



# 《多模态机器学习》

## 第五章 多模态表示

黄文炳

中国人民大学高瓴人工智能学院

[hwenbing@126.com](mailto:hwenbing@126.com)

2024年秋季

# 课程提纲

## 单模态表示

视觉模态

文本模态

三维点云

动作模态

## 基本概念

神经网络及其优化

## 经典多模态机器学习

多模态表示

多模态对齐

多模态推理

多模态生成

多模态迁移

## 通用多模态机器学习

通用多模态（大）模型

多模态预训练

多模态典型应用

# 内容提纲

---

- ① Cross-modal interactions
- ② Additive and multiplicative fusion
- ③ Gated fusion

# 内容提纲

---

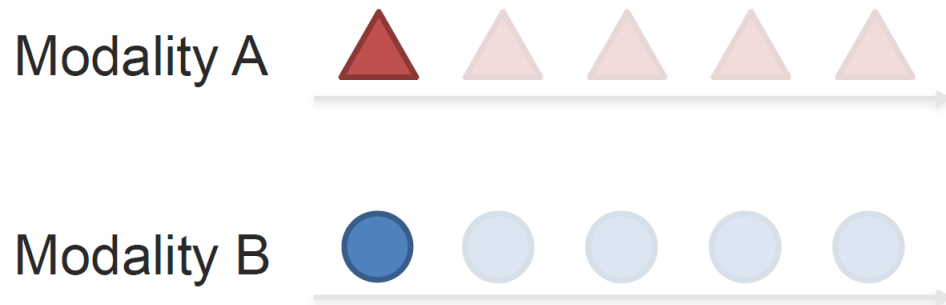
- ① Cross-modal interactions
- ② Additive and multiplicative fusion
- ③ Gated fusion

# Task 1: Representation (表示)

**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:**



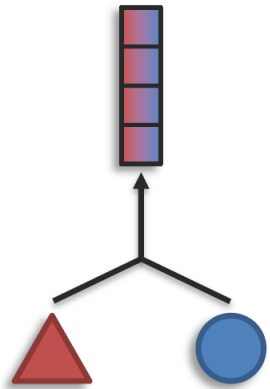
*It can be seen as a “local” representation  
or  
representation using holistic features*

# Task 1: Representation (表示)

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

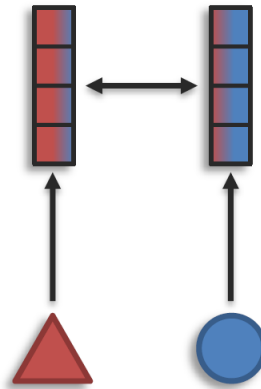
## Sub-challenges:

### Fusion



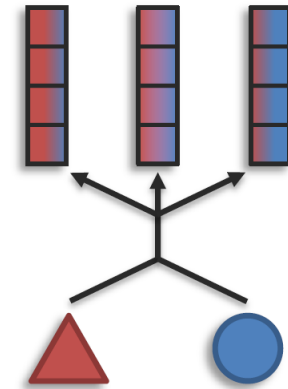
# modalities  $>$  # representations

### Coordination



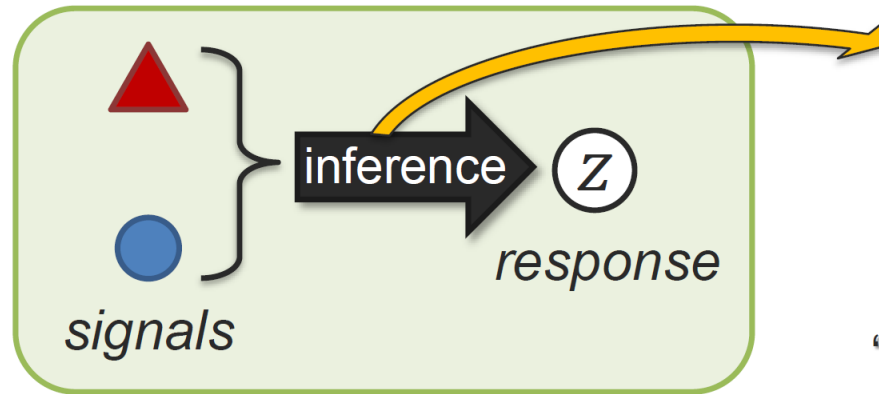
# modalities = # representations

### Fission

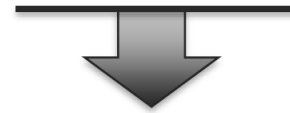


# modalities  $<$  # representations

# Cross-modal Interactions

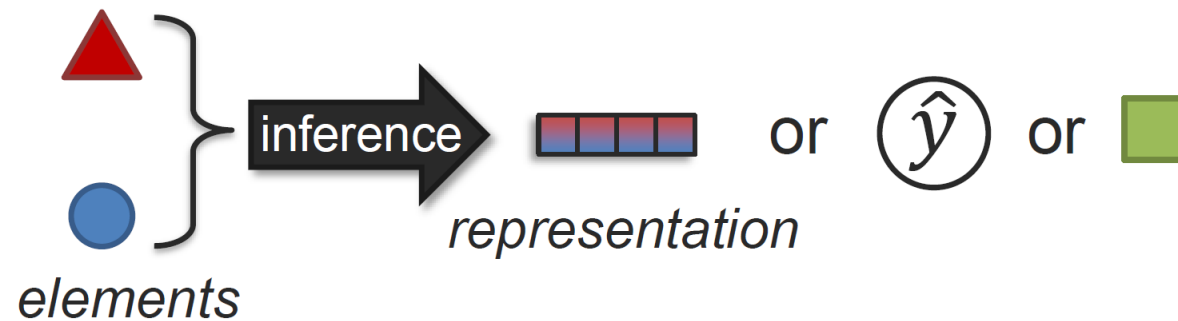


Interactions happen during inference!

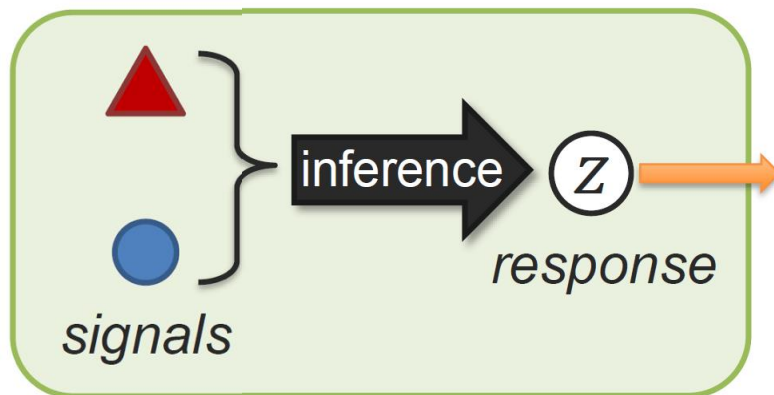


“Inference” examples:

- Representation fusion
- Prediction task
- Modality translation

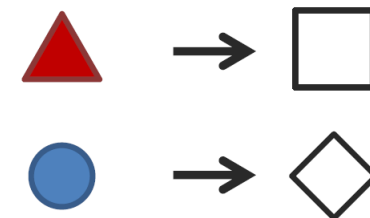


# Cross-modal Interactions



Types of interaction responses?  
(a taxonomy)

Unimodal  
Non-redundancy



***Is this a living room?***



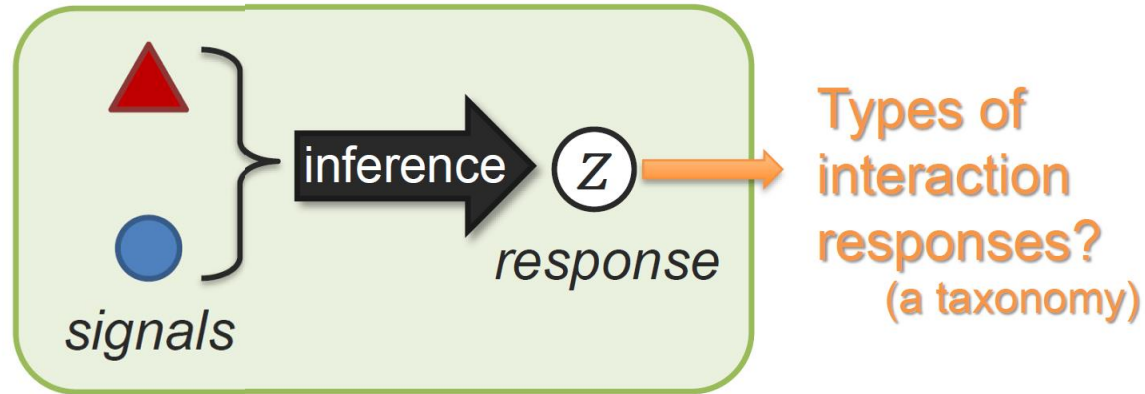
*A teacup on the right of a laptop in a clean room.*

**inference** → **Yes!**

**inference** → **No, probably study room.**



# Cross-modal Interactions



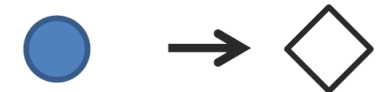
***Is this a living room?***



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

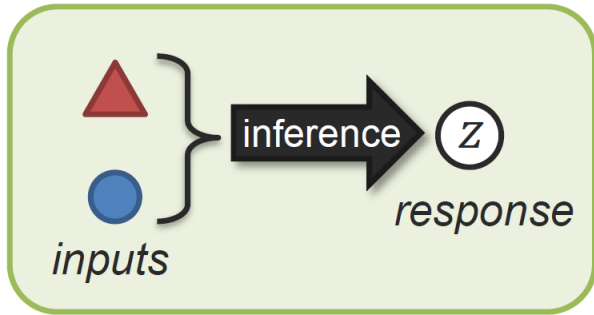
**inference** **Yes!**

Unimodal  
Non-redundancy



Multimodal  
dominance

# Taxonomy of Interaction Responses: A Behavioral Science View

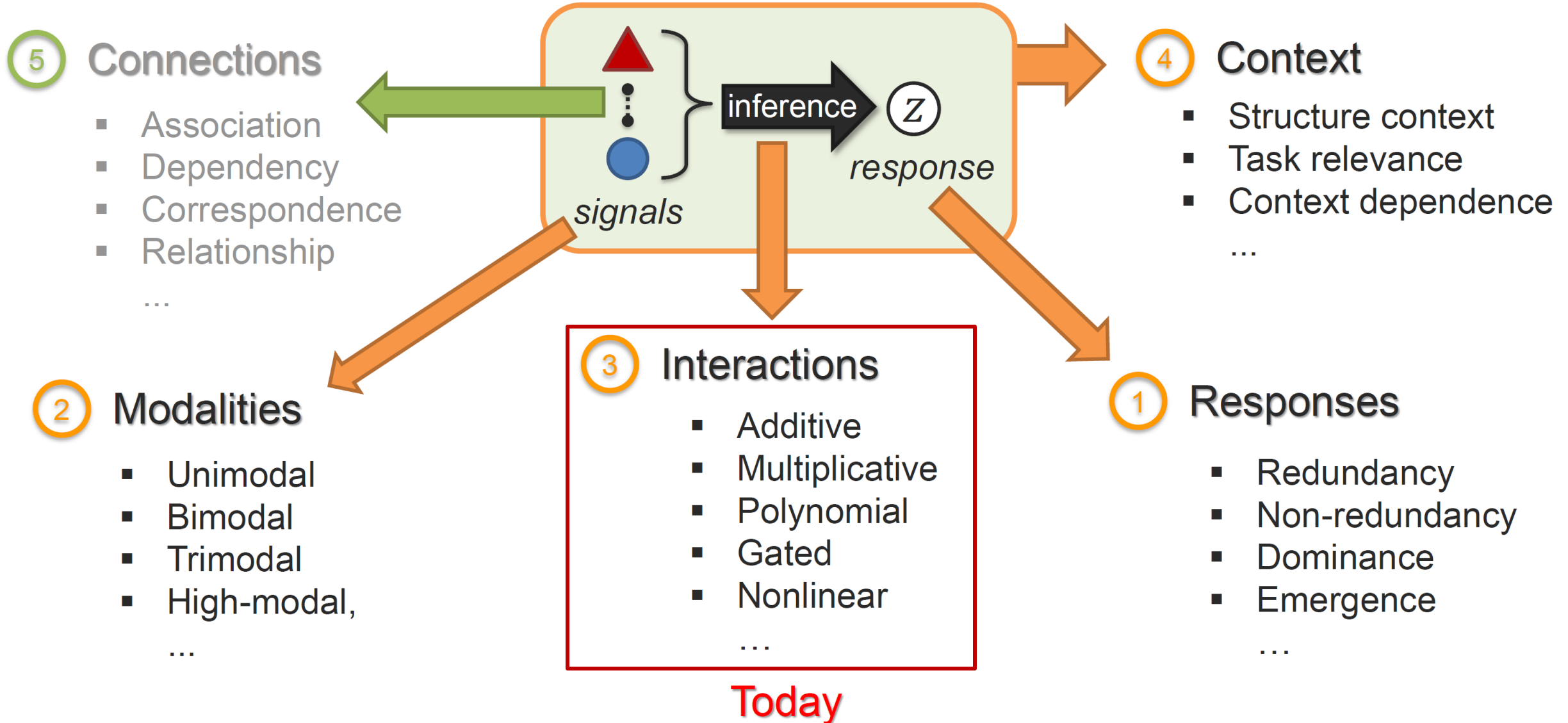


## Multimodal Communication



	signal	response	signal	response	
<b>Redundancy</b>	a	→ □	a+b	→ □	Equivalence
	b	→ □	a+b	→ □	Enhancement
<b>Nonredundancy</b>	a	→ □	a+b	→ □ and ○	Independence
	b	→ ○	a+b	→ □	Dominance
			a+b	→ □ (or □)	Modulation
			a+b	→ △	Emergence

# Cross-modal Interactions

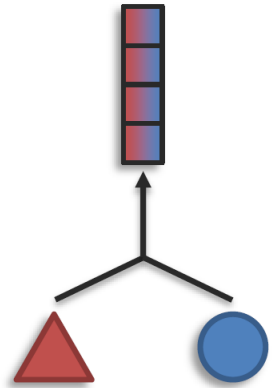


# 内容提纲

---

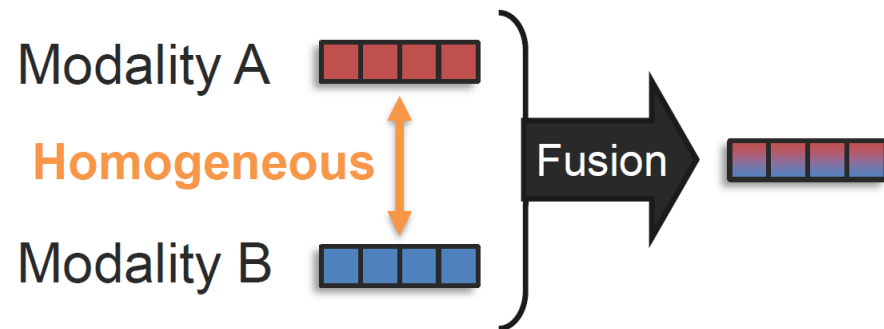
- ① Cross-modal interactions
- ② Additive and multiplicative fusion
- ③ Gated fusion

# Sub-Challenge 1a: Representation Fusion

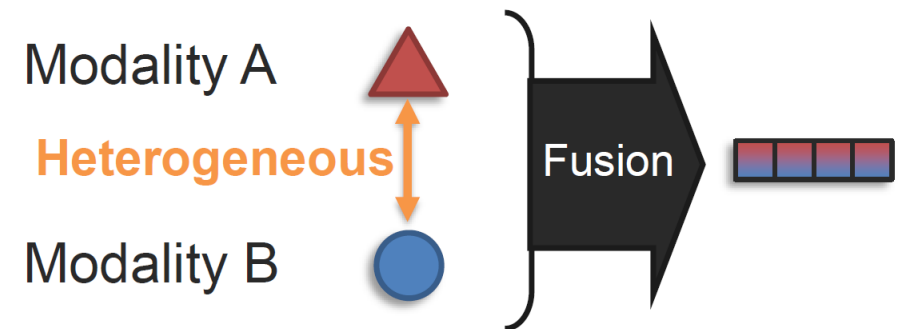


**Definition:** Learn a joint representation that models cross-modal interactions between individual elements of different modalities

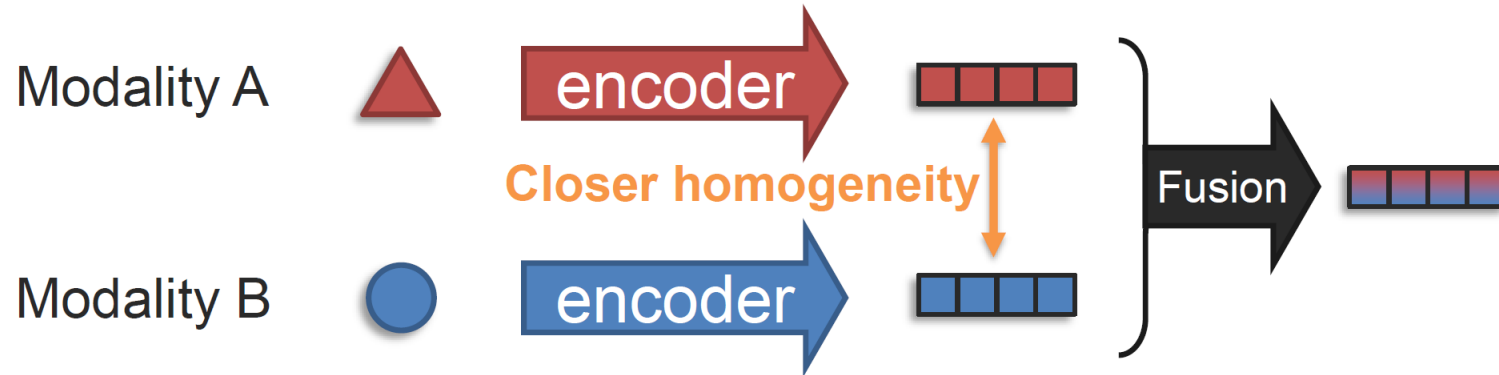
Basic fusion:



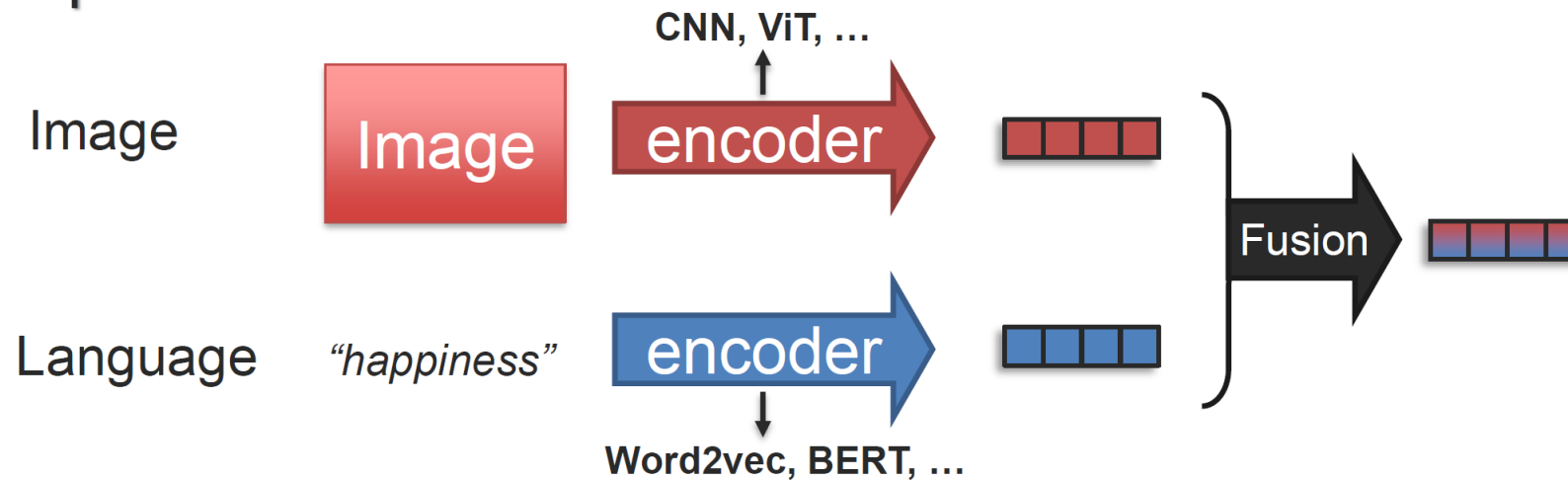
Raw-modality fusion:



# Fusion with Unimodal Encoders



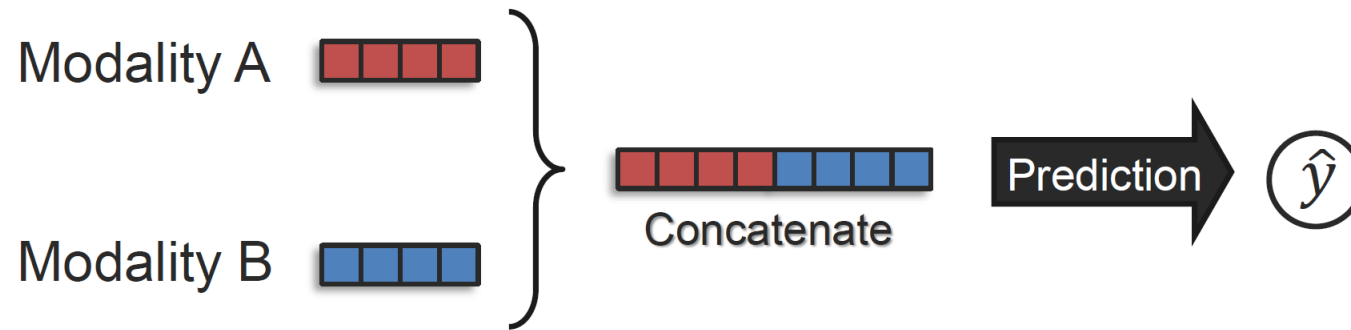
## Example:



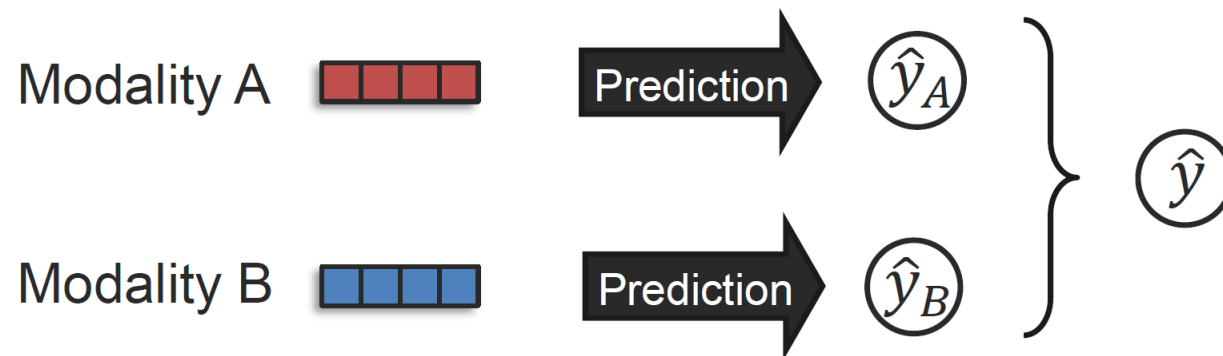
➡ Unimodal encoders can be jointly learned with fusion network, or pre-trained

# Early and Late Fusion – A historical View

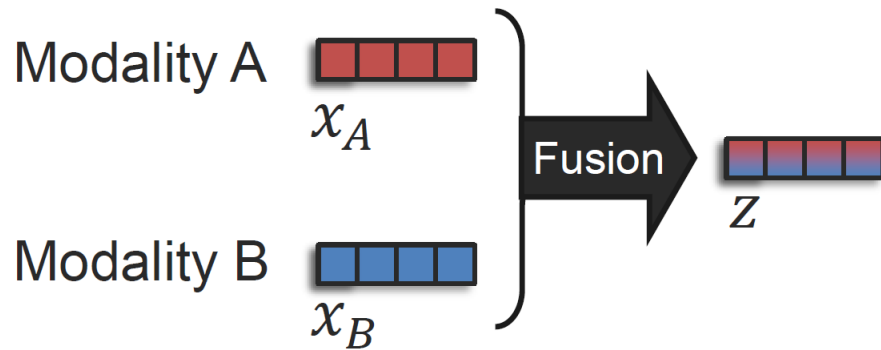
Early fusion:



Late fusion:



# Basic Concepts for Representation Fusion (aka, Basic Fusion)



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

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

Linear regression:

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

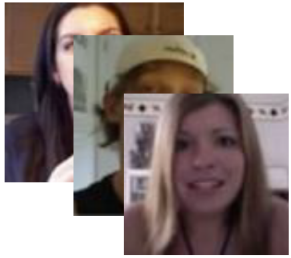
intercept (bias term)      Additive terms      Multiplicative term      error (residual term)



# Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

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:

$$y = w_0 + w_1x_A + w_2x_B + w_3(x_A \times x_B) + \epsilon$$

intercept (bias term)      Additive terms      Multiplicative term      error (residual term)

$w_0$ : average score when  $x_A$  and  $x_B$  are zero

$w_1$ : effect from  $x_A$  variable only

$w_2$ : effect from  $x_B$  variable only

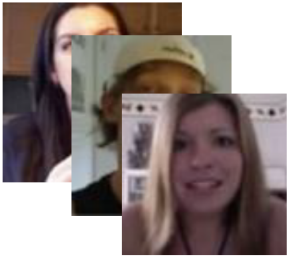
$w_3$ : effect from  $x_A$  and  $x_B$  interaction only

$\epsilon$ : residual not modeled by  $w_0$ ,  $w_1$ ,  $w_2$  or  $w_3$

# Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

300 book reviews



$y$ : audience score

$x_A$ : percentage of smiling

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

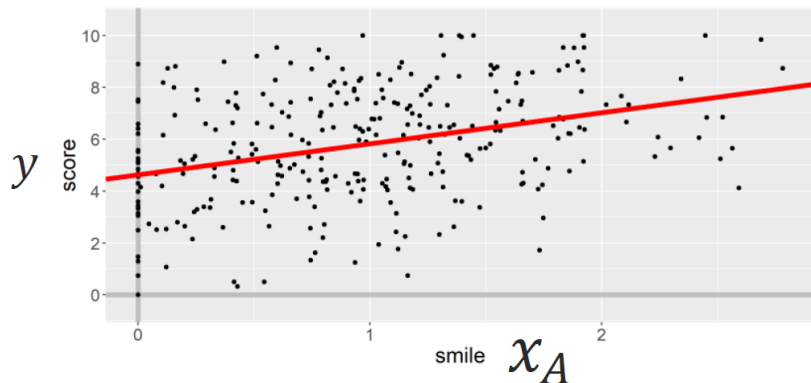
*H1: Does smiling reveal what the audience score was?*

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

Linear regression:

$$Z = W_0 + \boxed{W_1}x_A + \epsilon$$

slope



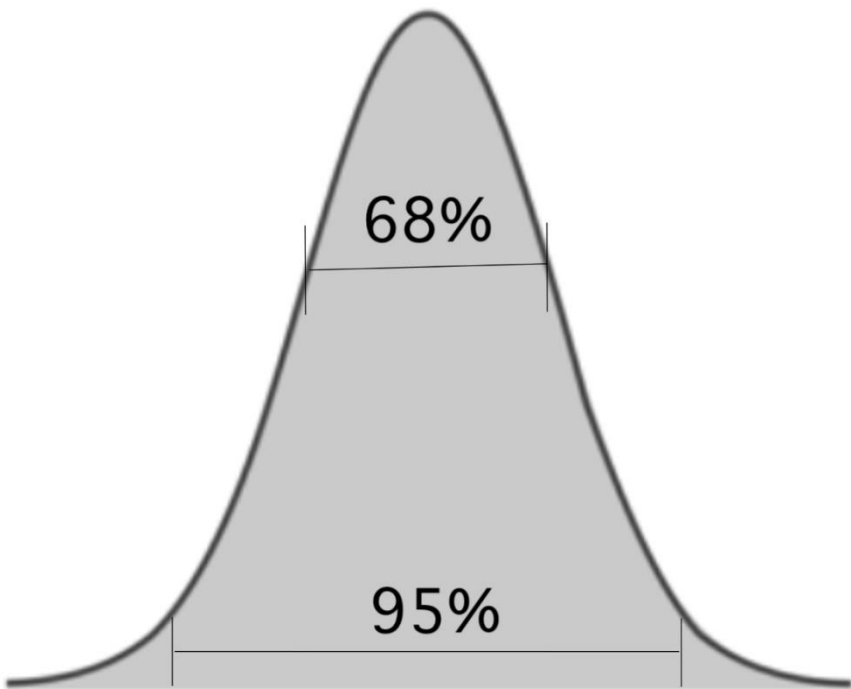
Confidence interval: “95% confident that  $w$  parameter is contained within this interval”

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

p-values would be another way to test hypothesis

Confidence interval does not contain 0, so effect is significant

# Confidence Interval



## 1. 已知总体标准差的情况

当总体标准差已知时，使用 **Z 分布** 来计算。

公式：

$$CI = \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

- $\bar{X}$ : 样本均值
- $Z_{\alpha/2}$ : Z 值，对应给定的置信水平（例如，95% 置信水平时， $Z_{\alpha/2} = 1.96$ ）
- $\sigma$ : 总体标准差
- $n$ : 样本大小

## 2. 未知总体标准差的情况

当总体标准差未知时，使用 **t 分布** 来计算。

公式：

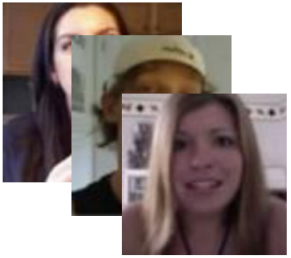
$$CI = \bar{X} \pm t_{\alpha/2, df} \cdot \frac{s}{\sqrt{n}}$$

- $\bar{X}$ : 样本均值
- $t_{\alpha/2, df}$ : t 值，基于样本大小  $n$  的自由度  $df = n - 1$  查表

# Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

300 book reviews



$y$ : audience score

$x_A$ : percentage of smiling

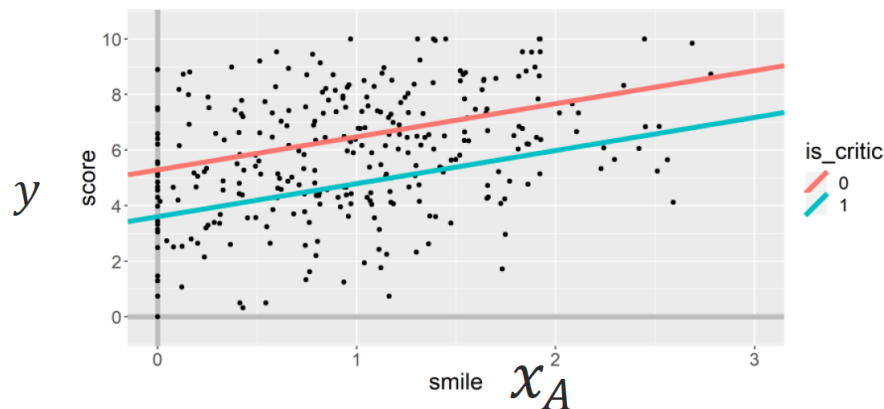
$x_B$ : professional status  
(0=non-critic, 1=critic)

*H1: Does smiling reveal what the audience score was?*

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

Linear regression:

$$Z = W_0 + W_1 x_A + W_2 x_B + \epsilon$$



	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]

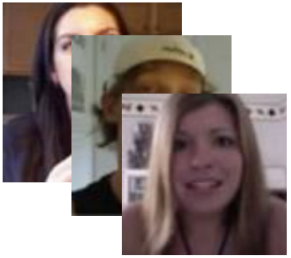
➔ Positive effect

➔ Negative effect

# Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

300 book reviews



$y$ : audience score

$x_A$ : percentage of smiling

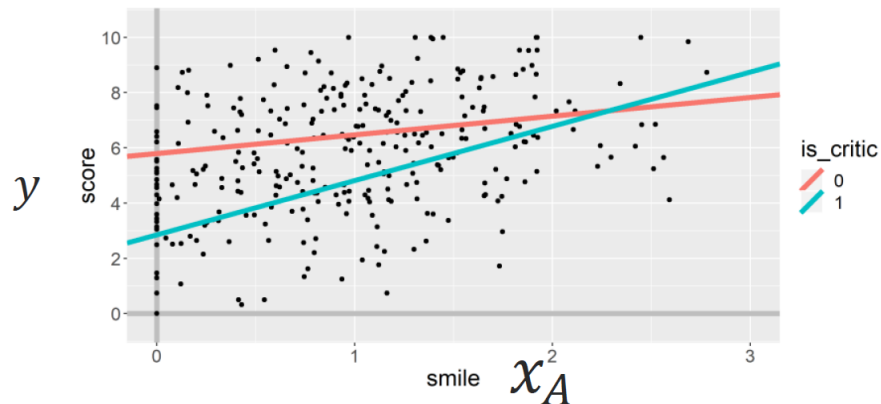
$x_B$ : professional status  
(0=non-critic, 1=critic)

*H1: Does smiling reveal what the audience score was?*

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

Linear regression:

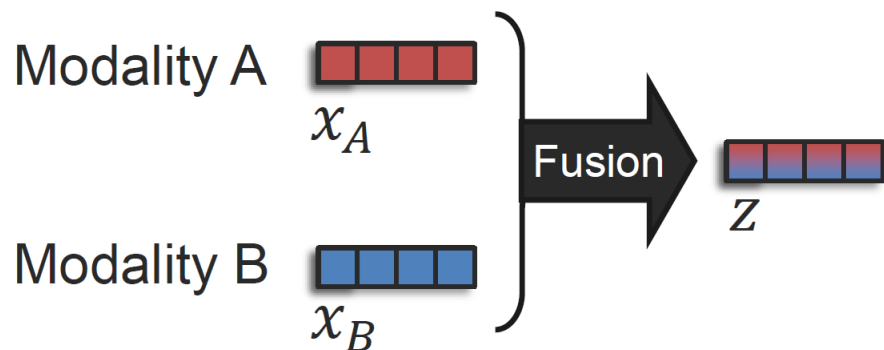
$$Z = W_0 + W_1 x_A + W_2 x_B + W_3 (x_A \times x_B) + \epsilon$$



	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 Concepts for Representation Fusion (aka, Basic Fusion)



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

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

Linear regression:

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

intercept (bias term)      Additive terms      Multiplicative term      error (residual term)

① Additive terms:

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

② Multiplicative “interaction” term:

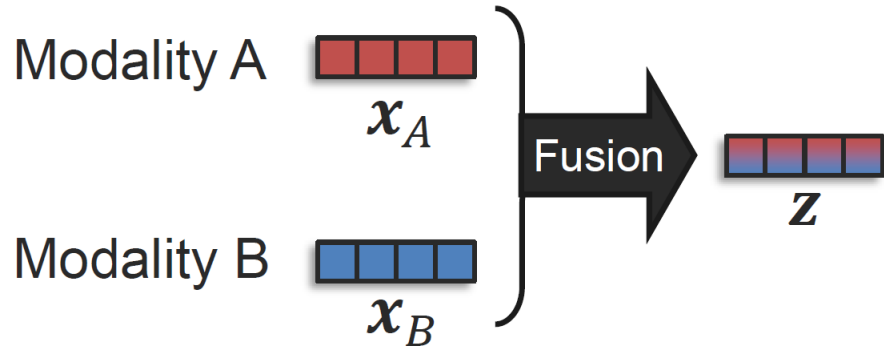
$$Z = w_3 (x_A \times x_B) + \epsilon$$

③ Additive and multiplicative terms:

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

# Additive Fusion

➔ Back to multivariate case!

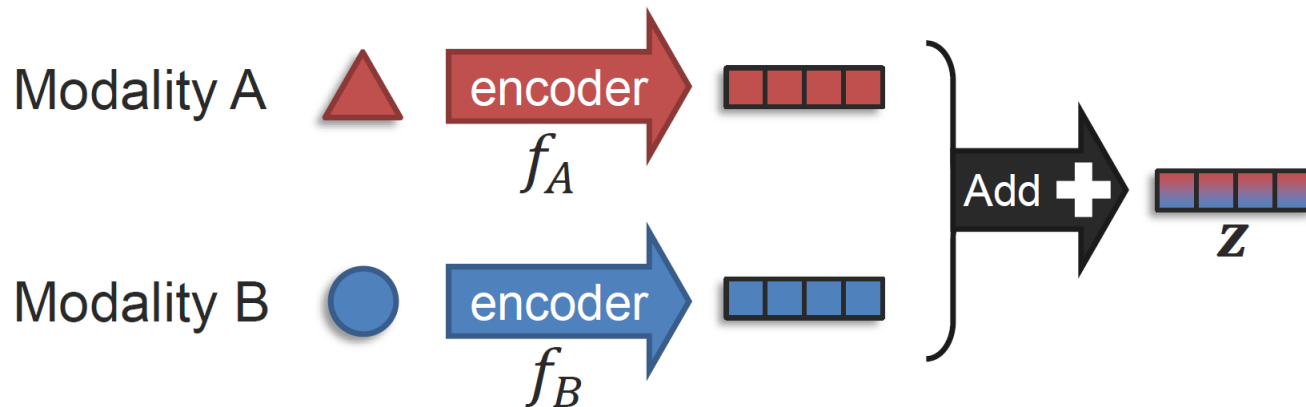


Additive fusion:

$$z = W_1 x_A + W_2 x_B$$

➔ 1-layer neural network can be seen as additive

With unimodal encoders:

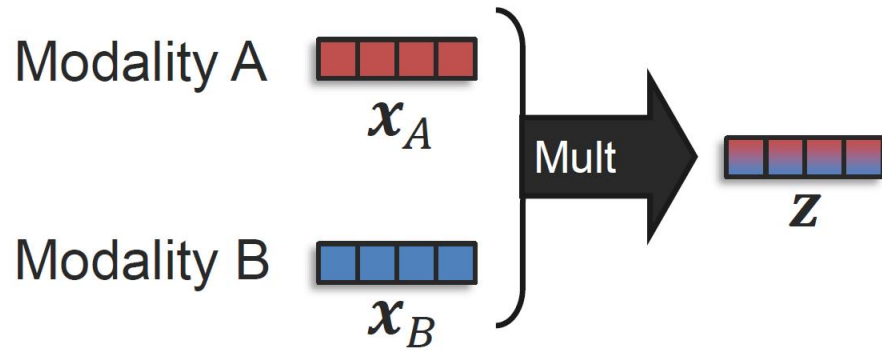


Additive fusion:

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

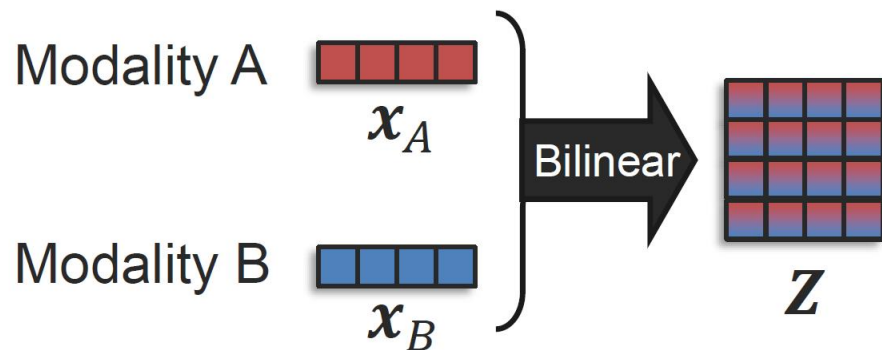
➔ It could be seen as an ensemble approach (late fusion)

# Multiplicative Fusion



Simple multiplicative fusion:

$$\mathbf{z} = \mathbf{x}_A \odot \mathbf{x}_B$$



Bilinear Fusion:

$$\mathbf{Z} = \mathbf{x}_A^\top \mathbf{x}_B$$

$$\text{vec}(\mathbf{Z}) = \mathbf{x}_A \otimes \mathbf{x}_B$$



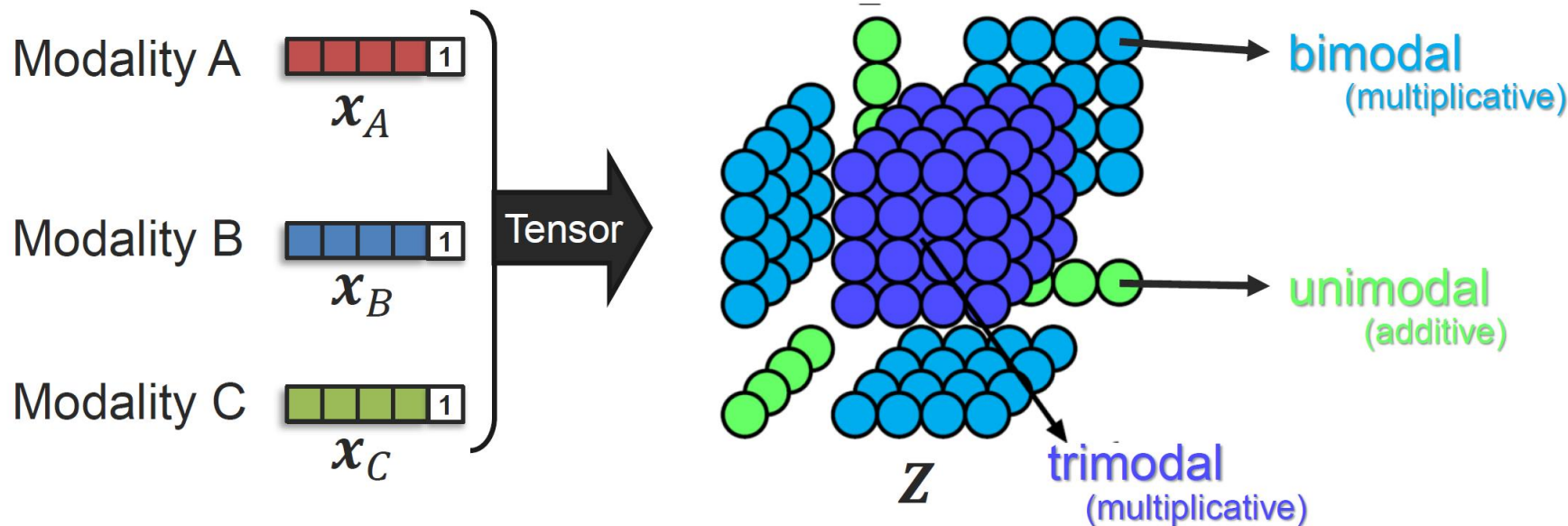
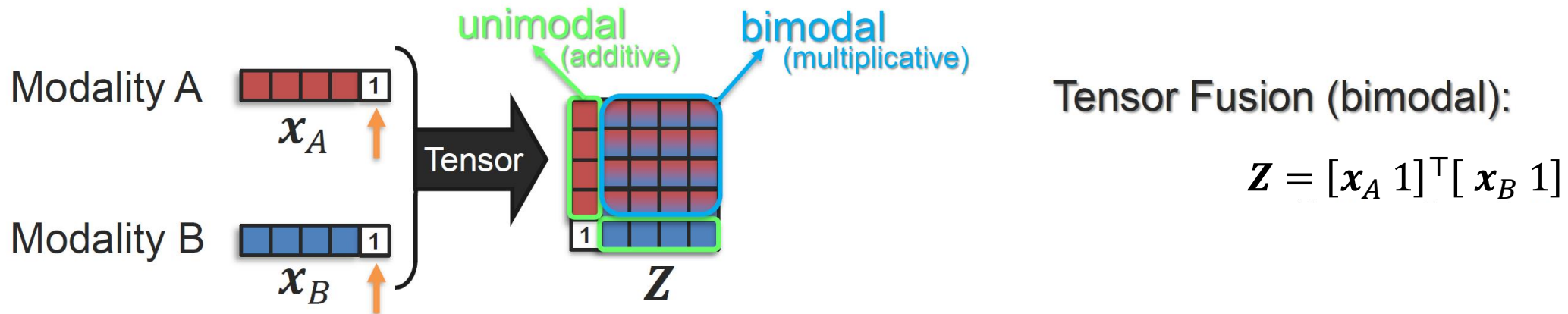
# Kronecker product

---

If  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{p \times q}$ , then the Kronecker product  $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{pm \times qn}$ .

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11} \mathbf{B} & \cdots & a_{1n} \mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1} \mathbf{B} & \cdots & a_{mn} \mathbf{B} \end{bmatrix}$$

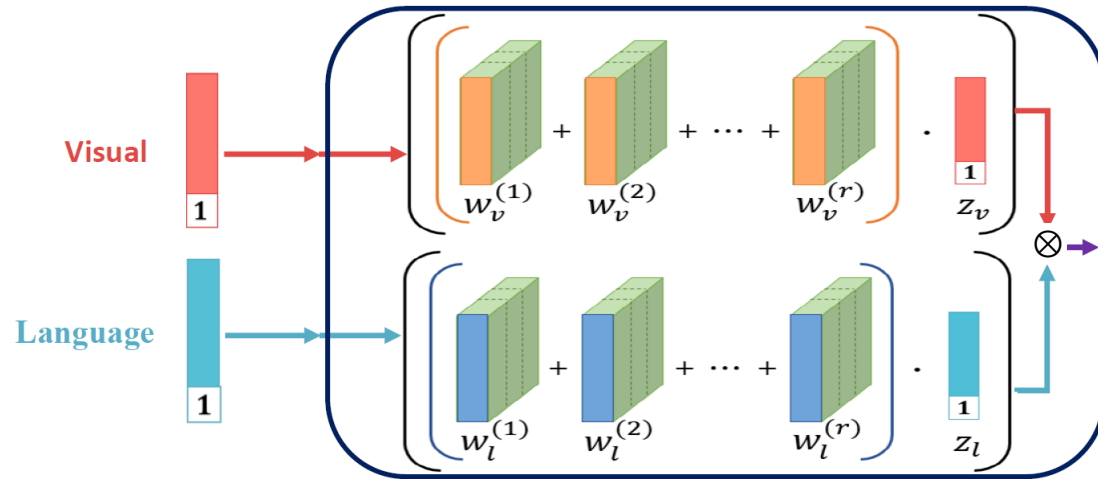
# Multiplicative Fusion



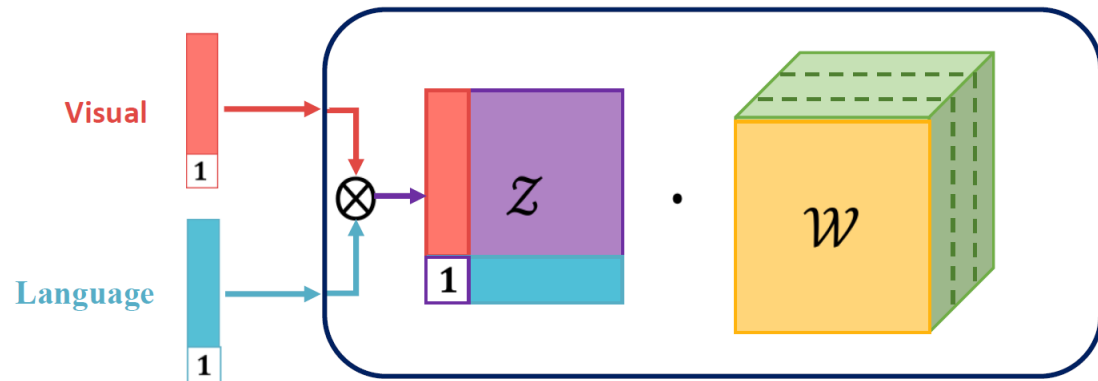
... but the weight matrix may end up quite large!

# Low-rank Fusion

## Low-rank Fusion



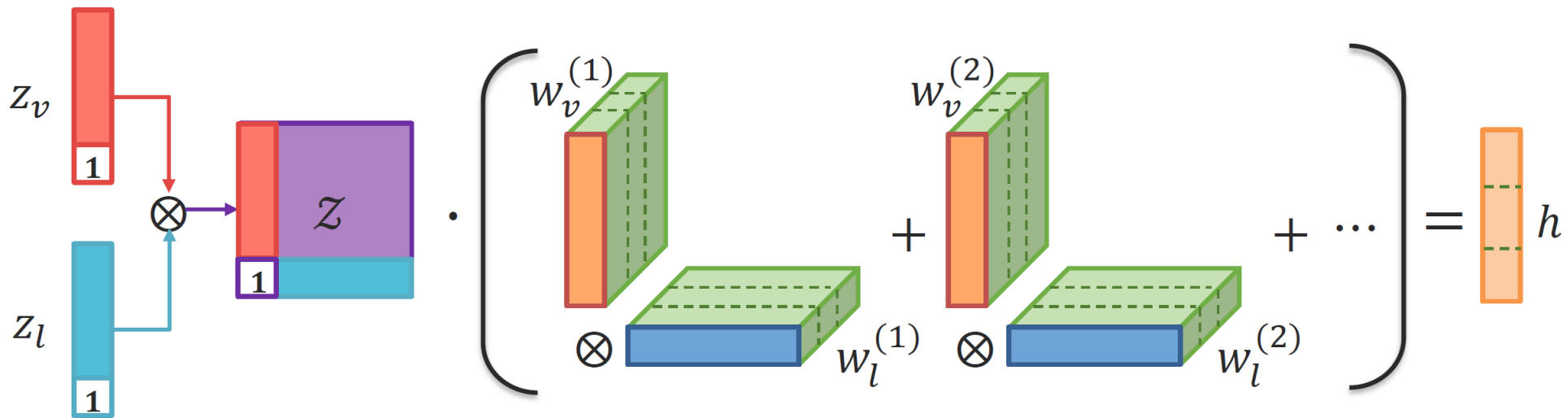
## Tensor Fusion



- ① Decomposition of weight  $W$ .
- ② Decomposition of input tensor  $Z$ .
- ③ Rearrange the computation of  $h$ .

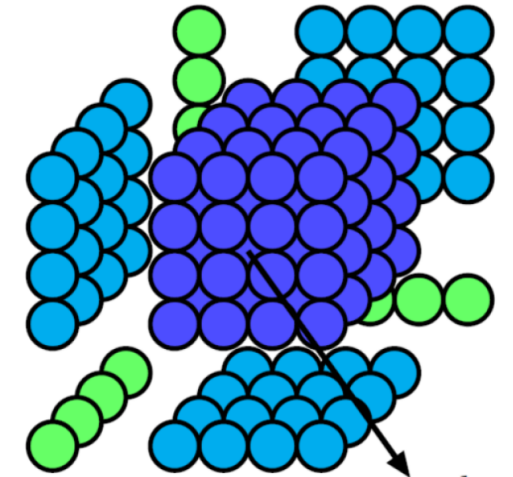
# Low-rank Fusion

$$\begin{aligned} & \text{vec}(\mathbf{Z}) \cdot \text{vec}(\mathbf{W}) \\ &= (\mathbf{z}_v \otimes \mathbf{z}_l) \cdot \left( \sum_{k=1}^r \mathbf{w}_v^{(k)} \otimes \mathbf{w}_l^{(k)} \right) \\ &= \sum_{k=1}^r \mathbf{z}_v \cdot \mathbf{w}_v^{(k)} \otimes \mathbf{z}_l \cdot \mathbf{w}_l^{(k)} \end{aligned}$$

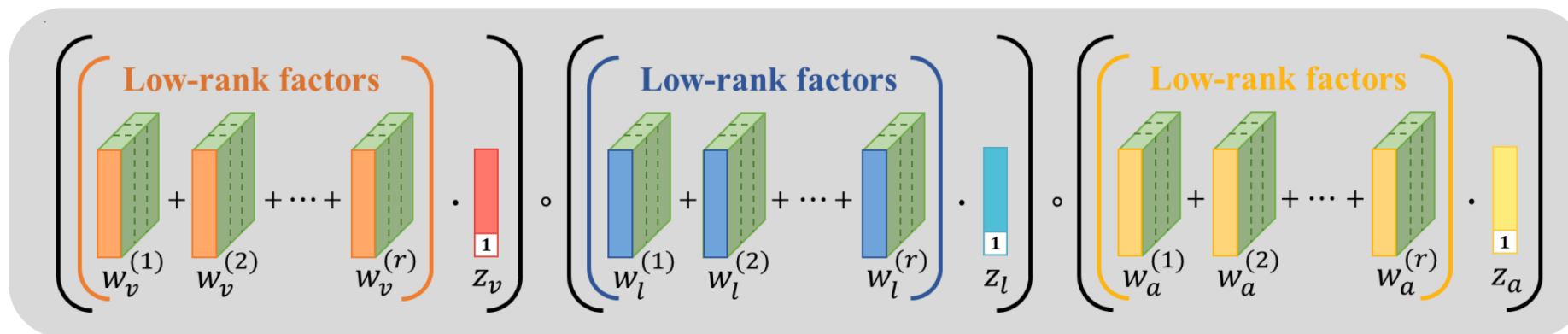


# Low-rank Fusion with Trimodal Input

Tensor Fusion



**Low-rank Fusion :**



Canonical Polyadic Decomposition

# Going Beyond Additive and Multiplicative Fusion

Additive interaction:

$$z = w_1 x_A + w_2 x_B$$

← First-order polynomial

Additive and multiplicative interaction:

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

← Second-order polynomial

Trimodal fusion (e.g., tensor fusion):

$$z = \underbrace{w_1 x_A + w_2 x_B + w_3 x_C}_{\substack{\text{Unimodal terms} \\ \text{(first-order)}}} + \underbrace{w_4 (x_A \times x_C) + w_5 (x_A \times x_C) + w_6 (x_B \times x_C)}_{\substack{\text{Bimodal terms} \\ \text{(second-order)}}} + \underbrace{w_7 (x_A \times x_B \times x_C)}_{\substack{\text{Trimodal terms} \\ \text{(third-order)}}$$

Can we add  
higher-order  
interaction terms?

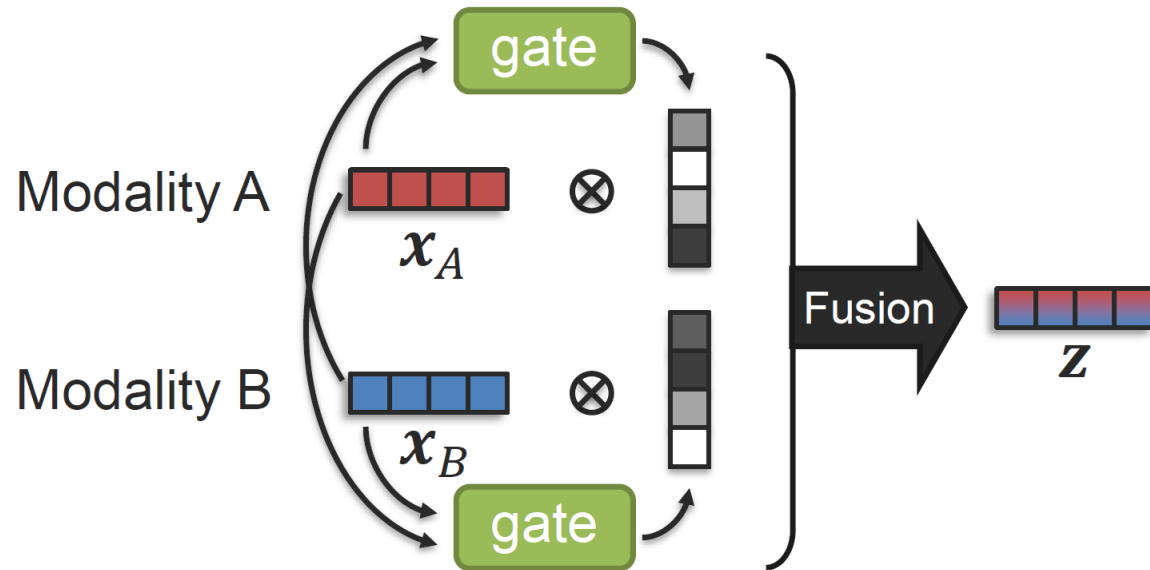
For example:

$$+w_8 (x_A^2 \times x_B^2 \times x_C^2)$$

$$+w_9 (x_A^3 \times x_B)$$

$$+w_{10} (x_B^3 \times x_C^3)$$

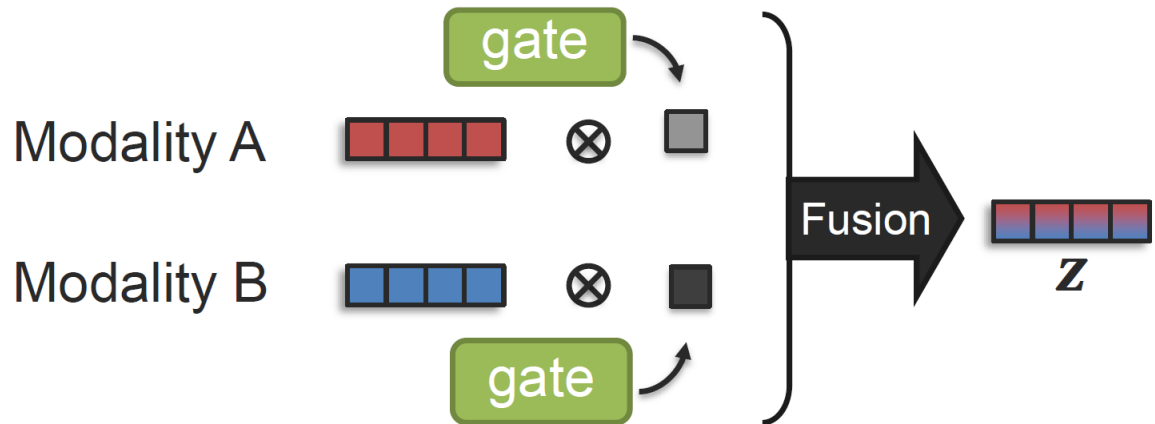
# Gated Fusion



Example with additive fusion:

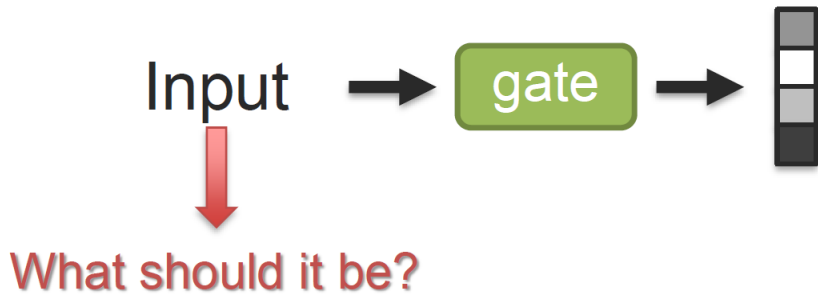
$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

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



$\Rightarrow$  Gating output can be one weight for the whole modality

# Gating Module (aka, attention module)



Target modality 

Other modality 

All modality   


*“Neural network designed to mask unwanted signal from propagating forward”* (gating)

...or with a more positive view:

*“Neural network designed to select preferable signal to move forward”* (attention)

Soft attention



Easier to compute derivative (gradient)

Hard attention



Derivative is harder (e.g., use reinforcement learning)

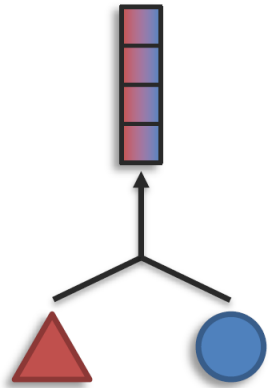


# Task 1: Representation (表示)

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

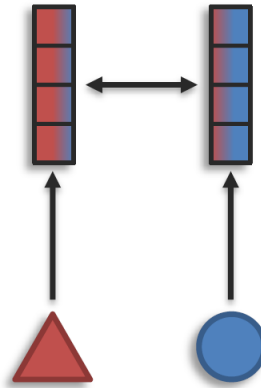
## Sub-challenges:

### Fusion



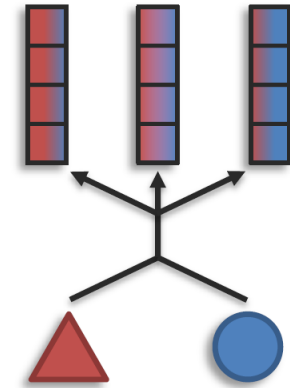
# modalities  $>$  # representations

### Coordination



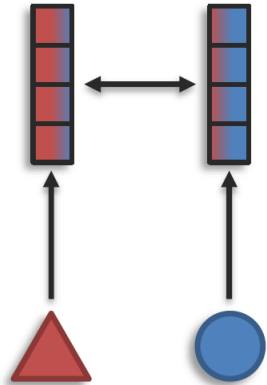
# modalities = # representations

### Fission



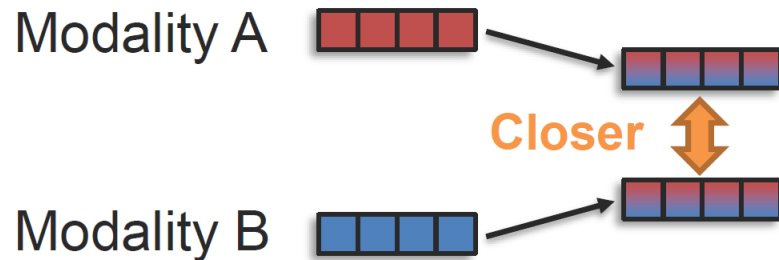
# modalities  $<$  # representations

# Sub-Challenge 1b: Representation Coordination

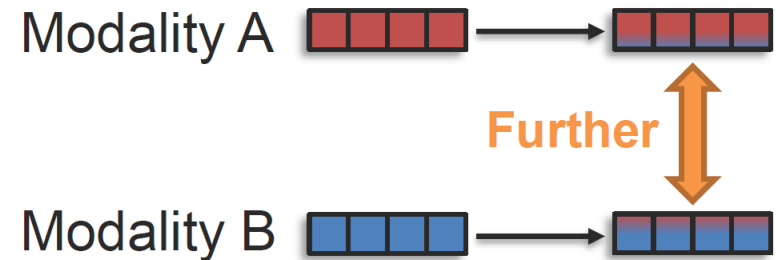


**Definition:** Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions

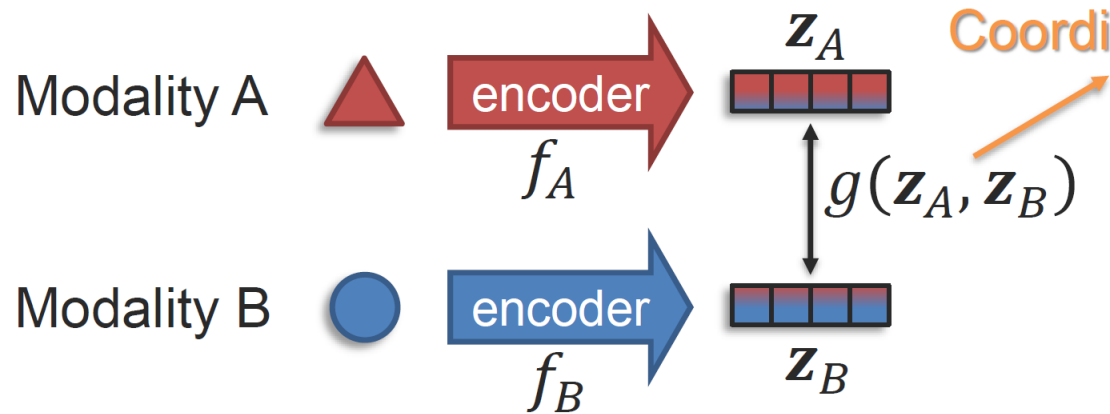
Strong Coordination:



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

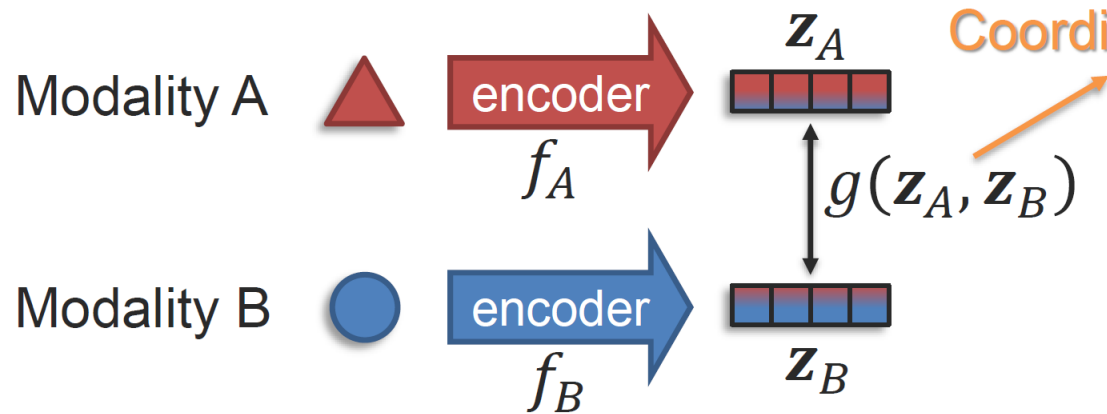
Examples of coordination function:

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  $\theta_{f_B}$

➡ Requires paired data

Examples of coordination function:

② Kernel similarity functions:

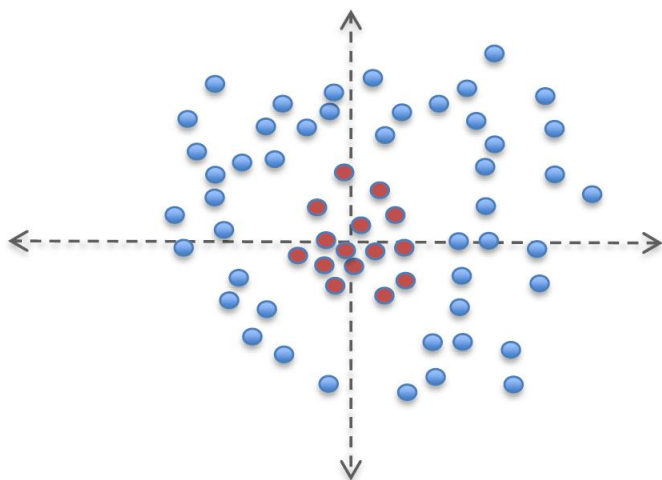
$$g(\mathbf{z}_A, \mathbf{z}_B) = k(\mathbf{z}_A, \mathbf{z}_B) \left\{ \begin{array}{l} \bullet \text{ Linear} \\ \bullet \text{ Polynomial} \\ \bullet \text{ Exponential} \\ \bullet \text{ RBF} \end{array} \right.$$

➡ All these examples bring relatively strong coordination between modalities

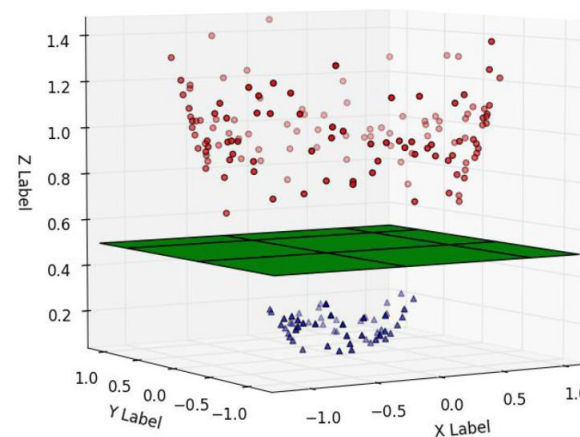
# Kernel Function

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

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



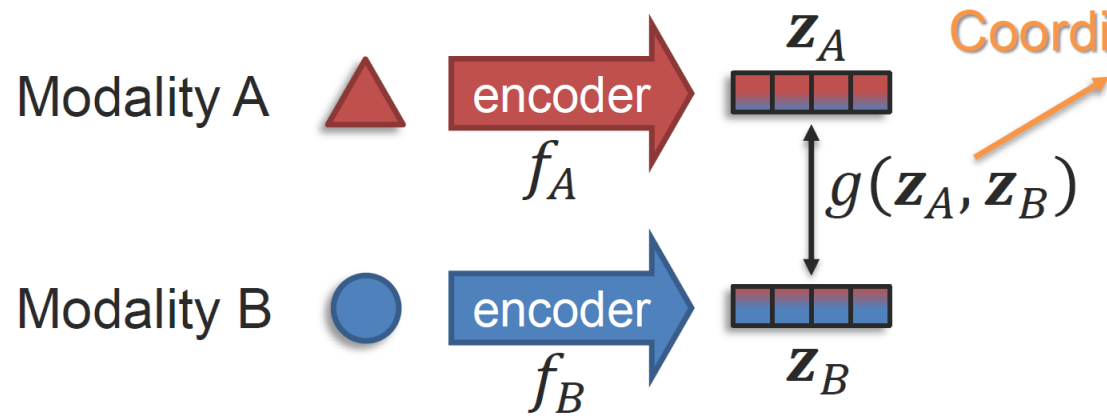
Not linearly separable in  $x$  space



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

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

# Coordination Function



Learning with coordination function:

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

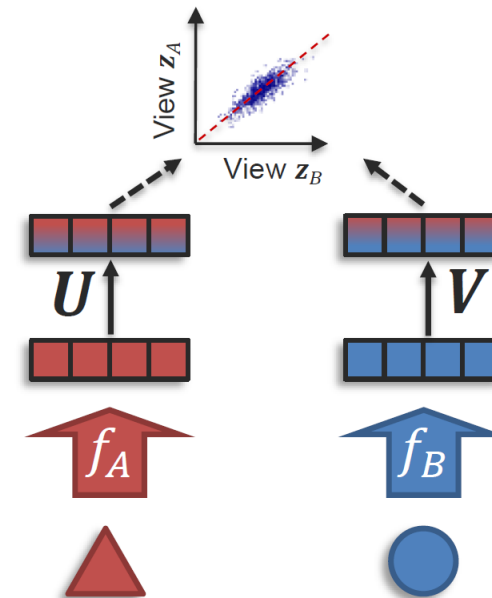
with model parameters  $\theta_g$ ,  $\theta_{f_A}$  and  $\theta_{f_B}$

## Examples of coordination function:

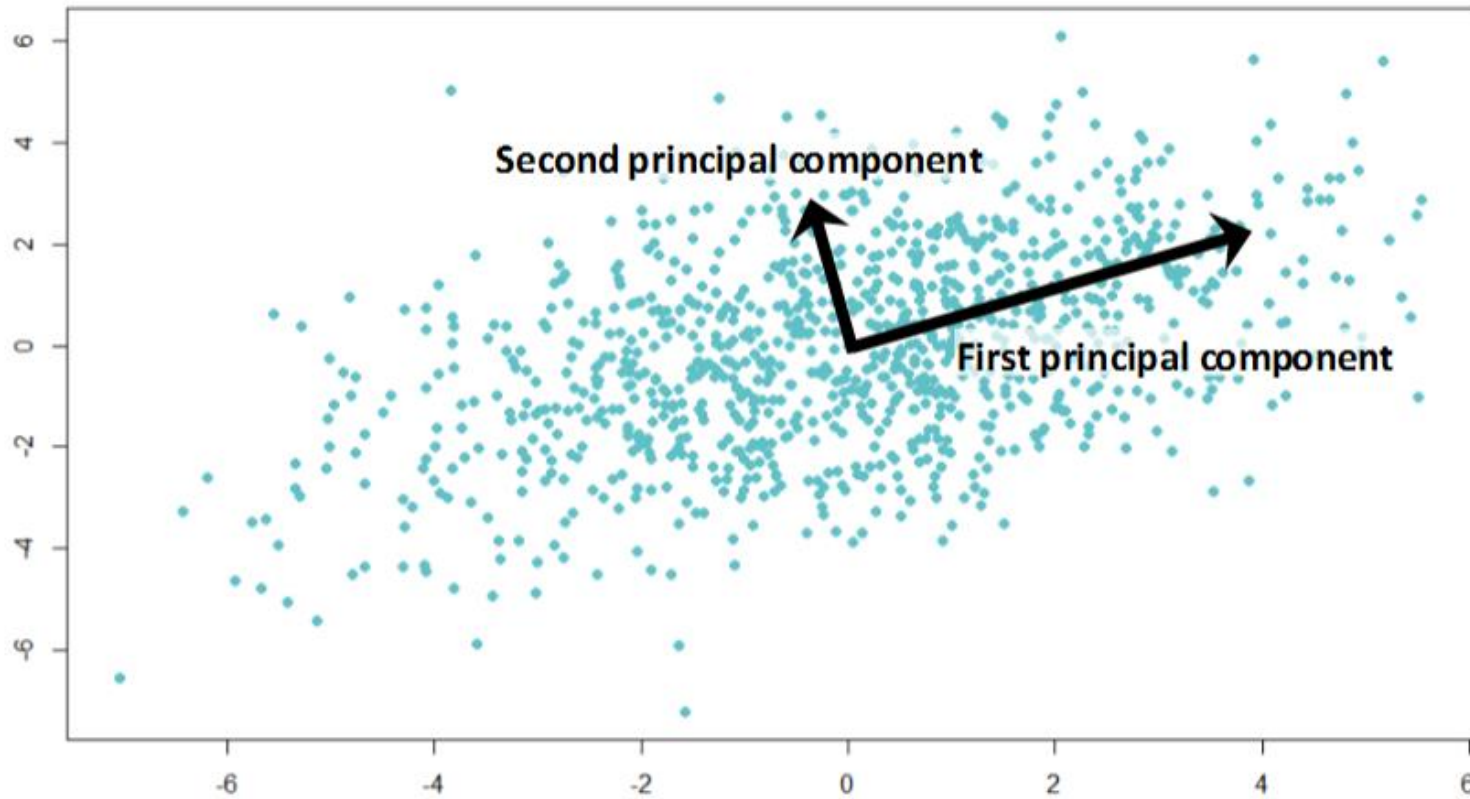
③ Canonical Correlation Analysis (CCA):

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

→ CCA includes multiple projections, all orthogonal with each others

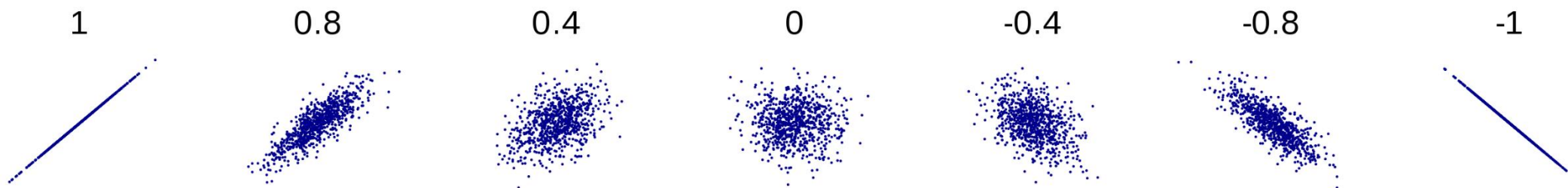


# Retrospect: Principal Component Analysis (PCA)



$$\arg \max \|\mathbf{u}^\top \mathbf{X}\|^2$$
$$\text{s. t. } \mathbf{u}^\top \mathbf{u} = 1$$

# Correlation



$$\rho_{X,Y} = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\text{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

$X, Y$  independent  $\Rightarrow \rho_{X,Y} = 0$  ( $X, Y$  uncorrelated)  
 $\rho_{X,Y} = 0$  ( $X, Y$  uncorrelated)  $\nRightarrow X, Y$  independent



# Correlation

## 1. 二次函数关系的例子

设随机变量  $X$  均匀分布在区间  $[-1, 1]$ , 定义  $Y = X^2$ 。虽然  $X$  和  $Y$  之间显然不是独立的, 因为  $Y$  完全由  $X$  决定, 但它们是 **不相关的**, 因为协方差为零。

- 证明不相关性:

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

由于  $Y = X^2$ , 我们可以计算:

$$\mathbb{E}[X] = 0 \quad (\text{因为 } X \text{ 是均匀分布的, 对称性导致均值为 } 0)$$

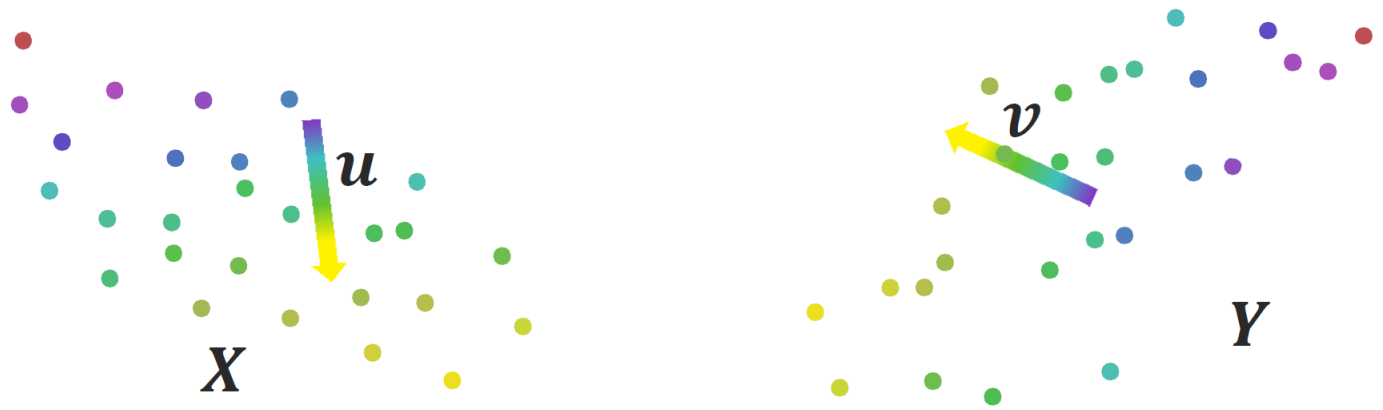
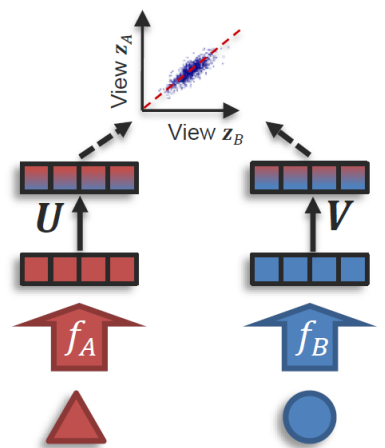
$$\mathbb{E}[X^3] = 0 \quad (\text{立方对称分布于 } [-1, 1], \text{ 故期望为 } 0)$$

因此,  $\text{Cov}(X, Y) = 0$ , 它们不相关, 但显然  $X$  和  $Y$  不是独立的。

# Correlated Projection

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

$$(u^*, v^*) = \arg \max \frac{u^T \Sigma_{XY} v}{\sqrt{u^T \Sigma_{XX} u} \sqrt{v^T \Sigma_{YY} v}}$$



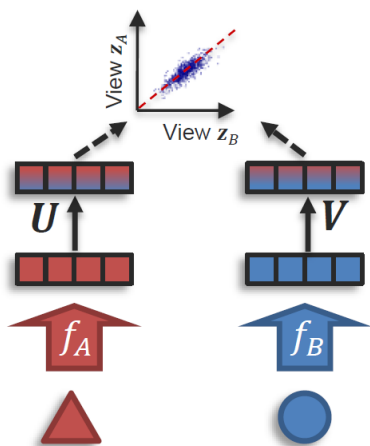
Two views  $X, Y$  where same instances have the same color

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

# Correlated Projection

The first pair of canonical variables:

$$(\mathbf{u}_1, \mathbf{v}_1) = \arg \max \frac{\mathbf{u}^\top \boldsymbol{\Sigma}_{XY} \mathbf{v}}{\sqrt{\mathbf{u}^\top \boldsymbol{\Sigma}_{XX} \mathbf{u}} \sqrt{\mathbf{v}^\top \boldsymbol{\Sigma}_{YY} \mathbf{v}}}$$



2

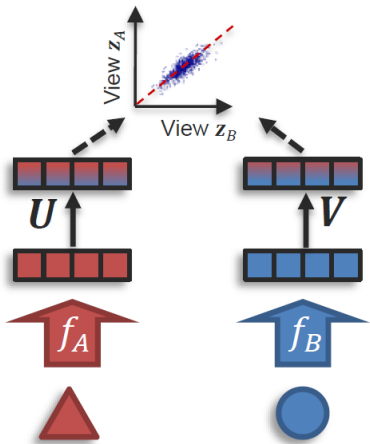
Since this objective function is invariant to scaling, we can constraint the projections to have unit variance:

$$\mathbf{u}_1^\top \boldsymbol{\Sigma}_{XX} \mathbf{u}_1 = \mathbf{v}_1^\top \boldsymbol{\Sigma}_{YY} \mathbf{v}_1 = 1$$

# Correlated Projection

The  $k$ -th pair of canonical variables:

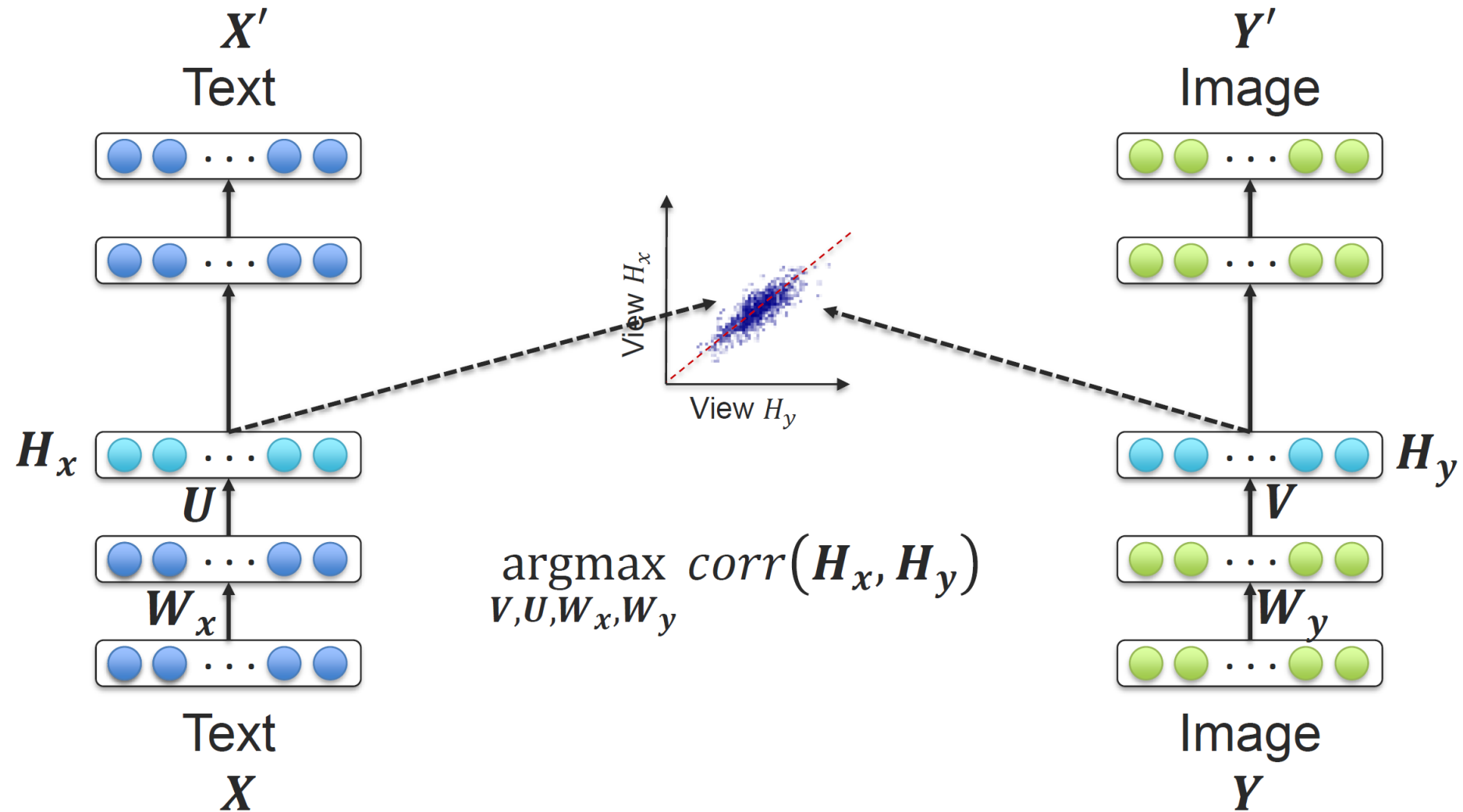
$$(\mathbf{u}_k, \mathbf{v}_k) = \arg \max \frac{\mathbf{u}^\top \boldsymbol{\Sigma}_{XY} \mathbf{v}}{\sqrt{\mathbf{u}^\top \boldsymbol{\Sigma}_{XX} \mathbf{u}} \sqrt{\mathbf{v}^\top \boldsymbol{\Sigma}_{YY} \mathbf{v}}}$$



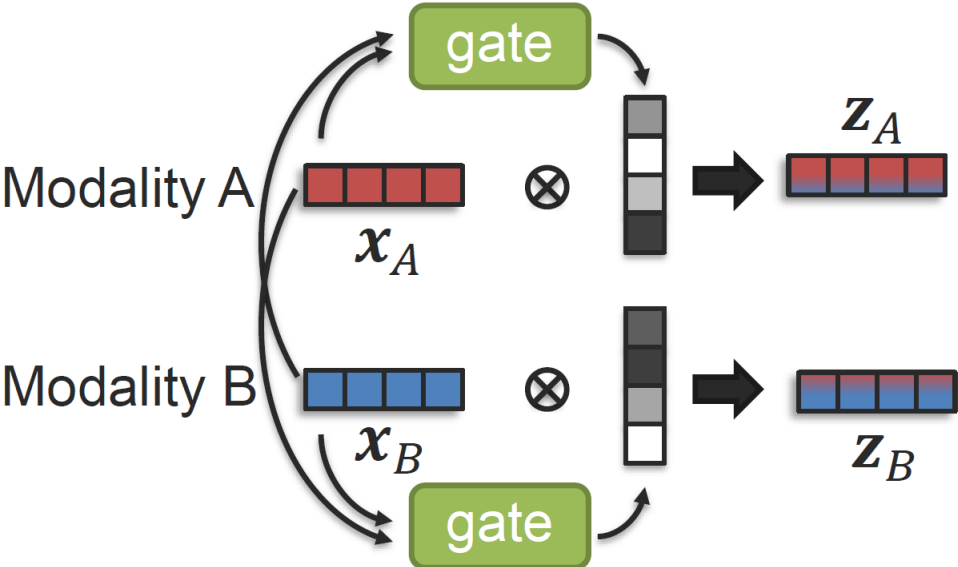
③ We want these multiple projection pairs to be orthogonal (“canonical”) to each other:

$$\mathbf{u}^\top \boldsymbol{\Sigma}_{XX} \mathbf{u}_j = \mathbf{v}^\top \boldsymbol{\Sigma}_{YY} \mathbf{v}_j = \mathbf{0}, \forall j = 1, \dots, k - 1$$

# Deep Canonically Correlated Autoencoders (DCCA)



# Gated Coordination



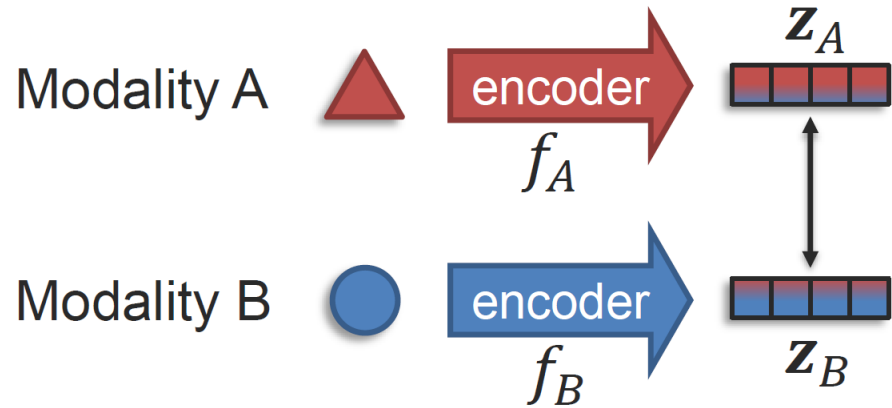
Gated coordination:

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

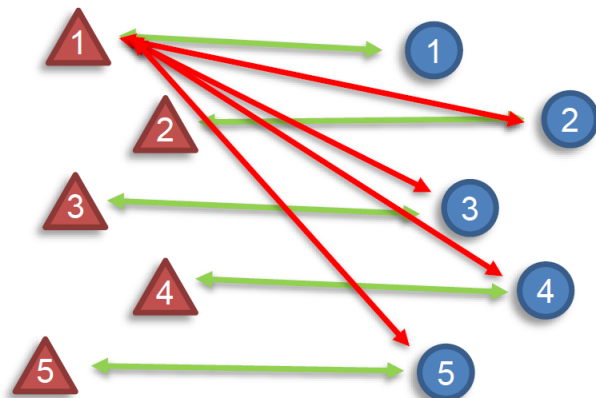
$$z_B = g_B(x_A, x_B) \cdot x_B$$

➡ Related to attention modules in transformers


# 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:

$$\max\{0, \alpha + \boxed{\text{sim}(z_A, z_B^+)} - \boxed{\text{sim}(z_A, z_B^-)}\}$$

 positive pairs

 negative pair

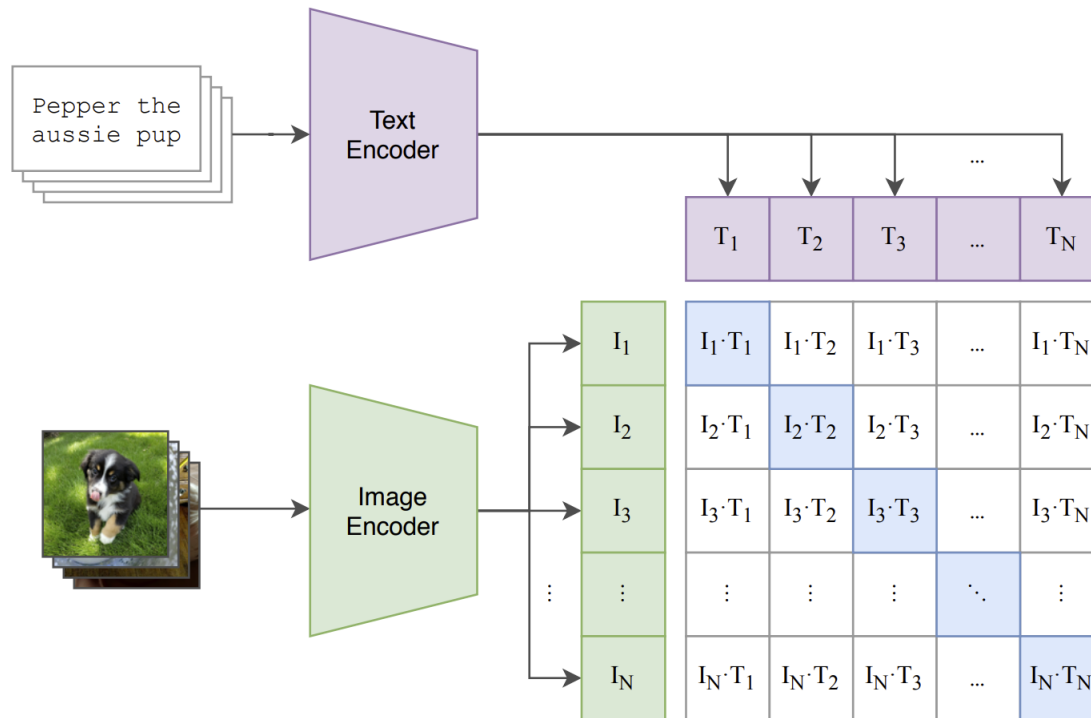
Similarity functions are often cosine similarity

 Similar to hinge loss

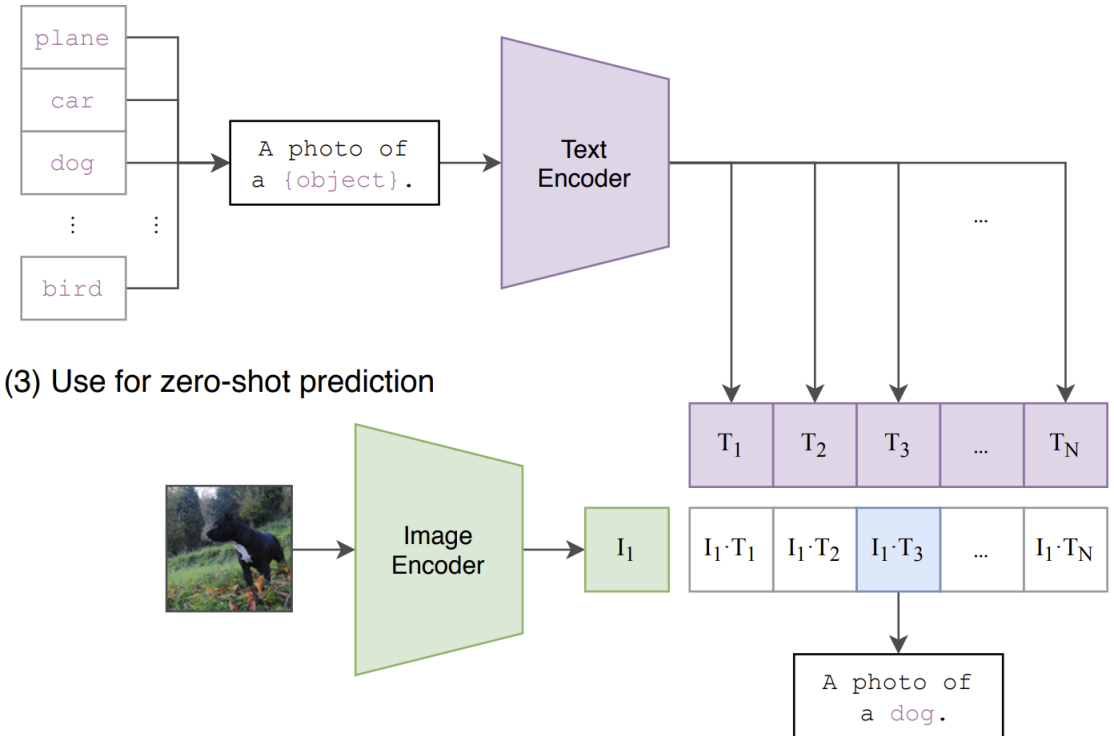
# Example – CLIP (Contrastive Language–Image Pre-training)

## Zero-Shot Image Classification

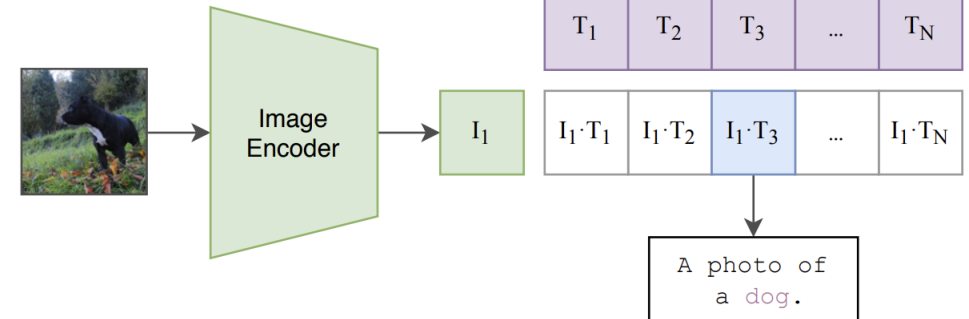
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction





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

---

## Pretrained Dataset (not open-sourced by openAI)

we constructed a new dataset of **400 million (image, text) pairs** collected from a variety of publicly available sources on the Internet. To attempt to cover as broad a set of visual concepts as possible, we search for (image, text) pairs as part of the construction process whose text includes one of a set of 500,000 queries. We approximately class balance the results by including up to 20,000 (image, text) pairs per query. The resulting dataset has a similar total word count as the WebText dataset used to train GPT-2. We refer to this dataset as **WIT** for **WebImageText**.

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

## Encoders

Model	Learning rate	Embedding dimension	Input resolution	ResNet blocks	width	Text Transformer layers	width	heads
RN50	$5 \times 10^{-4}$	1024	224	(3, 4, 6, 3)	2048	12	512	8
RN101	$5 \times 10^{-4}$	512	224	(3, 4, 23, 3)	2048	12	512	8
RN50x4	$5 \times 10^{-4}$	640	288	(4, 6, 10, 6)	2560	12	640	10
RN50x16	$4 \times 10^{-4}$	768	384	(6, 8, 18, 8)	3072	12	768	12
RN50x64	$3.6 \times 10^{-4}$	1024	448	(3, 15, 36, 10)	4096	12	1024	16

Table 19. CLIP-ResNet hyperparameters

Model	Learning rate	Embedding dimension	Input resolution	Vision Transformer layers	width	heads	Text Transformer layers	width	heads
ViT-B/32	$5 \times 10^{-4}$	512	224	12	768	12	12	512	8
ViT-B/16	$5 \times 10^{-4}$	512	224	12	768	12	12	512	8
ViT-L/14	$4 \times 10^{-4}$	768	224	24	1024	16	12	768	12
ViT-L/14-336px	$2 \times 10^{-5}$	768	336	24	1024	16	12	768	12

Table 20. CLIP-ViT hyperparameters

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

```
# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

## Symmetric InfoNCE (Noise Contrastive Estimation) loss

$$\mathcal{L} = \frac{1}{2} (\mathcal{L}_{\text{image}} + \mathcal{L}_{\text{text}})$$







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

**Positive pair**

**Negative pairs**

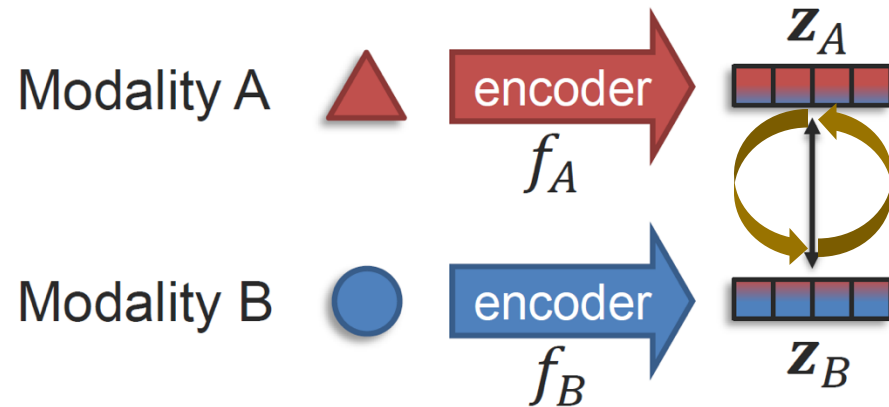
$$\mathcal{L}_{\text{text}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(\mathbf{z}_{i,\text{image}}, \mathbf{z}_{i,\text{text}})}{\sum_{j=1}^N \text{sim}(\mathbf{z}_{j,\text{image}}, \mathbf{z}_{i,\text{text}})}$$

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

	Dataset Examples	ImageNet ResNet101	Zero-Shot CLIP	$\Delta$ Score
ImageNet		76.2	76.2	0%
ImageNetV2		64.3	70.1	+5.8%
ImageNet-R		37.7	88.9	+51.2%
ObjectNet		32.6	72.3	+39.7%
ImageNet Sketch		25.2	60.2	+35.0%
ImageNet-A		2.7	77.1	+74.4%

# Homogeneous Coordination with Channel Exchanging

**Homogeneous Multimodal Learning:** The modalities to fuse are of the same shape; there is certain correspondence between their each element

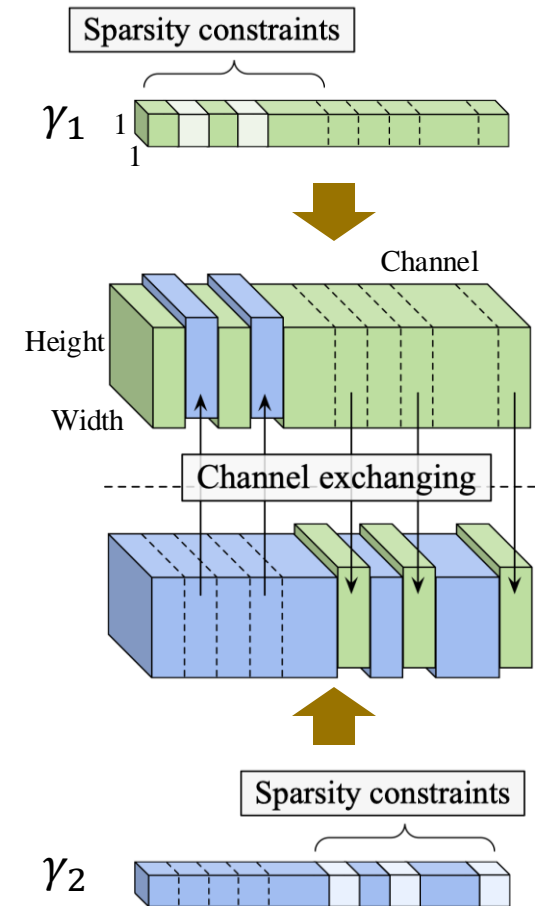


**Parameter-free, Self-adaptive**

Modality 1



Modality 2

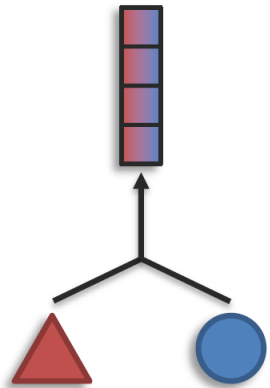


# Task 1: Representation (表示)

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

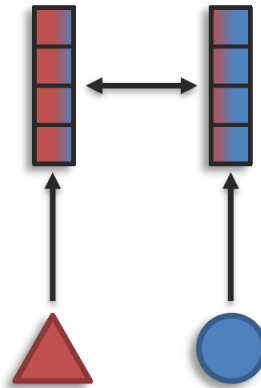
## Sub-challenges:

### Fusion



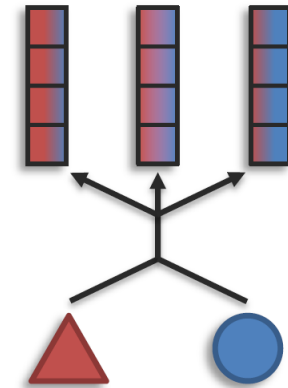
# modalities  $>$  # representations

### Coordination



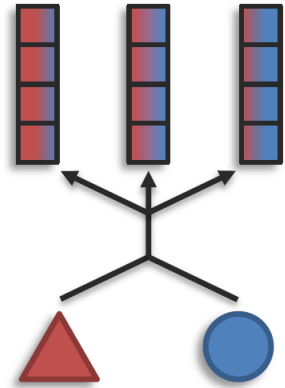
# modalities = # representations

### Fission



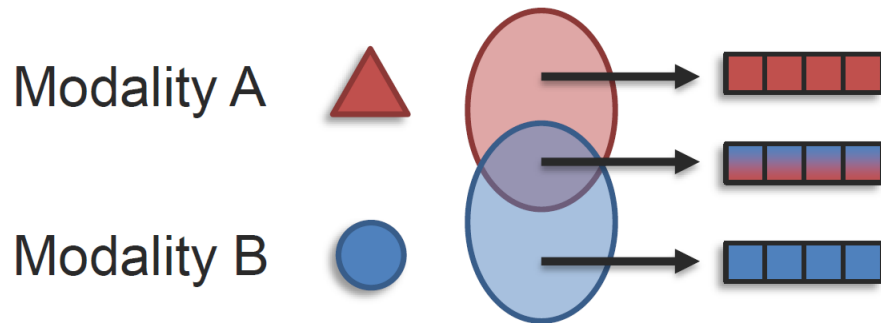
# modalities  $<$  # representations

# Sub-Challenge 1c: Representation Fission

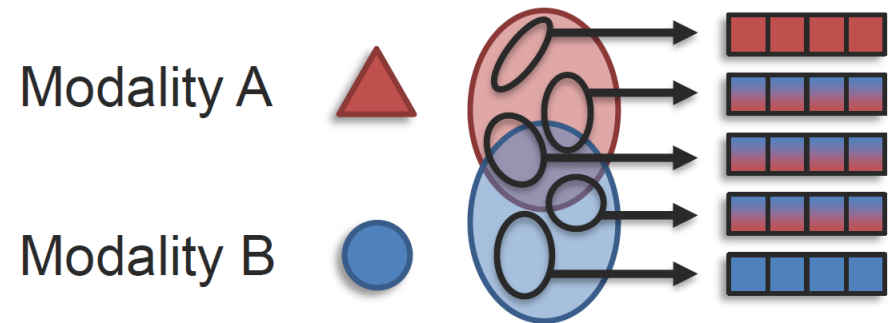


**Definition:** learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering

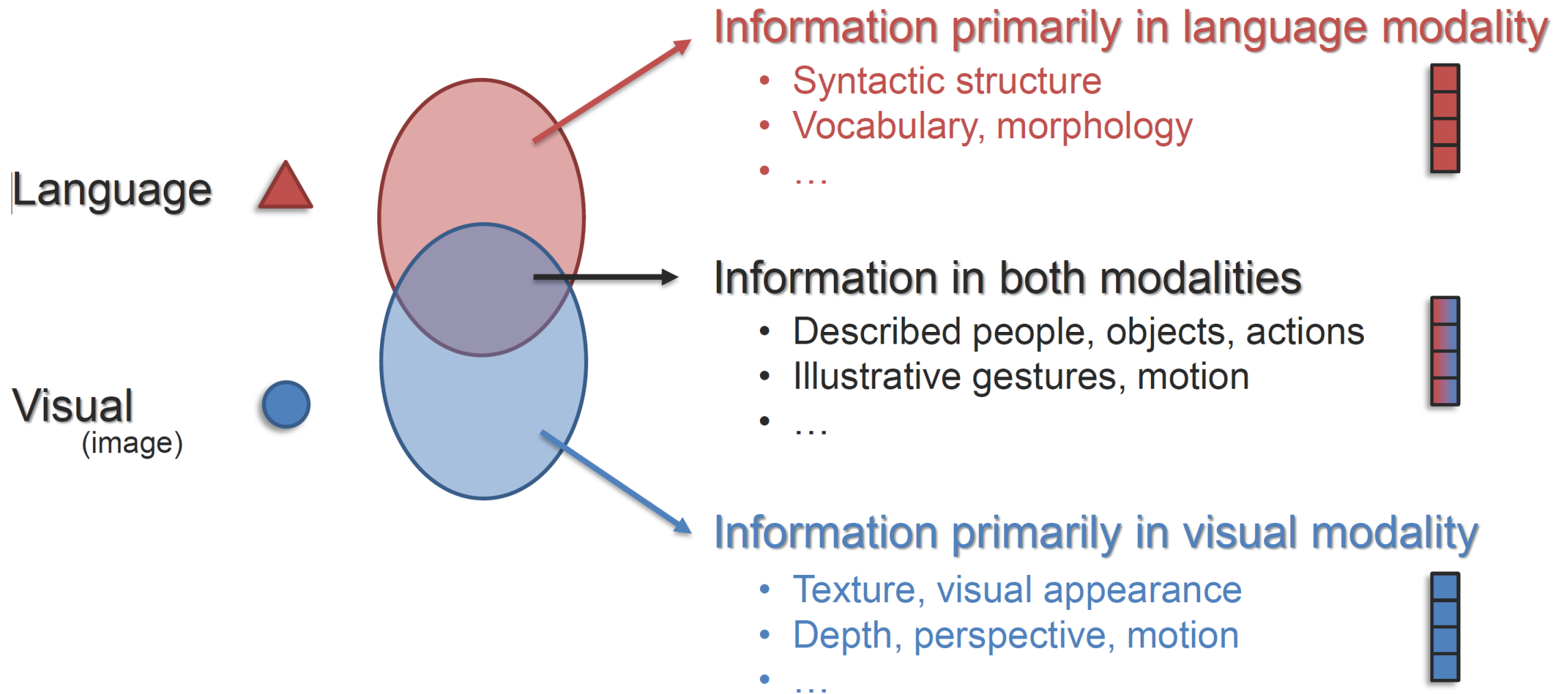
Modality-level fission:



Fine-grained fission:

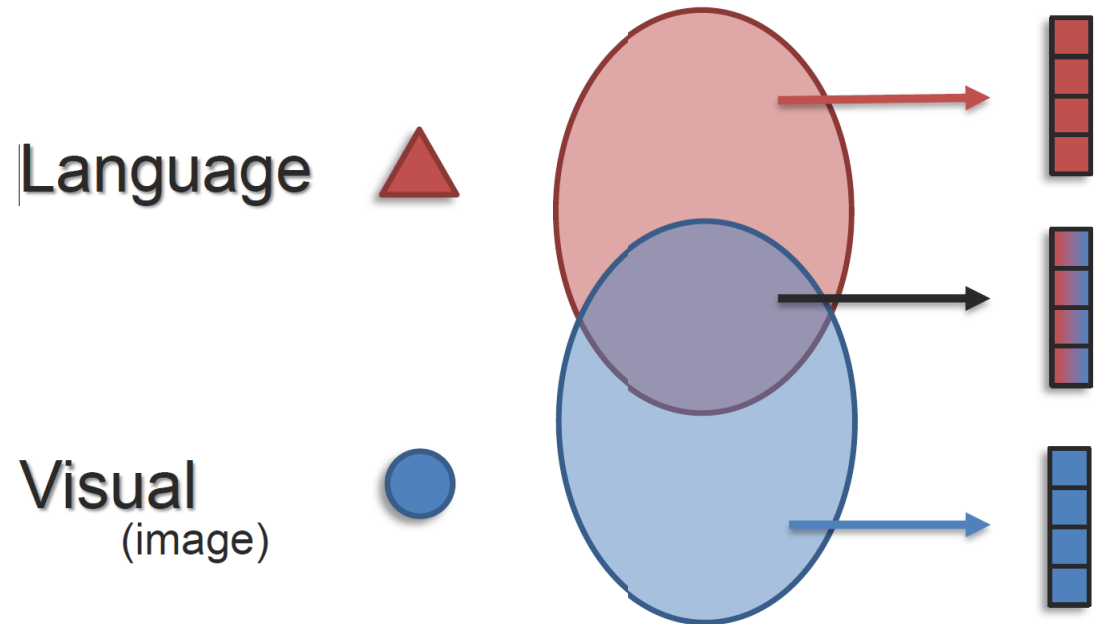


# Modality-Level Fission



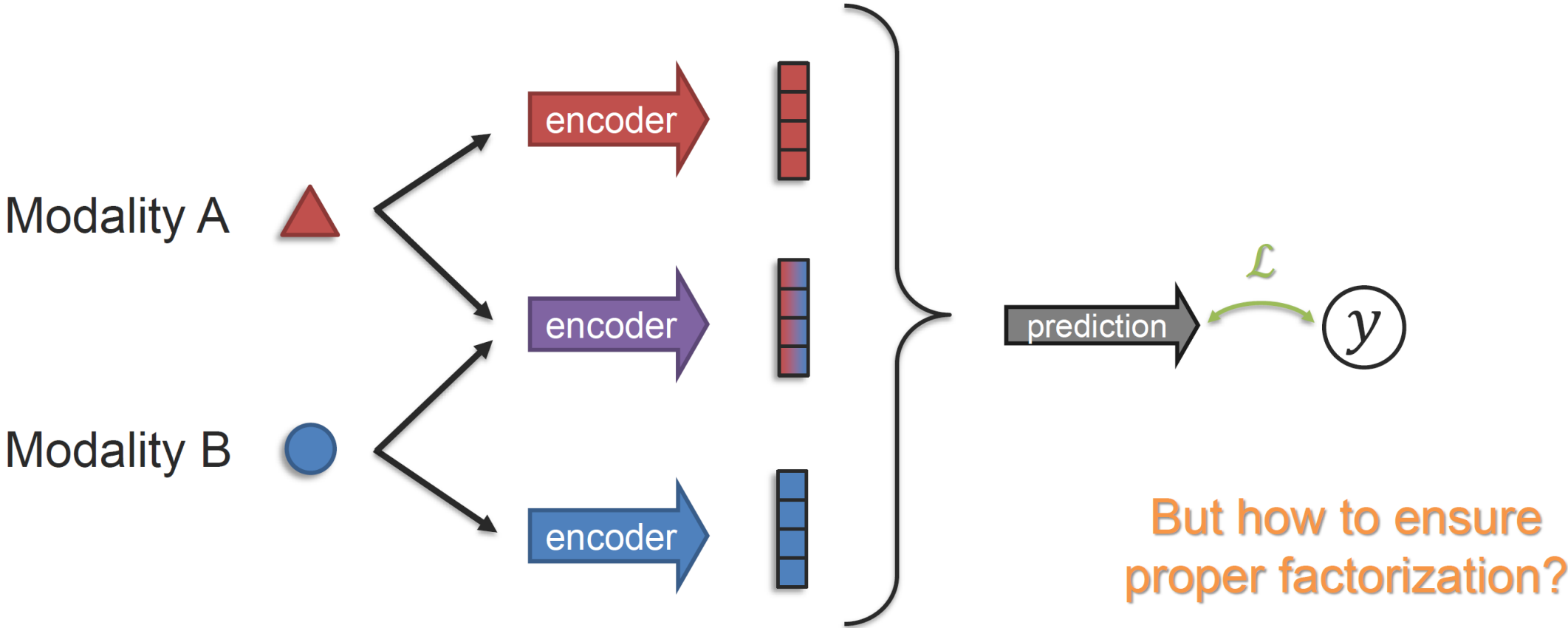


# Modality-Level Fission

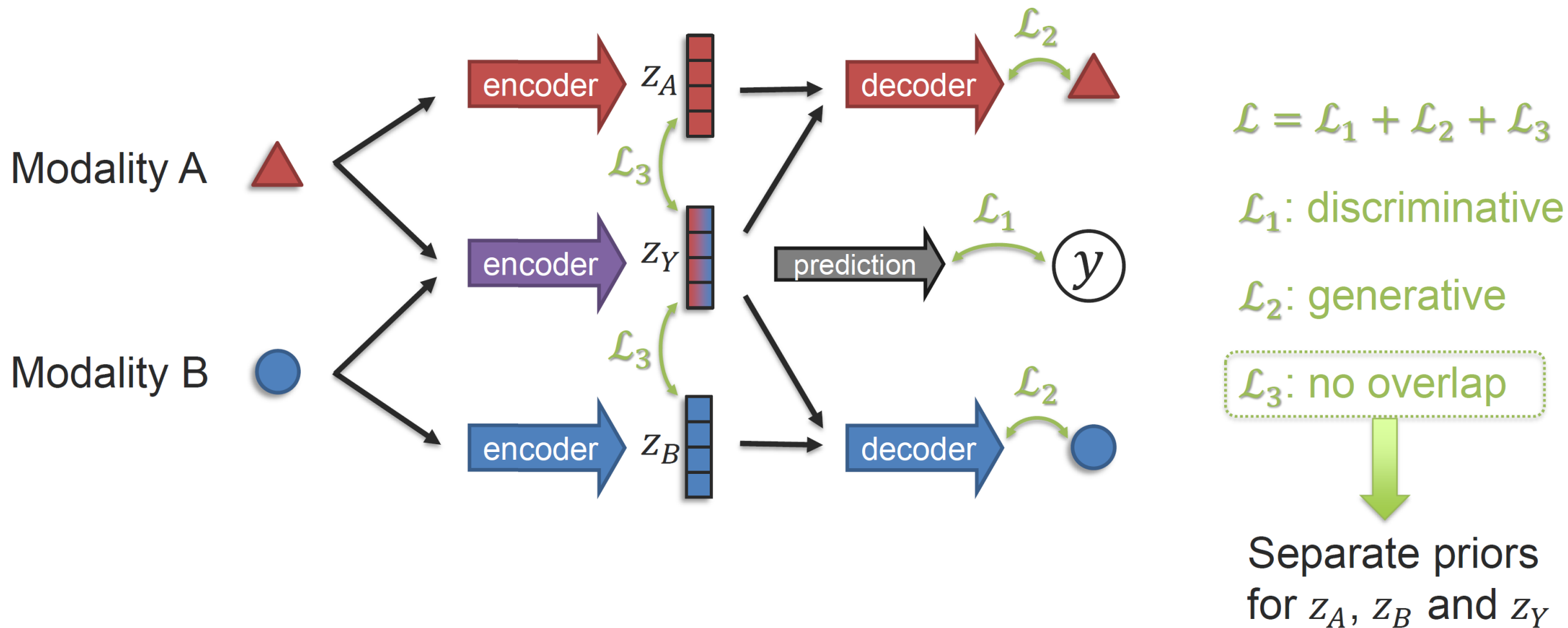


How to learn factorized multimodal representations?

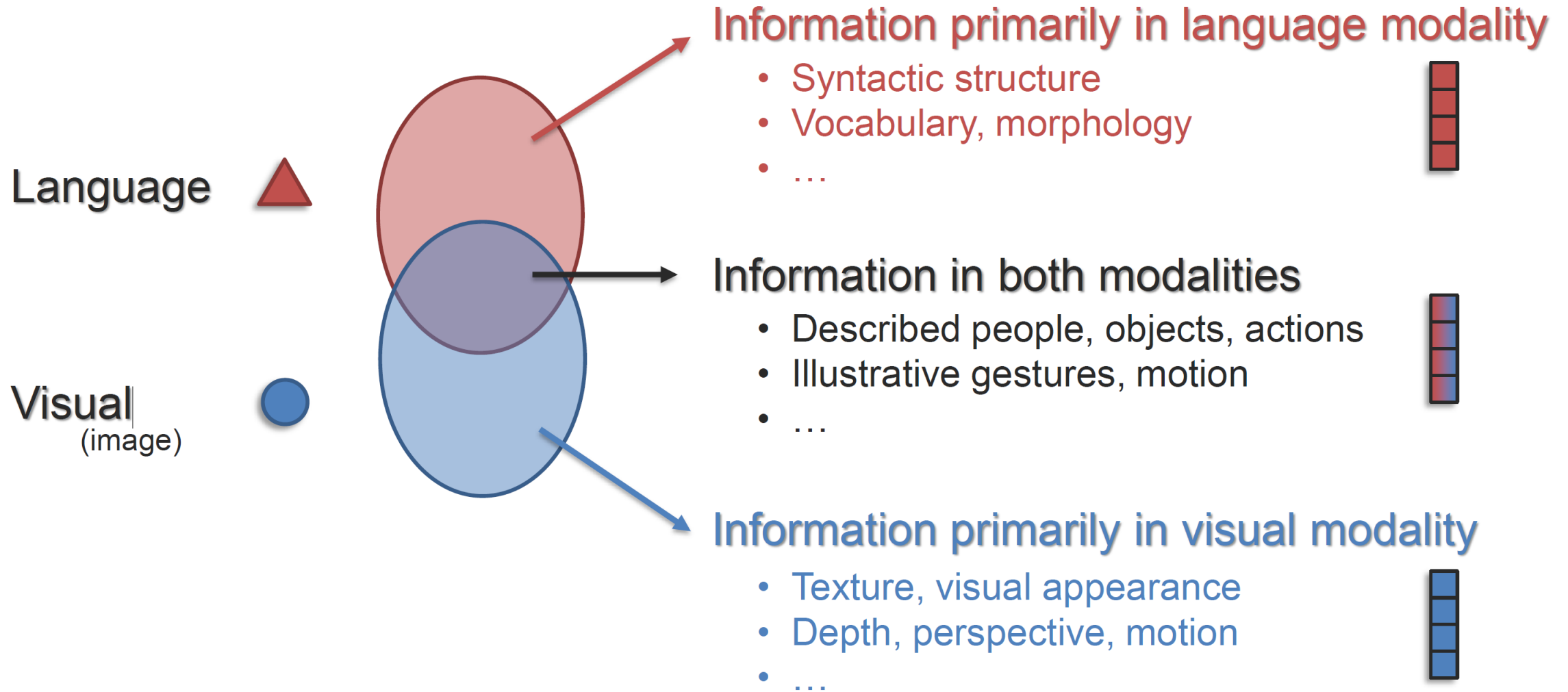
# A Discriminative Approach – Factorized Multimodal Representations



# A Generative-Discriminative Approach

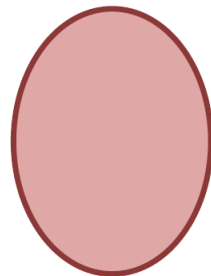


# Modality-Level Fission – Information Theory



# Information and Entropy – Information Theory

Language



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”

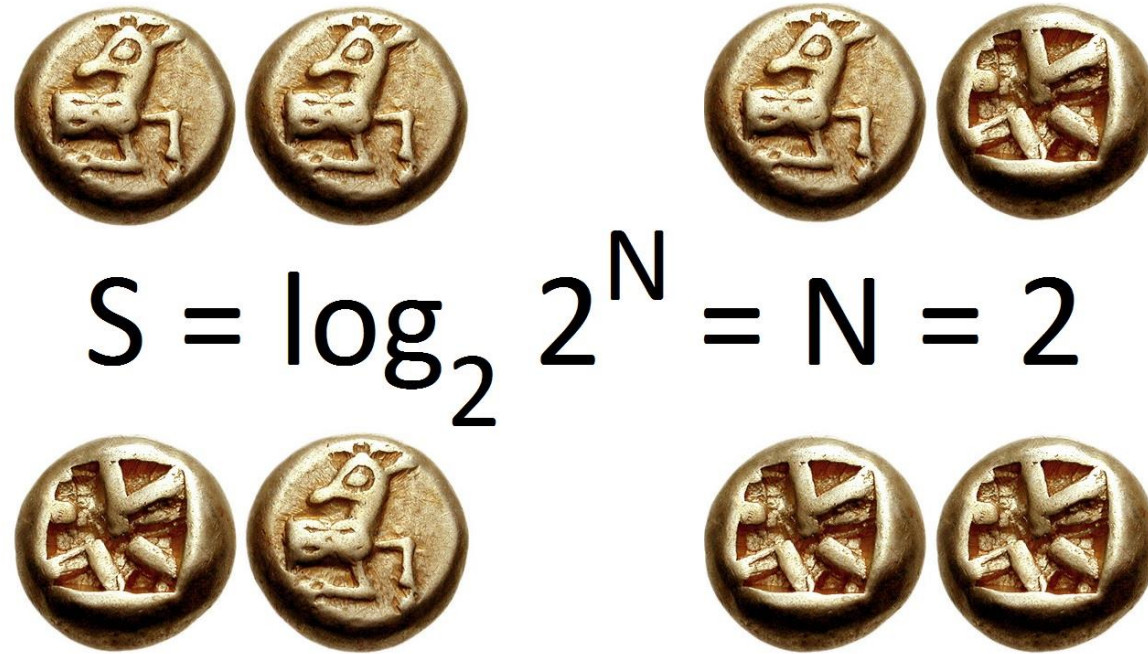
➔ Low chances... So, higher information

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

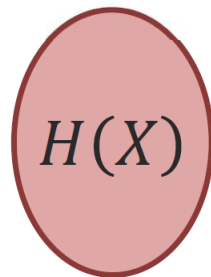


$$S = \log_2 2^N = N = 2$$

Information entropy (in bits) is the log-base-2 of the number of possible outcomes. With two coins there are four outcomes HH-HT-TH-TT, and the entropy is two bits.

# Information and Entropy – Information Theory

Language



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

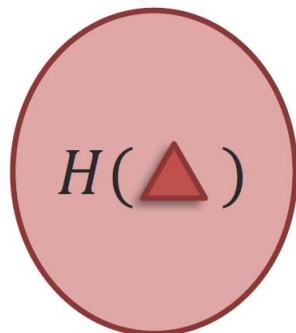
# Information and Entropy

If no overlapping information

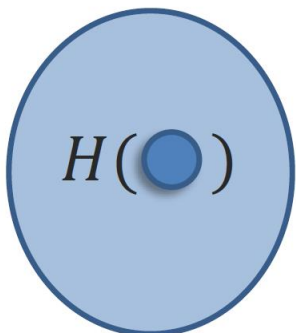


But in most real-world scenarios, modalities are *inter-connected*

Modality A



Modality B



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

Statistical

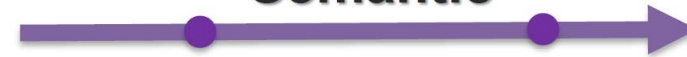


Association

Dependency



Semantic



Correspondence

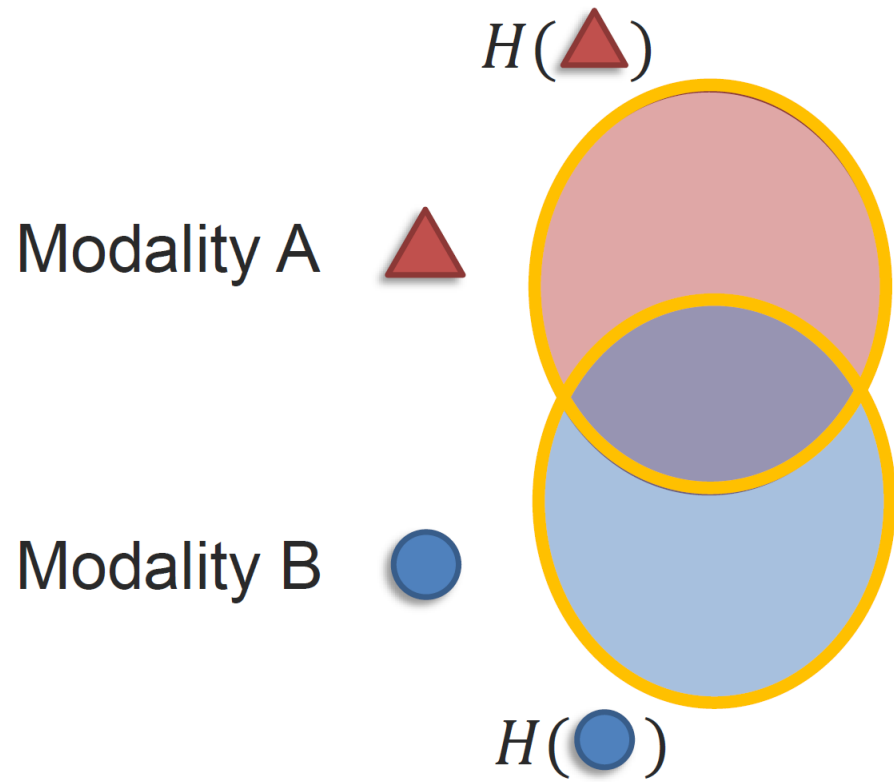
Relationship



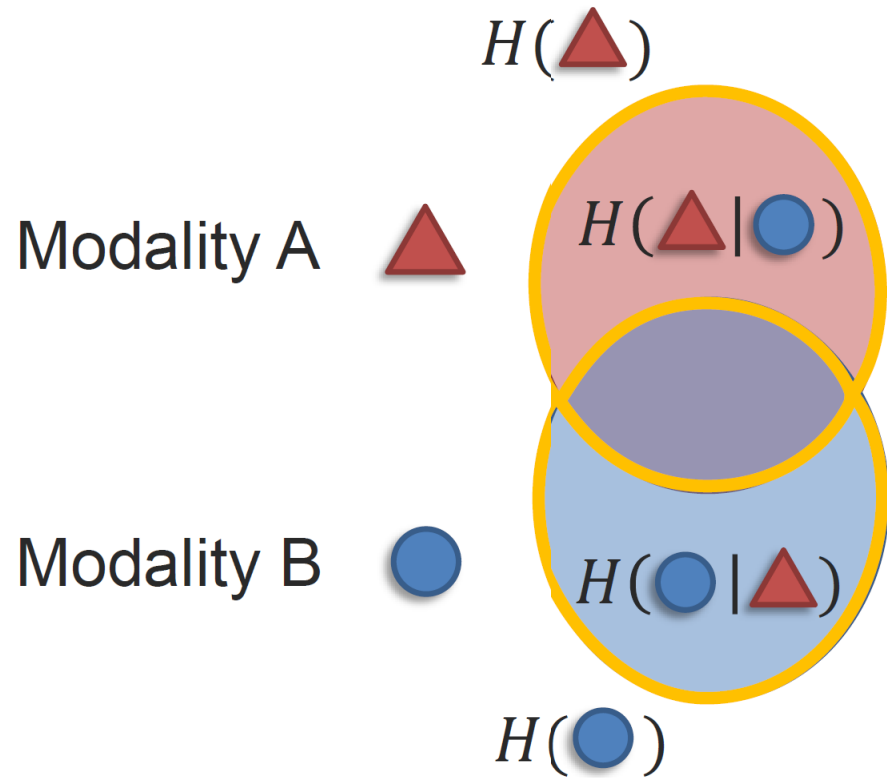


# Entropy with Two Modalities

---



# Entropy with Two Modalities

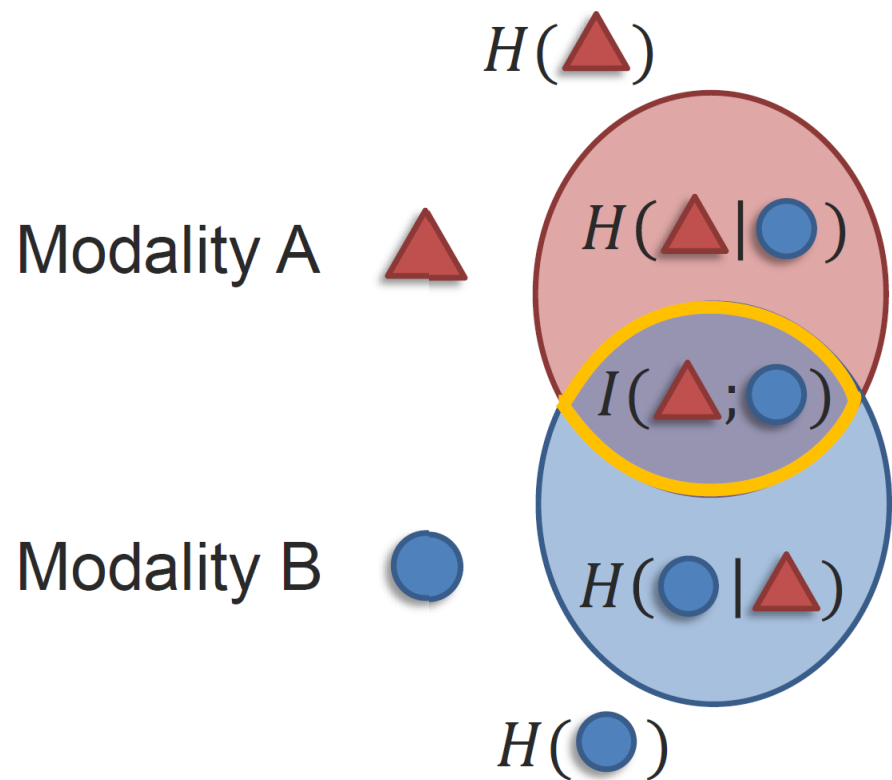


Conditional entropy  $H(Y|X)$

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

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

# Entropy with Two Modalities



## Mutual information $I(X; Y)$

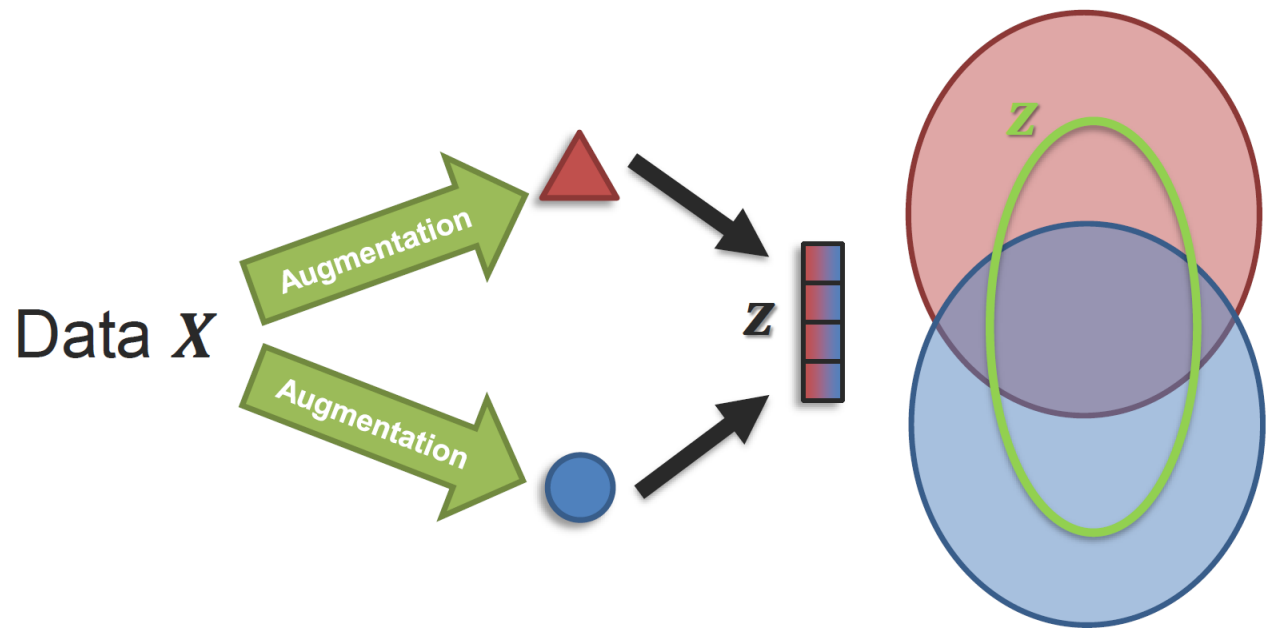
$$I(X; Y) = H(X) - H(X|Y)$$

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

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

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

# Link with Self-Supervised Learning



- 1 Maximize the mutual information

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

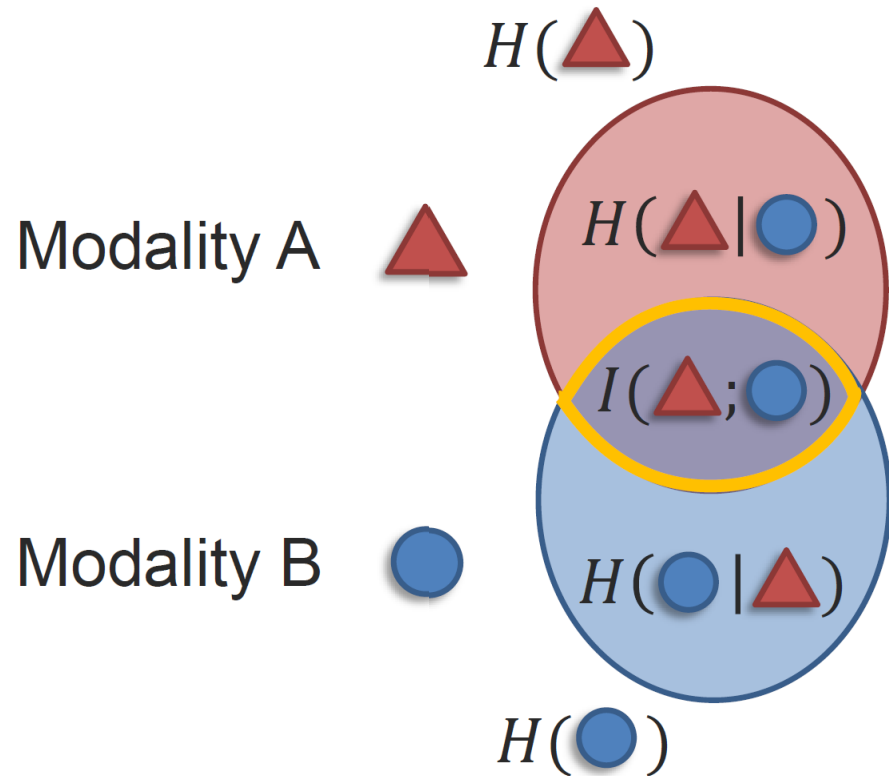
➔ Related to contrastive learning

- 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

# Some facts about information theory



## 1. Properties of mutual information:

$$I(X; Y) \geq 0, I(X; Y) = I(Y; X)$$

## 2. Subadditivity:

$$H(X) + H(Y) = H(X, Y) + I(X; Y) \geq H(X, Y)$$

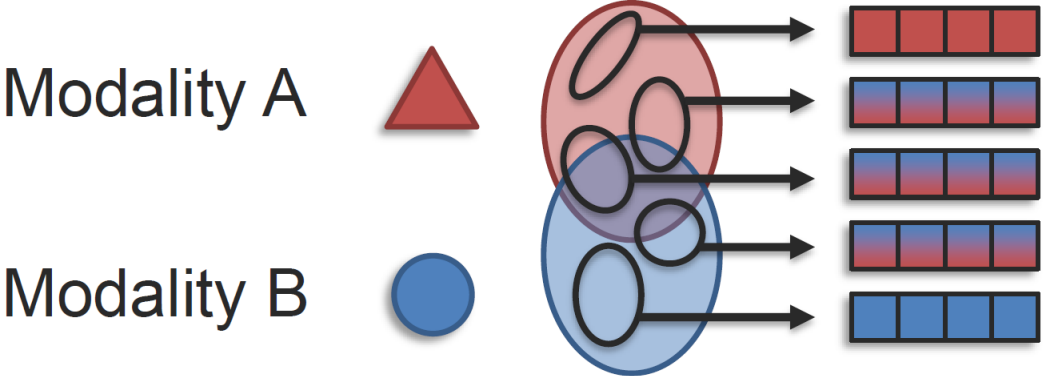
$$H(X) = H(X|Y) + I(X; Y) \geq H(X|Y)$$

3. The entropy or the amount of information revealed by evaluating X and Y simultaneously is equal to: first evaluating the value of Y, then revealing the value of X given that you know the value of Y.

$$H(X, Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$$

But, do we have  $H(X|Y) \geq 0$  ???

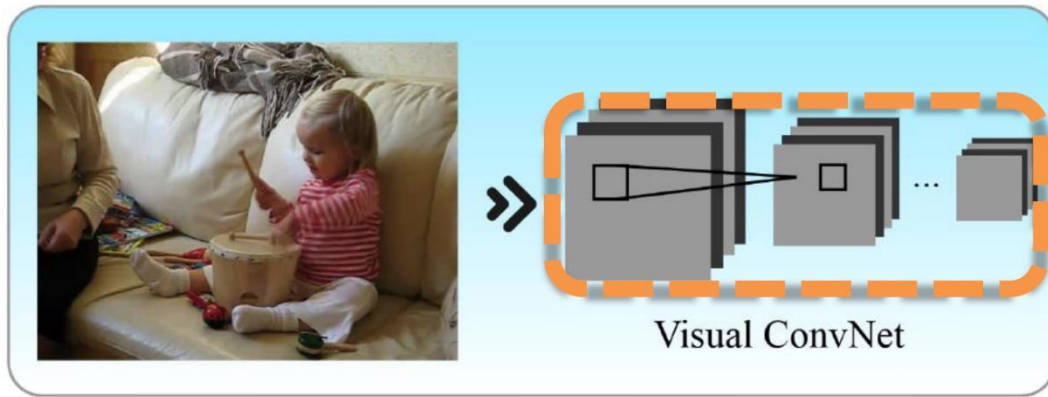
# Fine-Grained Fission



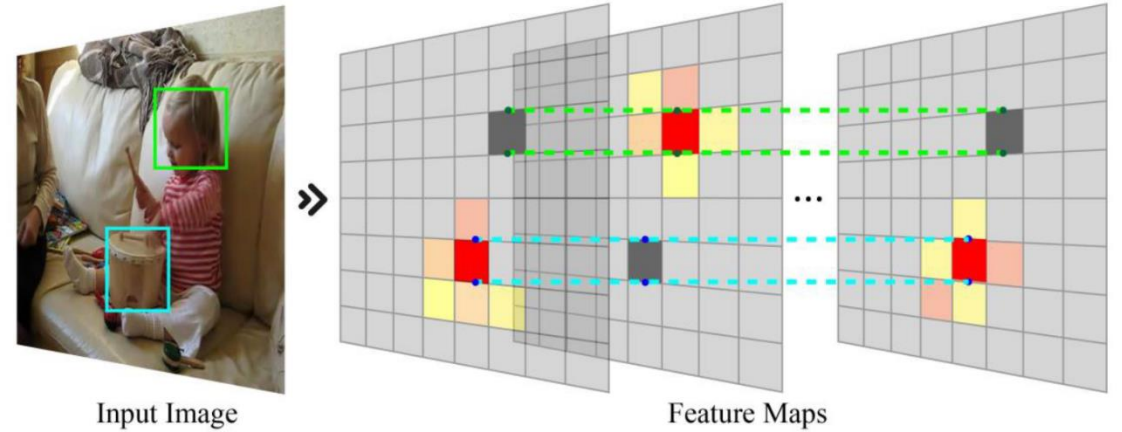
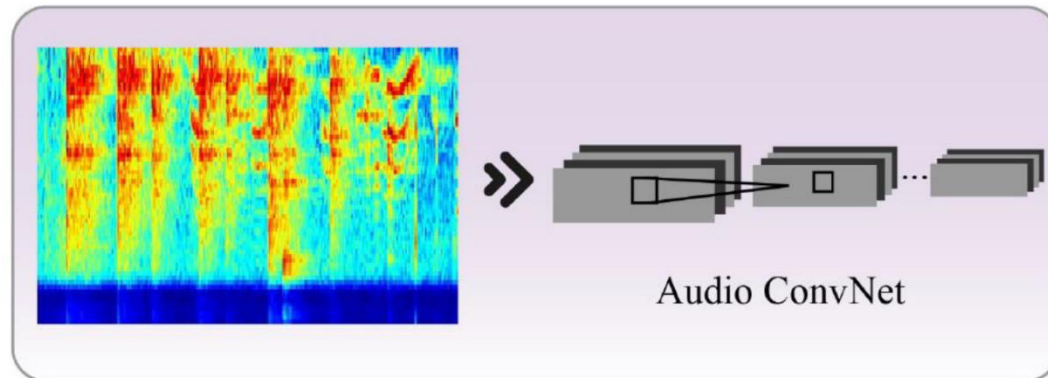
How to automatically discover these internal clusters, factors?

# Fine-Grained Fission – A Clustering Approach

## Unimodal Encoders

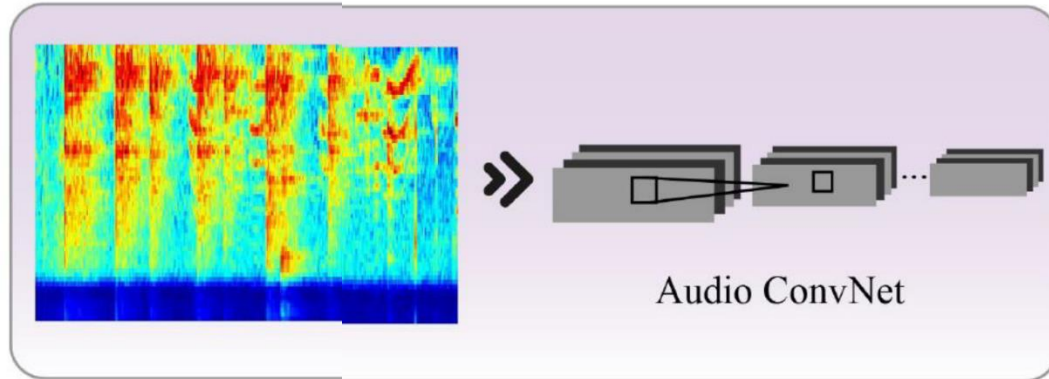
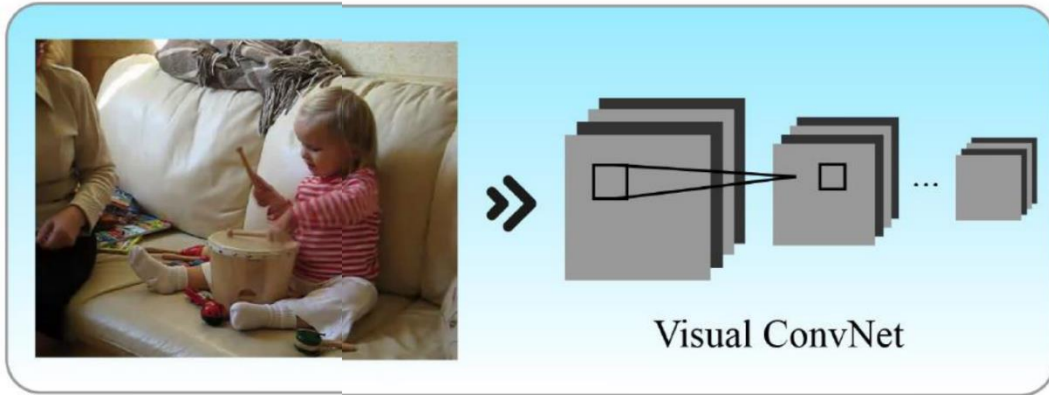


Localized activations for different objects

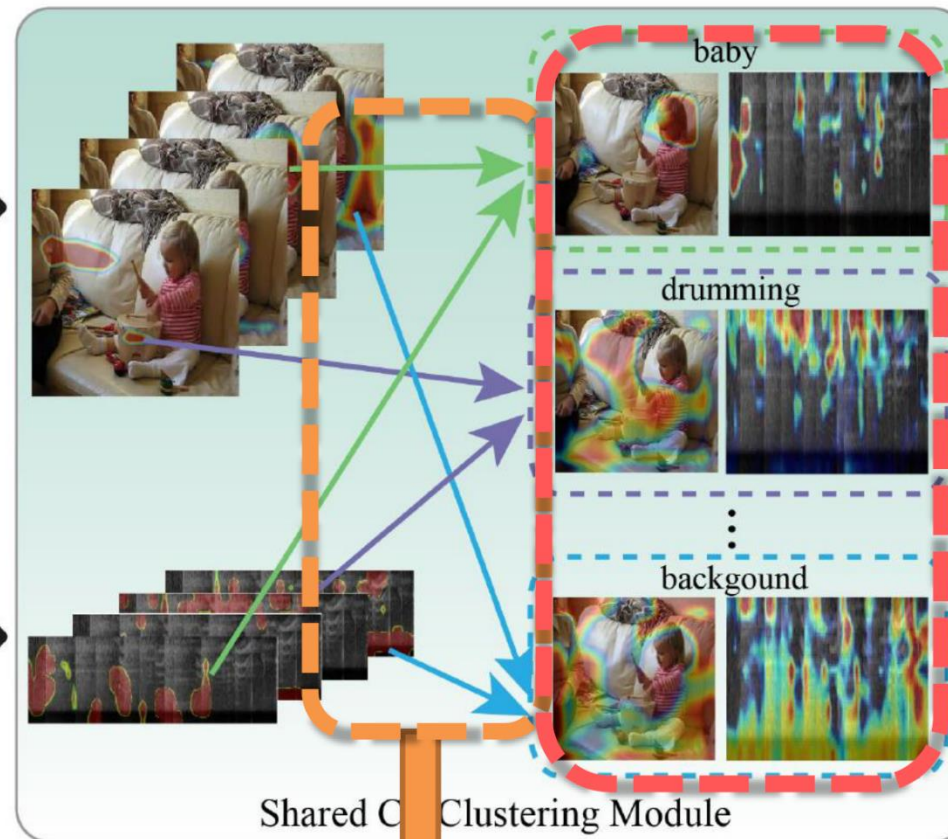


# Fine-Grained Fission – A Clustering Approach

## Unimodal Encoders



## Multimodal Fission



Discovers multiple audio-visual correspondences

Audiovisual Similarity

Explores different shared spaces (clusters)

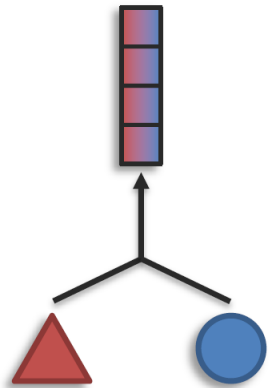


# Task 1: Representation (表示)

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

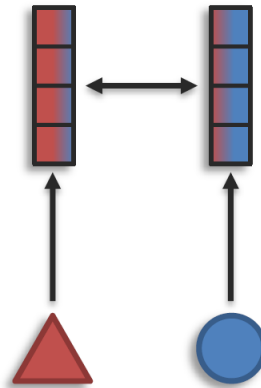
## Sub-challenges:

### Fusion



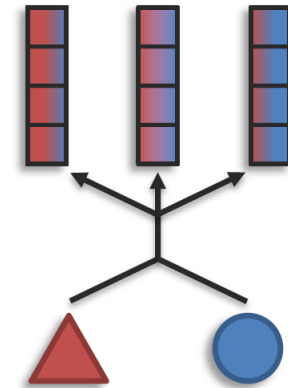
# modalities  $>$  # representations

### Coordination



# modalities = # representations

### Fission



# modalities  $<$  # representations