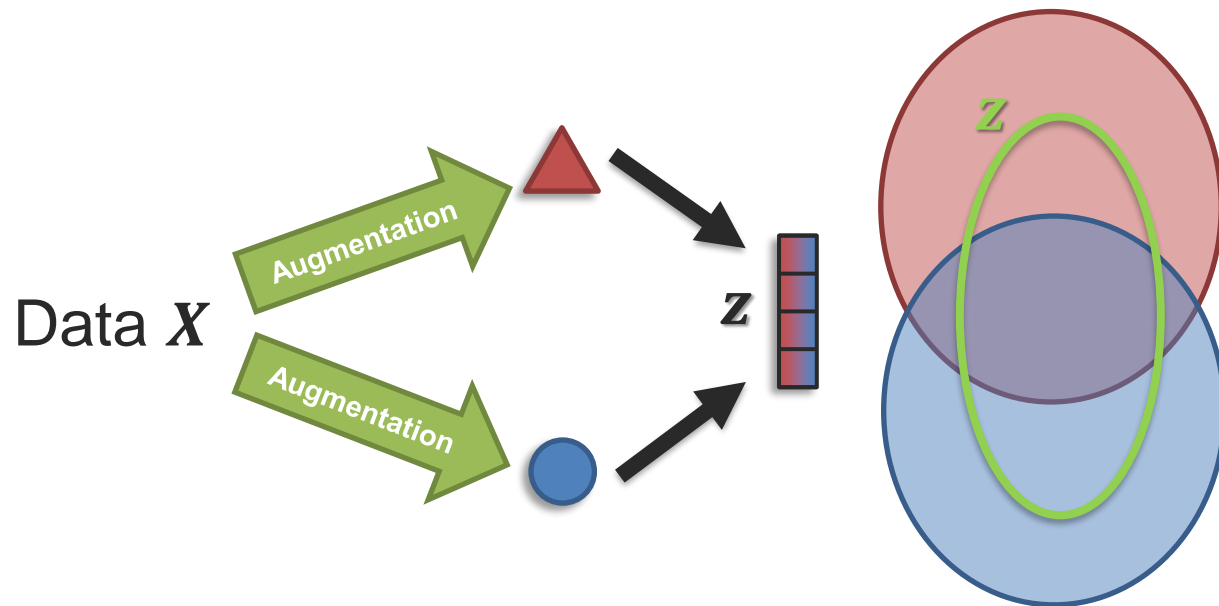
# Multimodal Fusion with Mutual Information



## Assumption?

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

# Link with Self-Supervised Learning



- 1 Maximize the mutual information

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

➔ Related to contrastive learning

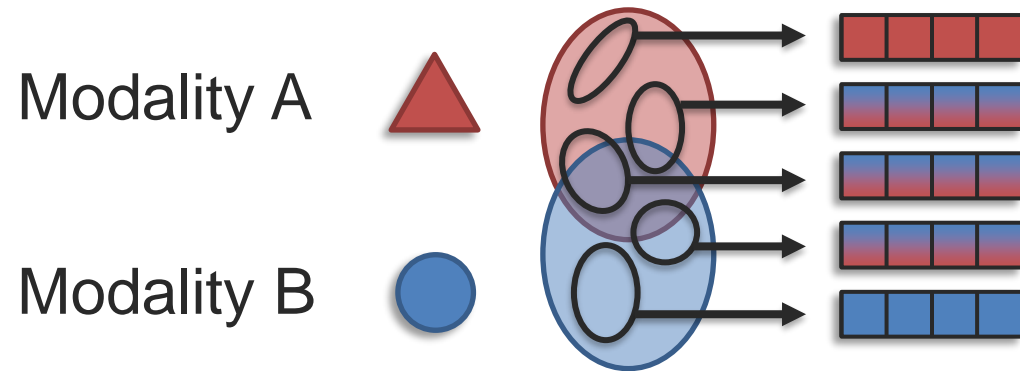
- 2 Minimize the conditional entropy

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

Information theory gives us a path towards  
disentangled representation learning

# Fine-Grained Fission

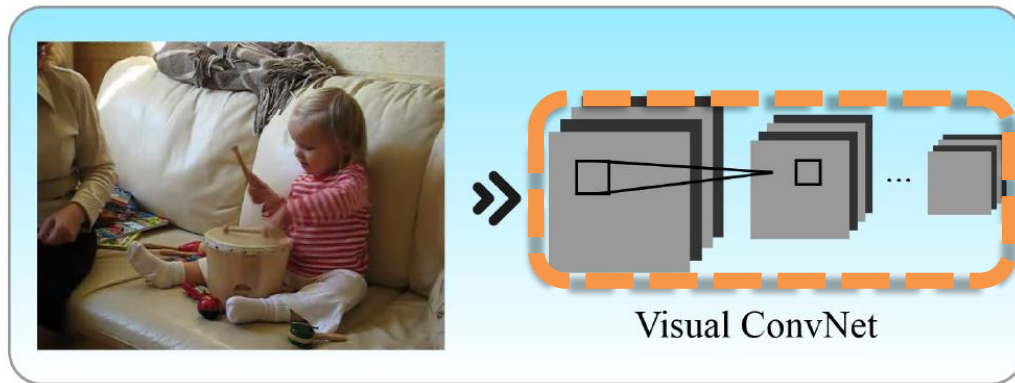
---



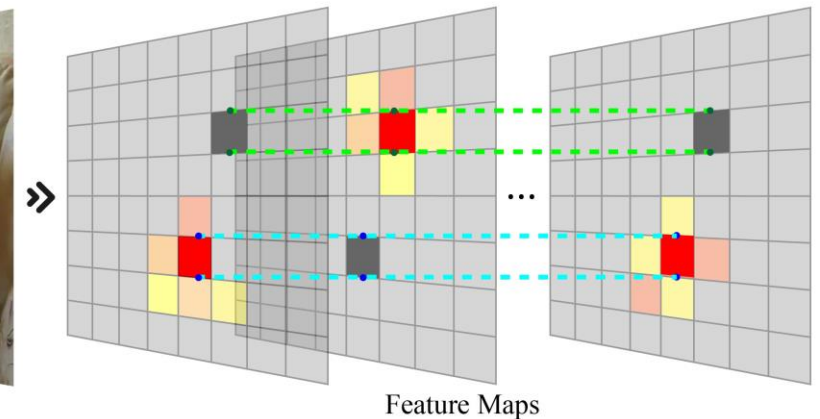
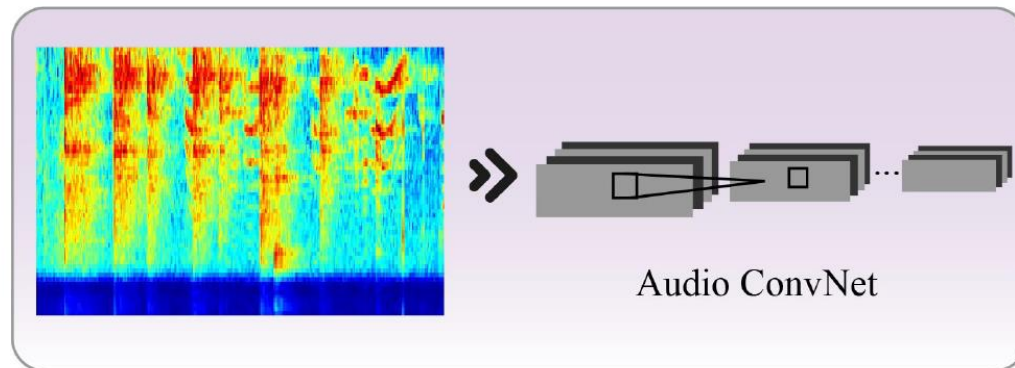
How to automatically discover these internal clusters, factors?

# Fine-Grained Fission – A Clustering Approach

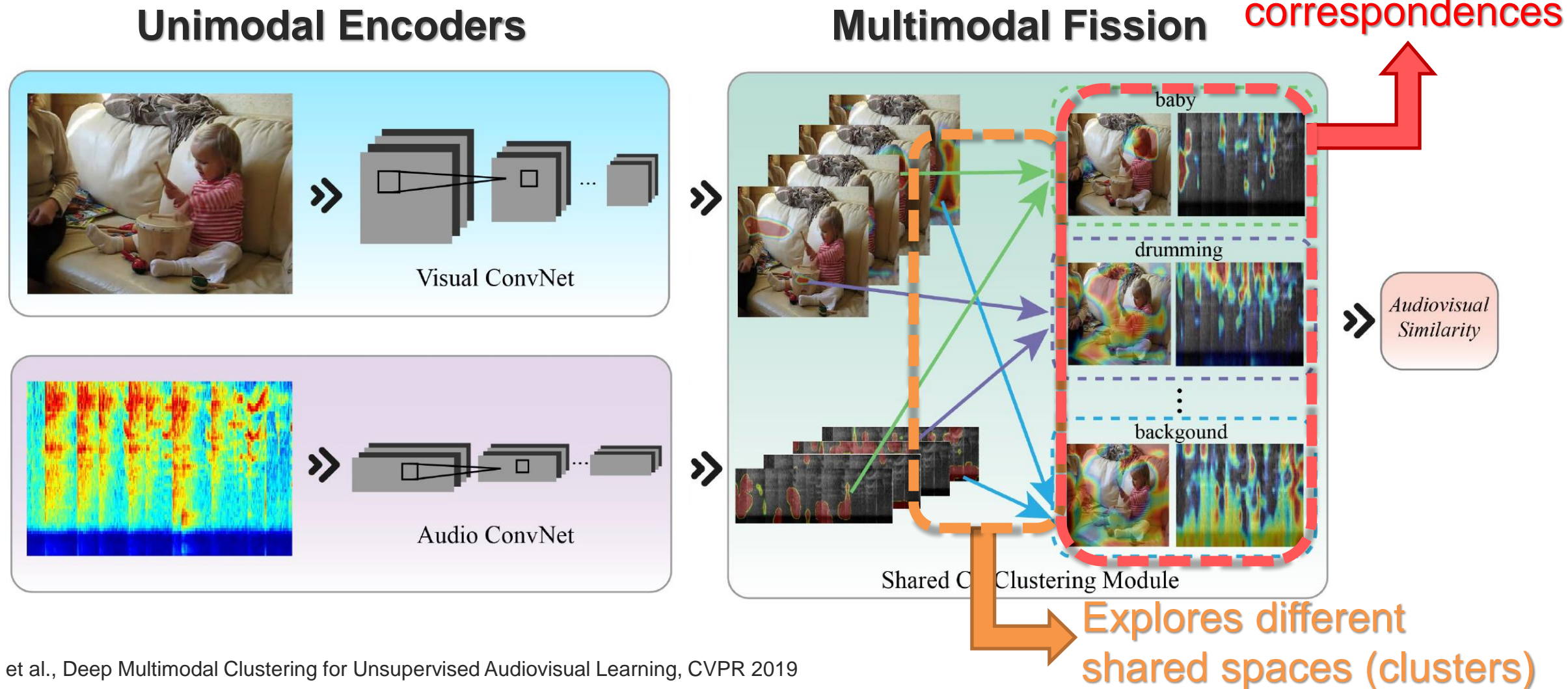
## Unimodal Encoders



Localized activations for different objects



# Fine-Grained Fission – A Clustering Approach



Hu et al., Deep Multimodal Clustering for Unsupervised Audiovisual Learning, CVPR 2019



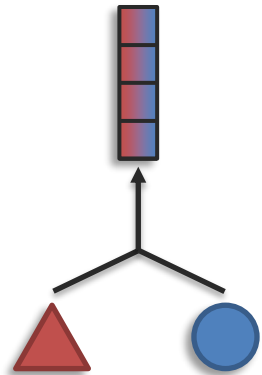
# Challenge 1: Representation

---

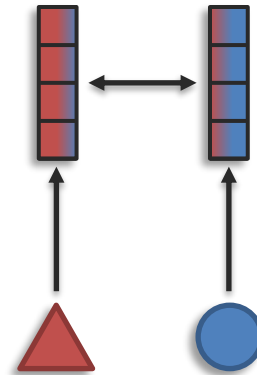
**Definition:** Learning representations that reflect cross-modal interactions between individual elements, across different modalities

## Sub-challenges:

### Fusion



### Coordination



### Fission

