

# HUST

**ĐẠI HỌC BÁCH KHOA HÀ NỘI**

HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

ONE LOVE. ONE FUTURE.

# Multimodal Machine Learning

# Multimodal Machine Learning

Week i: ...

ONE LOVE. ONE FUTURE.

# Lecture Objectives

- What is Multimodal?
  - Research-oriented definition
  - Dimensions of modality heterogeneity
  - Modality connections and interactions
- Core technical and conceptual challenges
  - Representation, alignment, reasoning, generation, transference and quantification

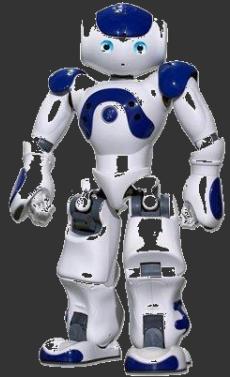


# What is Multimodal?

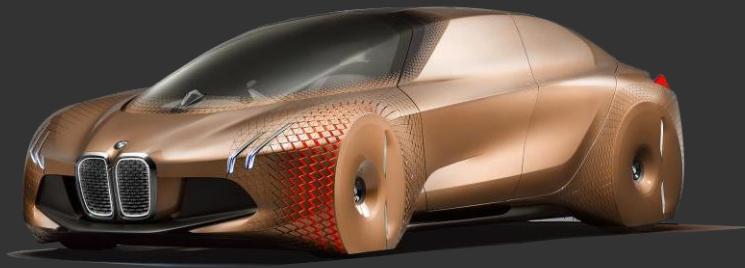


# Multimodal AI Technologies

## Robots



## Personal Vehicles



## Ubiquitous



## Mobile



## Wearable



# Multimodal AI Technologies

Robots



M



Personal Vehicles

Video Conferencing

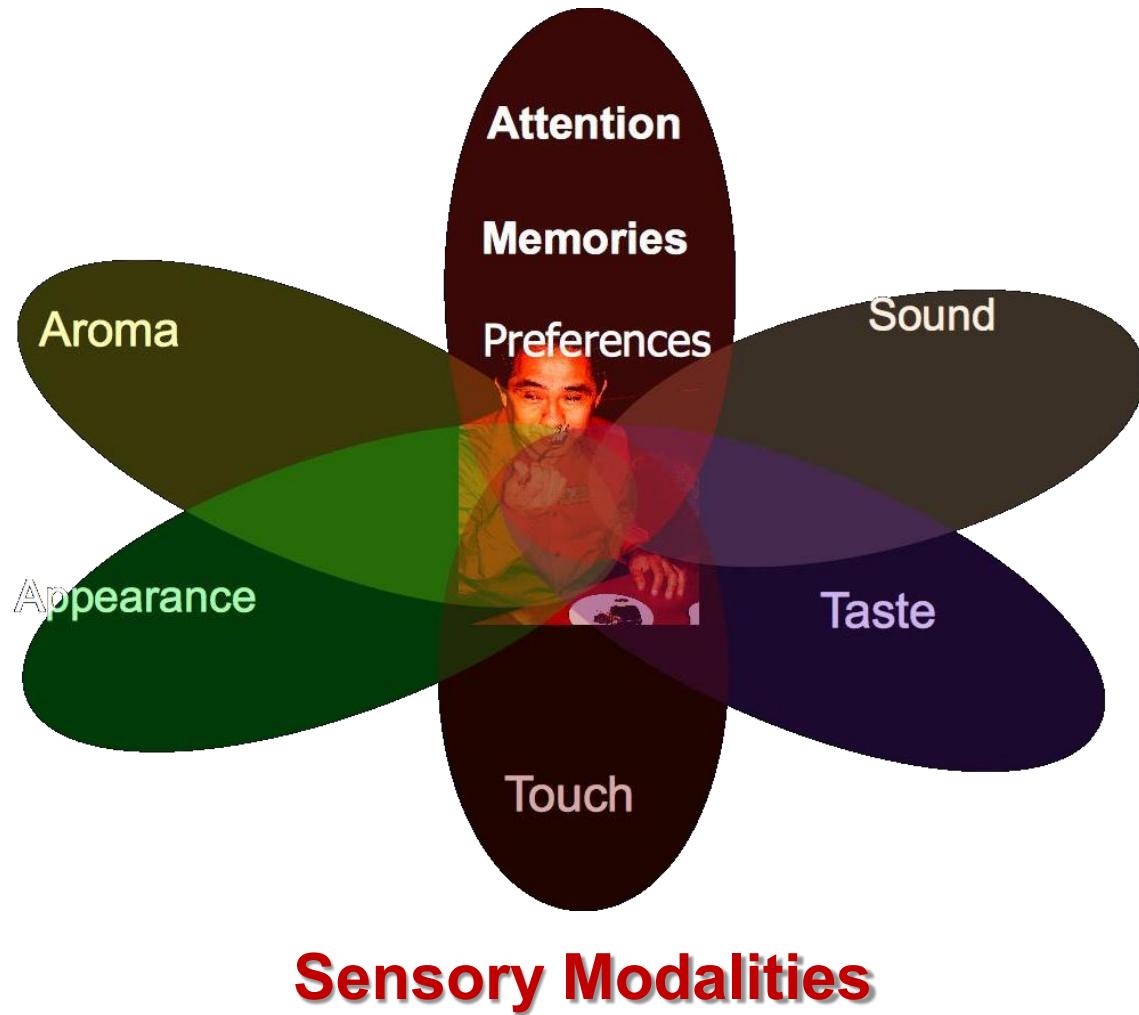


Ubiquitous



Wearable

# What is Multimodal?



# Multimodal Behaviors and Signals

## Language

- Lexicon
  - Words
- Syntax
  - Part-of-speech
  - Dependencies
- Pragmatics
  - Discourse acts

## Acoustic

- Prosody
  - Intonation
  - Voice quality
- Vocal expressions
  - Laughter, moans

## Visual

- Gestures
  - Head gestures
  - Eye gestures
  - Arm gestures
- Body language
  - Body posture
  - Proxemics
- Eye contact
  - Head gaze
  - Eye gaze
- Facial expressions
  - FACS action units
  - Smile, frowning

## Touch

- Haptics
- Motion

## Physiological

- Skin conductance
- Electrocardiogram

## Mobile

- GPS location
- Accelerometer
- Light sensors



# Multimodal Behaviors and Signals

High-Modality  
Multimodal  
Transformer:  
Quantifying  
Modality  
&Interaction  
Heterogeneity for  
High-Modality  
Representation  
Learning Liang et  
al 2023

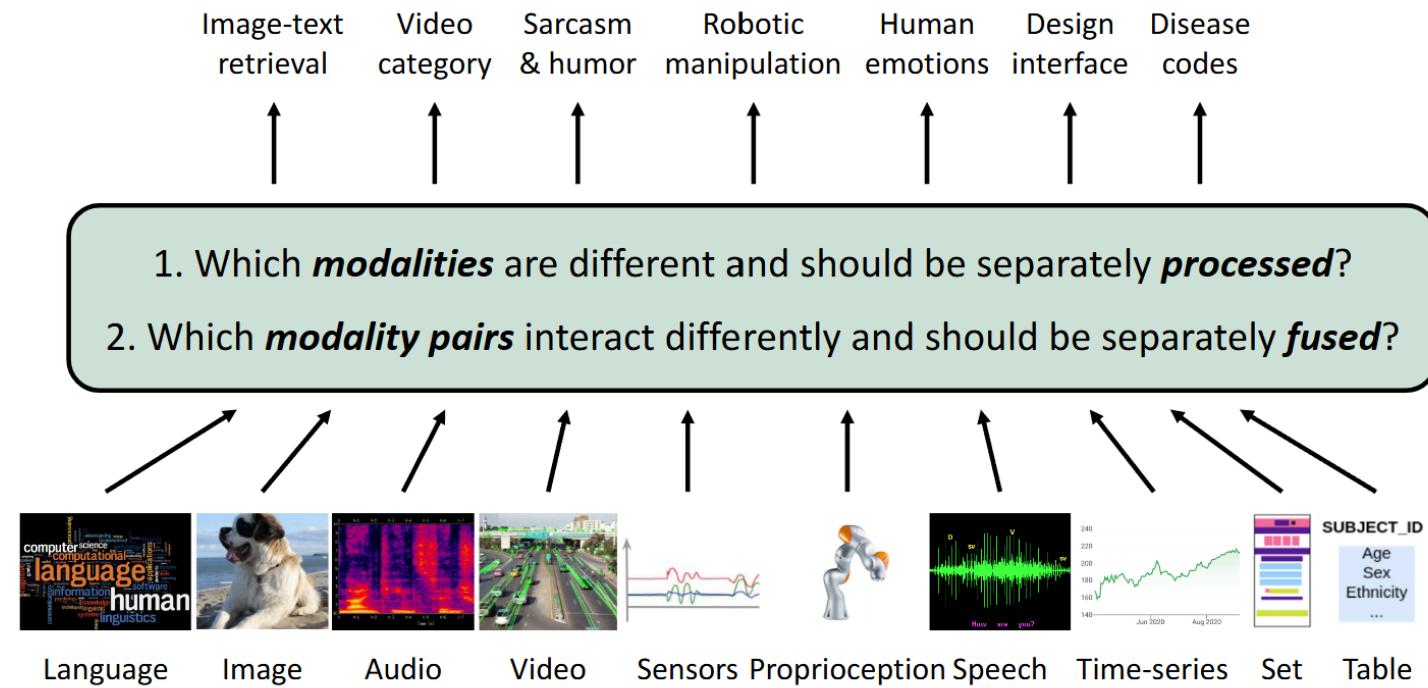
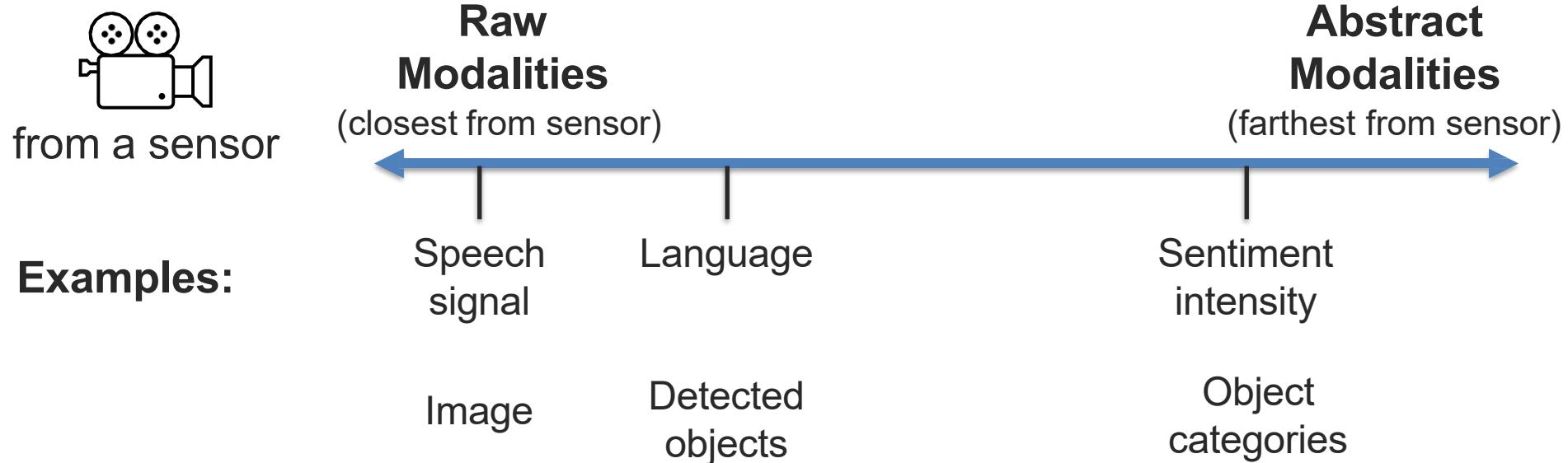


Figure 1: **Heterogeneity quantification:** Efficiently learning from many modalities requires measuring (1) *modality heterogeneity*: which modalities are different and should be separately processed, and (2) *interaction heterogeneity*: which modality pairs interact differently and should be separately fused. HIGHMMT uses these measurements to dynamically group parameters balancing performance and efficiency.

# What is a Modality?

## Modality

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



# What is Multimodal?

A dictionary definition...

**Multimodal:** with multiple modalities

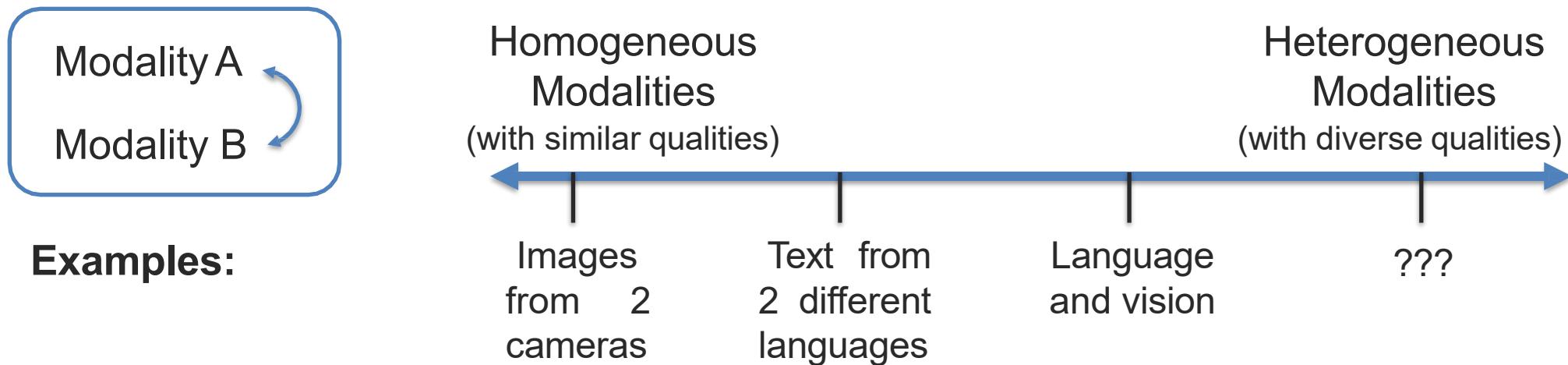
A research-oriented definition...

*Multimodal is the science of  
heterogeneous and interconnected data*



# Heterogeneous Modalities

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



Abstract modalities are more likely to be heterogeneous



# Dimensions of Heterogeneity

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



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

# Dimensions of Heterogeneity

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



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

①

**Element representations:** discrete, continuous, granularity



{*teacup, right, laptop, clean, room*}

# Dimensions of Heterogeneity

"Cross-Modal Discrete Representation Learning" by Alexander H. Liu et al

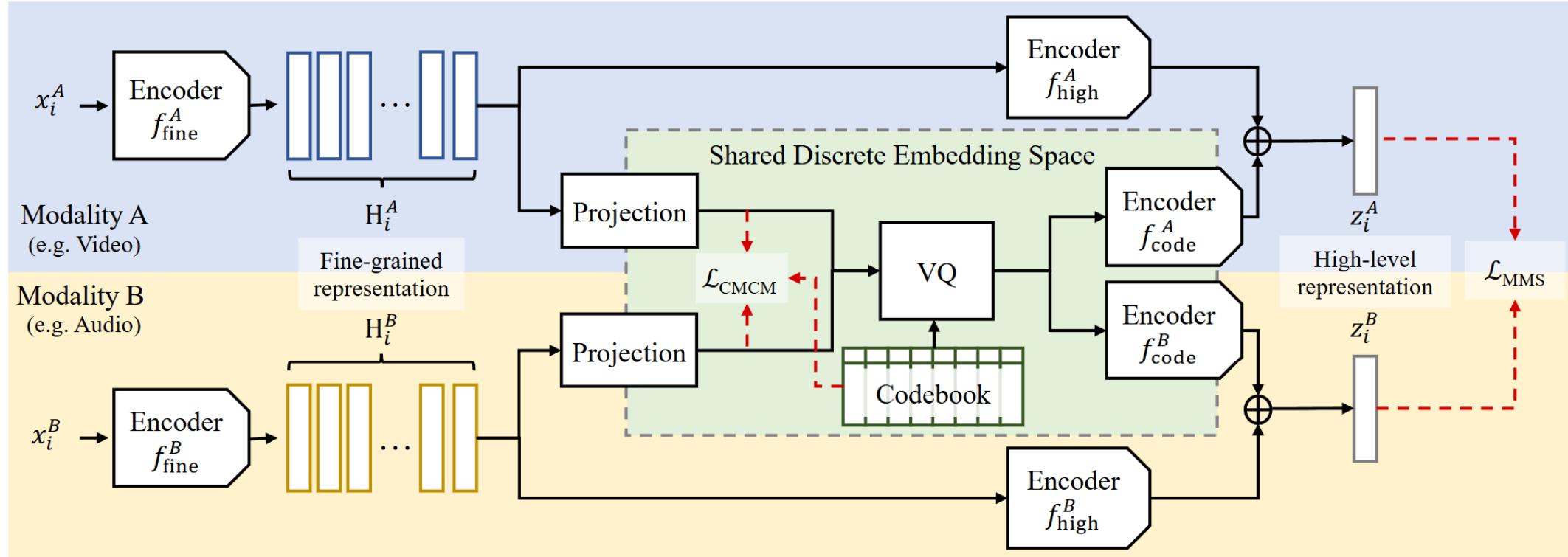


Figure 1: An overview of the proposed framework. The proposed shared discrete embedding space (green region, described in Section 2.2) is based on a cross-modal representation learning paradigm (blue and yellow regions, described in Section 2.1). The proposed Cross-Modal Code Matching  $\mathcal{L}_{\text{CMCM}}$  objective is detailed in Section 2.3 and Figure 2.

# Dimensions of Heterogeneity

②

## Element distributions: density, frequency



objects per image



words per minute

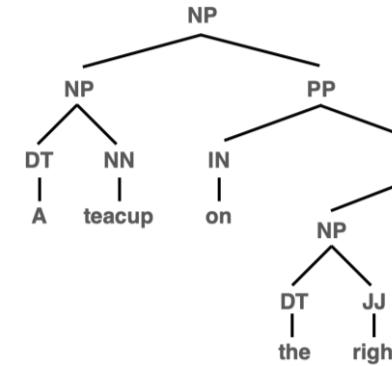
- Different modalities have different **densities and frequencies** of information.
  - Text: Words per minute.
  - Vision: Objects per image.
  - Audio: Frequencies of sound waves.

### Example

- **Speech signals** are **dense** (continuous waveforms).
- **Text documents** are **sparse** (symbolic, limited vocabulary).

# Dimensions of Heterogeneity

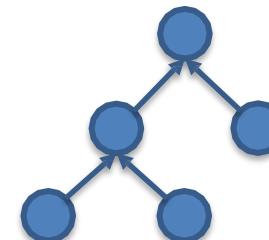
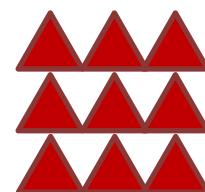
Information present in different modalities will often show diverse qualities, structures, and representations.



*A teacup on the right ...*

... } Latent (implicit)  
} Explicit (observable)

③ **Structure:** temporal, spatial, hierarchical, latent, explicit



# Dimensions of Heterogeneity

## Structural Differences

- Data can be **temporal, spatial, hierarchical, latent, or explicit**:
  - **Temporal**: Speech signals evolve over time.
  - **Spatial**: Objects in an image have spatial relationships.
  - **Hierarchical**: Sentences follow syntax trees.
  - **Latent**: Hidden features extracted by deep learning models.
  - **Explicit**: Directly observable structures.

## Example

- **Speech**: Time-series structure (temporal).
- **Image**: Object layout in a scene (spatial).
- **Text**: Syntactic dependency trees (hierarchical).



*A teacup on the right ...*



# Dimensions of Heterogeneity

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



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

④

**Information:** abstraction, entropy

# Dimensions of Heterogeneity

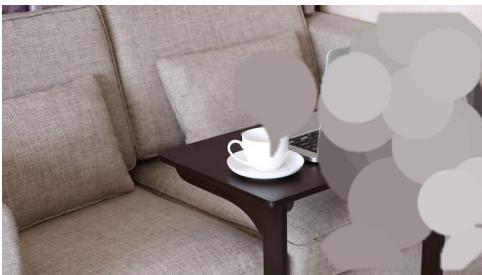
Information present in different modalities will often show diverse qualities, structures, and representations.



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

⑤

**Noise:** uncertainty, signal-to-noise ratio, missing data



teacup → **teacip**  
right → **rihjt**

# Dimensions of Heterogeneity

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



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

⑥

**Relevance:** task relevance, context dependence



- recreational
- living room
- right-handed

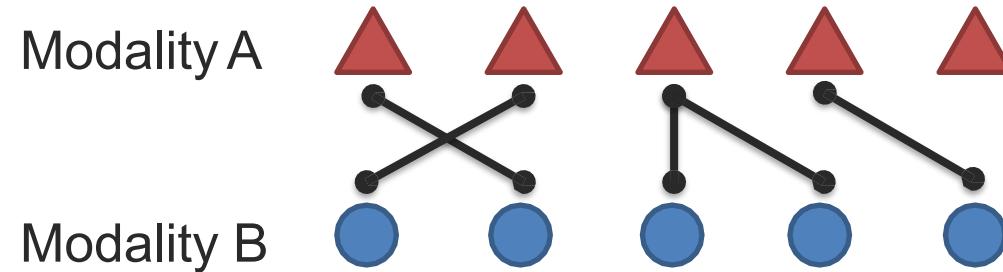
*A teacup on the  
right of a laptop  
in a clean room.*

- workspace
- study room

# Interconnected Modalities

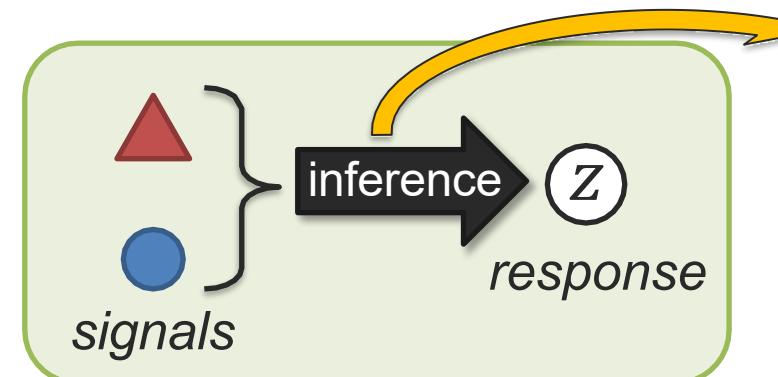
## ① Modality connections

*Modalities are often related and share commonality*



## ② Modality interactions

*Modality elements often interact during inference*



**Interactions happen during inference!**

“Inference” examples:

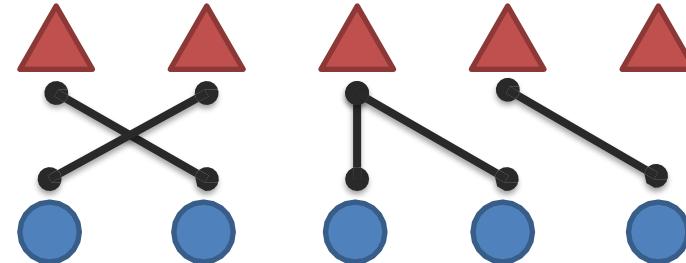
- Behavior perception
- Recognition task
- Modality translation

# Interconnected Modalities

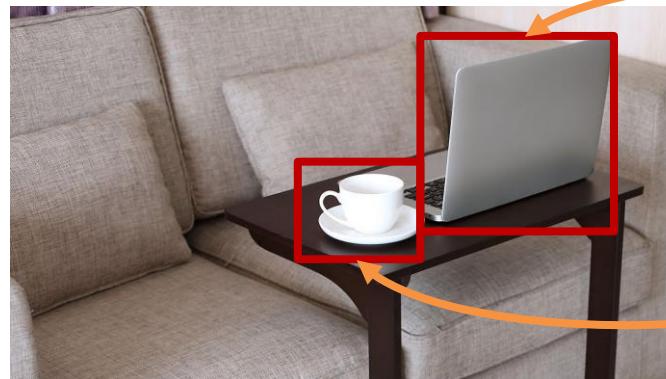
## ① Modality connections

*Modalities are often related and share commonality*

Modality A



Modality B



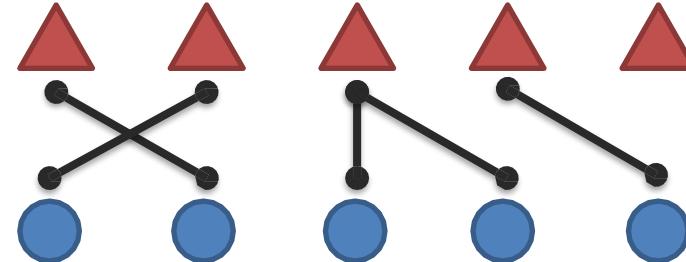
A **teacup** on the right of a **laptop** in a clean room.

# Interconnected Modalities

## ① Modality connections

*Modalities are often related and share commonality*

Modality A



Modality B

**Statistical**

**Semantic**



Association



e.g., correlation,

co-occurrence

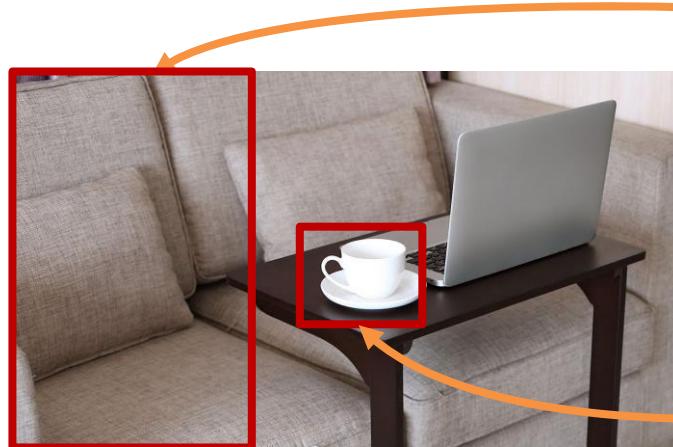
Correspondence



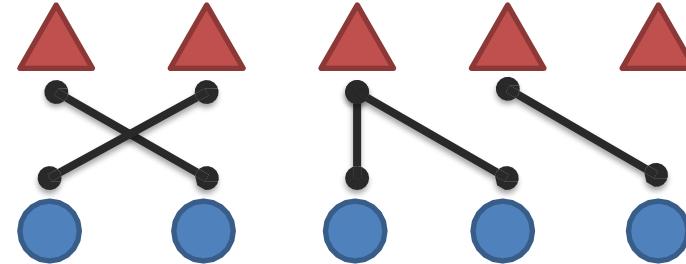
e.g., grounding

## ① Modality connections

*Modalities are often related and share commonality*



Modality A



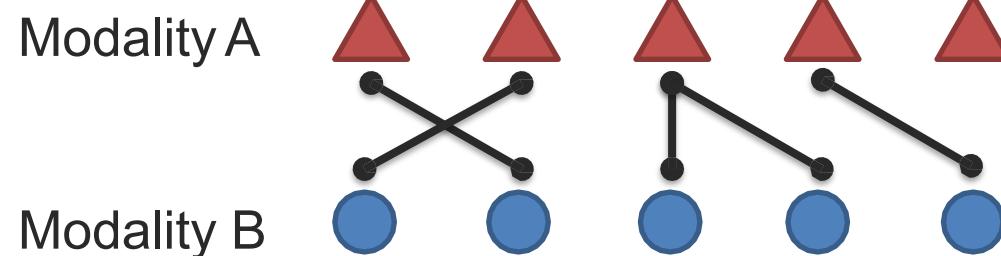
Modality B

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

# Interconnected Modalities

## ① Modality connections

Modalities are often related  
and share commonality



### Statistical



#### Association



e.g., correlation,  
co-occurrence

#### Dependency



e.g., causal,  
temporal

### Semantic



#### Correspondence



e.g., grounding

#### Relationship



e.g., function

# Interconnected Modalities

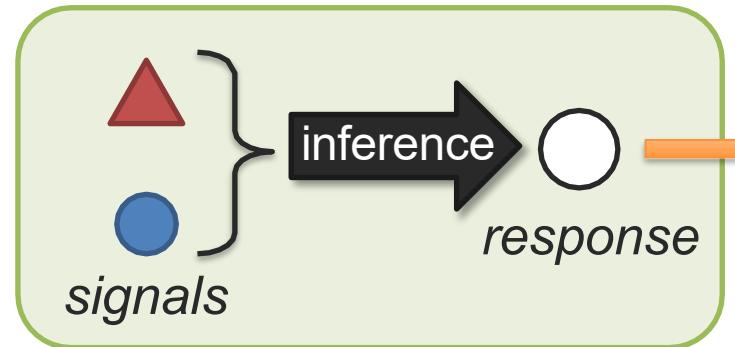
## ② Modality interactions

*Modality elements often interact during inference*

**Is this  
indoors?**



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

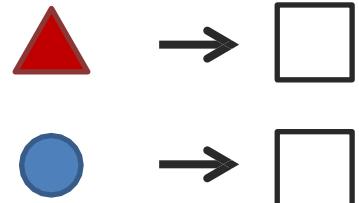


**Types of  
interaction  
responses?  
(a taxonomy)**

inference → Yes!

inference → Yes!

**Unimodal  
redundancy**



# Interconnected Modalities

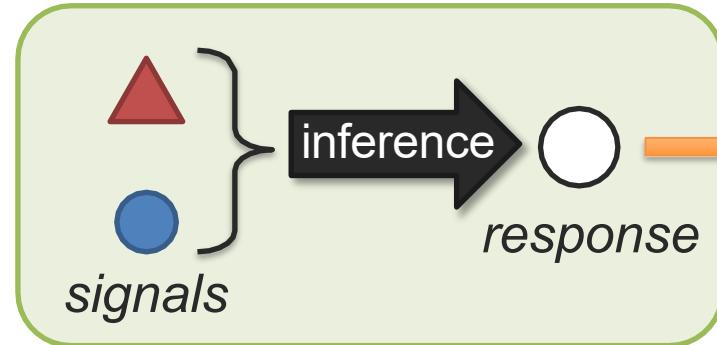
## ② Modality interactions

*Modality elements often interact during inference*

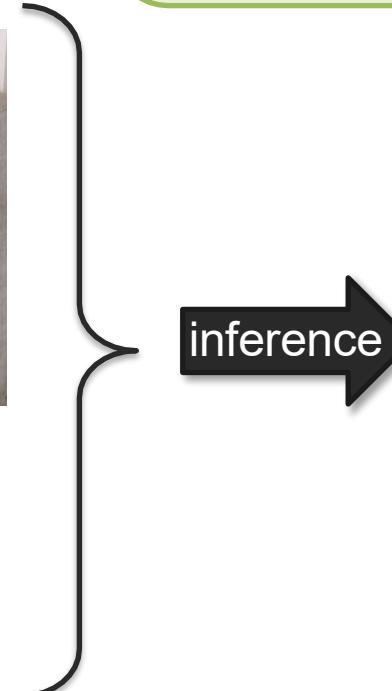
**Is this  
indoors**

?

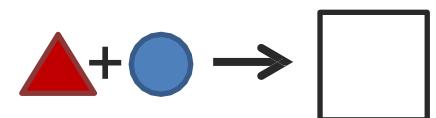
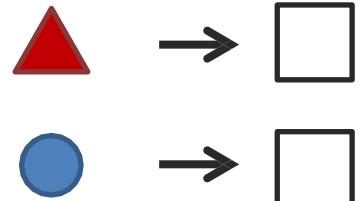
*A teacup on the right of a laptop in a clean room.*



**Yes!**



**Unimodal redundancy**



**Multimodal enhancement**

# Interconnected Modalities

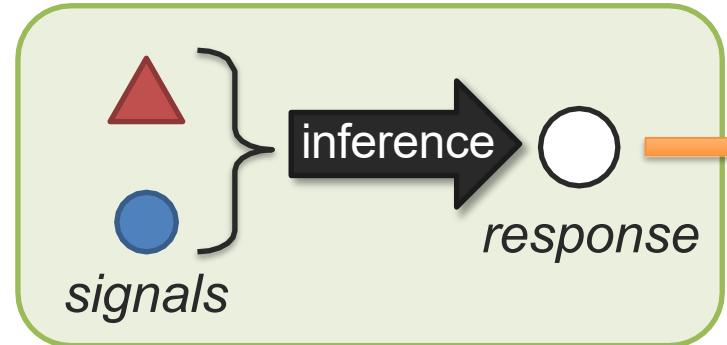
## ② Modality interactions

*Modality elements often interact during inference*

*Is this  
a living  
room?*



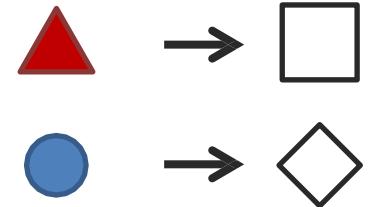
*A teacup on the right of a laptop in a clean room.*



inference → Yes!

inference → No, probably  
study room.

Unimodal  
Non-redundancy



# Interconnected Modalities

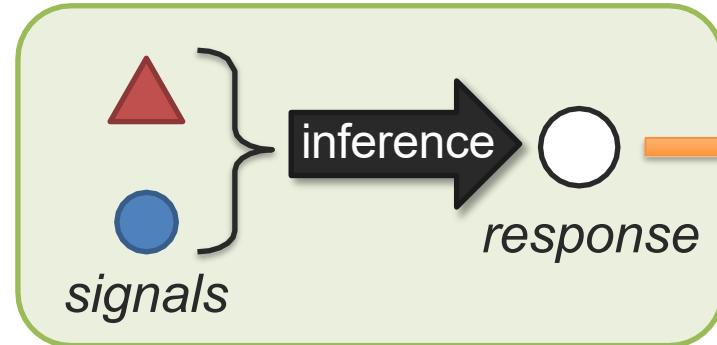
## ② Modality interactions

*Modality elements often interact during inference*

**Is this  
a  
living  
room?**

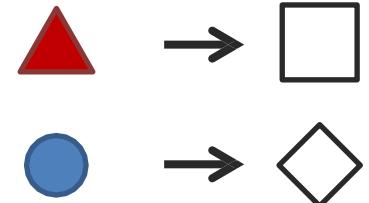


*A teacup on the right of a laptop in a clean room.*



**Types of  
interaction  
responses?**

**Unimodal  
Non-redundancy**

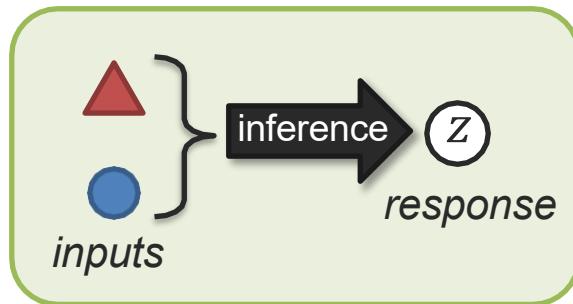


**Yes!**



**Multimodal  
dominance**

# Taxonomy of Interaction Responses – A Behavioral Science View



## Multimodal Communication

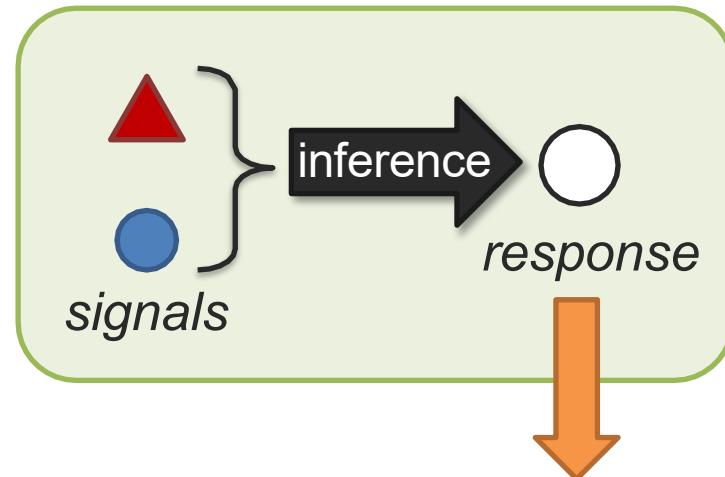


Redundancy	signal	response	
Redundancy	a. $\rightarrow$	<input type="square"/>	
	b. $\rightarrow$	<input type="square"/>	
Nonredundancy	$a+b \rightarrow$	<input type="square"/>	Equivalence
	$a+b \rightarrow$	<input type="square"/>	Enhancement
Nonredundancy	$a \rightarrow$	<input type="square"/>	
	$b \rightarrow$	<input type="circle"/>	
Nonredundancy	$a+b \rightarrow$	<input type="square"/> and <input type="circle"/>	Independence
	$a+b \rightarrow$	<input type="square"/>	Dominance
Nonredundancy	$a+b \rightarrow$	<input type="square"/> (or <input type="square"/> )	Modulation
	$a+b \rightarrow$	<input type="triangle"/>	Emergence



# Dimensions of Modality Interactions

What are the dimensions  
for **digitally-represented**  
modalities?



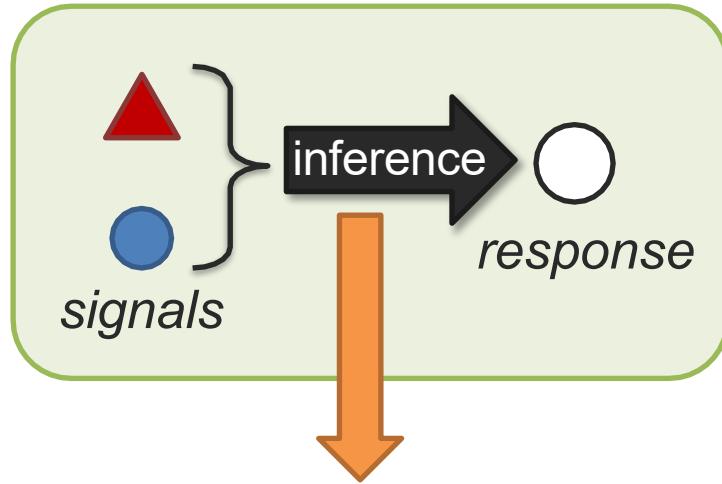
## ① Interaction Responses:

- Redundancy
- Non-redundancy
- Dominance
- Emergence...



# Dimensions of Modality Interactions

What are the dimensions  
for **digitally-represented**  
modalities?



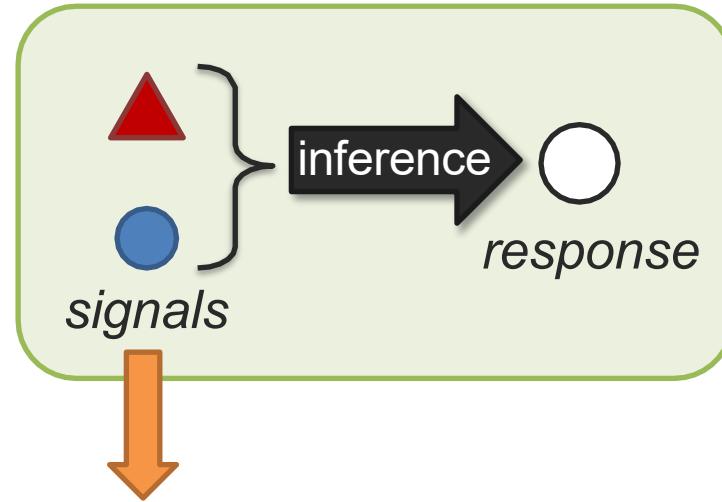
## ② Interaction Mechanics:

- Additive
- multiplicative
- Nonlinear
- Causal,
- Logical, ...



# Dimensions of Modality Interactions

What are the dimensions for **digitally-represented** modalities?

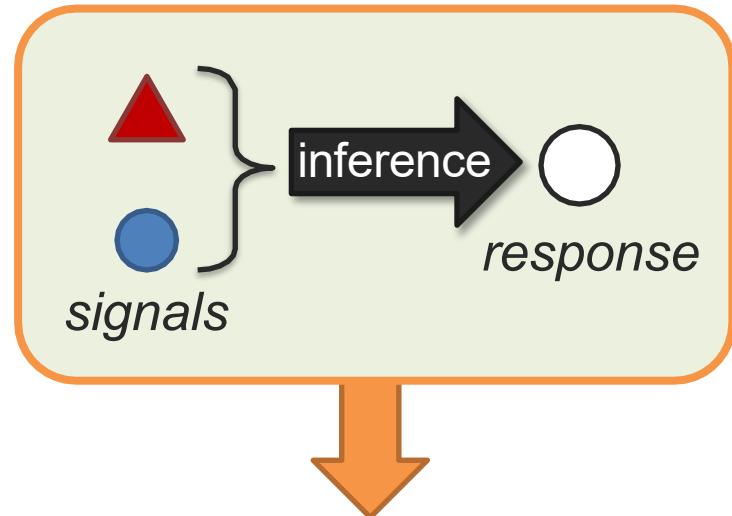


## ③ Input modalities:

- Unimodal
- Bimodal
- Trimodal
- High-modal, ...

# Dimensions of Modality Interactions

What are the dimensions  
for **digitally-represented**  
modalities?



## ④ Context:

- Structure context
- Task relevance
- Context dependence
- High-modal, ...

# What is Multimodal?

*Multimodal* is the science of  
**heterogeneous** and **interconnected** data 😊



# Multimodal Machine Learning



# What is Multimodal Machine Learning?

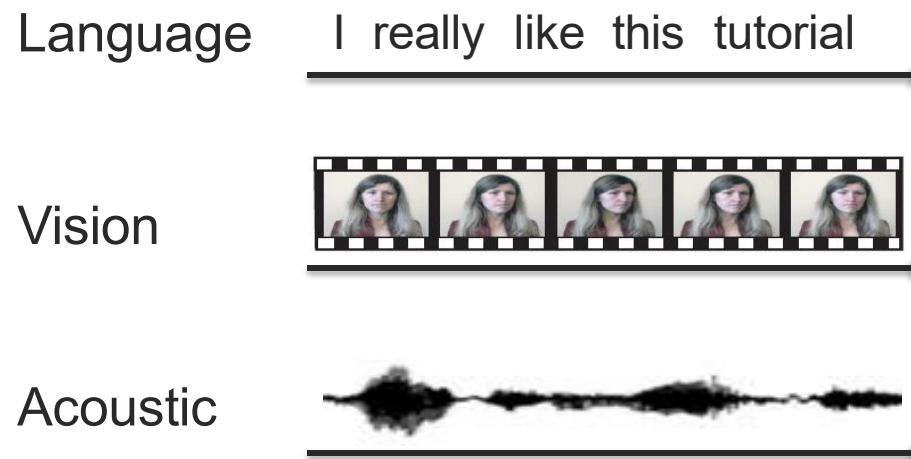
***Multimodal Machine Learning (ML)*** is the study of computer algorithms that learn and improve through the use and experience of data from multiple modalities

***Multimodal Artificial Intelligence (AI)*** studies computer agents able to demonstrate intelligence capabilities such as understanding, reasoning and planning, through multimodal experiences, and data

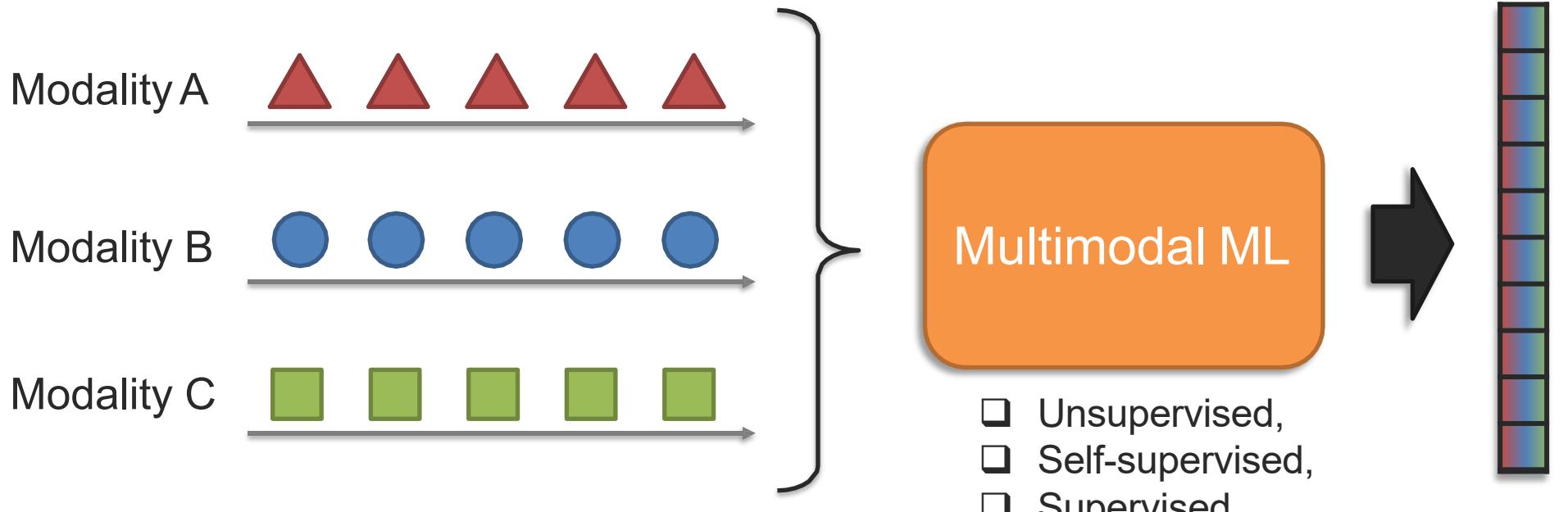
**Multimodal AI is a superset of Multimodal ML**



# Multimodal Machine Learning



# Multimodal Machine Learning



*What are the core multimodal technical challenges,*

## **Multimodal Machine Learning**

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



# Multimodal Technical Challenges – Surveys, Tutorials and Courses

## 2016

### Multimodal Machine Learning: A Survey and Taxonomy

*Tadas Baltrusaitis, Chaitanya Ahuja and Louis-Philippe Morency*  
(Arxiv 2017, IEEE TPAMI journal, February 2019)

<https://arxiv.org/abs/1705.09406>

**Tutorials:** CVPR 2016, ACL 2016, ICMI 2016, ...

### Graduate-level courses:

#### Multimodal Machine learning (11<sup>th</sup> edition)

<https://cmu-multicomp-lab.github.io/mmml-course/fall2020/>

#### Advanced Topics in Multimodal Machine learning

<https://cmu-multicomp-lab.github.io/adv-mmml-course/spring2022/>

## 2022

### Fundamentals of Multimodal ML: A Taxonomy & Open Challenges

*Paul Liang, Amir Zadeh and Louis-Philippe Morency*

- 6 core challenges
- 50+ taxonomic classes
- 600+ referenced papers

**Tutorials:** CVPR 2022, NAACL 2022, ...

### Updated graduate-level course:

**Multimodal Machine learning (12<sup>th</sup> edition)**  
Fall 2022 semester



# Challenge 1: Representation

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

This is a **core challenge** in multimodal learning because different modalities (text, images, audio, etc.) have **diverse structures** and need to be unified in a meaningful way.

**Individual elements:**

Modality A    

*It can be seen as a “local” representation*

or

Modality B    

*representation using holistic features*

# Challenge 1: Representation

## Example 1: Vision-Language Representation (CLIP)

- **Task:** Match images with text descriptions.
- **Problem:** Text and images have different feature spaces.
- **Solution:** CLIP (Contrastive Language–Image Pretraining) learns a **shared latent space** where related text and images are close together.

Learning  
Transferable  
Visual  
Models From  
Natural  
Language  
Supervision  
Radford  
2021

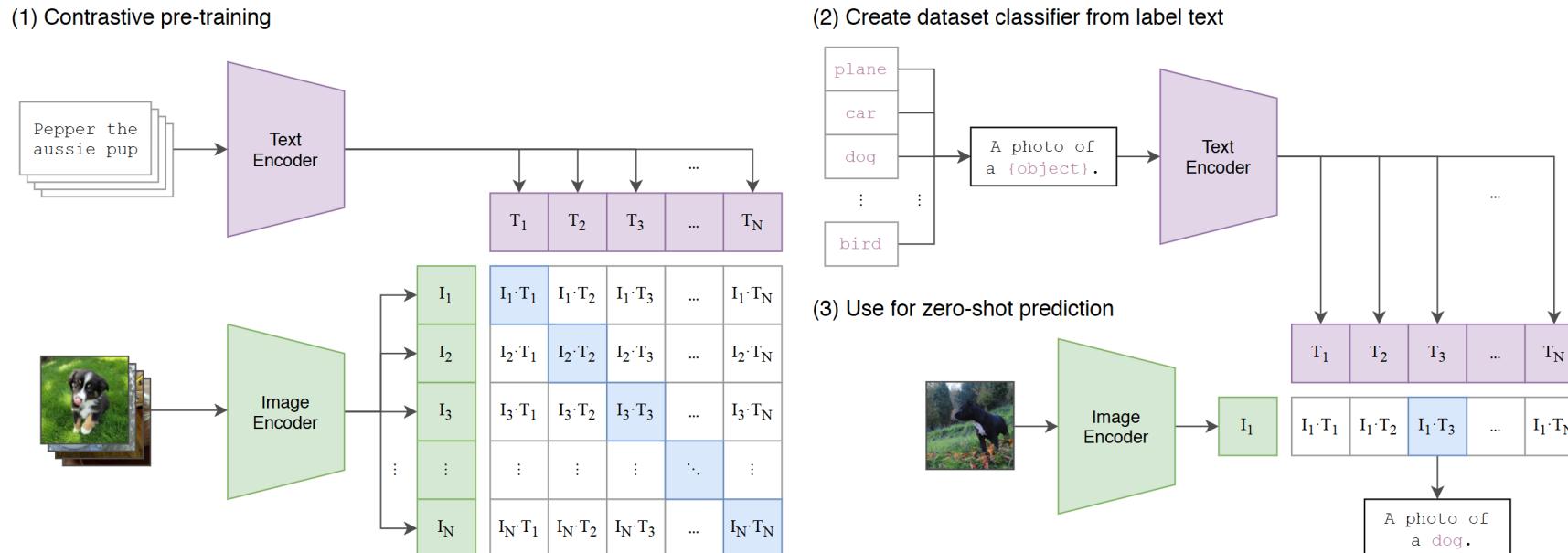


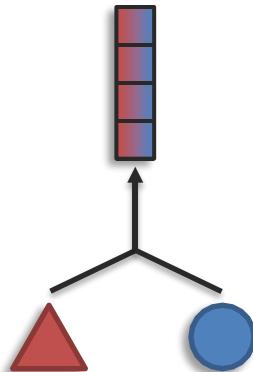
Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

# Challenge 1: Representation

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

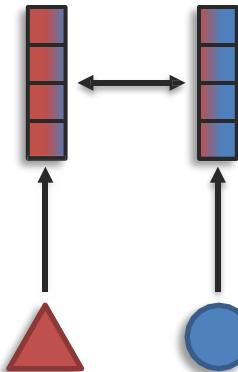
## Sub-challenges:

### Fusion



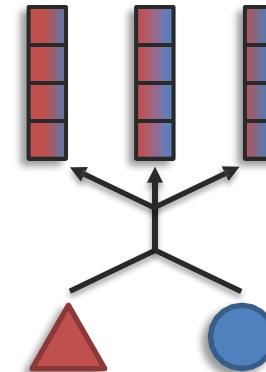
# modalities > # representations

### Coordination



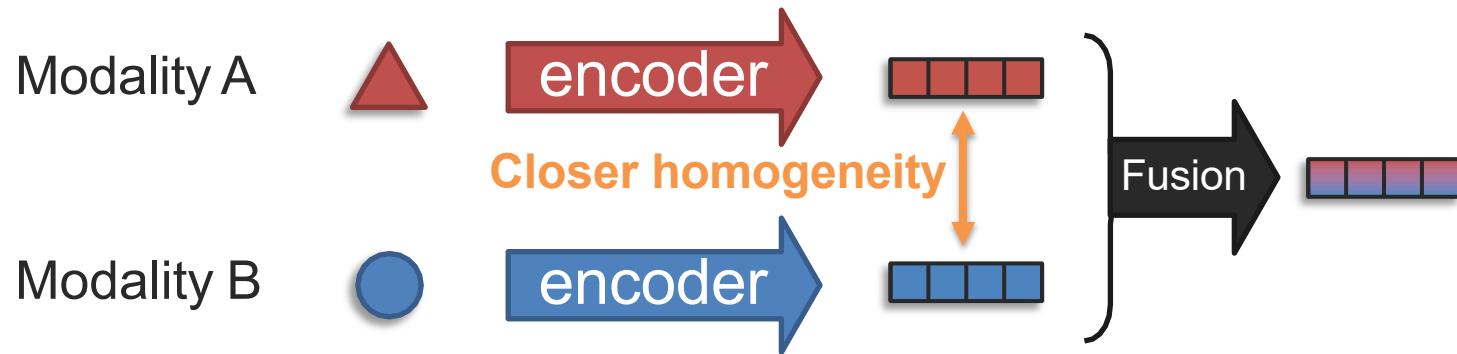
# modalities = # representations

### Fission

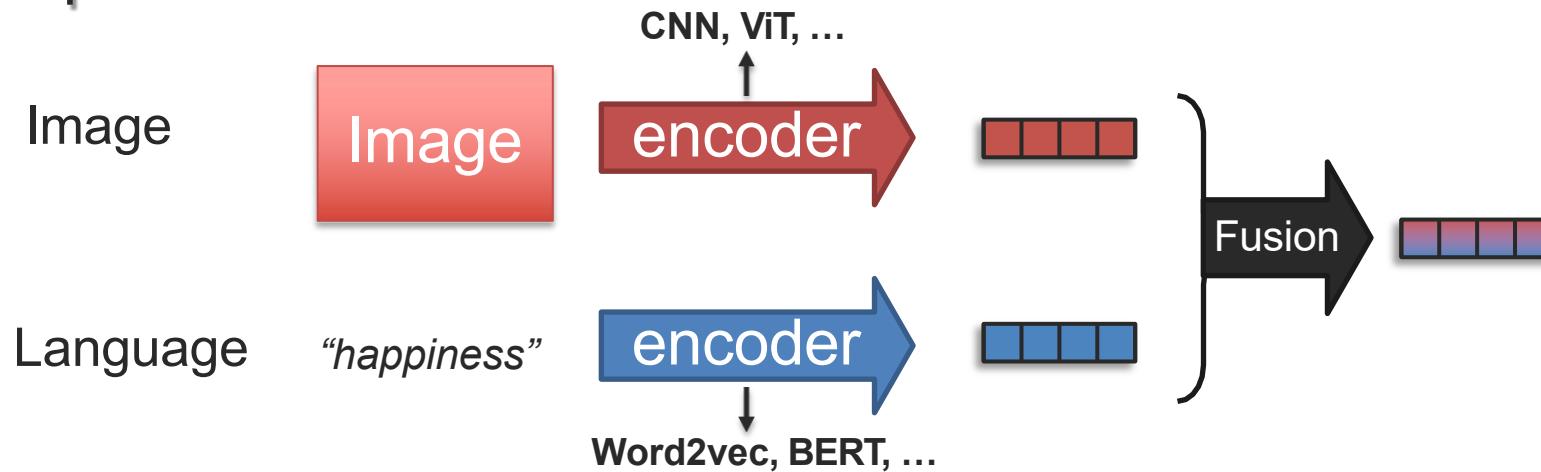


# modalities < # representations

# Fusion with Unimodal Encoders



Example:

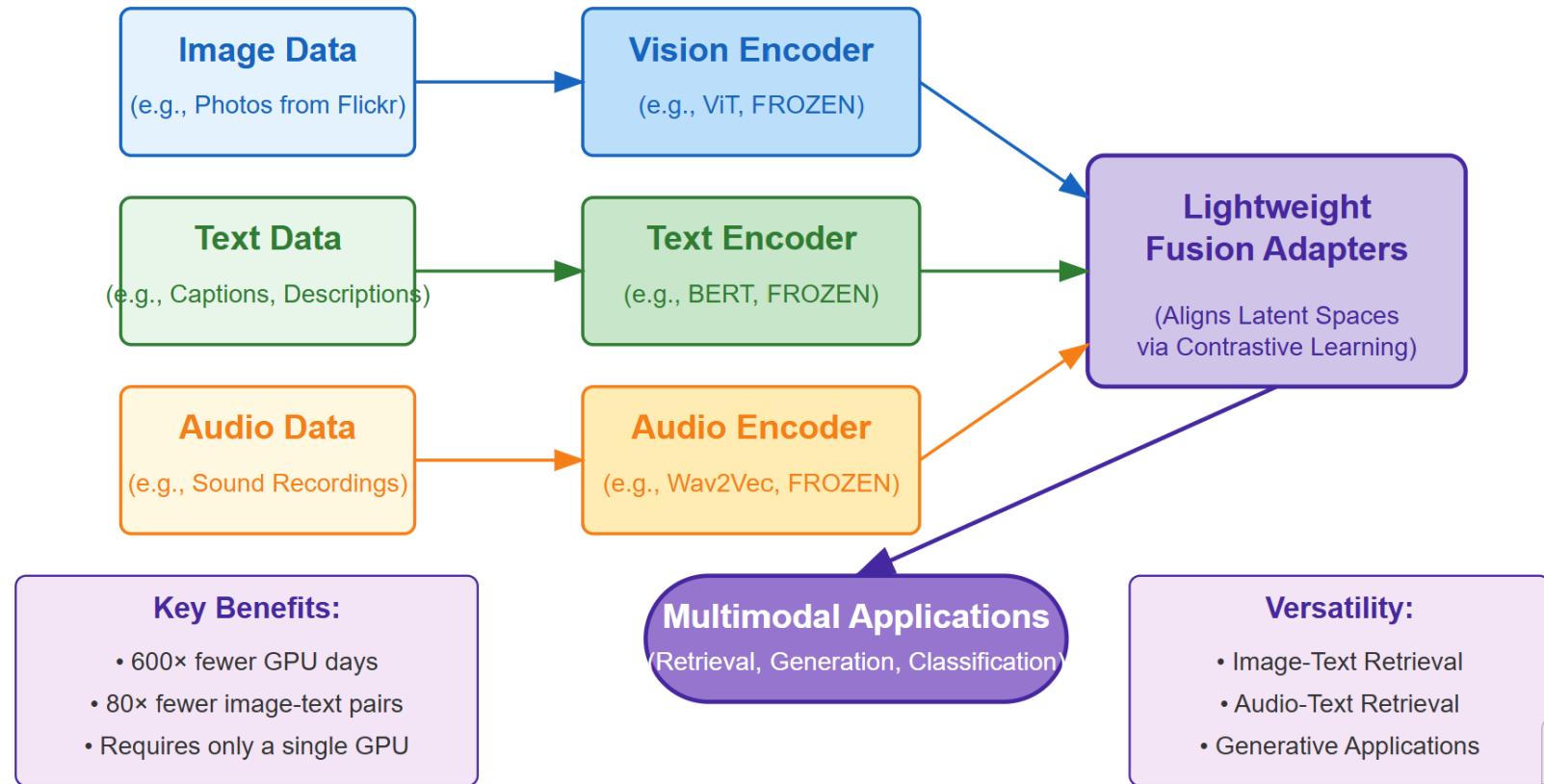


# Example paper

## FuseMix Framework Overview

Data-Efficient Multimodal Fusion on a Single GPU

Ideal Example Paper: "***Data-Efficient Multimodal Fusion on a Single GPU***"

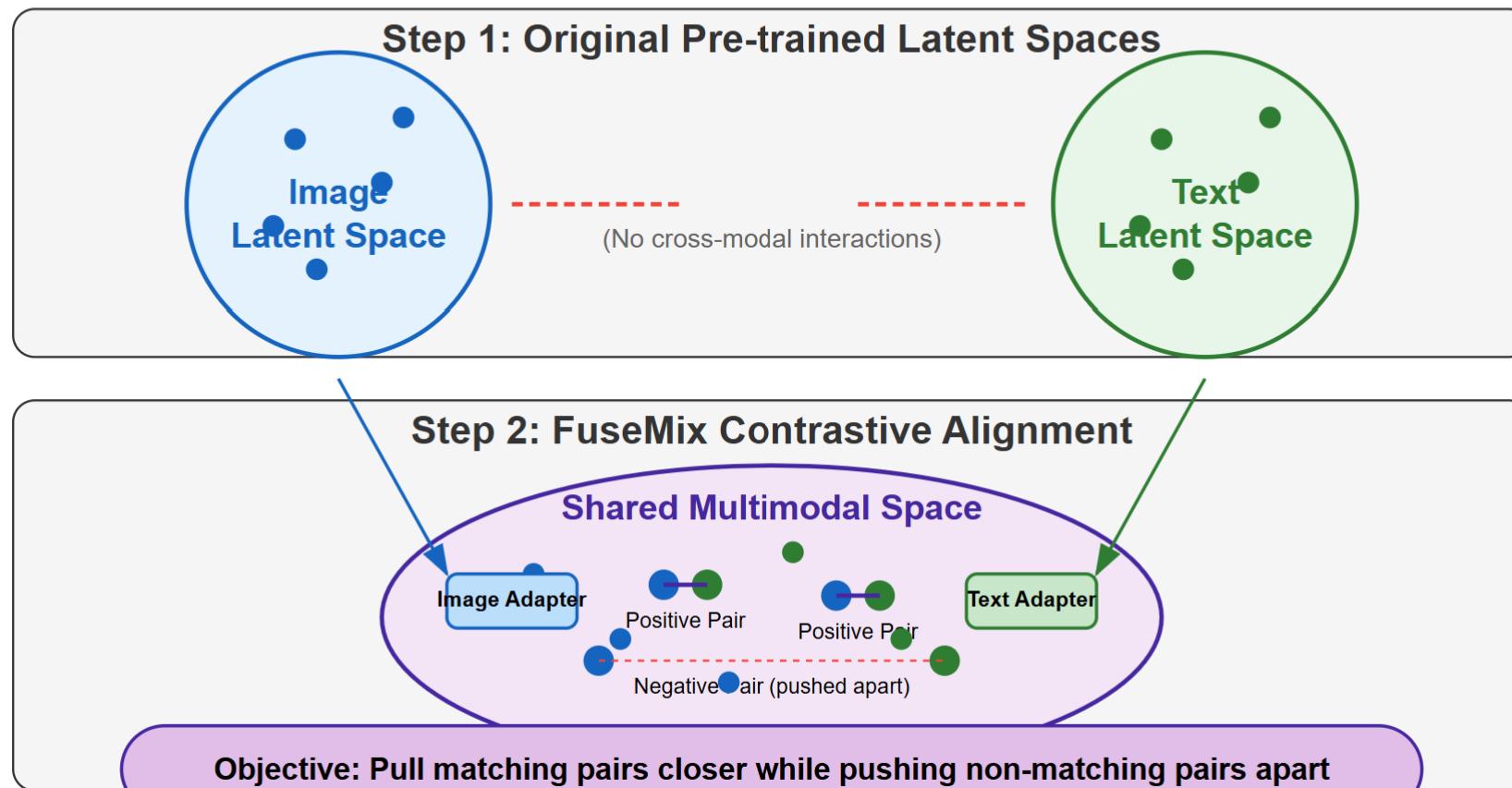


Efficient multimodal fusion using frozen pre-trained encoders and lightweight adapters

# Example paper

## Contrastive Latent Space Alignment

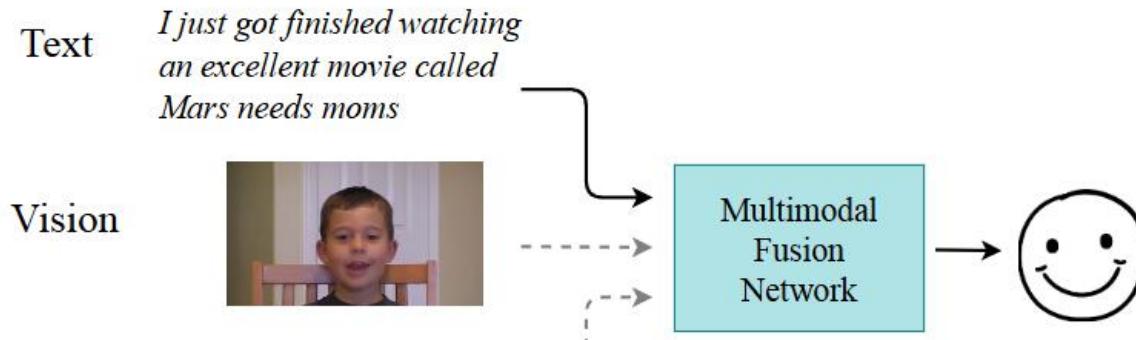
How FuseMix Aligns Modality-Specific Representations



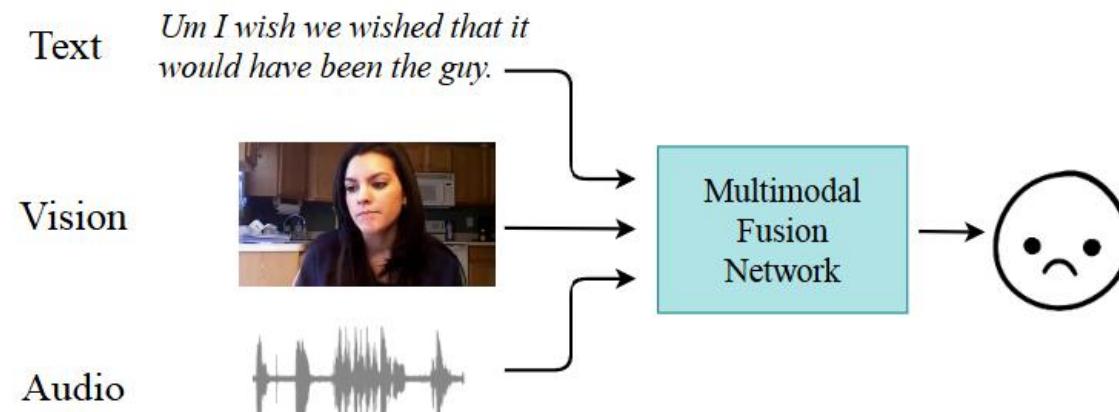
Ideal Example Paper: "Data-Efficient Multimodal Fusion on a Single GPU"

# Challenge 1: Representation

Dynamic  
Multimod  
al Fusion  
2023



(a)



(b)

Figure 1. Two examples in CMU-MOSEI [51] for emotion recognition. Figure (a) shows an “easy” multimodal instance as using textual information is sufficient to predict emotions correctly (this is a positive emotion). Figure (b) shows a “hard” example where all three modalities are required to make correct predictions (this is a negative emotion). While static multimodal fusion networks process “hard” and “easy” inputs identically, we propose *dynamic instance-wise inference* that can achieve computational savings for “easy” examples and preserve representation power for “hard” instances. For (a), DynMM only activates the text path and skips paths corresponding to the other two modalities, thus leading to computational efficiency.

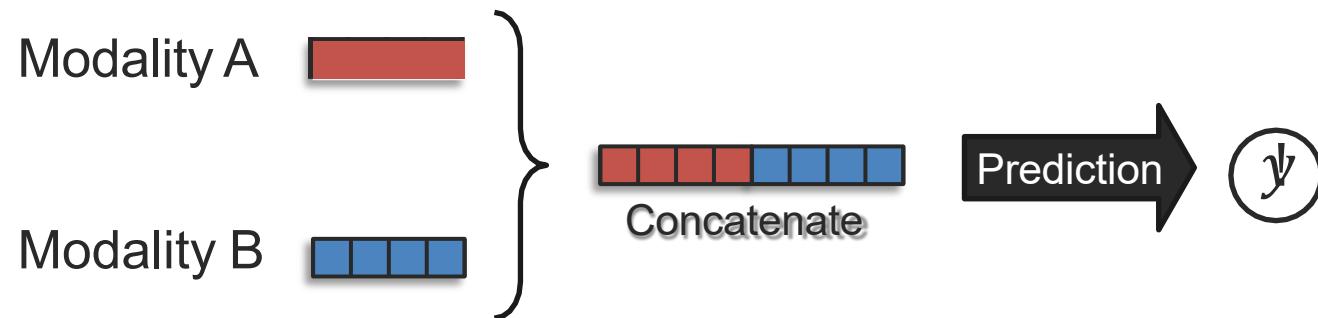


ĐẠI HỌC BÁCH

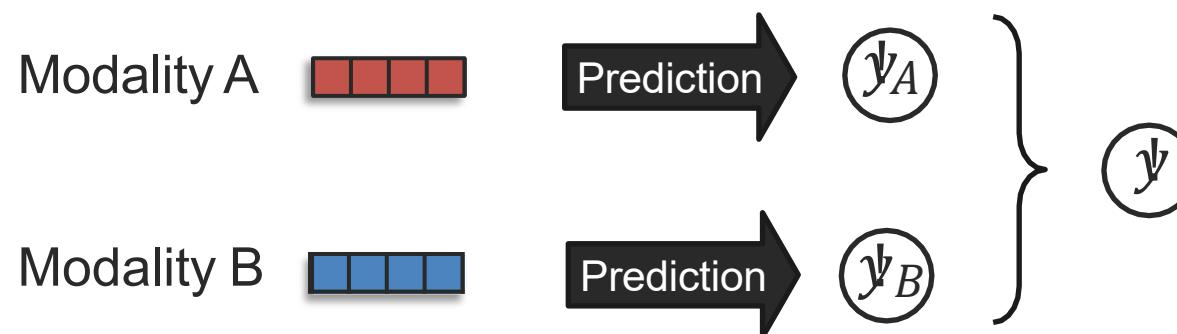
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

# Early and Late Fusion – A historical View

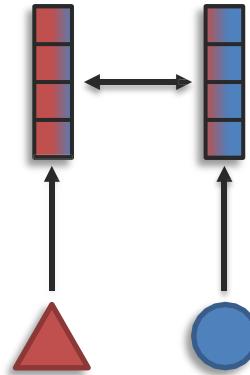
Early fusion:



Late fusion:

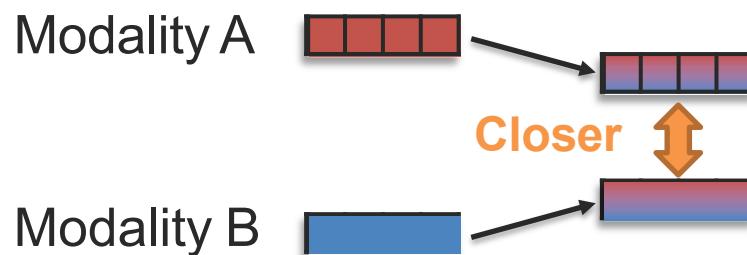


# Sub-Challenge 1b: Representation Coordination

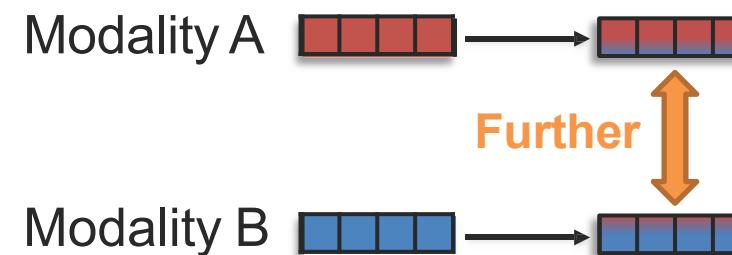


**Definition:** Learn multimodally-contextualized representations that are coordinated through their cross-modal interactions

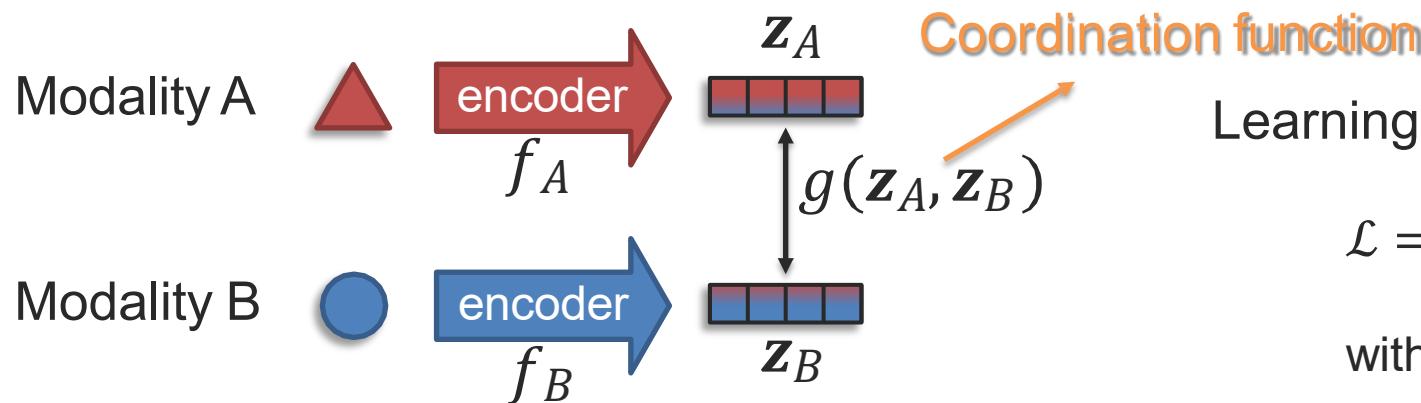
Strong Coordination:



Partial Coordination:



# Coordination Function



Learning with coordination function:

$$\mathcal{L} = g(f_A(\text{red triangle}), f_B(\text{blue circle}))$$

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

→ Requires paired data

Examples of coordination function:

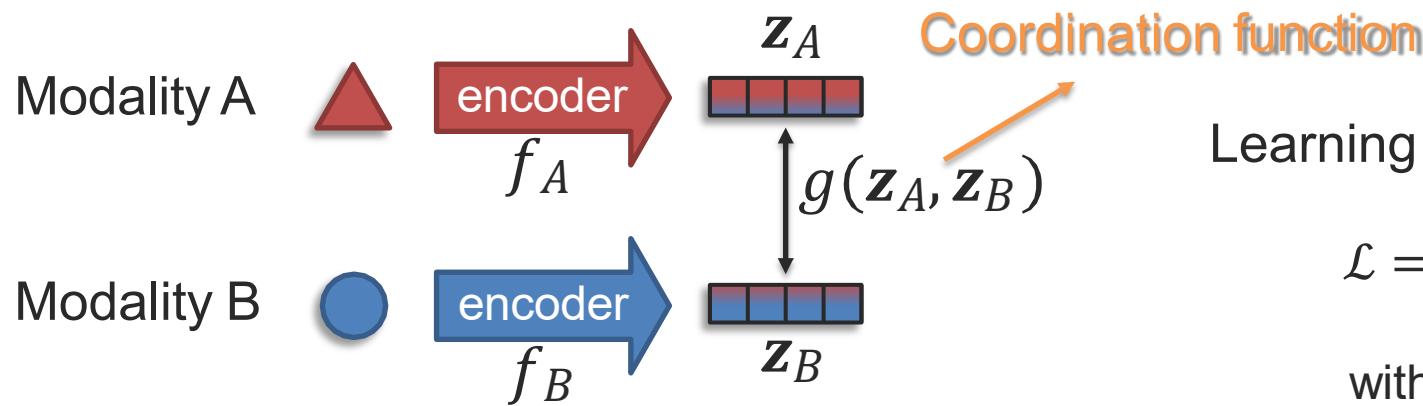
1 Cosine similarity:

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

Strong coordination!

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

# Coordination Function



Learning with coordination function:

$$\mathcal{L} = g(f_A(\text{red triangle}), f_B(\text{blue circle}))$$

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

Examples of coordination function:

2 Kernel similarity functions:

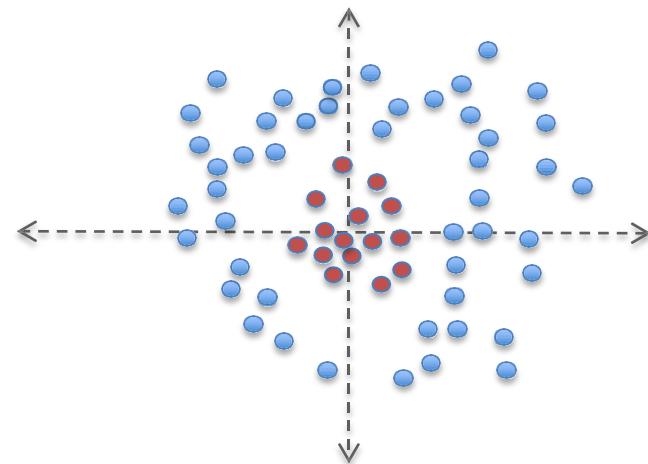
$$g(\mathbf{z}_A, \mathbf{z}_B) = k(\mathbf{z}_A, \mathbf{z}_B) \quad \left\{ \begin{array}{l} \cdot \text{ Linear} \\ \cdot \text{ Polynomial} \\ \cdot \text{ Exponential} \\ \cdot \text{ RBF} \end{array} \right.$$

All these examples bring relatively strong coordination between modalities

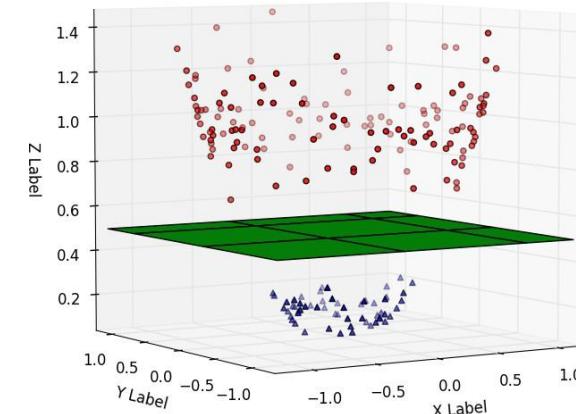
# Kernel Function

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

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

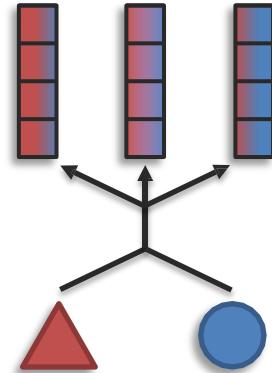


Not linearly separable in  $x$  space



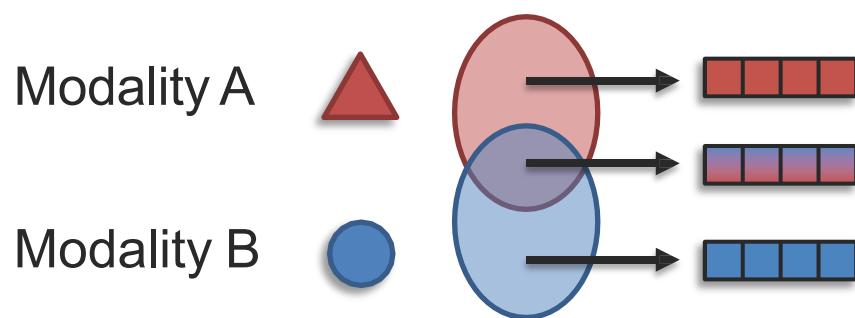
Same data, but now linearly separable in  $\phi(x)$  space

# Sub-Challenge 1c: Representation Fission

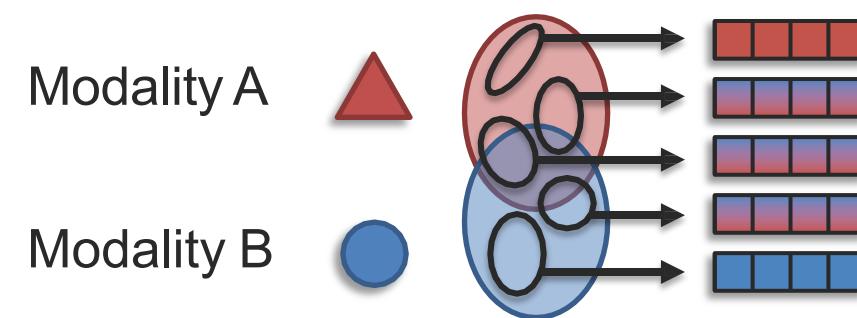


**Definition:** Learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering

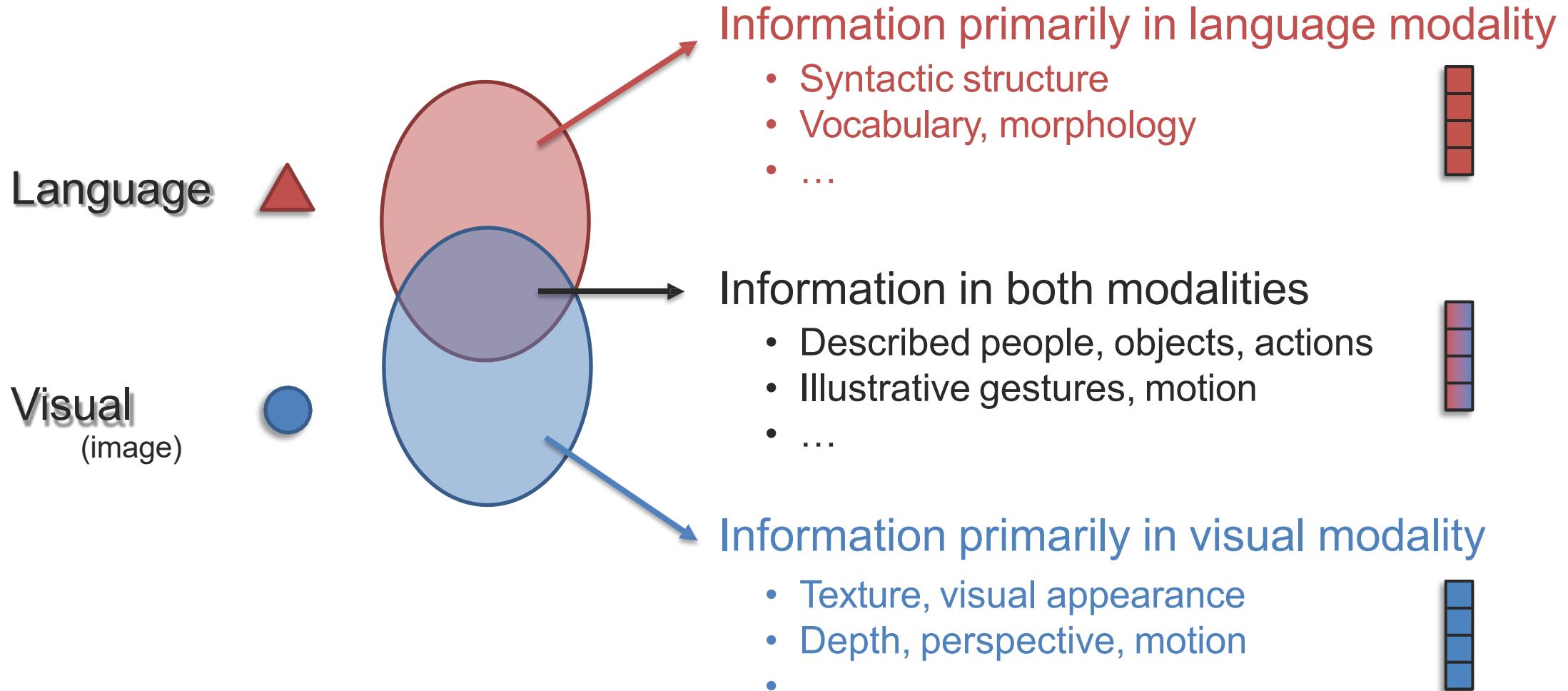
Modality-level fission:



Fine-grained fission:



# Modality-Level Fission

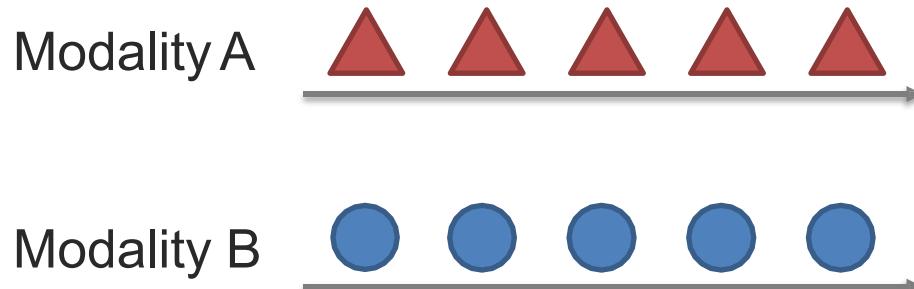


# Challenge 2: Alignment

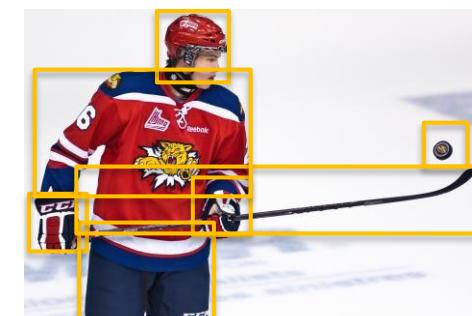
**Definition:** Identifying and modeling cross-modal connections between all elements of multiple modalities, building from the data structure

→ Most modalities have internal structure with multiple elements

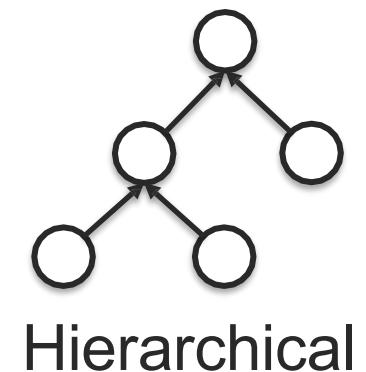
Elements with temporal structure:



Other structured examples:



Spatial



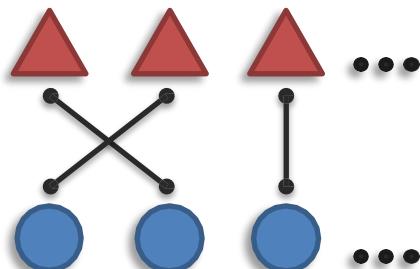
Hierarchical

# Challenge 2: Alignment

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

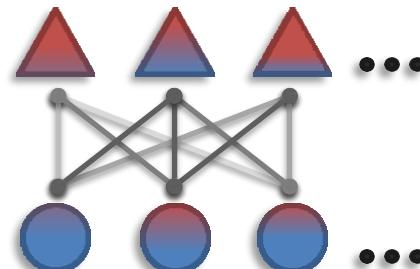
## Sub-challenges:

### Connections



Explicit alignment  
(e.g., grounding)

### Aligned Representation



Alignment + representation  
(aka, contextualized representation)

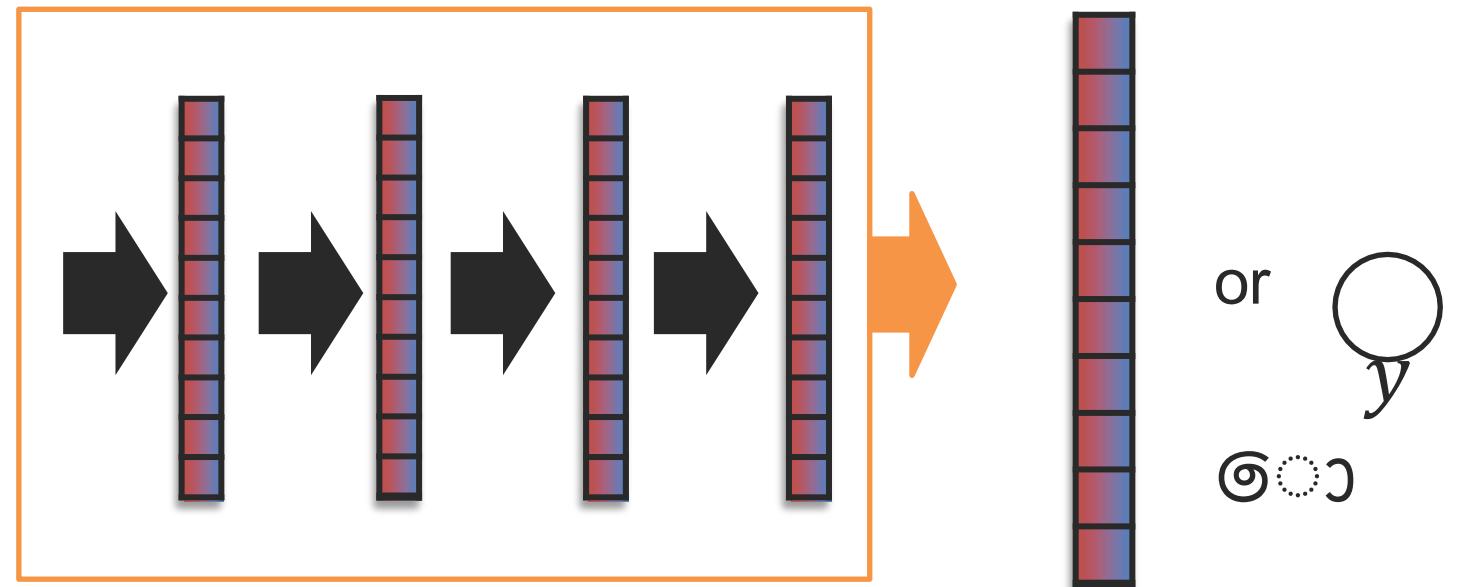
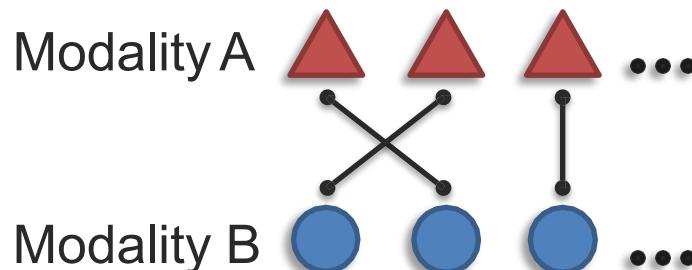
### Elements



Segmentation of  
individual elements

# Challenge 3: Reasoning

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure to derive meaningful conclusions.



# Challenge 3: Reasoning

"Learning to  
Reason: End-to-End  
Module Networks  
for Visual Question  
Answering"

 *Andreas et al.,  
NeurIPS 2016*

There is a shiny object that is right of the gray metallic cylinder;  
does it have the same size as the large rubber sphere?

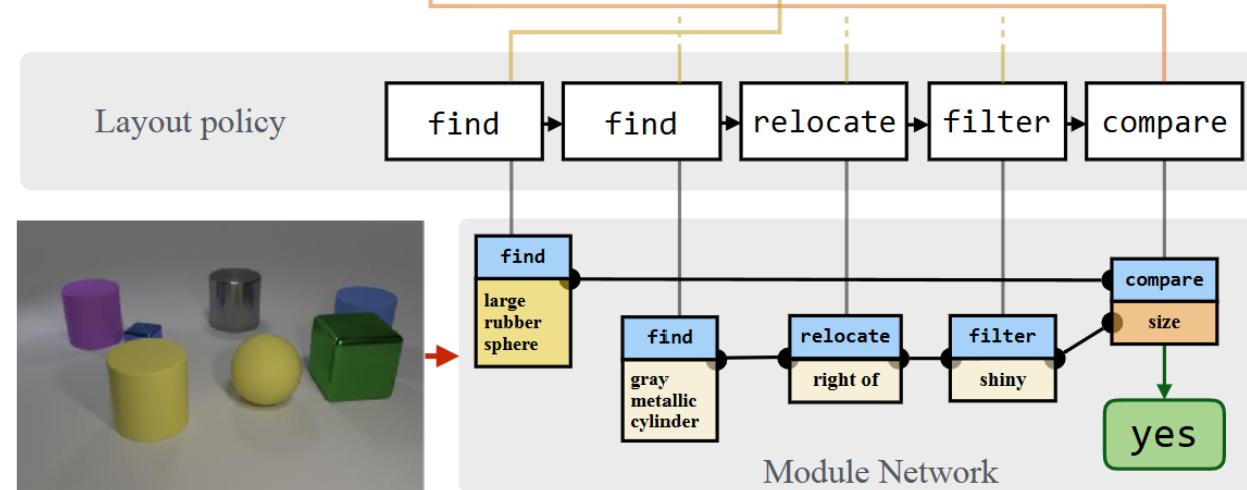
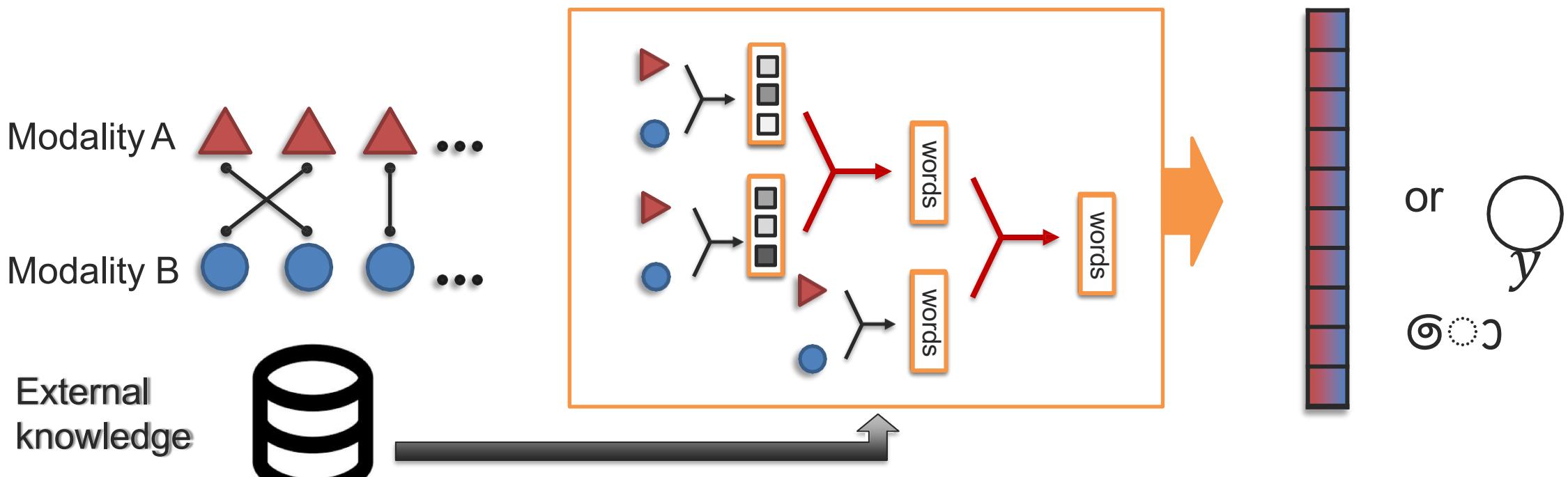


Figure 1: For each instance, our model predicts a computational expression and a sequence of attentive module parameterizations. It uses these to assemble a concrete network architecture, and then executes the assembled neural module network to output an answer for visual question answering. (The example shows a real structure predicted by our model, with text attention maps simplified for clarity.)

# Challenge 3: Reasoning

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure to derive meaningful conclusions.



# Challenge 3: Reasoning



**Q: Which American president is associated with the stuffed animal seen here?**

**A: Teddy Roosevelt**

## Outside Knowledge

Another lasting, popular legacy of Roosevelt is the stuffed toy bears—teddy bears—named after him following an incident on a hunting trip in Mississippi in 1902.

Developed apparently simultaneously by toymakers ... and named after President Theodore "Teddy" Roosevelt, the teddy bear became an iconic children's toy, celebrated in story, song, and film.

At the same time in the USA, Morris Michtom created the first teddy bear, after being inspired by a drawing of Theodore "Teddy" Roosevelt with a bear cub.

**"OK-VQA: A Visual Question Answering Benchmark Requiring External Knowledge"**  
 *Marino et al., CVPR 2019*

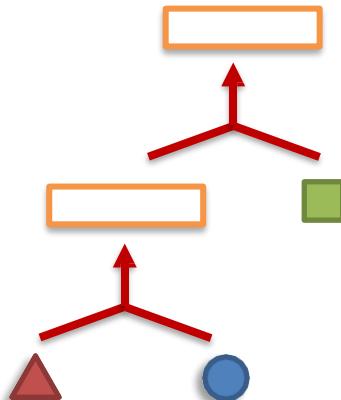
Figure 1: We propose a novel dataset for visual question answering, where the questions require external knowledge resources to be answered. In this example, the visual content of the image is not sufficient to answer the question. A set of facts about teddy bears makes the connection between teddy bear and the American president, which enables answering the question.

# Challenge 3: Reasoning

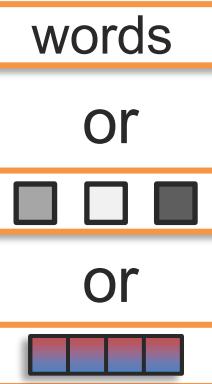
**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure

## Sub-challenges:

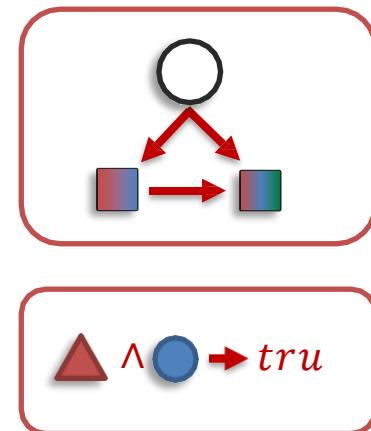
### Structure Modeling



### Intermediate concepts



### Inference Paradigm



### External Knowledge



# The Challenge of Compositionality

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



(a) some plants surrounding a lightbulb



(b) a lightbulb surrounding some plants

CLIP, ViLT, ViLBERT, etc.  
**All random chance**

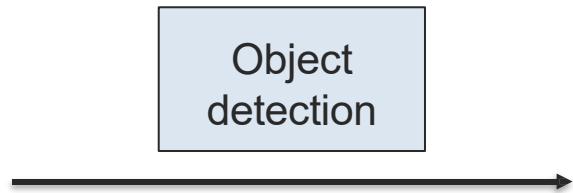
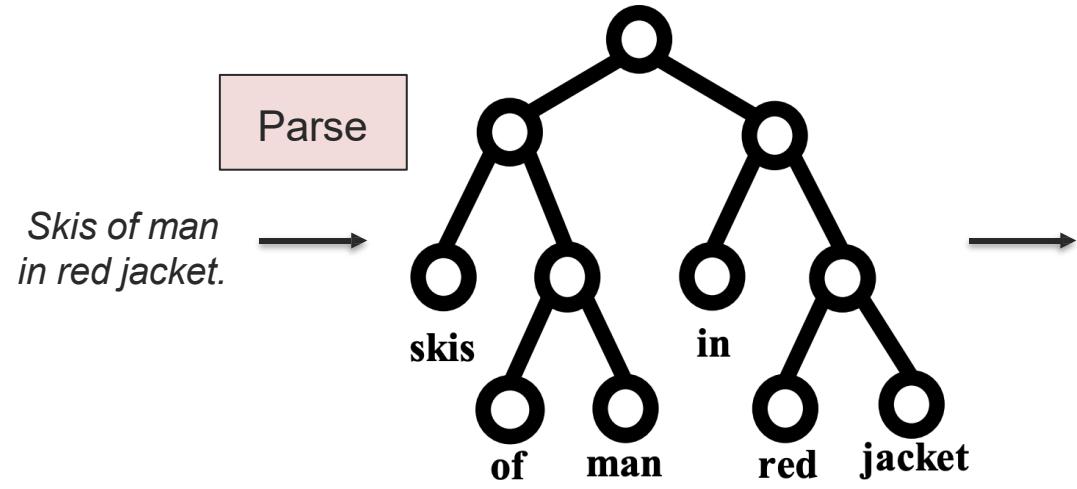
Compositional Generalization  
**to novel combinations outside of training data**

1. Structure: <subject> <verb> <object>
2. Concepts: ‘plants’, ‘lightbulb’
3. Inference: ‘surrounding’ – spatial relation
4. Knowledge: from humans!

[Thrush et al., Winoground: Probing Vision and Language Models for Visio-Linguistic Compositionality. CVPR 2022]

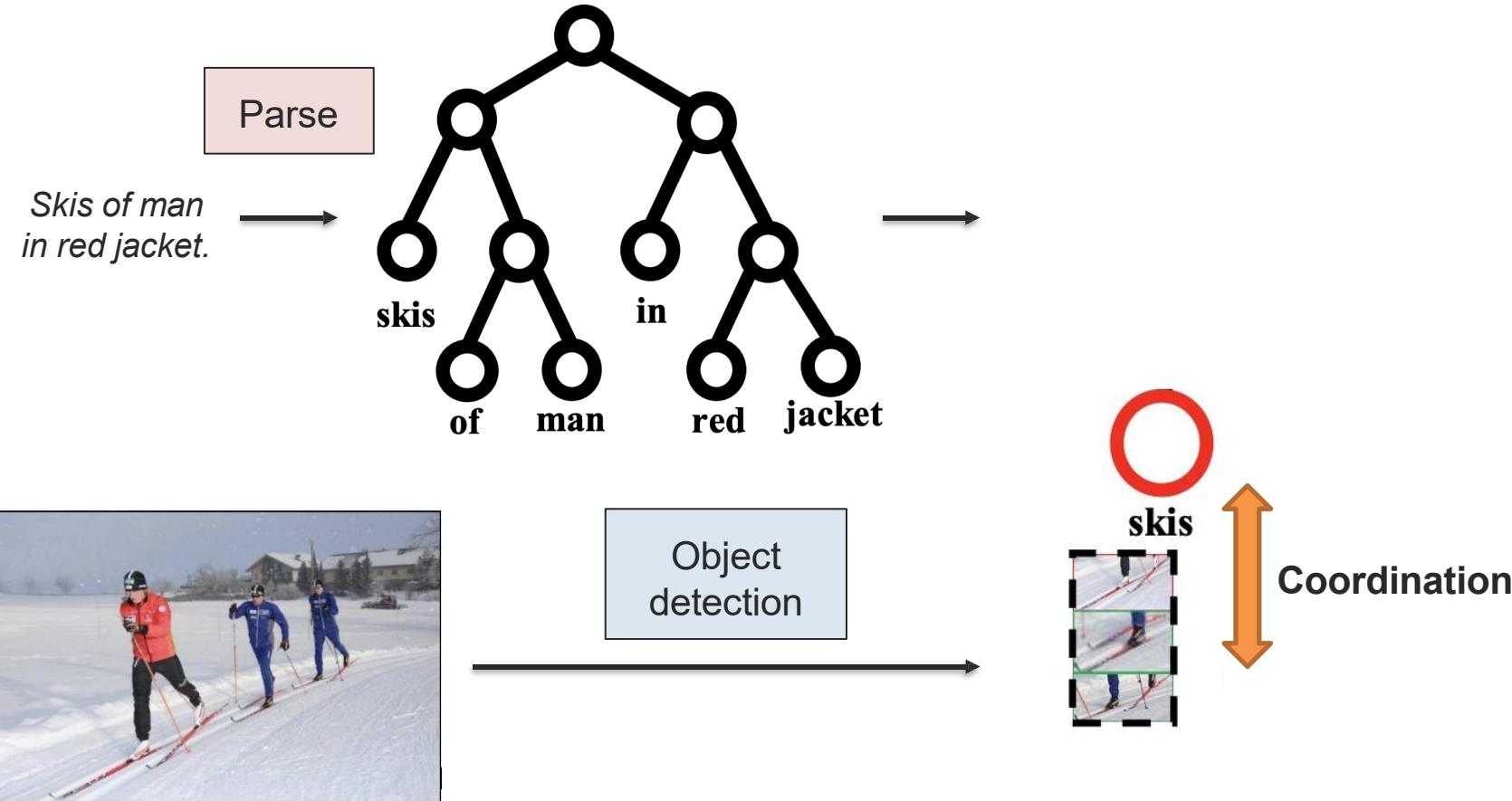
# Hierarchical Structure

# Leverage syntactic structure of language



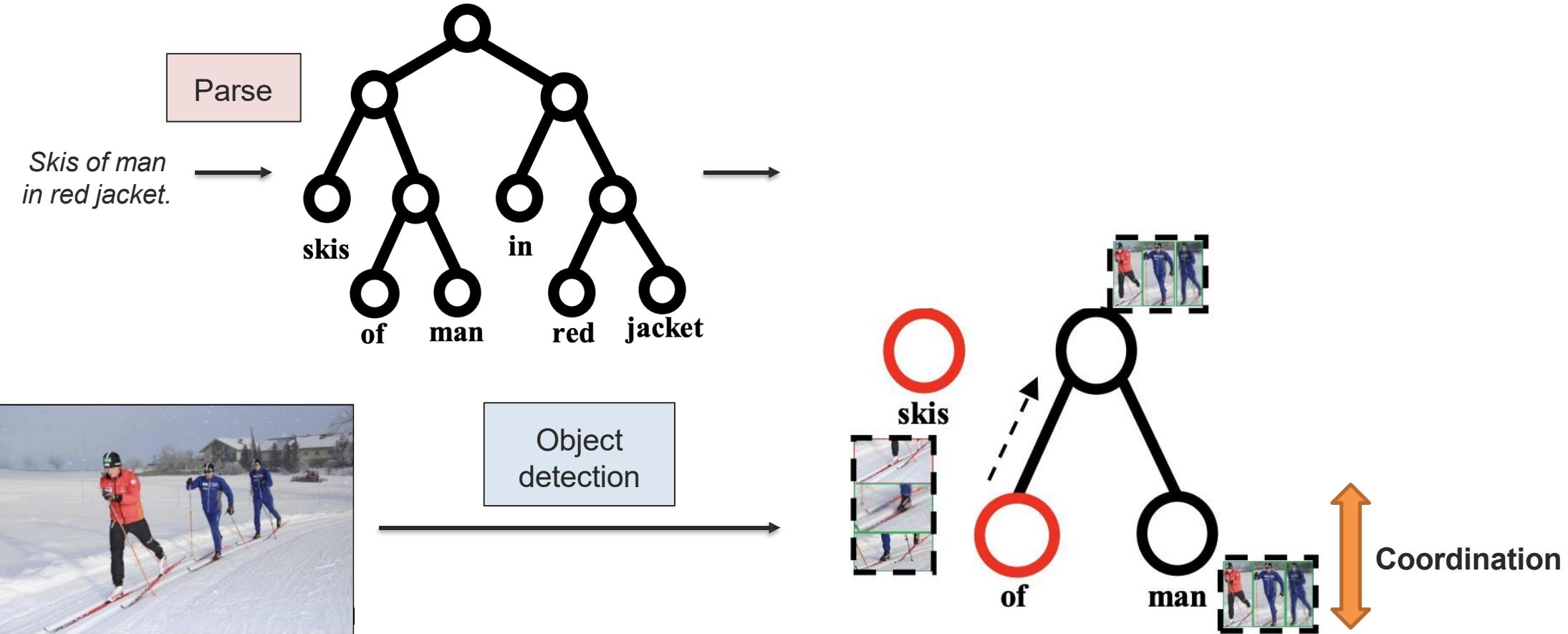
# Hierarchical Structure

Leverage syntactic structure of language



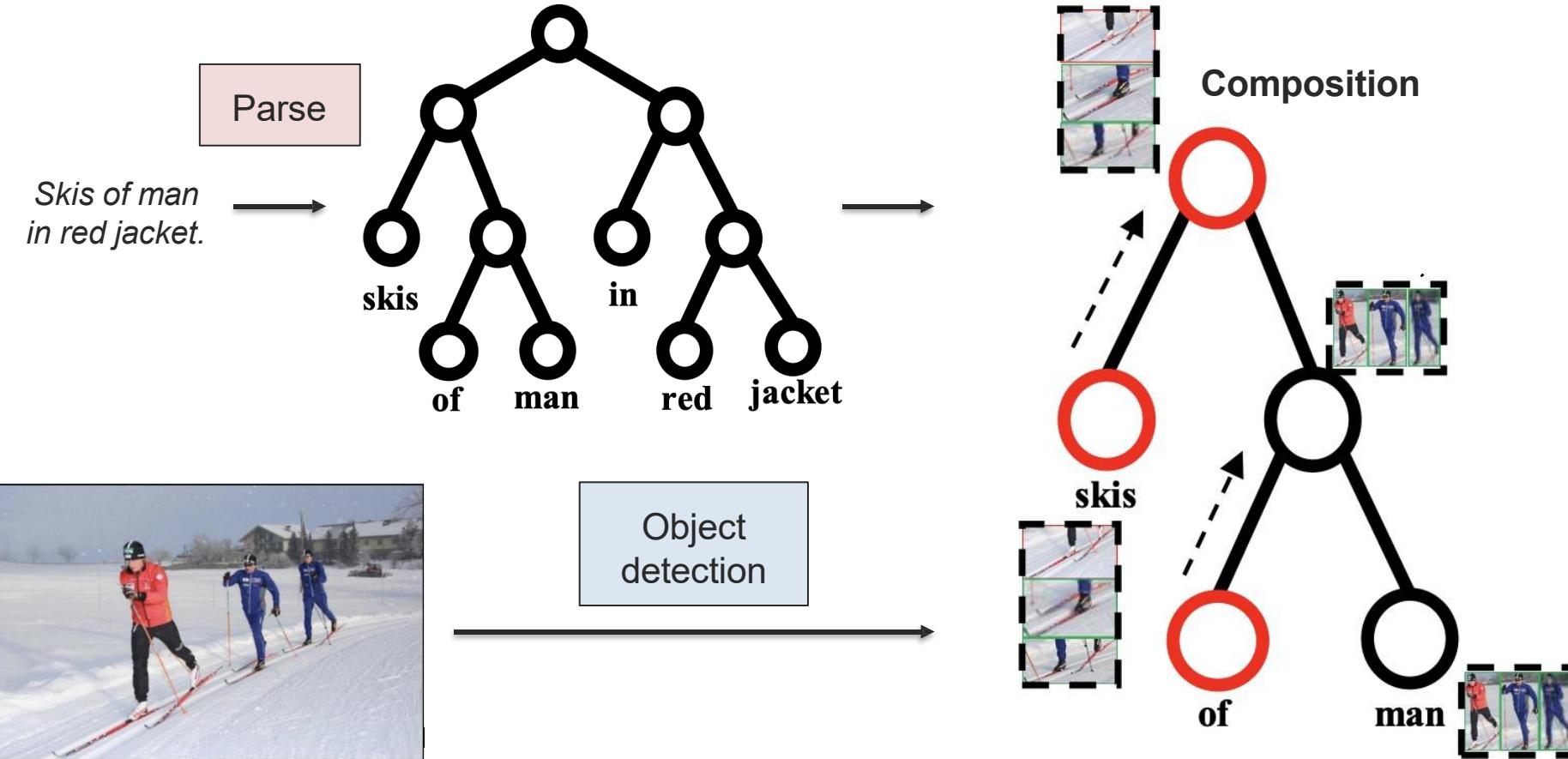
# Hierarchical Structure

Leverage syntactic structure of language



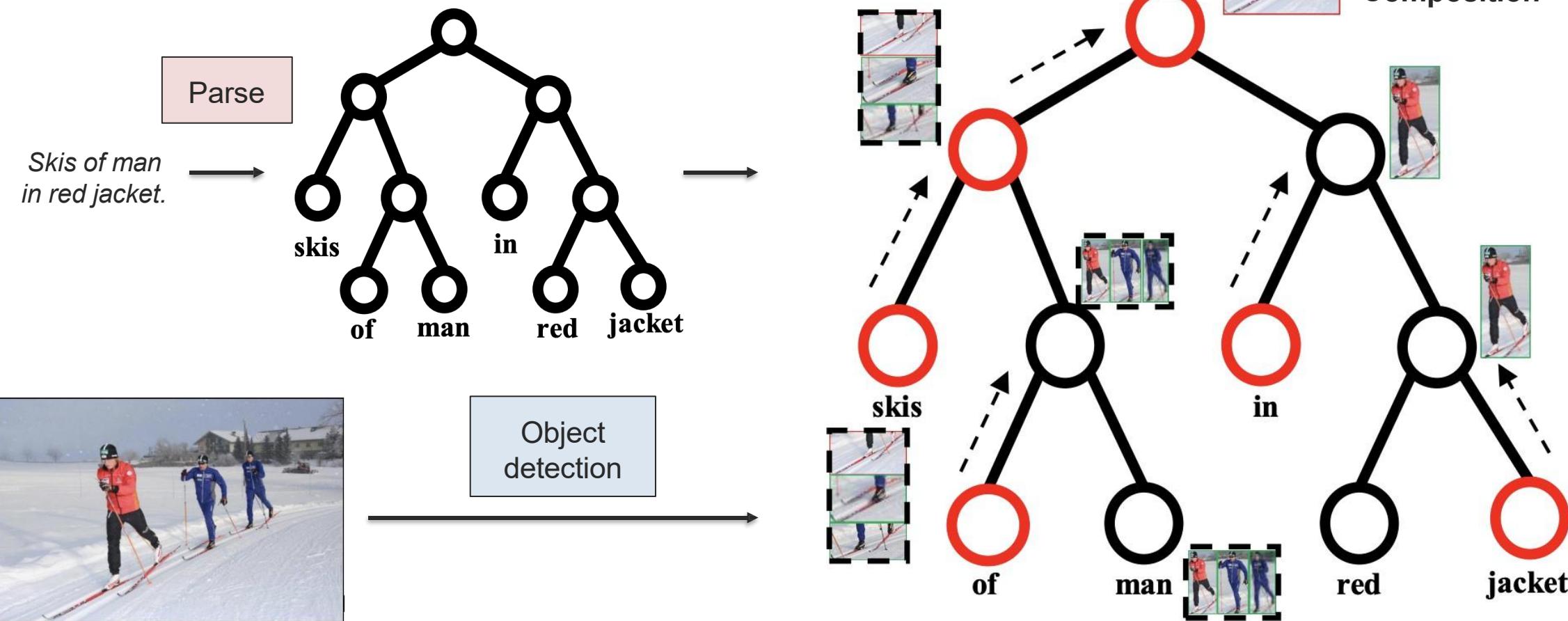
# Hierarchical Structure

Leverage syntactic structure of language



# Hierarchical Structure

Leverage syntactic structure of language



# Hierarchical Structure

## Leverage syntactic structure of language

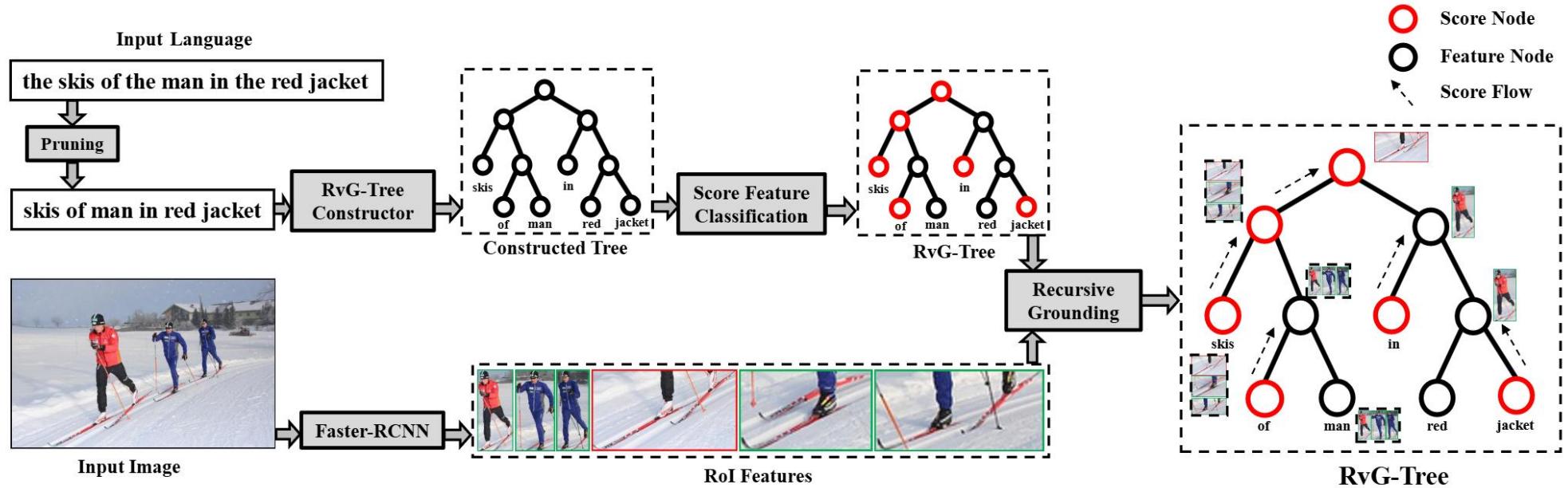


Fig. 2: The overview of using RvG-TREE for natural language grounding. Given a natural language sentence as input, we first prune the sentence and then construct RvG-TREE (Section 3.2). Then, we have a score-feature classifier (Section 3.3) to determine each node as the “score node” or “feature node”, where the score node returns the recursive score and the feature node returns the feature (Section 3.3). The final score of the root node is accumulated recursively in a bottom-up fashion (Section 3.3) and the visual region with the highest score is considered as the grounding result. Note that all the nodes can be visualized by the corresponding confidence scores and only qualitative regions are visualized.

# Module Network V2: End-to-End Learning

There is a shiny object that is right of the gray metallic cylinder;  
does it have the same size as the large rubber sphere?

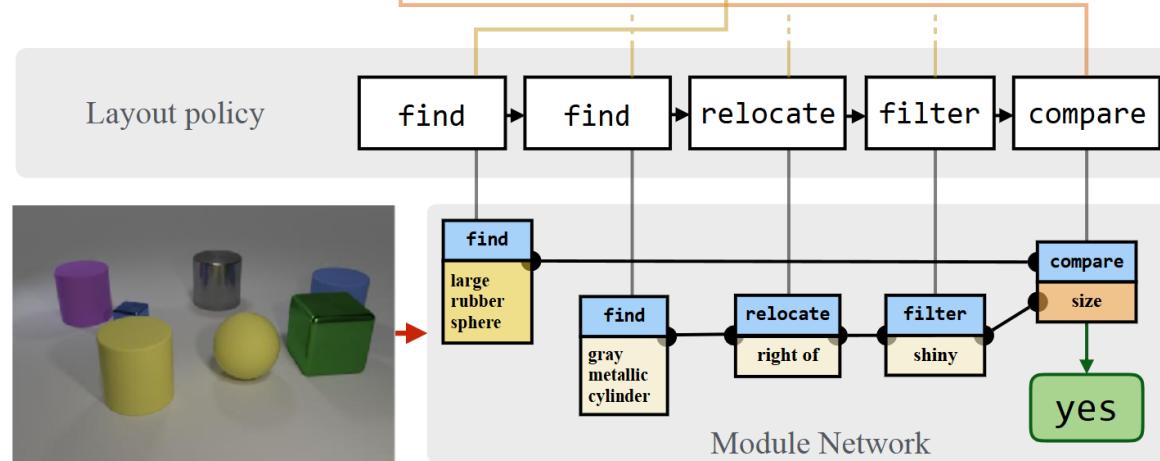
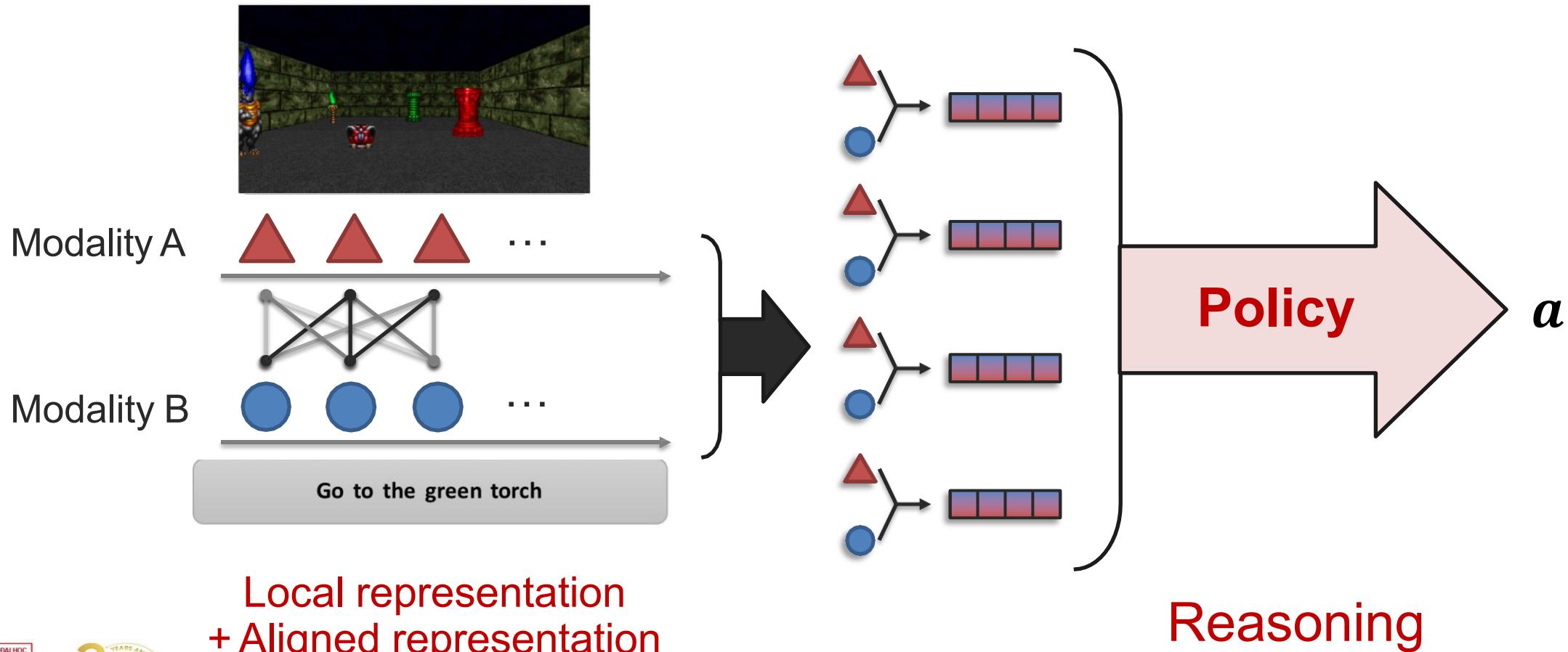


Figure 1: For each instance, our model predicts a computational expression and a sequence of attentive module parameterizations. It uses these to assemble a concrete network architecture, and then executes the assembled neural module network to output an answer for visual question answering. (The example shows a real structure predicted by our model, with text attention maps simplified for clarity.)

# Interactive Structure

**Structure defined through interactive environment**

Main difference from temporal - actions taken at previous time steps affect future states



# Challenge 3: Reasoning

## Final Summary

### Sub-Challenge

### Structure Modeling

### Intermediate Concepts

### Inference Paradigm

### External Knowledge

### Definition

Learning structured relationships

Extracting feature representations

Logical or probabilistic reasoning over modalities

Using external databases or knowledge graphs

### Concrete Example

Scene graphs for VQA,  
Graph-based summarization

Sentiment analysis, Lip-reading phonemes

VQA (Vision-Text fusion),  
Commonsense AI

Wikipedia-augmented  
Chatbots, COMET-based  
causal reasoning



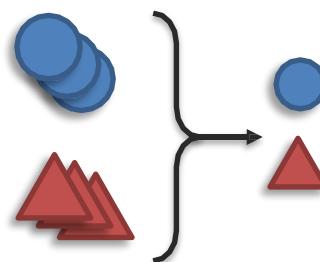
# Challenge 4: Generation

**Definition:** Learning a generative process to produce new data while maintaining cross-modal interactions, structure, and coherence.

Covers tasks such as **summarization**, **translation**, and **creative generation** across multiple modalities.

## Sub-challenges:

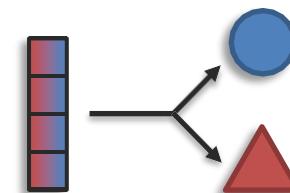
### Summarization



### Translation



### Creation



**Information:**  
(content)

Reduction

$$\square > \square$$

Maintenance

$$\square = \square$$

Expansion

$$\square < \square$$

# Sub-challenge 4a: Summarization

**Definition:** Summarizing multimodal data to reduce information content while highlighting the most salient parts of the input.

Transcript

today we are going to show you how to make spanish omelet . i 'm going to dice a little bit of peppers here . i 'm not going to use a lot , i 'm going to use very very little . a little bit more then this maybe . you can use red peppers if you like to get a little bit color in your omelet . some people do and some people do n't .... t is the way they make there spanish omelets that is what she says . i loved it , it actually tasted really good . you are going to take the onion also and dice it really small . you do n't want big chunks of onion in there cause it is just pops out of the omelet . so we are going to dice the up also very very small . so we have small pieces of onions and peppers ready to go .

Video



How2 video dataset

Complementary  
cross-modal  
interactions

Summary

Cuban breakfast  
Free cooking video

(not present in text)

how to cut peppers to make a spanish omelette; get expert tips and advice on making cuban breakfast recipes in this free cooking video .

# Sub-challenge 4a: Summarization

**Definition:** Summarizing multimodal data to reduce information content while highlighting the most salient parts of the input.

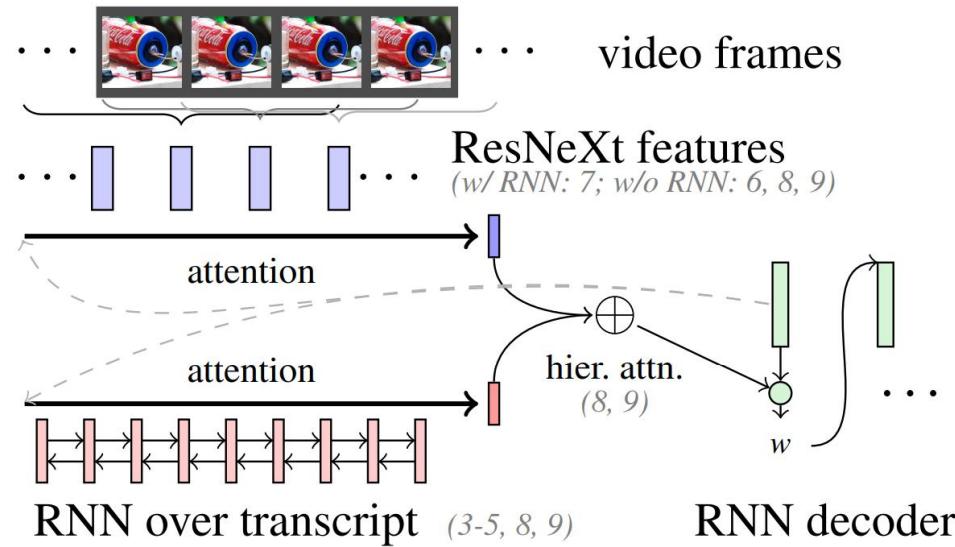


Figure 2: Building blocks of the sequence-to-sequence models, gray numbers in brackets indicate which components are utilized in which experiments.

# Sub-challenge 4a: Summarization

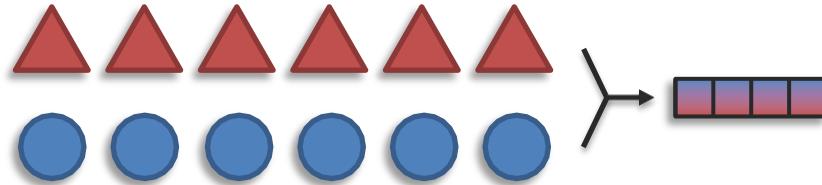
## Video summarization

A

Content

Fusion via  
joint representation

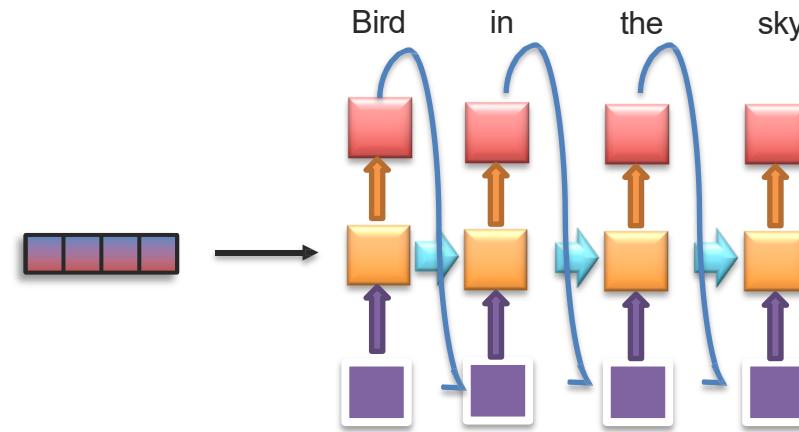
Capture complementary  
cross-modal interactions



B

Generation

Generative  $\approx$  abstractive summarization  
Exemplar  $\approx$  extractive summarization



# Sub-challenge 4b: Translation

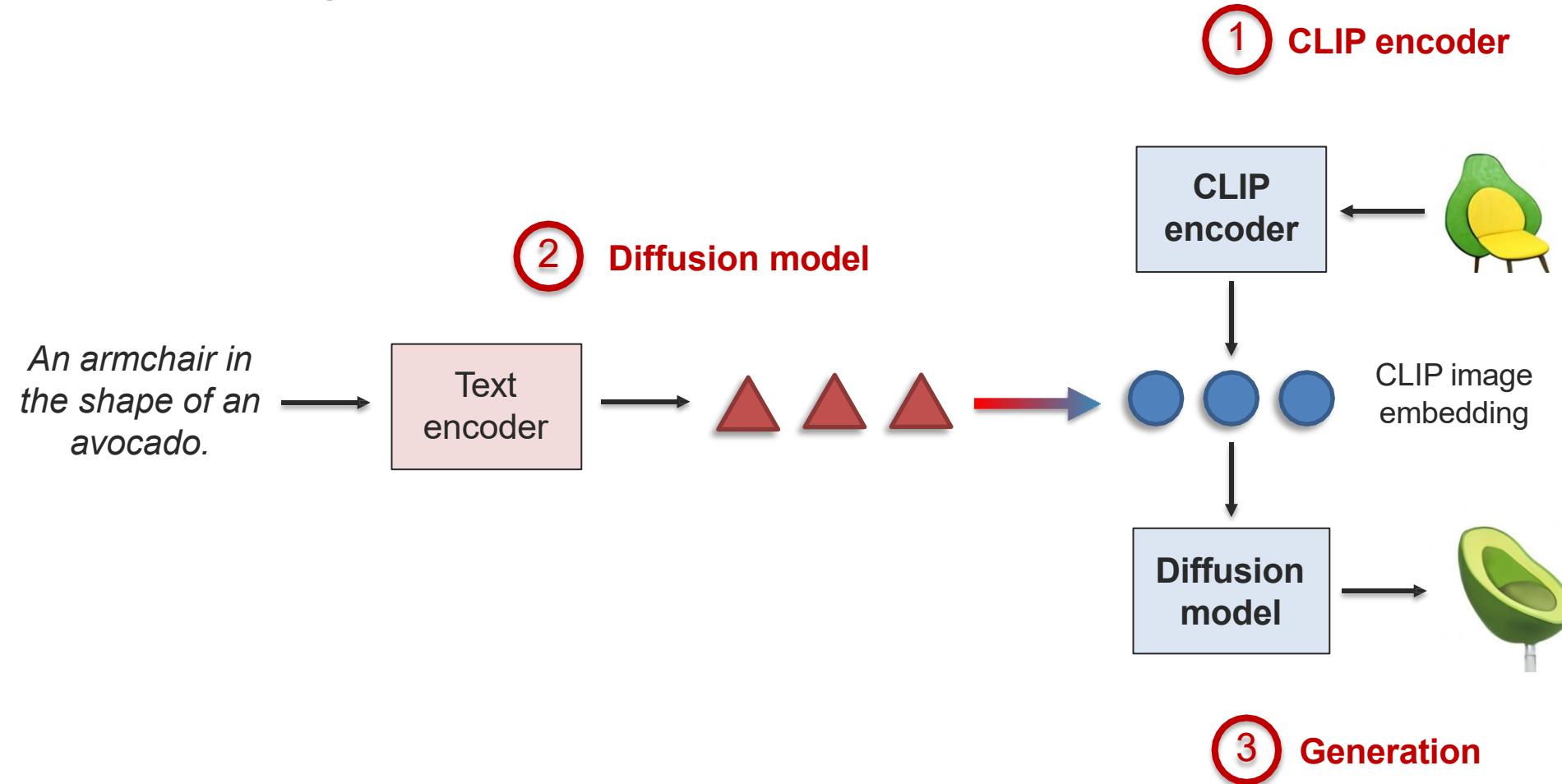
**Definition:** Translating from one modality to another and keeping information content while being consistent with cross-modal interactions.

*An armchair in the shape of an avocado*



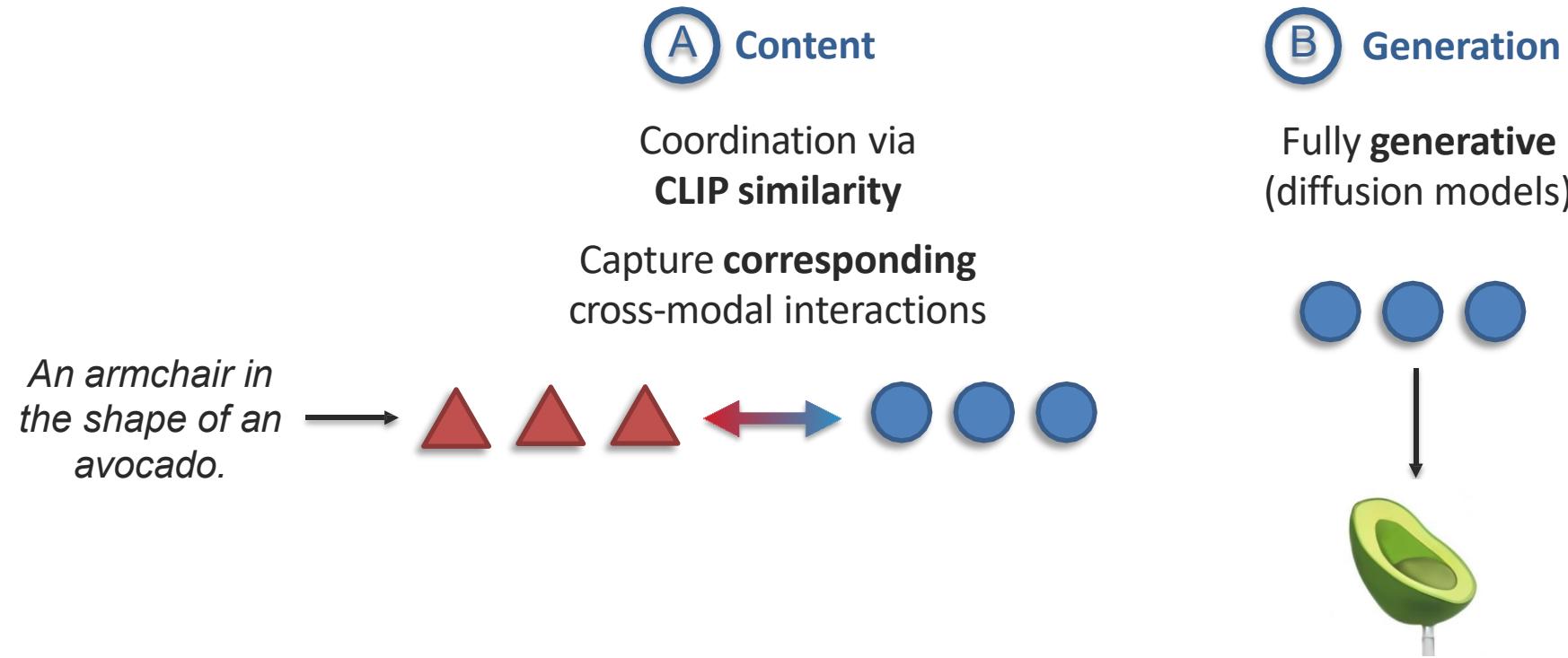
# Sub-challenge 4b: Translation

## DALL-E 2: Combining with CLIP, diffusion models



# Sub-challenge 4b: Translation

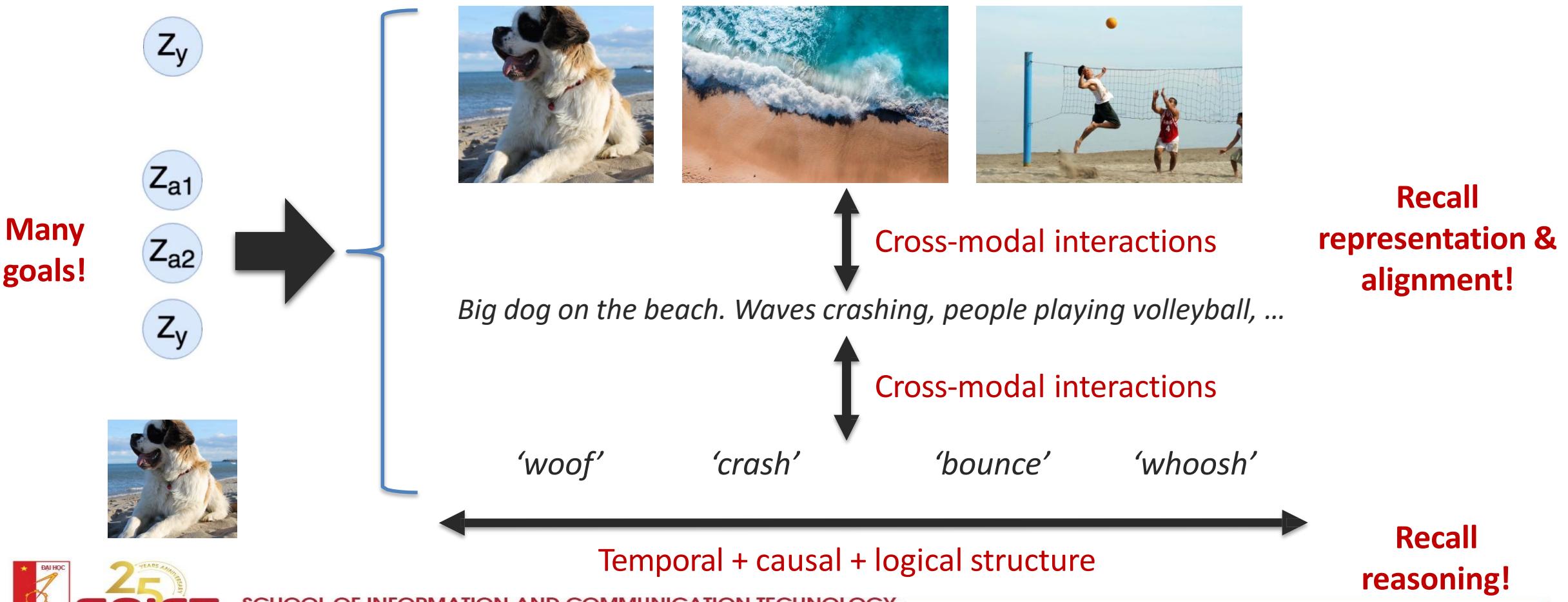
## DALL·E 2: Combining with CLIP, diffusion models



## Sub-challenge 4c: Creation

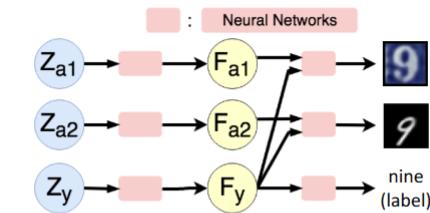
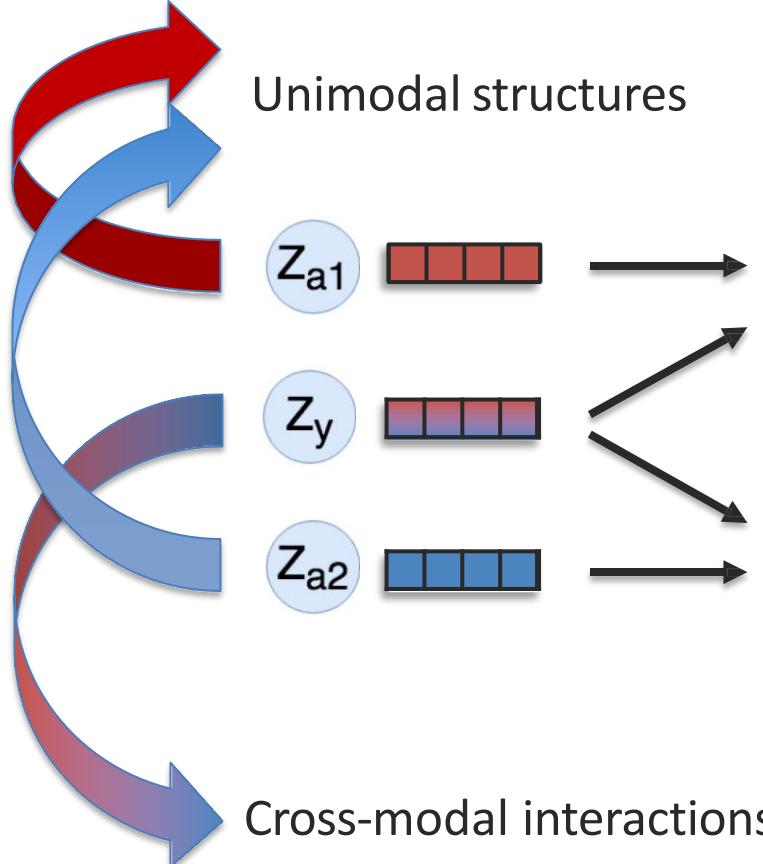
Open challenges

**Definition:** Simultaneously generating multiple modalities to increase information content while maintaining coherence within and across modalities.



# Sub-challenge 4c: Creation

Some initial attempts: factorized generation



(a)

Method	UM(SVHN)	UM(MNIST)	MM	MFM
Acc.	91.84	99.01	99.20	<b>99.36</b>

(b)



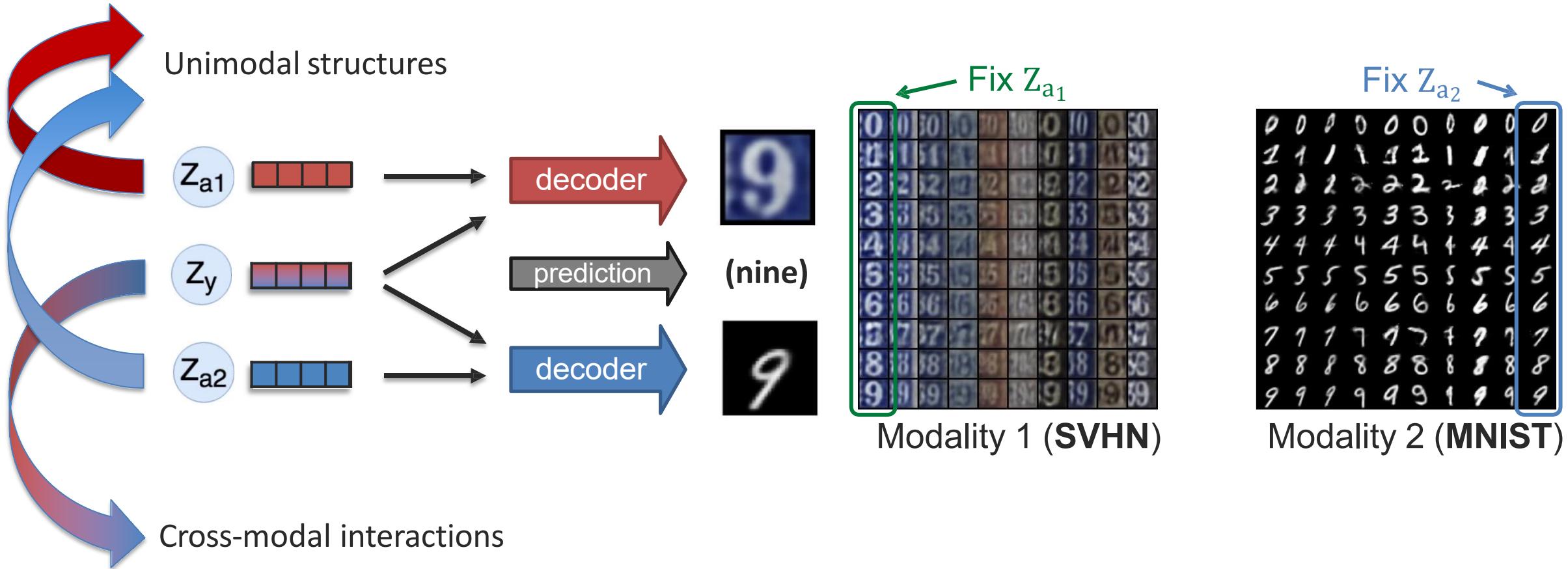
Modality 1 (SVHN)



Modality 2 (MNIST)

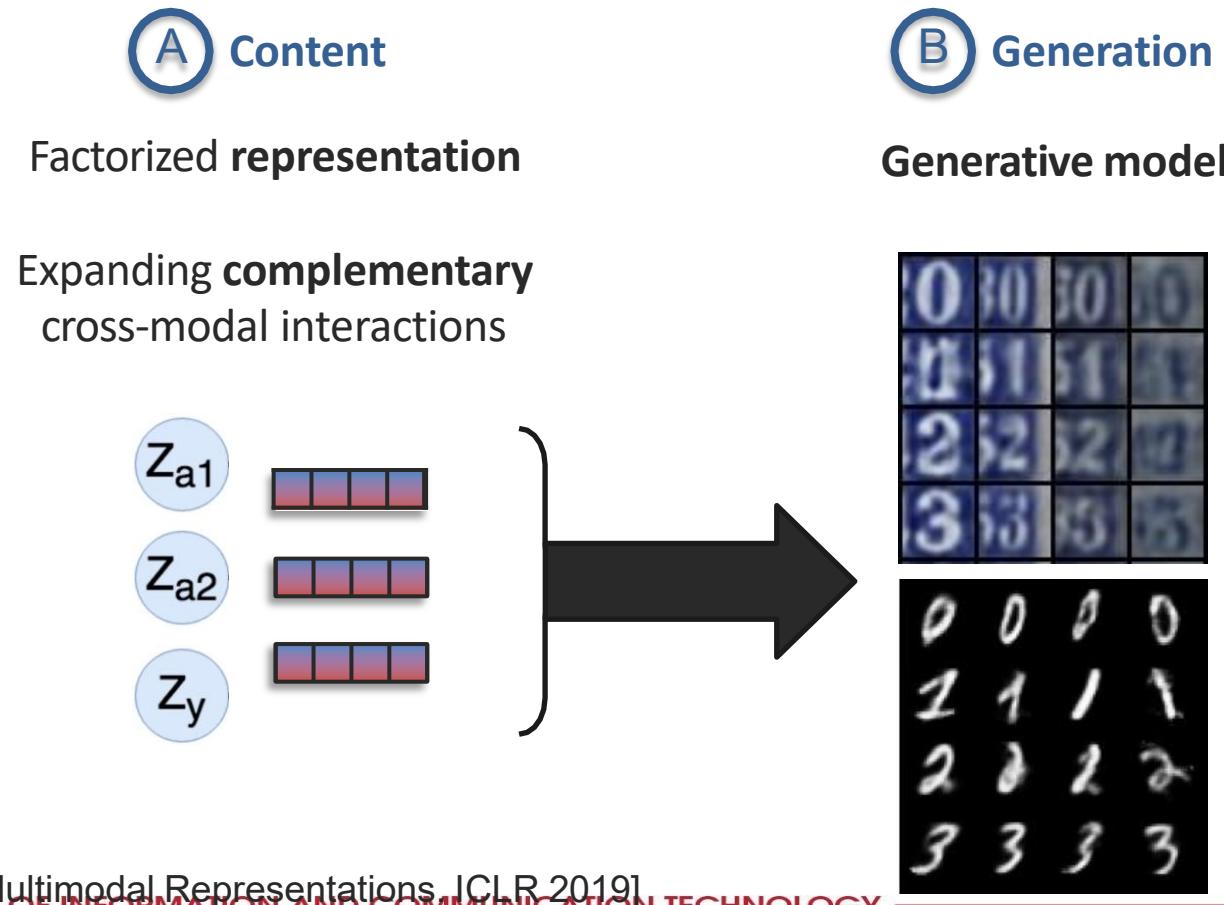
# Sub-challenge 4c: Creation

Some initial attempts: factorized generation



# Sub-challenge 4c: Creation

Some initial attempts: factorized generation



# Challenge 4: Generation

## Final Summary

### Sub-Challenge

#### Summarization (Reduction)

#### Translation (Preserving Meaning)

#### Creation (Expansion)

### Definition

Extract key information while removing redundancy

Convert data between modalities while retaining meaning

Generate new content beyond the given data

### Concrete Example

Video summarization, News article condensation

Speech-to-text, Image captioning

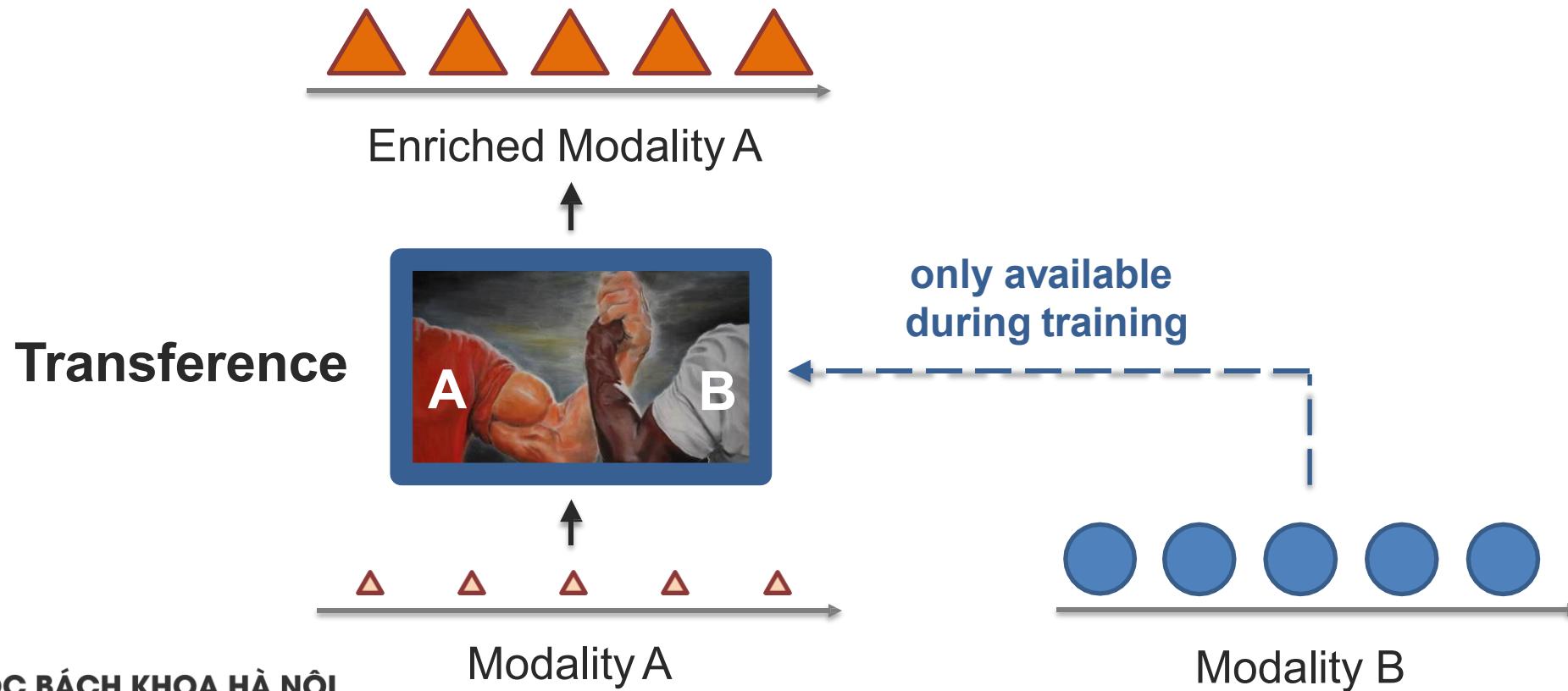
Text-to-image (DALL-E), Story generation



# Challenge 5: Transference

**Definition:** Transferring knowledge between different modalities to help a **weaker modality** that might be **noisy, incomplete, or resource-limited**.

Used in **low-resource learning scenarios, data-efficient AI, and multimodal adaptation**.



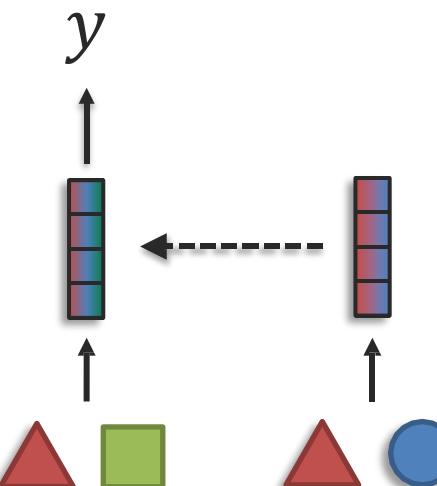
# Challenge 5: Transference

**Definition:** Transferring knowledge between different modalities to help a **weaker modality** that might be **noisy, incomplete, or resource-limited**.

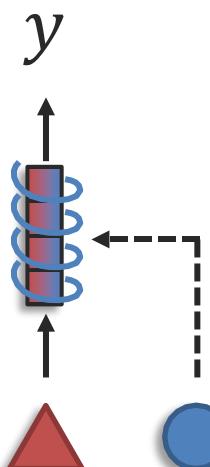
Used in **low-resource learning scenarios, data-efficient AI, and multimodal adaptation**.

## Sub-challenges:

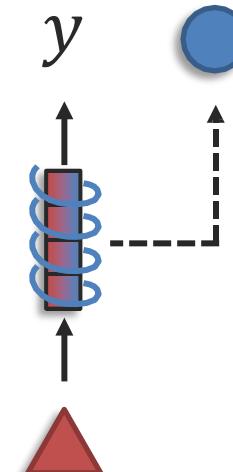
### Transfer



### Co-learning via representation

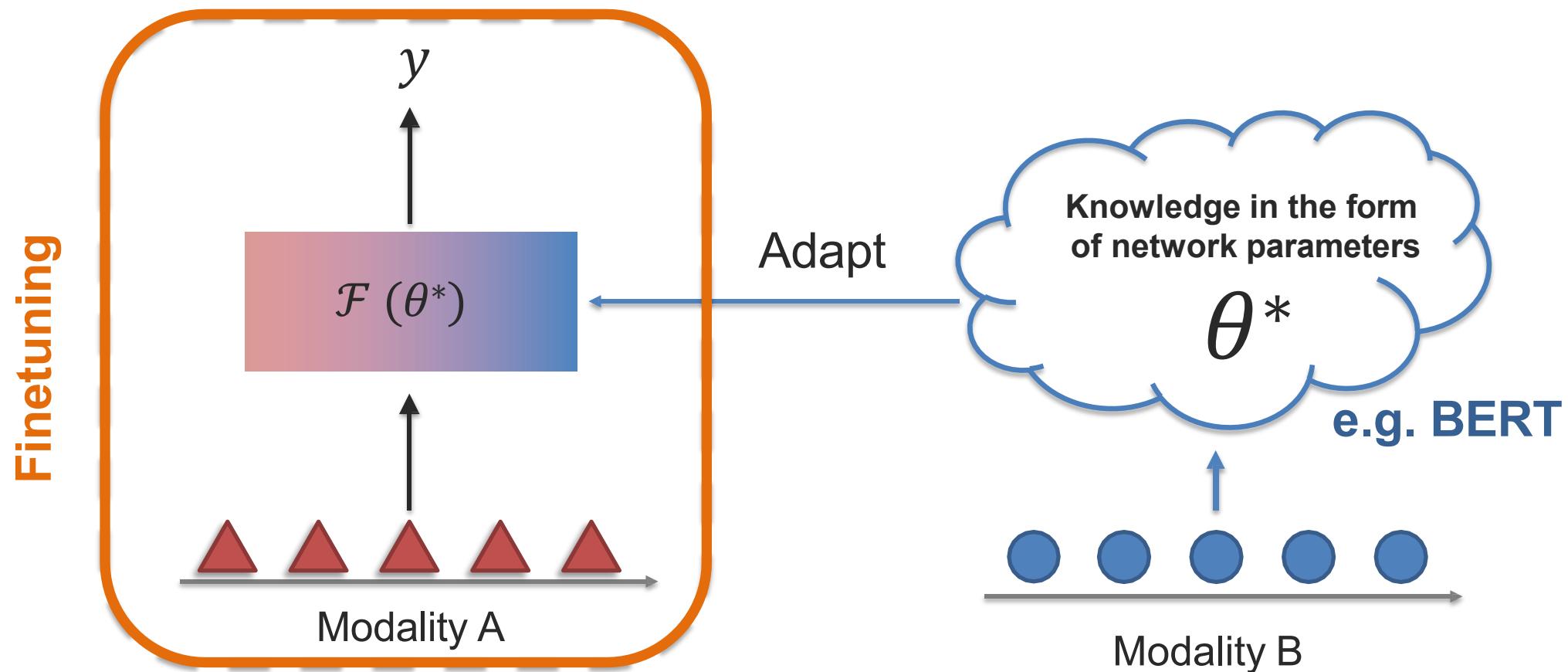


### Co-learning via generation



## Sub-Challenge 5a: Transfer via Pretrained Models

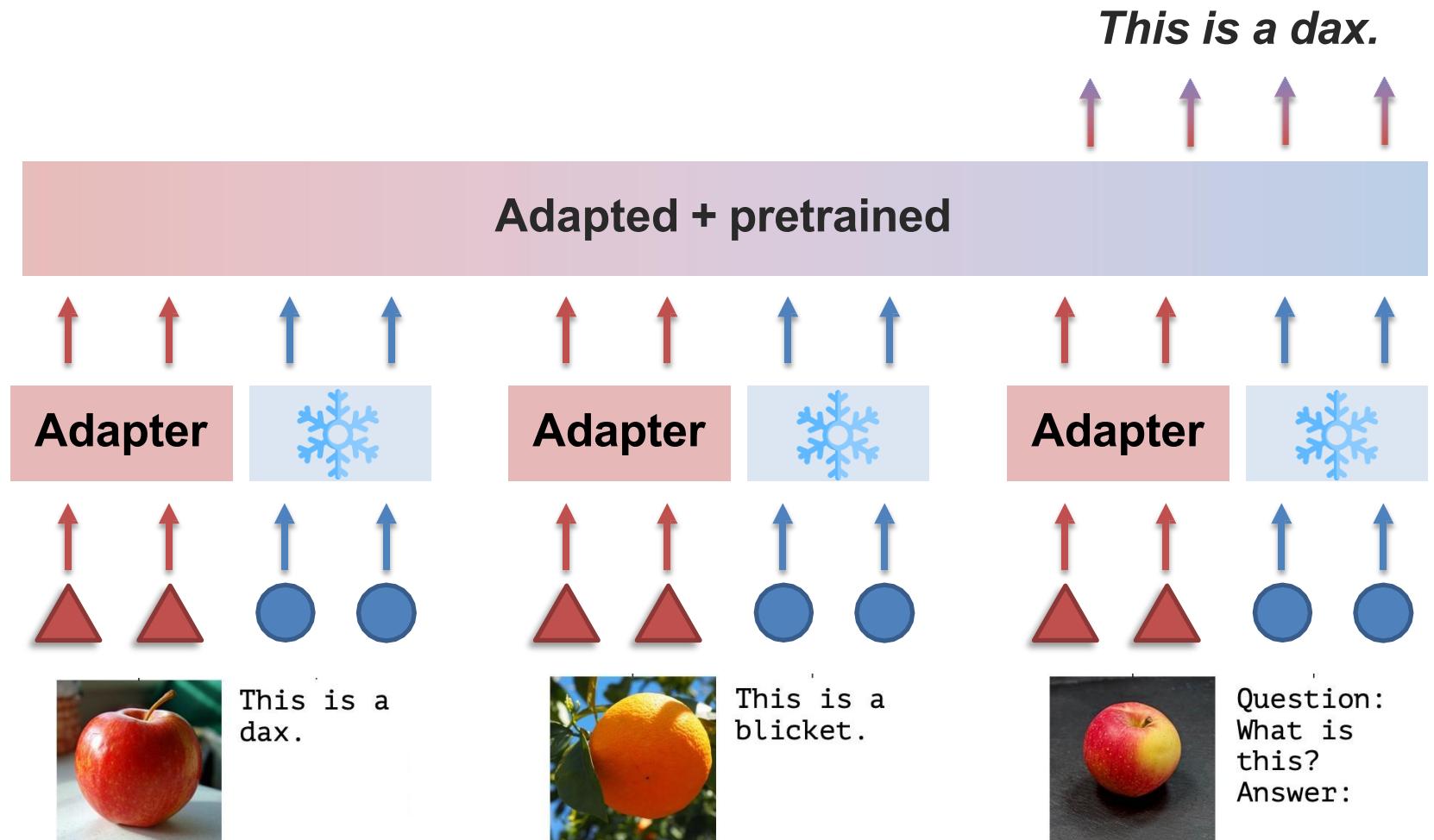
**Definition:** Transferring knowledge from large-scale pretrained models to downstream tasks involving the primary modality.



# Sub-Challenge 5a: Transfer via Pretrained Models

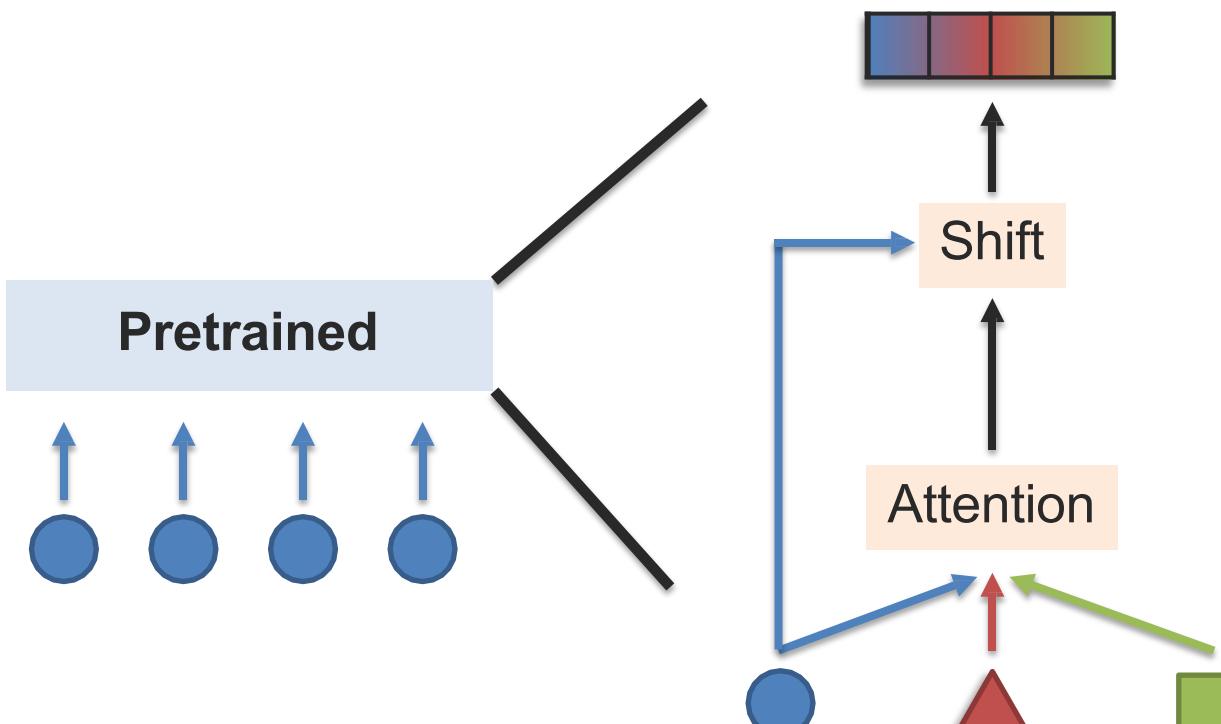
Transfer via prefix tuning

Few-shot image classification:

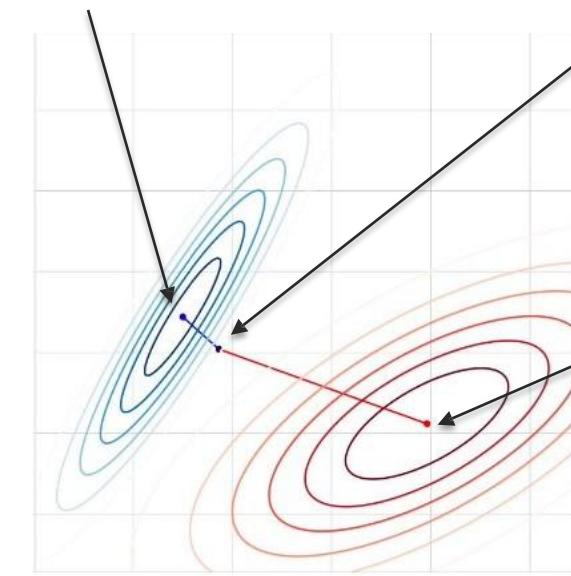


# Sub-Challenge 5a: Transfer via Pretrained Models

Transfer via representation tuning



With positive AV



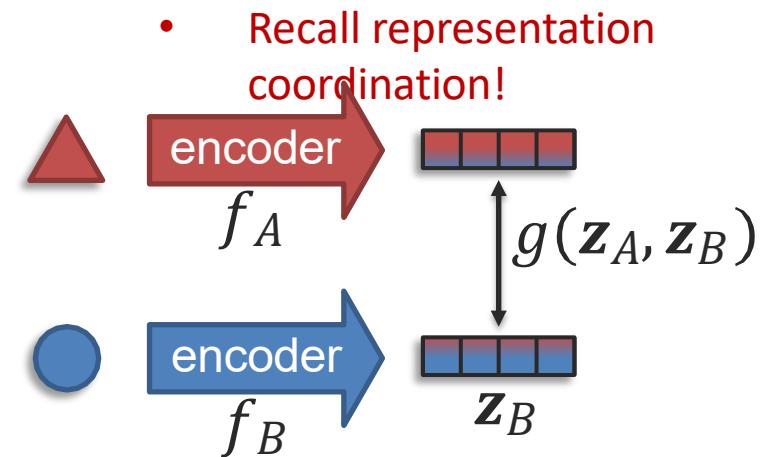
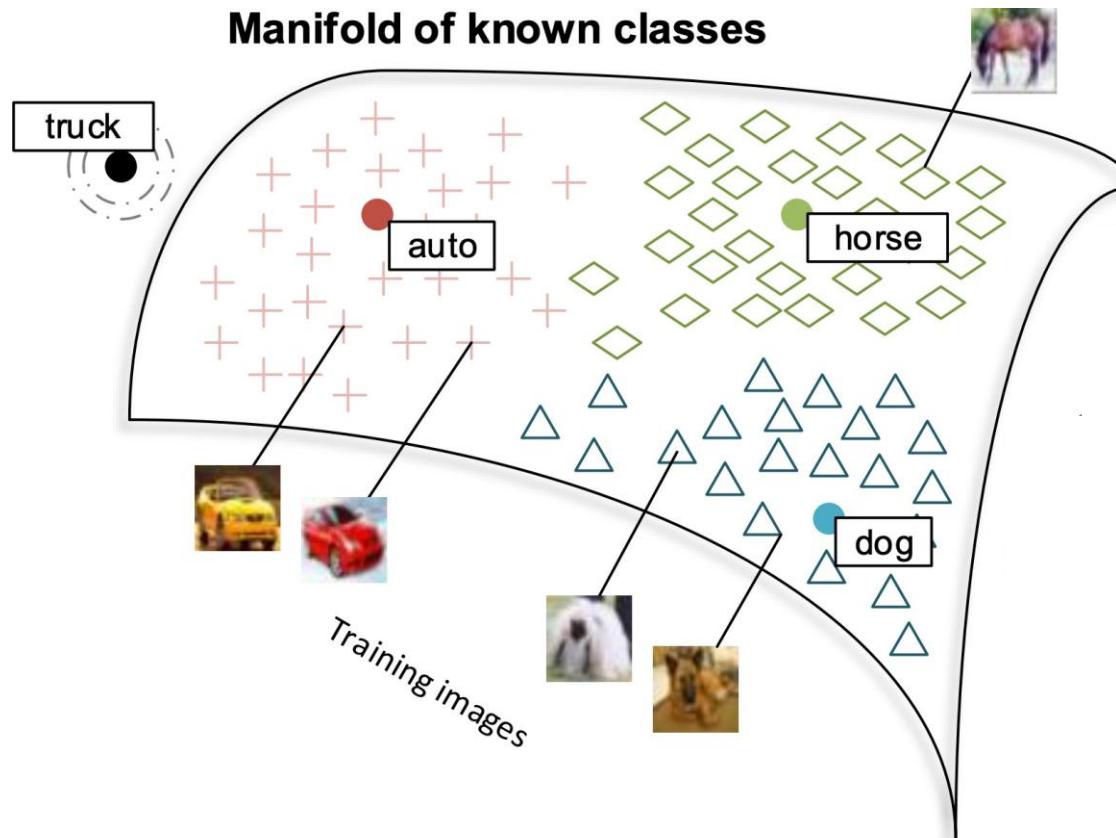
Without AV

With negative AV

Lexical Space

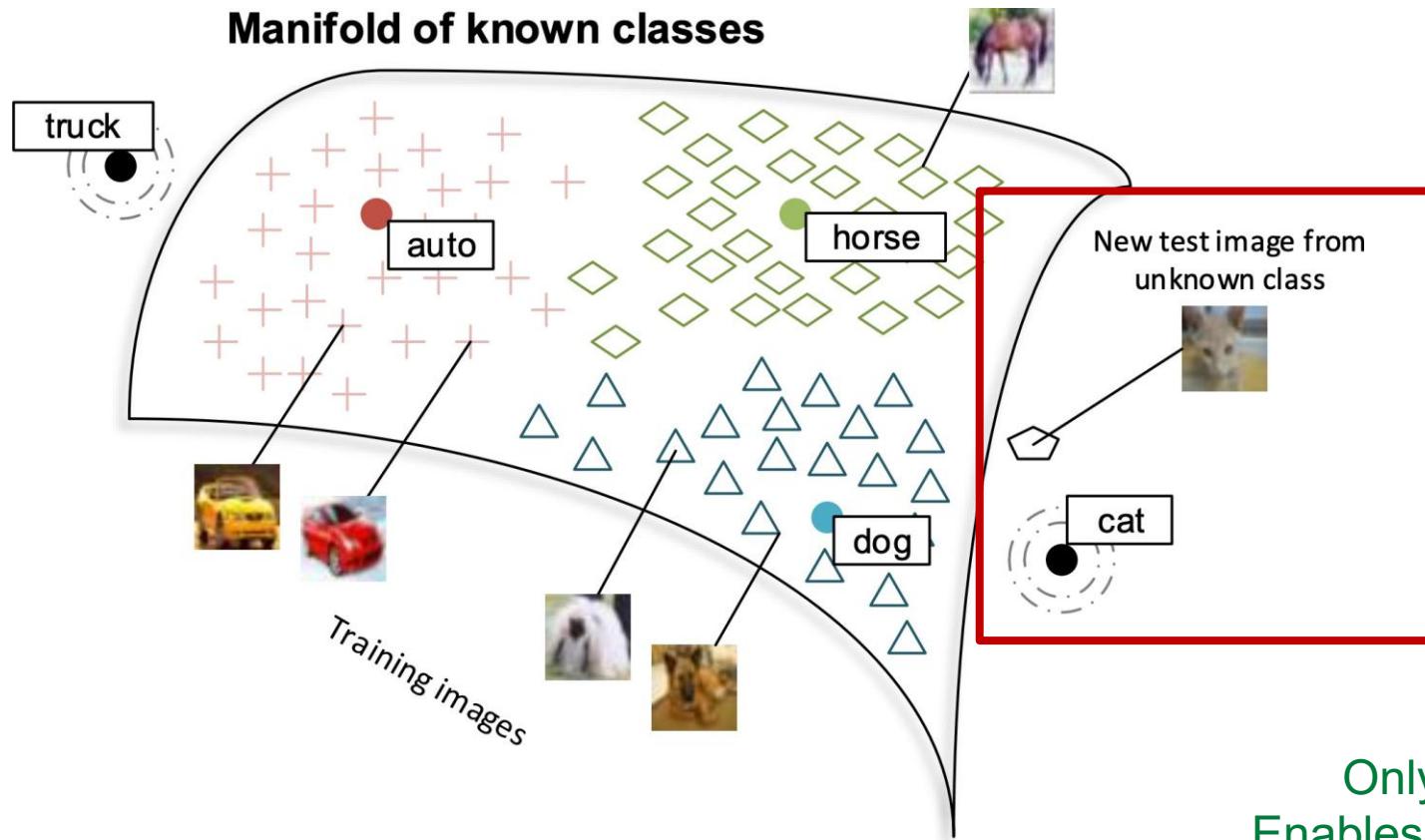
# Co-learning via Representation

- Representation coordination: word embedding space for zero-shot visual classification

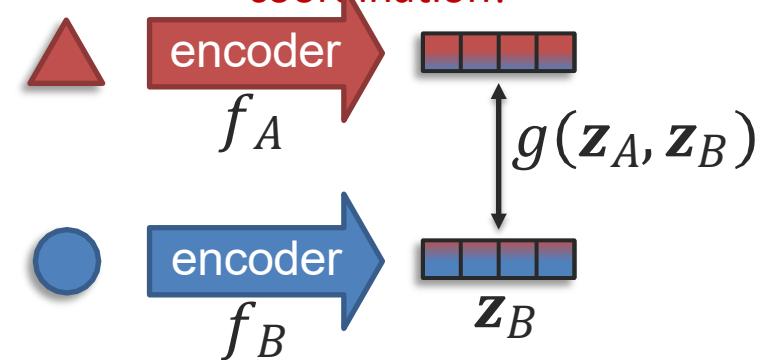


# Co-learning via Representation

- Representation coordination: word embedding space for zero-shot visual classification



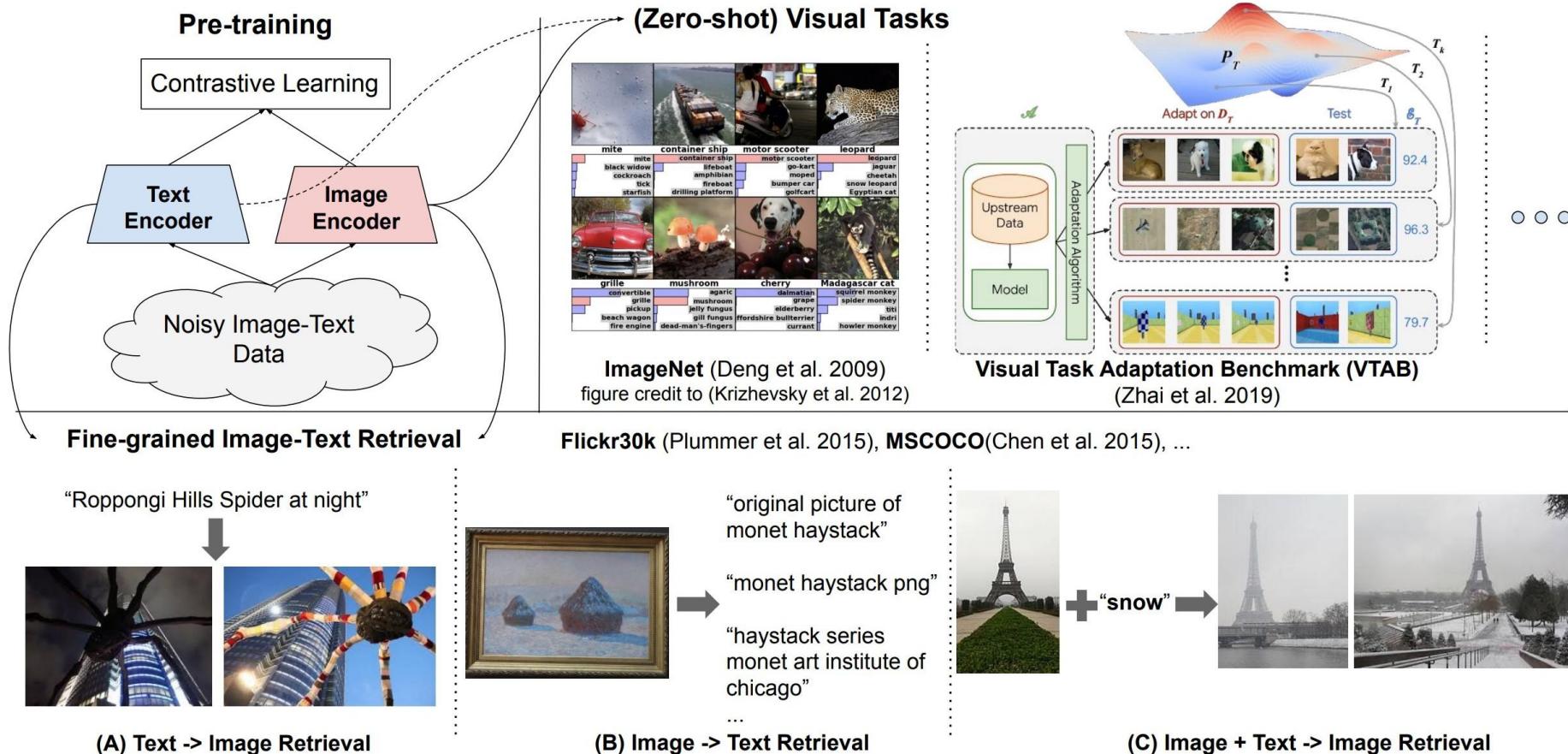
- Recall representation coordination!



Only images used at test-time  
Enables zero-shot image classification

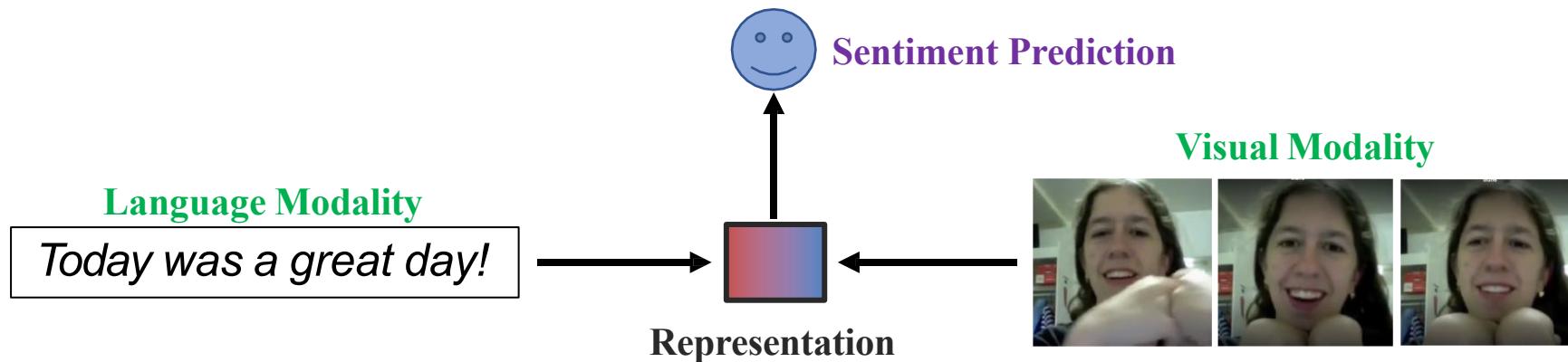
# Co-learning via Representation

## Representation coordination at scale



# Co-learning via Generation

## Bimodal translations



Both modalities required at test time!  
Sensitive to noisy/missing visual modality.

We want to leverage information from visual modality  
while being robust to it during test-time.

# Challenge 5: Transference

## Final Summary

Sub-Challenge	Definition	Concrete Example
Transfer	One modality transfers knowledge to another	Text-trained models help low-resource speech recognition
Co-Learning via Representation	Learning shared representations for robustness	Lip-reading assists noisy speech recognition
Co-Learning via Generation	One modality generates missing data for another	Text-to-speech, synthetic image-text augmentation

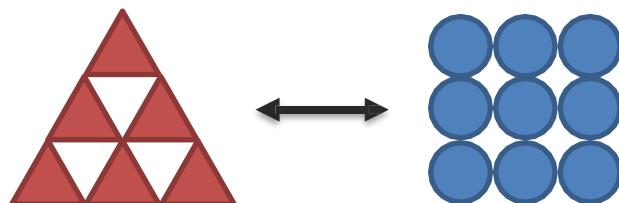


# Challenge 6: Quantification

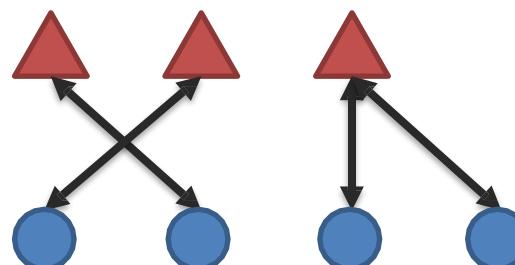
**Definition:** Empirical and theoretical study to improve understanding of **heterogeneity**, **cross-modal interactions**, and **learning dynamics** in multimodal AI.  
Involves **metrics**, **evaluation techniques**, and **optimization strategies** to enhance multimodal learning.

## Sub-challenges:

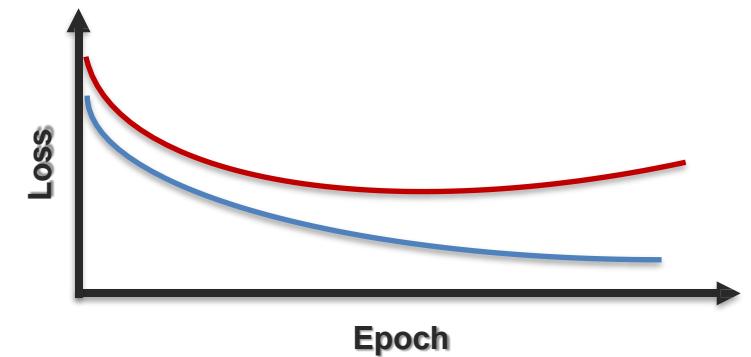
### Heterogeneity



### Interactions

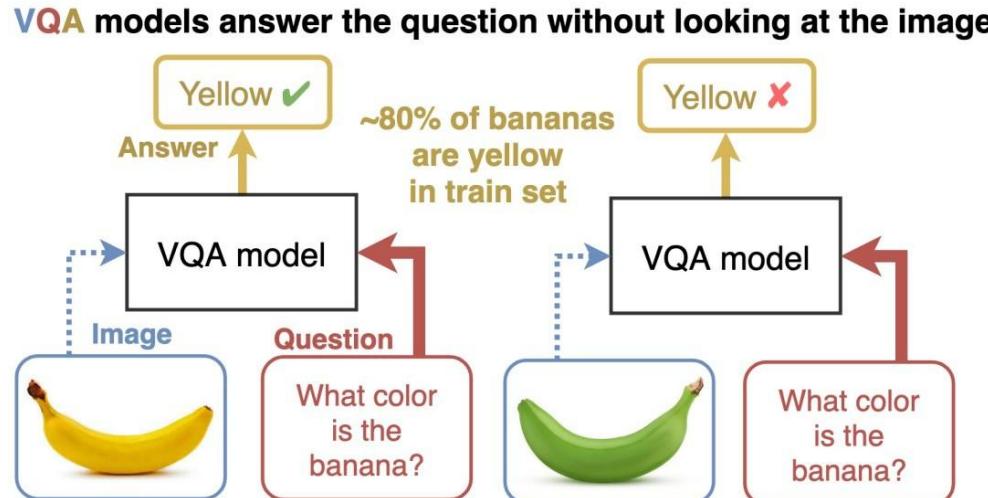


### Learning



# Modality Biases

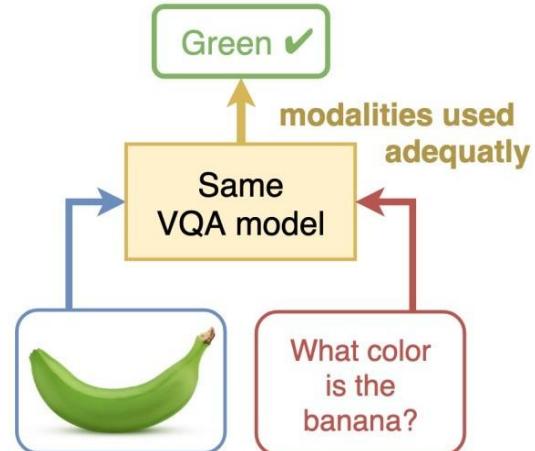
**Heterogeneity in information and relevance**  
Unimodal biases and modality collapse



Balancing modalities  
Balancing training



Not the case when trained with RUBi



[Wu et al., Characterizing and Overcoming the Greedy Nature of Learning in Multi-modal Deep Neural Networks. ICML 2022]  
[Javaloy et al., Mitigating Modality Collapse in Multimodal VAEs via Impartial Optimization. ICML 2022]

# Modality Biases

## Heterogeneity in information and relevance

Fairness and social biases – unimodal social biases

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

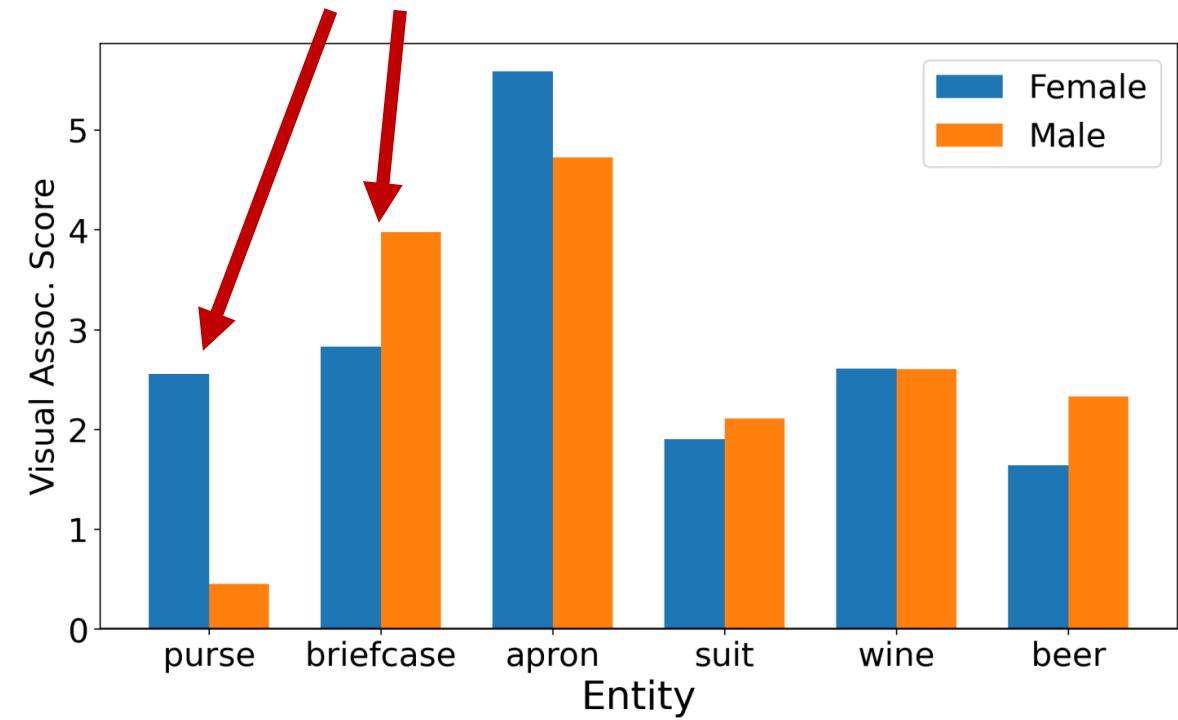
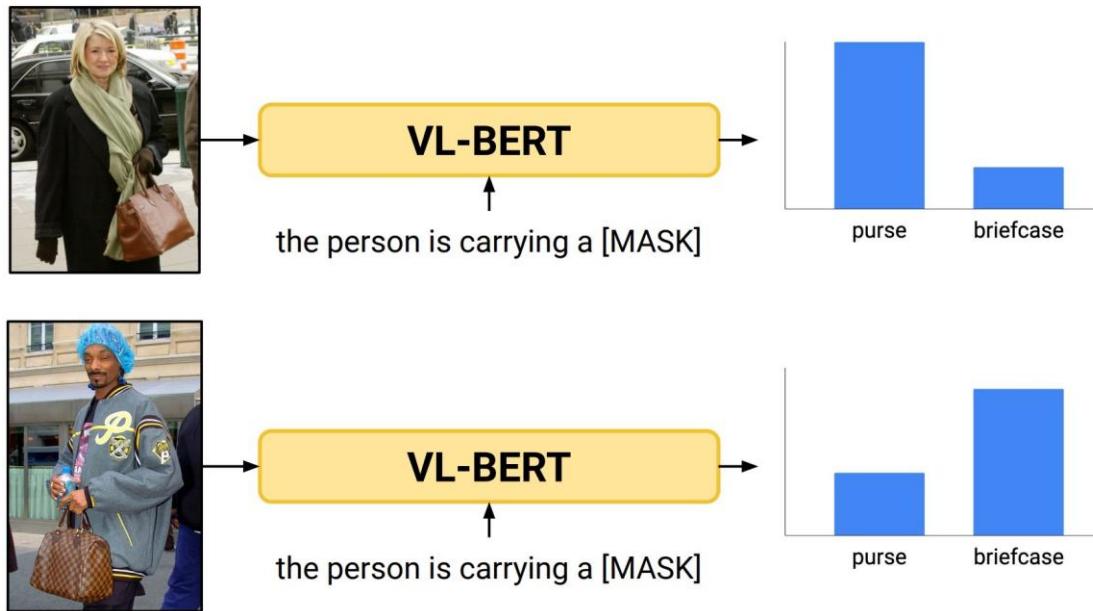


# Modality Biases

## Heterogeneity in information and relevance

Fairness and social biases – cross-modal interactions worsen social biases

Visual information makes model more confident  
in reinforcing gender stereotypes



# Challenge 6: Quantification

## Final Summary

### Sub-Challenge

### Heterogeneity

### Interactions

### Learning

### Definition

Quantifying differences across modalities

Evaluating cross-modal dependencies

Studying optimization and generalization

### Concrete Example

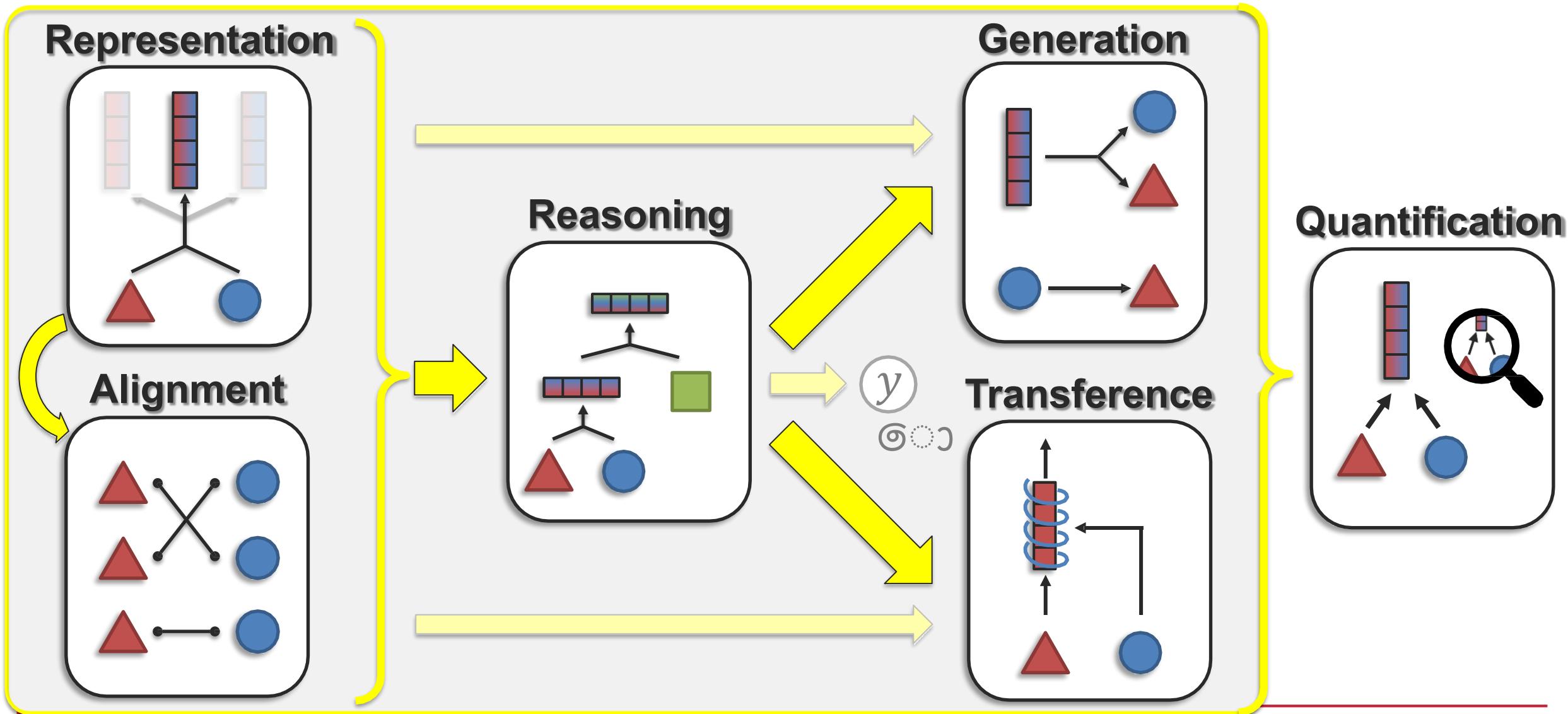
Measuring **modality importance** in model performance

Attention weights in vision-language transformers

Tracking **loss convergence for each modality**



# Core Multimodal Challenges



# Multimodal Machine Learning Models

## Models for Multimodal Learning

- 1. Concatenation-based Models** (e.g., Early Fusion)
- 2. Attention-based Models** (e.g., Cross-modal transformers like CLIP, ViLBERT)
- 3. Graph Neural Networks (GNNs)** (e.g., Graph-based multimodal learning)
- 4. Contrastive Learning** (e.g., Aligning image-text embeddings)



# Future of Multimodal Learning

## Research Challenges

### 1. Multimodal Explainability

1. Understanding how models integrate different modalities.

### 2. Efficient & Scalable Fusion

1. Reducing computation for real-time multimodal applications.

### 3. Few-Shot & Zero-Shot Learning

1. Generalizing to unseen modalities without retraining.

## Emerging Trends

- GPT-4V (multimodal ChatGPT).
- Self-Supervised Multimodal Learning (reducing reliance on labeled data).
- Embodied AI (robots using multimodal perception for interaction).



# References

- Multi Modal Machine Learning, 11-777 • Fall 2023 • Carnegie Mellon University





**HUST**

Thanks