



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Multimodal Machine Learning

## Lecture 4.1: Multimodal alignment

Louis-Philippe Morency

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

# Administrative Stuff

## Primary TAs

---

- Each team will have one primary TA
- Contact your primary TA anytime
  - Groups were created in Piazza for each team
- Some projects may have a secondary TA, with complementary expertise

**Schedule a meeting with your Primary TA this week!**

# First Project Assignment

---

Due date: Sunday 9/24 at 8pm

Four main sections:

- Introduction
- Related work
- Experimental setup
- Research ideas



The two main sections are related work and research ideas



# teammates = # research ideas



Page limit depends on team size:

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

Follows ICML paper format

## Team Meetings with Instructor

---

- No lecture on Tuesday 10/3
- 15-mins meeting with instructor
  - Optional, but highly suggested
  - Not all teammates are required to attend
- Meetings next week: Wednesday 9/27 until Friday 9/29
- Signup form will be shared via Piazza



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Multimodal Machine Learning

## Lecture 4.1: Multimodal alignment

Louis-Philippe Morency

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

## Lecture objectives

---

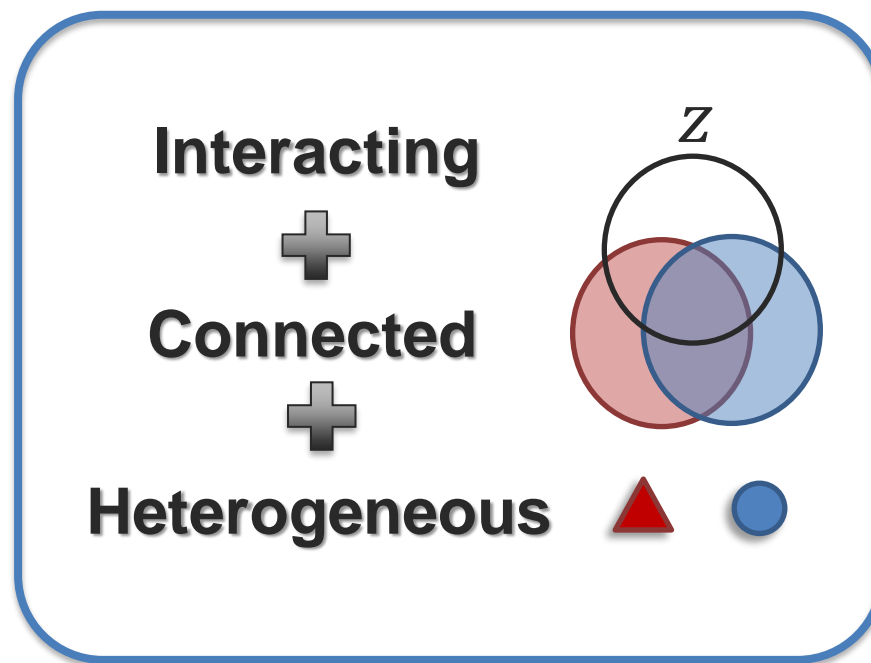
- A quick review
  - Connections, coordinated representations and mutual information
  - Modality interactions and factorized representations
- Discrete alignment
  - Local alignment
    - Coordinated representations; hard and soft attention
  - Global alignment
    - Assignment problem and optimal transport
- Continuous alignment
  - Continuous warping
    - Dynamic time warping
  - Discretization and segmentation

# A Quick Review

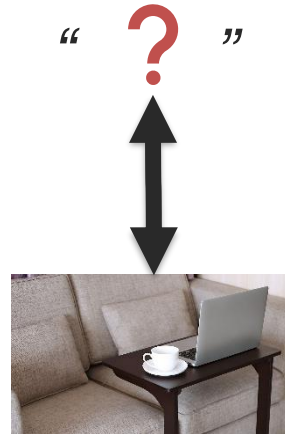
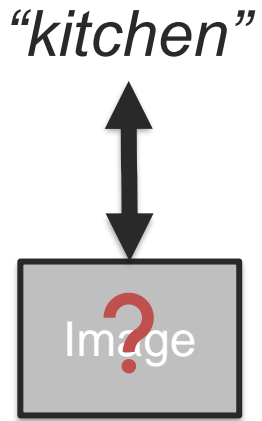
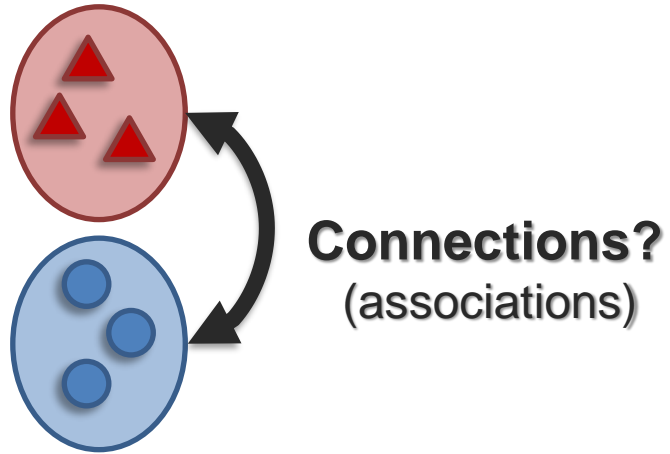
---



**Multimodal is the scientific study of**  
**heterogeneous, connected and interacting data**



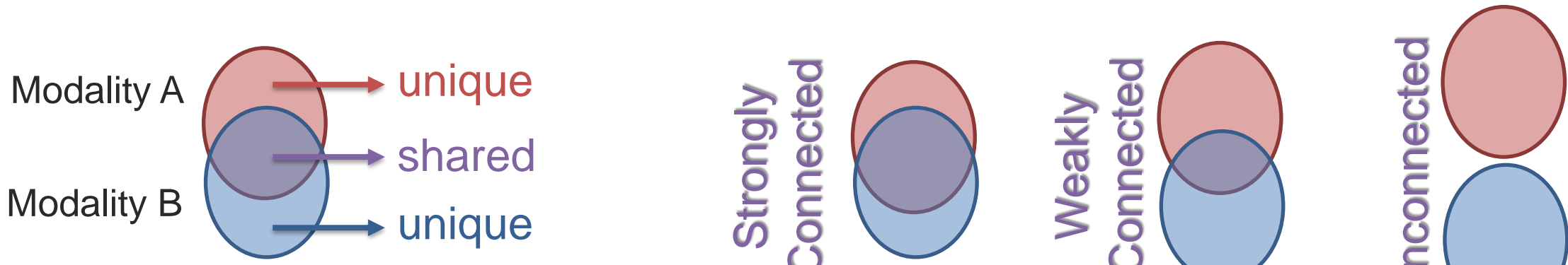
# Connections between Modalities



## Connection types:

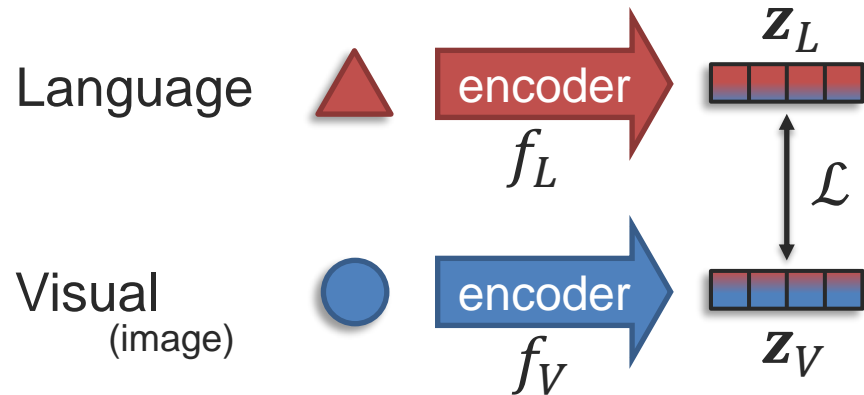
- Co-occurrence
- Correlation
- Causality
- Relationship
- ⋮

➔ knowledge of one modality provides information about the other modality



Connections are part of the data... and models will try to learn them.

# Coordinated Representations – Example: CLIP

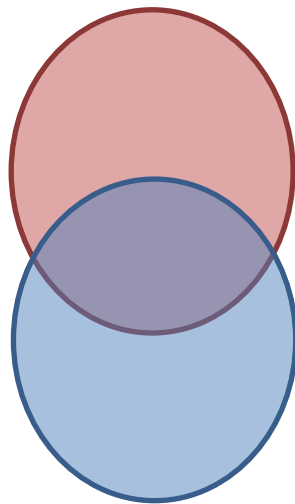
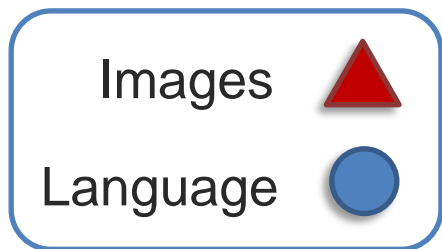


## Popular contrastive loss: InfoNCE

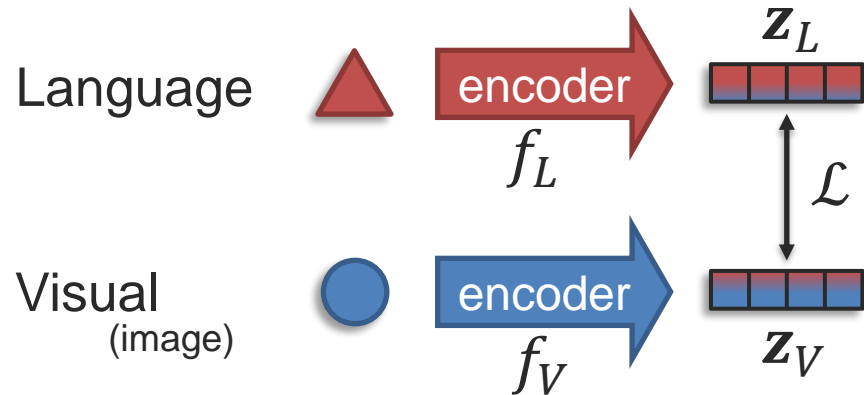
$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(\mathbf{z}_A^i, \mathbf{z}_B^i)}{\sum_{j=1}^N \text{sim}(\mathbf{z}_A^i, \mathbf{z}_B^j)}$$

Annotations for the equation:  
- A green box highlights  $\text{sim}(\mathbf{z}_A^i, \mathbf{z}_B^i)$  with an arrow pointing to the text 'positive pairs'.  
- A red box highlights the denominator  $\sum_{j=1}^N \text{sim}(\mathbf{z}_A^i, \mathbf{z}_B^j)$  with an arrow pointing to the text 'negative pairs and positive pairs'.  
- An orange arrow points from the text 'Similarity function can be cosine similarity' to the  $\text{sim}$  function in the numerator.

## Connected modalities:



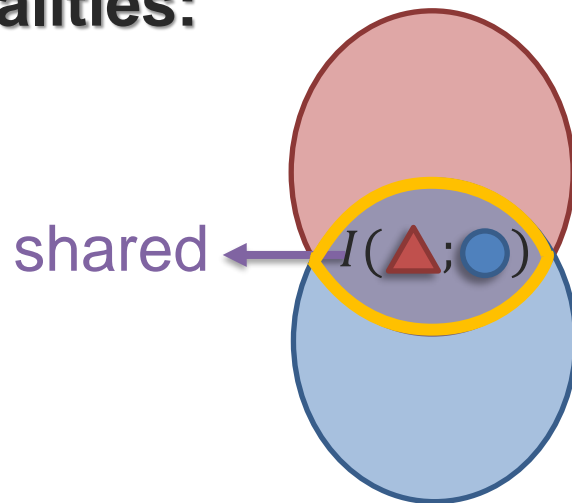
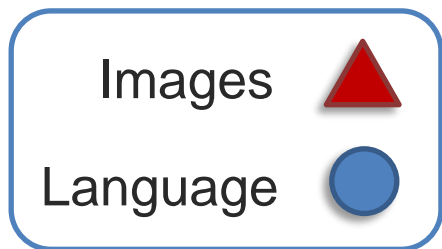
# Coordinated Representations – Example: CLIP



Popular contrastive loss: InfoNCE

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

Connected modalities:



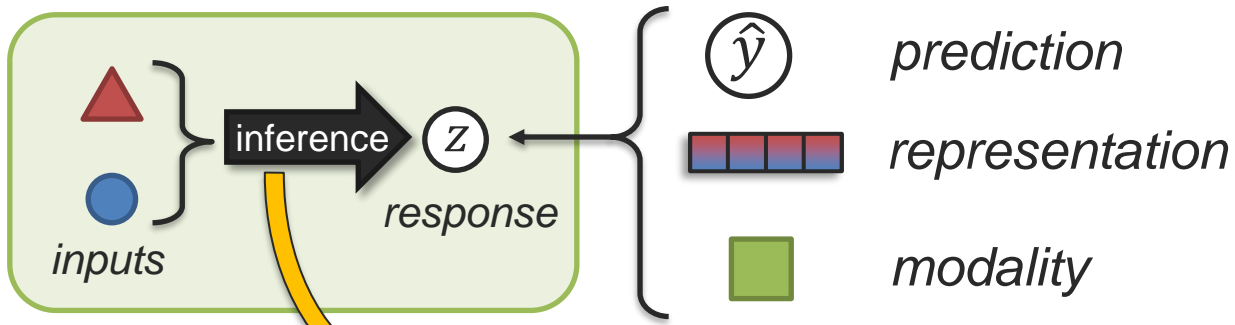
Mutual information  $I(X; Y)$

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



CLIP focuses on shared connections

# Modality Interactions

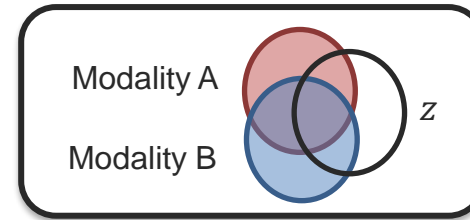


**Interactions happen during inference!  
(from human or model)**

Interactions require more than the input modalities!

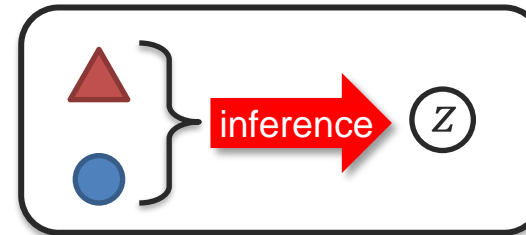
## Interactions taxonomy:

### Level 1: Response(s) and Input Modalities



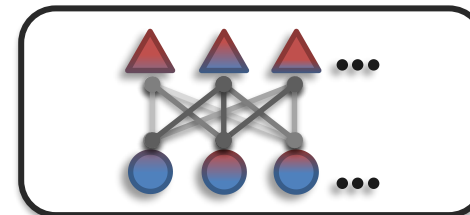
- Co-occurrence
- Redundancy
- Dominance
- Emergence
- ...

### Level 2: Interactions – Internal Mechanics



- Additive
- Multiplicative
- Polynomial
- Nonlinear
- ...

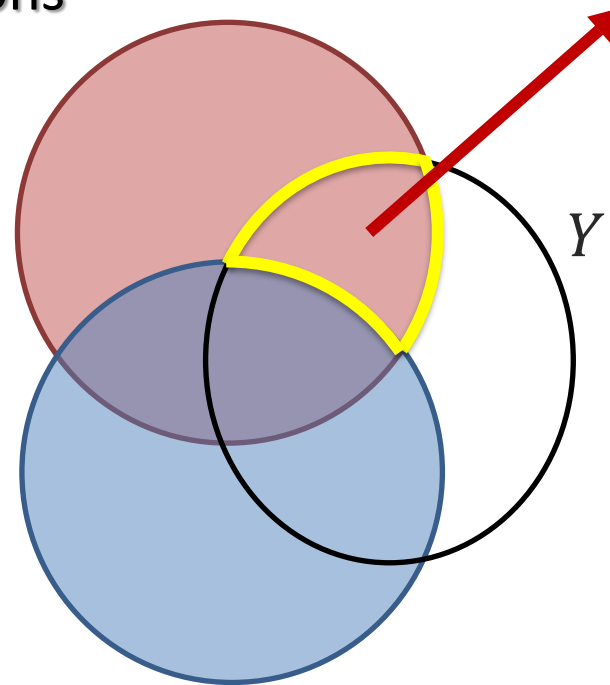
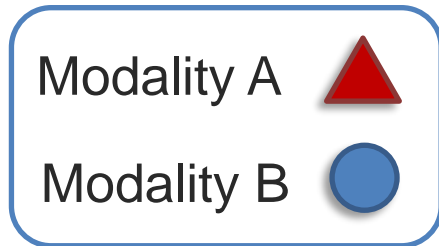
### Level 3: Contextualized Interactions



- Temporal
- Hierarchy
- Multimodal
- ...

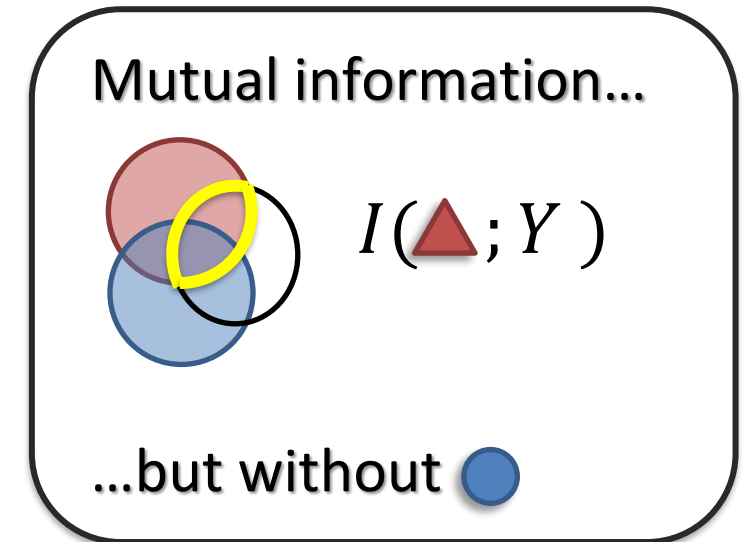
# Modeling Interactions – Response and Input Modalities (Level 1)

**Information theory** as a framework for modeling Level 1 interactions between input modalities and responses

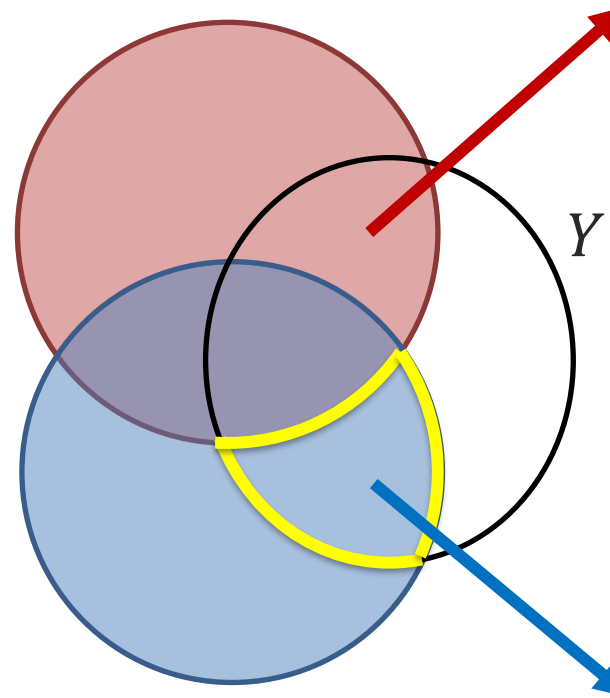
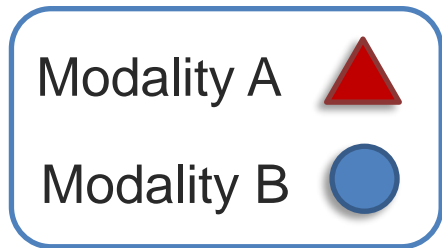


Task-relevant **unique** information

$$I(\triangle; Y | \bullet)$$



# Modeling Interactions – Response and Input Modalities (Level 1)



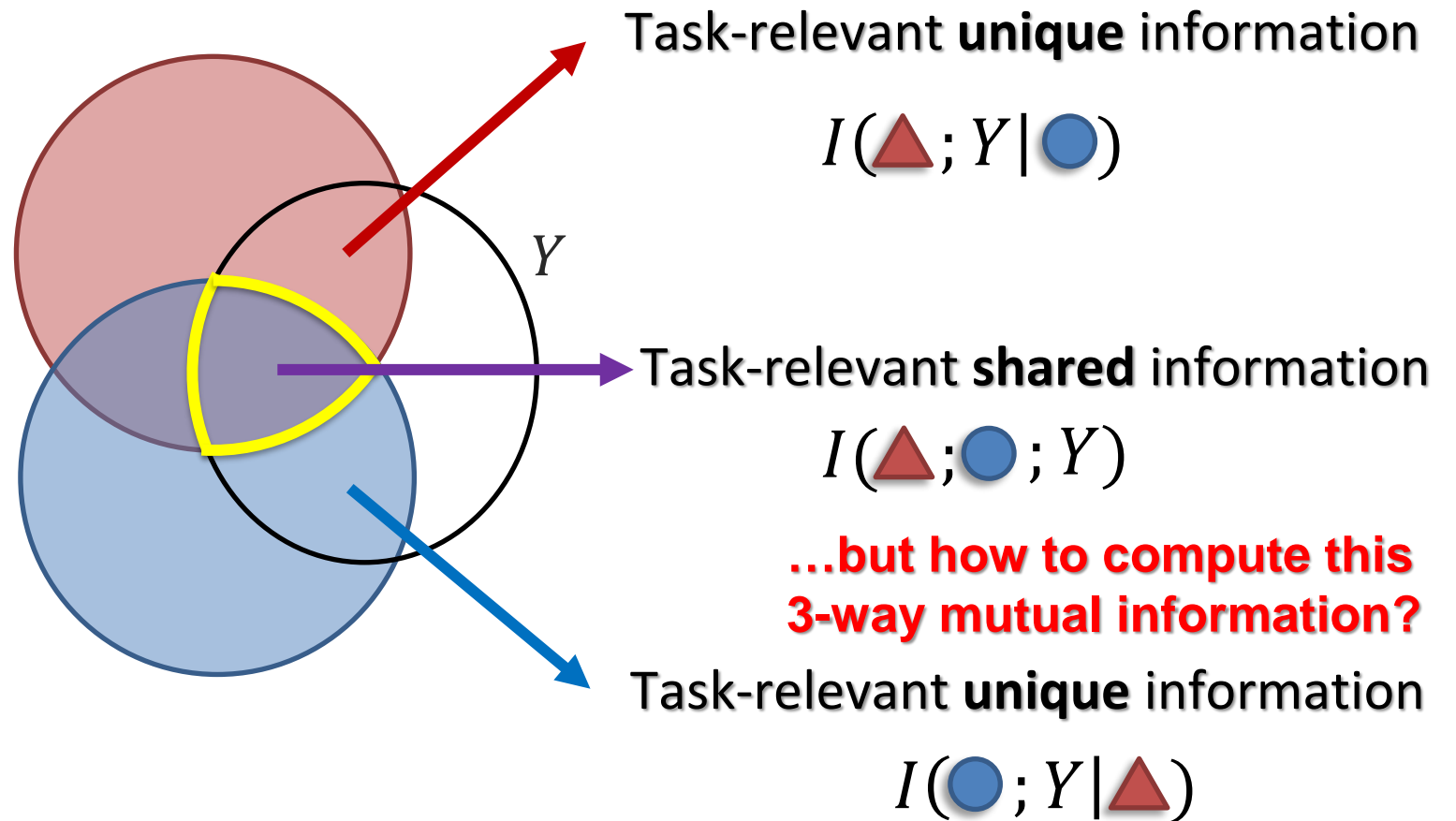
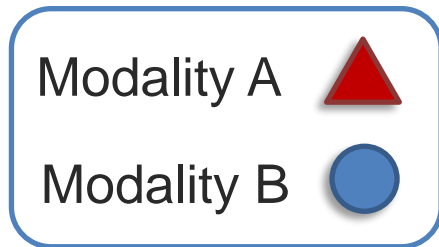
Task-relevant **unique** information

$$I(\blacktriangle; Y | \bullet)$$

Task-relevant **unique** information

$$I(\bullet; Y | \blacktriangle)$$

# Modeling Interactions – Response and Input Modalities (Level 1)





# Modeling Interactions – Response and Input Modalities (Level 1)

## Partial Information Decomposition (PID)

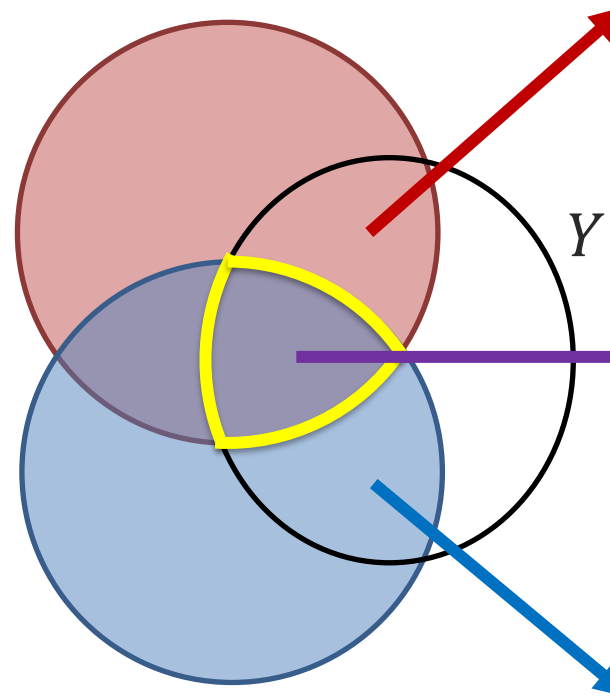
$$I(\triangle; \bullet; Y) = R - S$$

factorizes 3-way mutual information into:

R: redundancy

S: Synergy

More about PID in future lectures  
and reading assignments



Task-relevant **unique** information

$$I(\triangle; Y | \bullet)$$

Task-relevant **shared** information

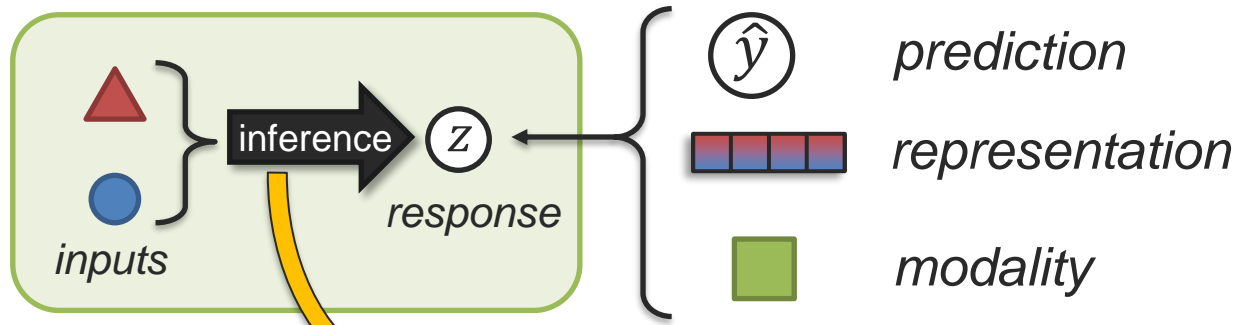
$$I(\triangle; \bullet; Y)$$

...but how to compute  
3-way mutual information?

Task-relevant **unique** information

$$I(\bullet; Y | \triangle)$$

# Modality Interactions

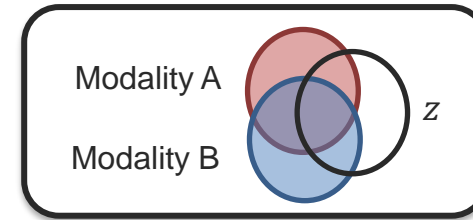


**Interactions happen during inference!  
(from human or model)**

Interactions require more than the input modalities!

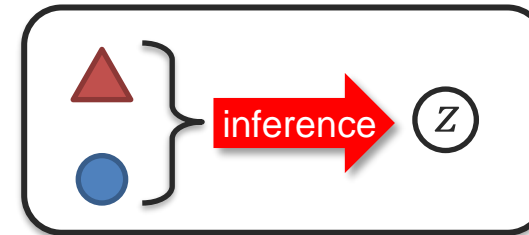
## Interactions taxonomy:

### Level 1: Response(s) and Input Modalities



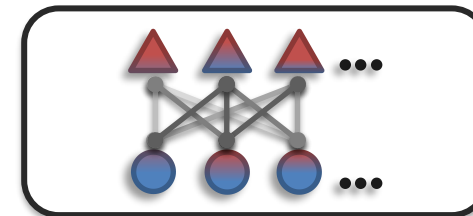
- Co-occurrence
- Redundancy
- Dominance
- Emergence
- ...

### Level 2: Interactions – Internal Mechanics



- Additive
- Multiplicative
- Polynomial
- Nonlinear
- ...

### Level 3: Contextualized Interactions



- Temporal
- Hierarchy
- Multimodal
- ...

# Challenge 2: Alignment

---

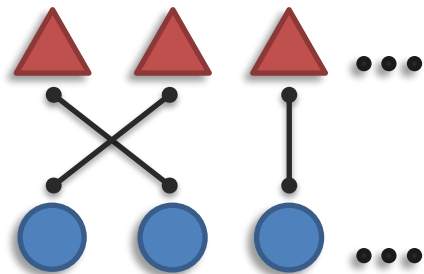
## Challenge 2: Alignment

---

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

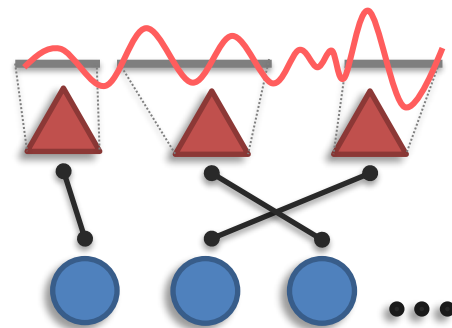
### Sub-challenges:

#### Discrete Alignment



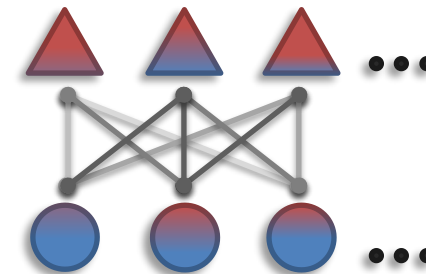
Discrete elements and connections

#### Continuous Alignment



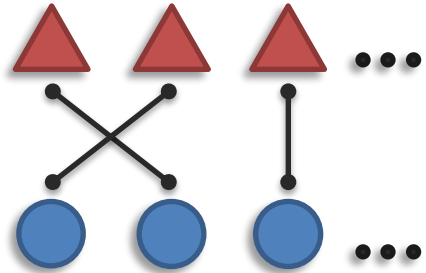
Segmentation and continuous warping

#### Contextualized Representation

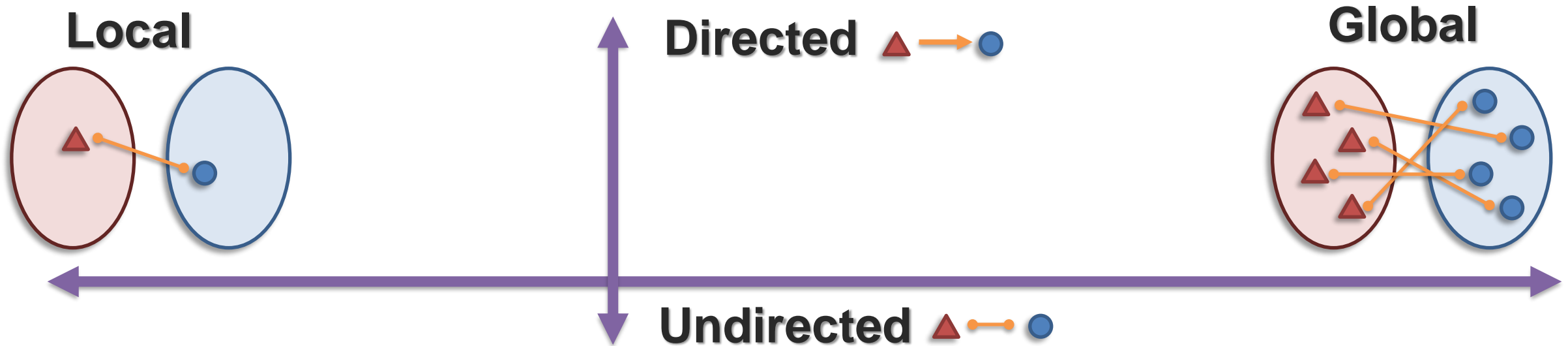


Alignment + representation

## Sub-Challenge 2a: Discrete Alignment



**Definition:** Identify and model discrete connections between elements of multiple modalities



# Language Grounding

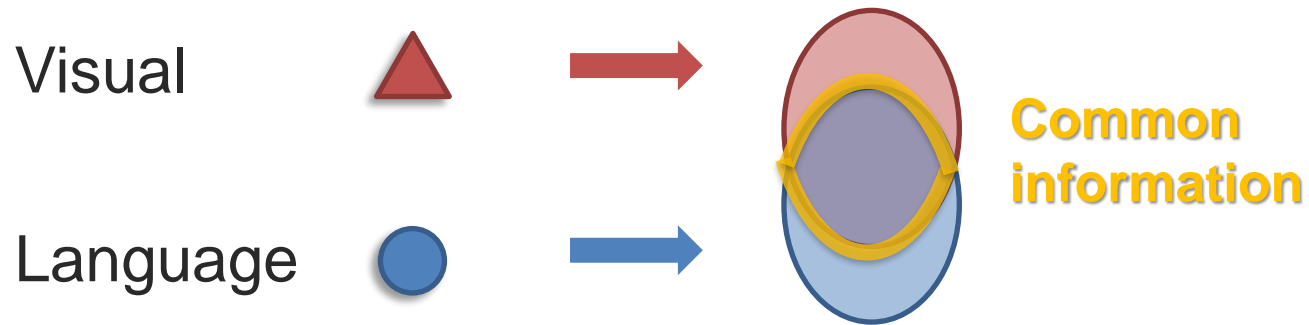
---

**Definition:** Tying language (words, phrases,...) to non-linguistic elements, such as the visual world (objects, people, ...)



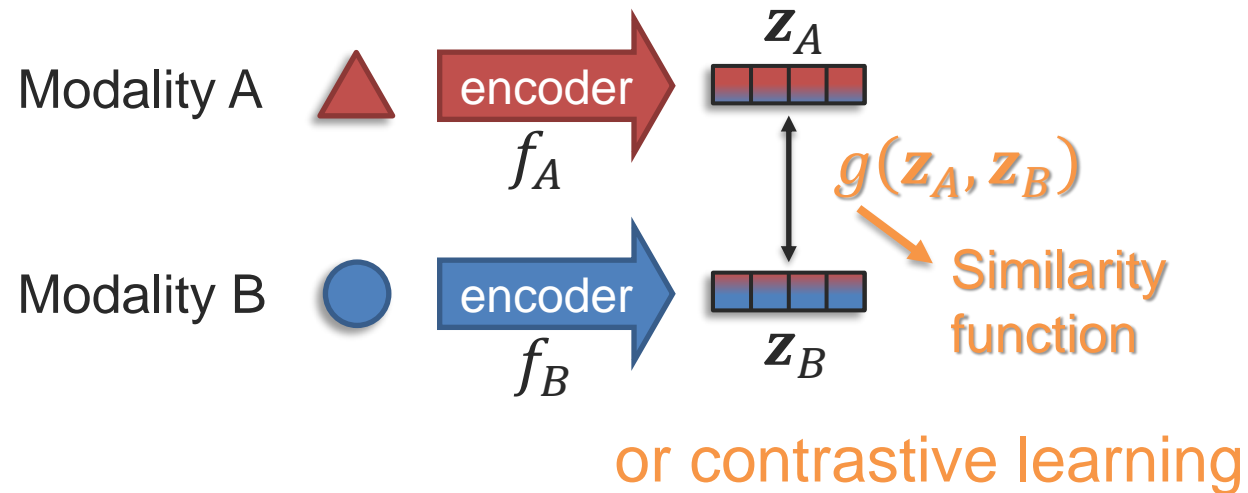
A woman reading newspaper

# Local Alignment – Coordinated Representations

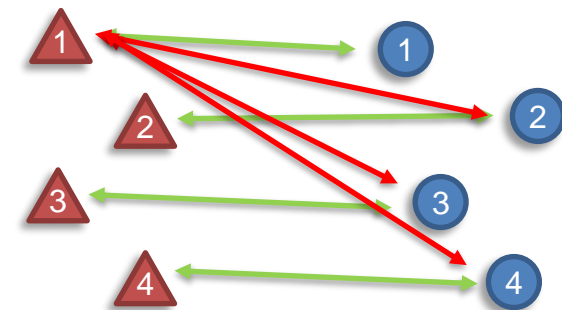


A **woman** reading **newspaper**

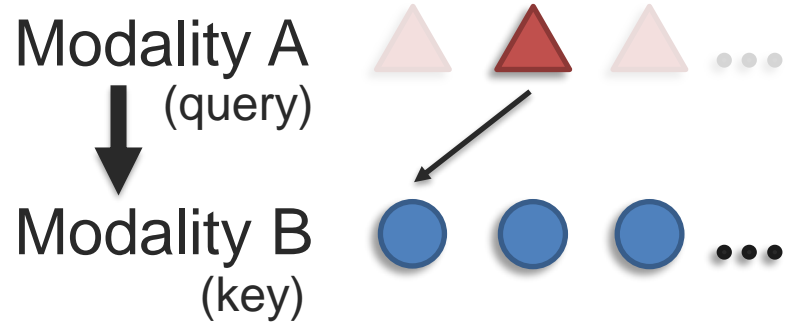
Learning coordinated representations:



Supervision: Paired data



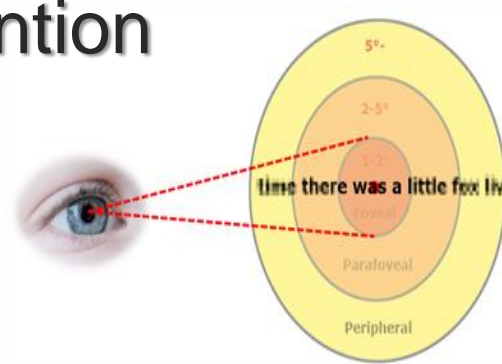
# Directed Alignment



A woman is throwing a frisbee

Which object?

## Attention



1 Soft attention

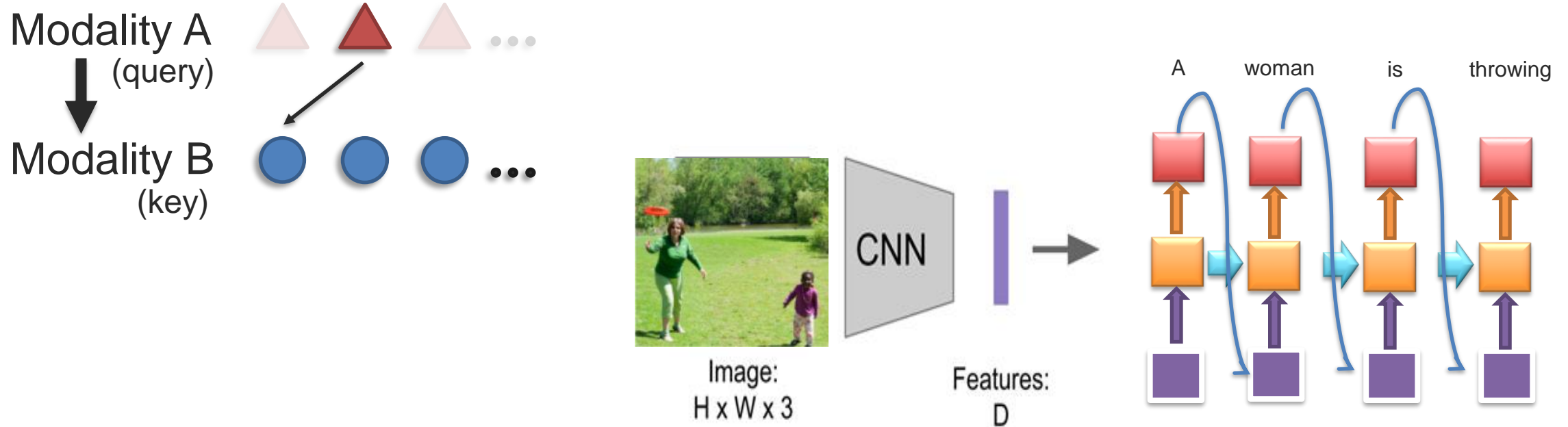


2 Hard attention



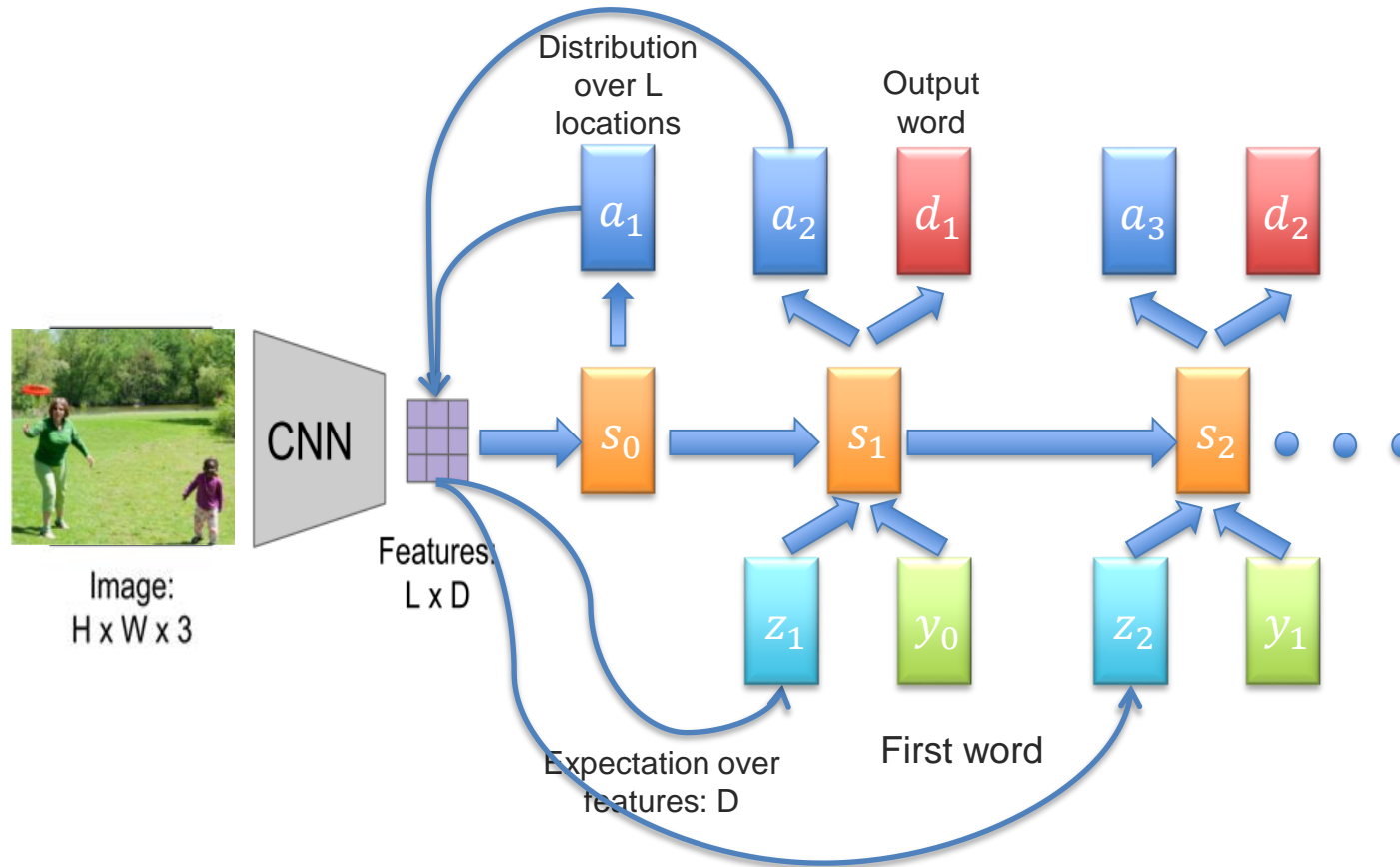


# Directed Alignment – Image Captioning



Should we always use the final layer of the CNN for all generated words?

# Directed Alignment – Image Captioning



# Attention Gates

---

Before:

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, \mathbf{s}_i, \mathbf{z}),$$

where  $\mathbf{z} = \mathbf{h}_T$ , last encoder state and  $\mathbf{s}_i$  is the current state of the decoder

Now:

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, \mathbf{s}_i, \mathbf{z}_i)$$

Have an attention “gate”

- A different context  $\mathbf{z}_i$  used at each time step!

- $\mathbf{z}_i = \sum_{j=1}^{T_x} \alpha_{ij} \mathbf{h}_j$

$\alpha_{ij}$  is the (scalar) attention for word  $j$  at generation step  $i$

# Attention Gates

---

So how do we determine  $\alpha_{ij}$ ?

$$\alpha_{i,j} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad \Rightarrow \text{softmax, making sure they sum to 1}$$

where:

$$e_{ij} = \mathbf{v}^T \sigma(W \mathbf{s}_{i-1} + U \mathbf{h}_j)$$

a feedforward network that can tell us how important the current encoding is

$\mathbf{v}$ ,  $W$ ,  $U$ — learnable weights

$$\mathbf{z}_i = \sum_{j=1}^{T_x} \alpha_{ij} \mathbf{h}_j$$

← expectation of the context (a fancy way to say it's a weighted average)

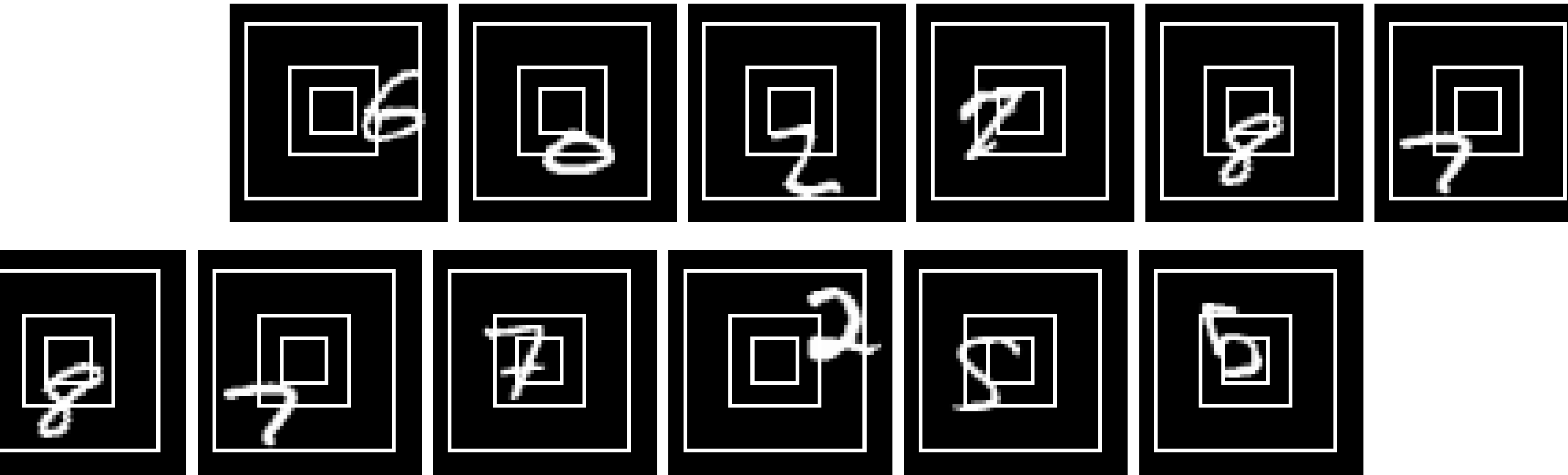
# Example – Image Captioning



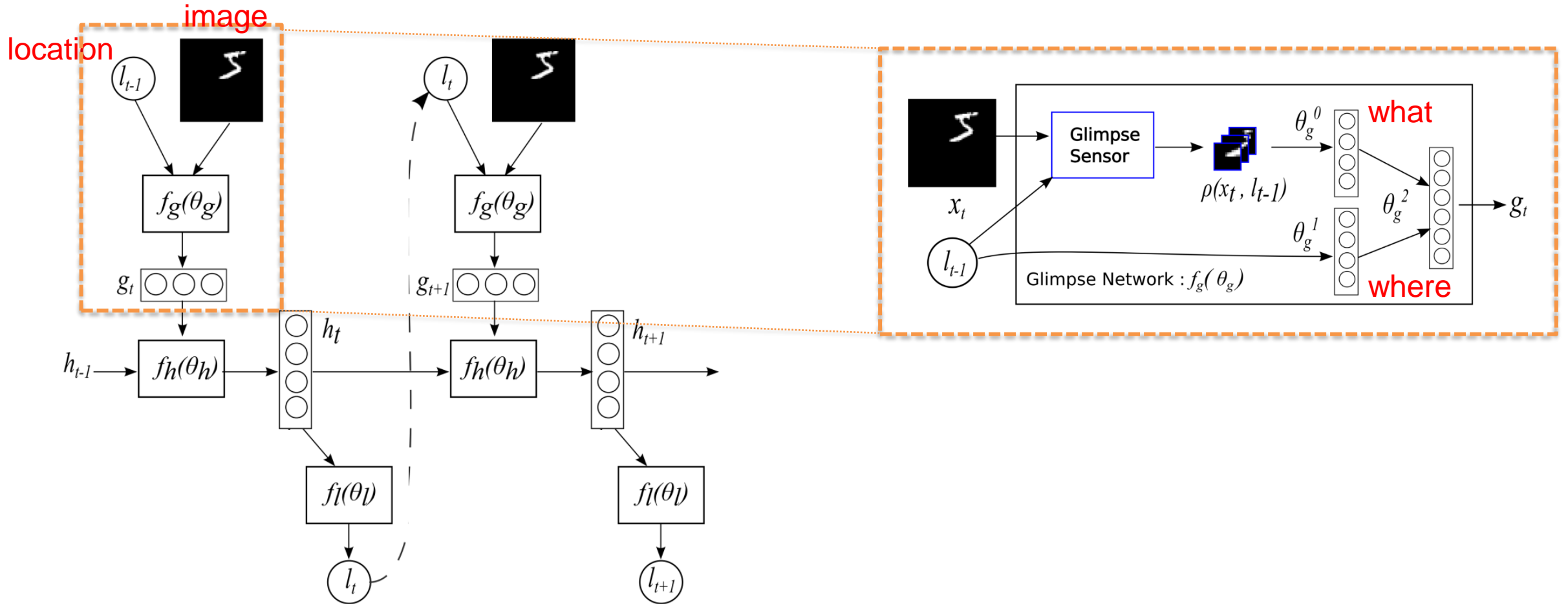
[Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Xu et al., 2015]

# Hard attention - Example

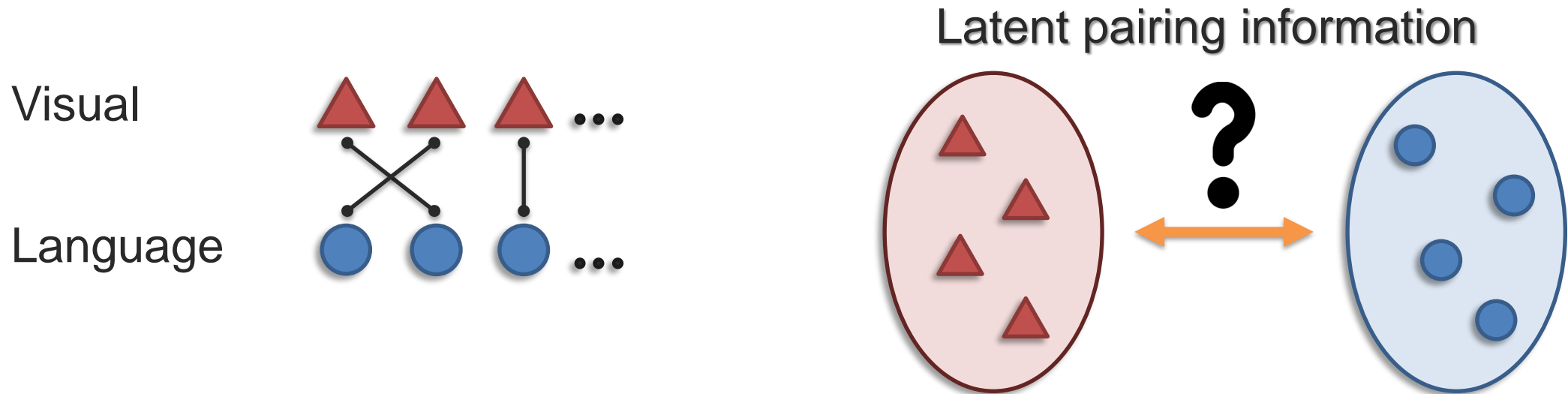
---



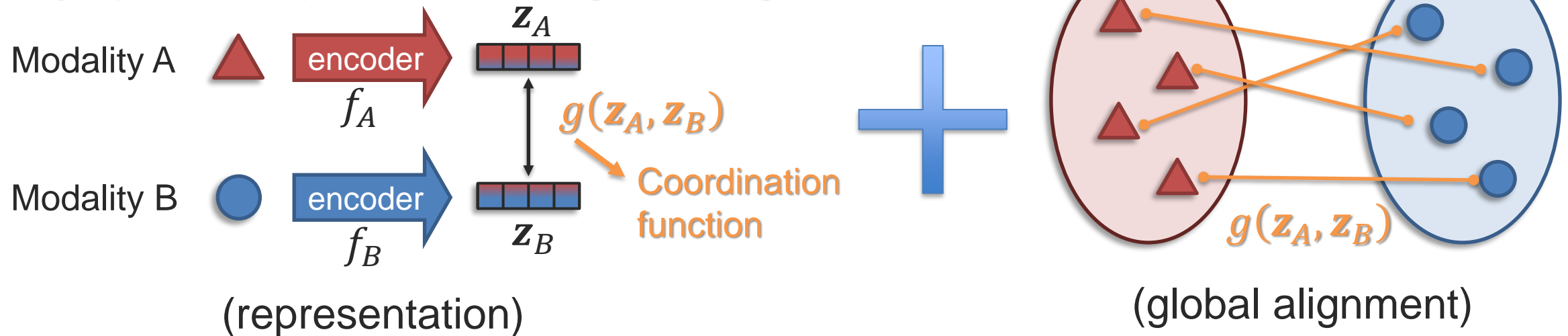
# Hard Attention – Recurrent Model of Visual Attention



# Global Alignment

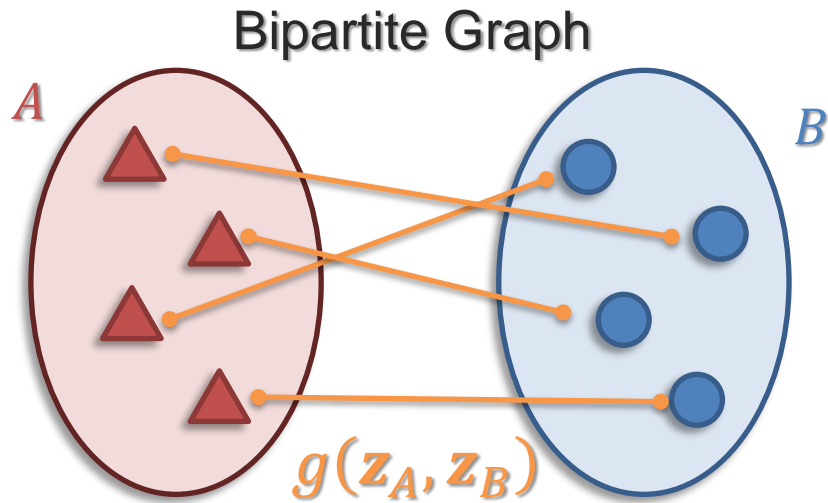


Jointly optimize representation + global alignment:





# Assignment Problem



## Initial assumptions:

- Same number of elements in A and B modalities
- 1-to-1 “hard” alignment between elements
- All elements assigned (aka “perfect matching”)

➔ How to solve?

Naive solution: check all assignments

Better solution: Linear Programming

Assignment:  ~~$f: A \rightarrow B$~~   
(vector of indices)

$x_{ij} = 1$  when matching connection, otherwise 0

Similarity weights:  ~~$w_{(l, f(l))} = g(z_A^l, z_B^{f(l)})$~~

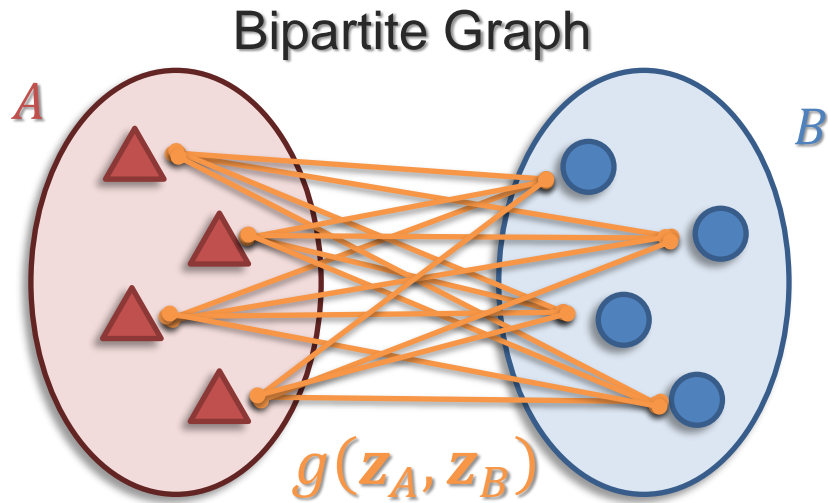
$w_{(i,j)} = g(z_A^i, z_B^j)$

Maximize:  ~~$\max_{f \in \text{Perm}(N)} \sum_{i=1}^N w_{i, f(i)}$~~

$\max_{\{x_{ij}\}} \sum_{(i,j) \in A \times B} w_{i,j} x_{ij}$

➔ Can be solved with simplex algorithm

# Optimal transport



## New assumptions:

- Different number of elements in A and B modalities
- Many-to-many “soft” alignment between elements

➔ It can be seen as “transporting” elements from modality A to modality B (and vice-versa)

Assignments:  $x_{(i,j)}$ : soft alignment between  $\mathbf{z}_A^i$  and  $\mathbf{z}_B^j$

Similarity weights:  $w_{(i,j)} = g(\mathbf{z}_A^i, \mathbf{z}_B^j)$

Maximize:  $\max_{\{x_{ij}\}} \sum_{(i,j) \in A \times B} w_{i,j} x_{ij}$

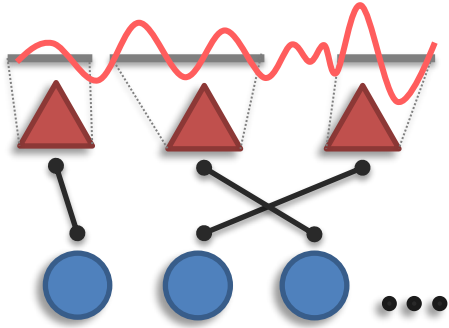
➔ Wasserstein distance gives optimal transport

# Continuous Alignment



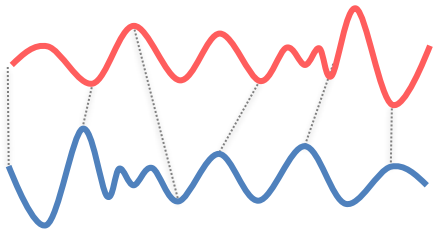
## Challenge 2b: Continuous Alignment

---

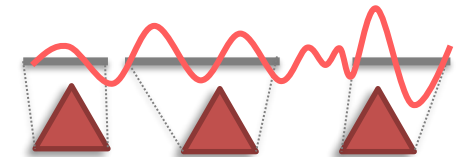


**Definition:** Model alignment between modalities with continuous signals and no explicit elements

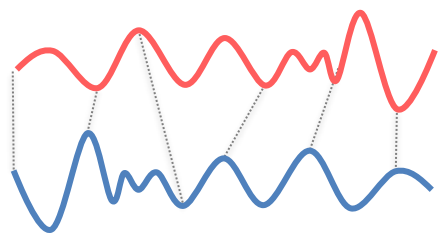
Continuous  
warping



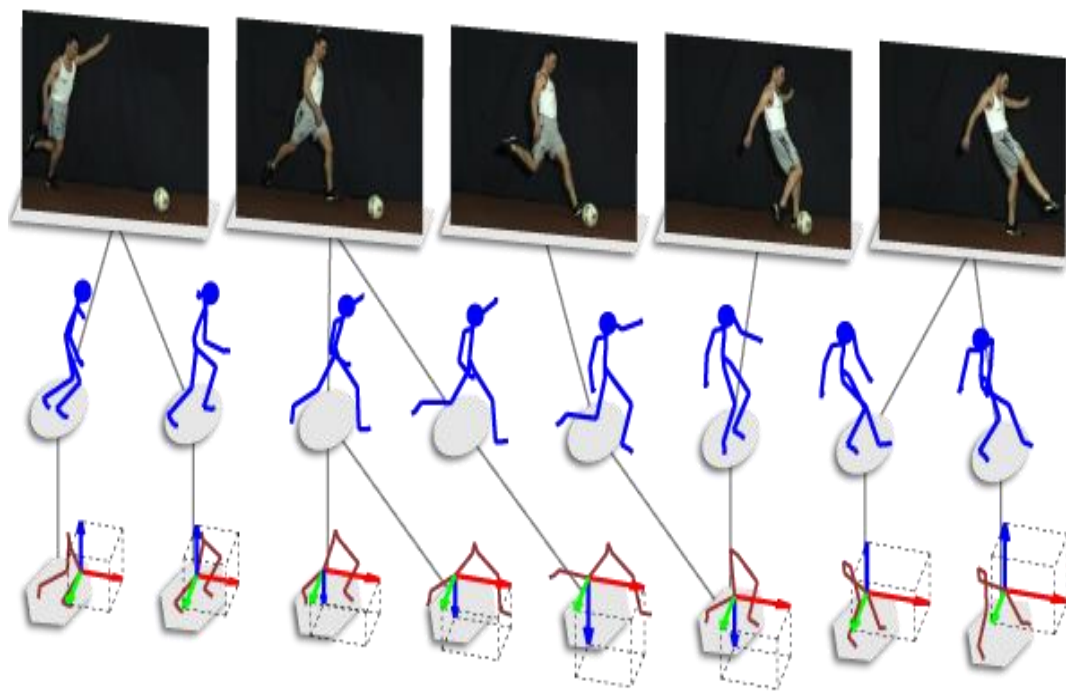
Discretization  
(segmentation)



# Continuous Warping – Example



➔ Aligning video sequences



# Dynamic Time Warping (DTW)

We have two unaligned temporal unimodal signals

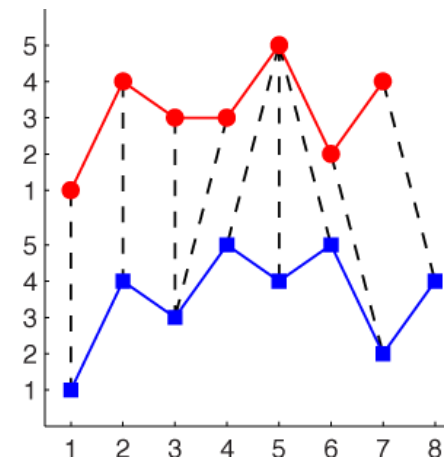
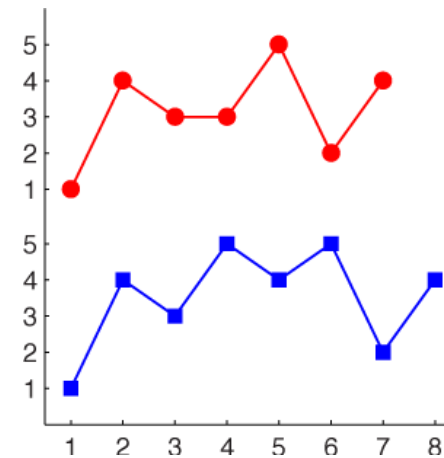
- $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_x}] \in \mathbb{R}^{d \times n_x}$
- $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{n_y}] \in \mathbb{R}^{d \times n_y}$

Find set of indices to minimize the alignment difference:

$$L(\mathbf{p}^x, \mathbf{p}^y) = \sum_{t=1}^l \left\| \mathbf{x}_{p_t^x} - \mathbf{y}_{p_t^y} \right\|_2^2$$

where  $\mathbf{p}^x$  and  $\mathbf{p}^y$  are index vectors of same length

Dynamic Time Warping is designed to find these index vectors!

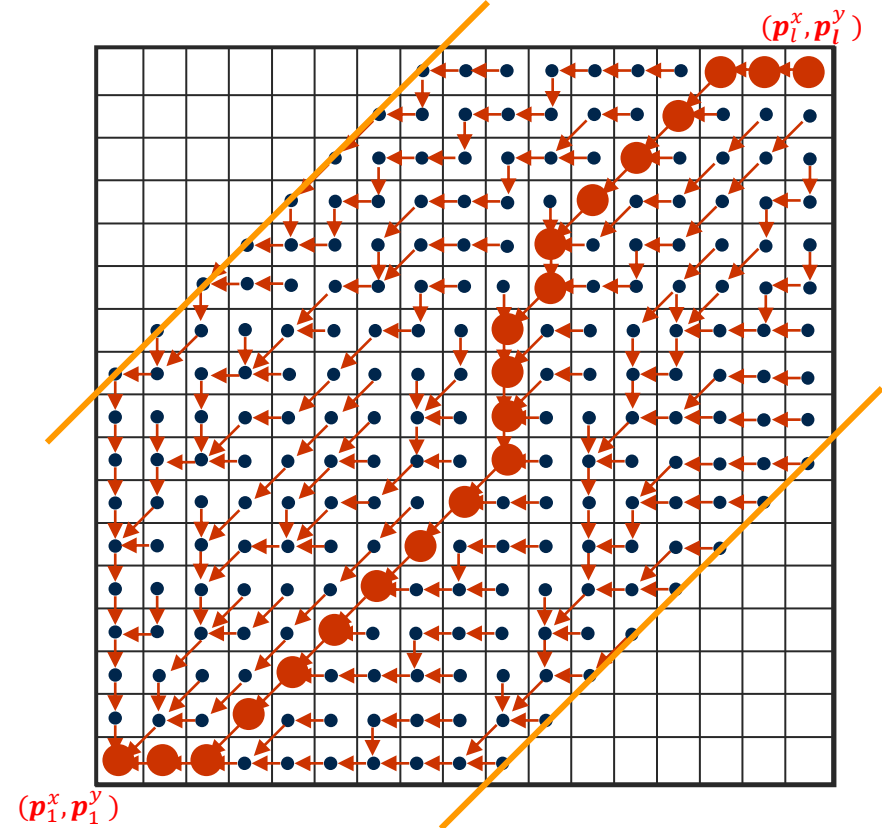


# Dynamic Time Warping (DTW)

Lowest cost path in a cost matrix

- Restrictions?
  - Monotonicity – no going back in time
  - Continuity - no gaps
  - Boundary conditions - start and end at the same points
  - Warping window - don't get too far from diagonal
  - Slope constraint – do not insert or skip too much

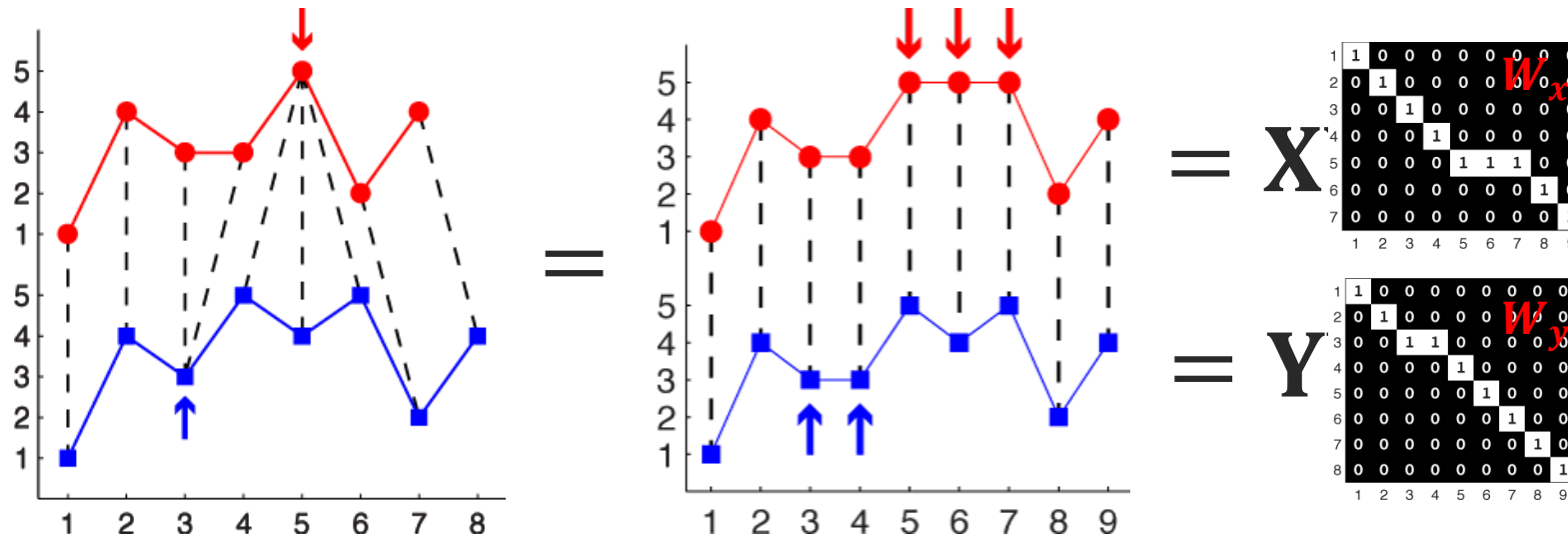
Solved using dynamic programming while respecting the restrictions



# DTW alternative formulation

$$L(\mathbf{p}^x, \mathbf{p}^y) = \sum_{t=1}^l \left\| \mathbf{x}_{p_t^x} - \mathbf{y}_{p_t^y} \right\|_2^2$$

Replication doesn't change the objective!



Alternative objective:

$$L(\mathbf{W}_x, \mathbf{W}_y) = \left\| \mathbf{X}\mathbf{W}_x - \mathbf{Y}\mathbf{W}_y \right\|_F^2$$

Frobenius norm  $\|\mathbf{A}\|_F^2 = \sum_i \sum_j |a_{i,j}|^2$

$\mathbf{X}, \mathbf{Y}$  – original signals (same #rows, possibly different #columns)

$\mathbf{W}_x, \mathbf{W}_y$  - alignment matrices

A differentiable version of DTW also exists...

<https://arxiv.org/pdf/1703.01541.pdf>



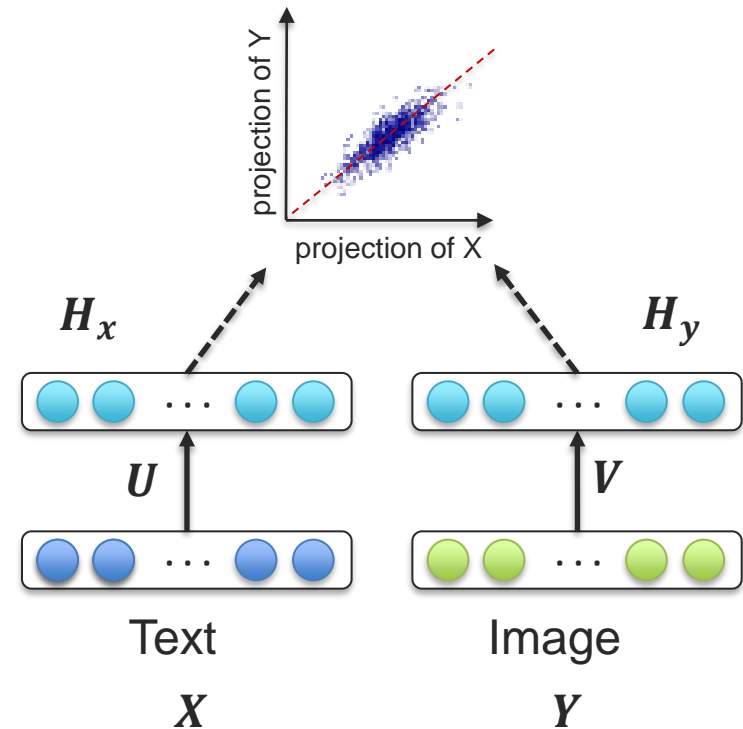
# Canonical Correlation Analysis – Reminder

CCA loss can also be re-written as:

$$L(U, V) = \|\mathbf{U}^T \mathbf{X} - \mathbf{V}^T \mathbf{Y}\|_F^2$$

subject to:

$$\mathbf{U}^T \boldsymbol{\Sigma}_{YY} \mathbf{U} = \mathbf{V}^T \boldsymbol{\Sigma}_{YY} \mathbf{V} = \mathbf{I}, \mathbf{u}_{(j)}^T \boldsymbol{\Sigma}_{XY} \mathbf{v}_{(i)} = 0$$



# Canonical Time Warping

---

Dynamic Time Warping + Canonical Correlation Analysis = Canonical Time Warping

$$L(\mathbf{U}, \mathbf{V}, \mathbf{W}_x, \mathbf{W}_y) = \|\mathbf{U}^T \mathbf{X} \mathbf{W}_x - \mathbf{V}^T \mathbf{Y} \mathbf{W}_y\|_F^2$$

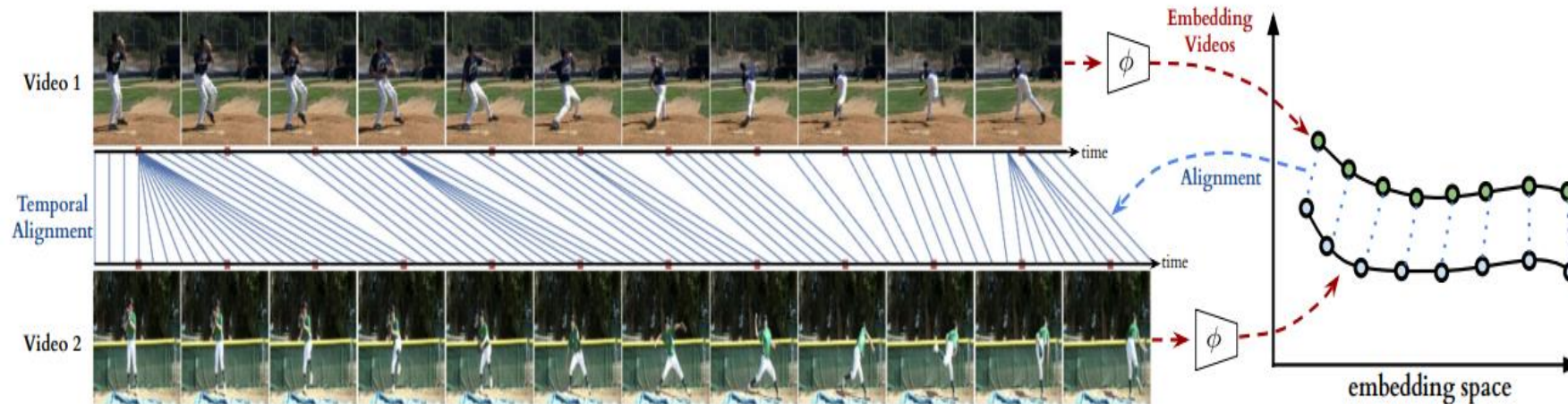
Allows to align multi-modal or multi-view (same modality but from a different point of view)

- $\mathbf{W}_x, \mathbf{W}_y$  – temporal alignment
- $\mathbf{U}, \mathbf{V}$  – cross-modal (spatial) alignment

[Canonical Time Warping for Alignment of Human Behavior, Zhou and De la Torre, 2009]

# Temporal Alignment and Neural Representation Learning

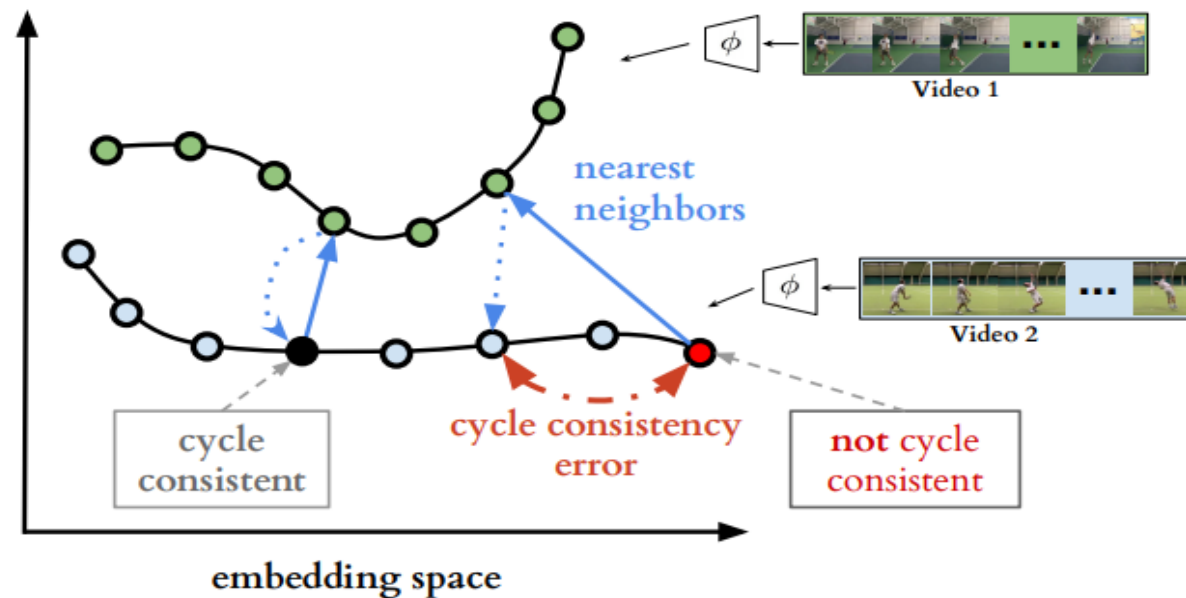
**Premise:** we have paired video sequences that can be temporally aligned



How can we define a loss function to enforce the alignment between sequences while at the same time learning good representations?

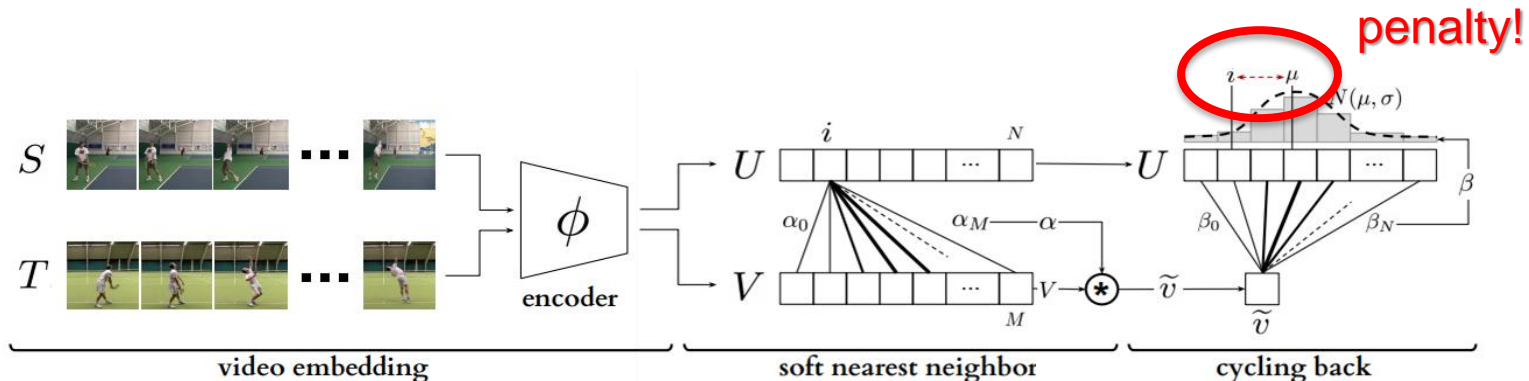
# Temporal Cycle-Consistency Learning

Solution: Representation learning by enforcing **Cycle consistency**



**Main idea:** My closest neighbor also views me as their closest neighbor

# Temporal Cycle-Consistency Learning



Compute “soft” / “weighted” nearest neighbour:

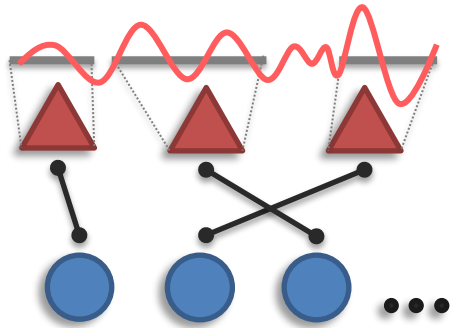
distances:  $\alpha_j = \frac{e^{-\|u_i - v_j\|^2}}{\sum_k^M e^{-\|u_i - v_k\|^2}}$       Soft nearest neighbor:  $\tilde{v} = \sum_j^M \alpha_j v_j$

Find the nearest neighbor the other way and then penalize the distance:

$$\beta_k = \frac{e^{-\|\tilde{v} - u_k\|^2}}{\sum_j^N e^{-\|\tilde{v} - u_j\|^2}} \quad L_{cbr} = \frac{|i - \mu|^2}{\sigma^2} + \lambda \log(\sigma)$$

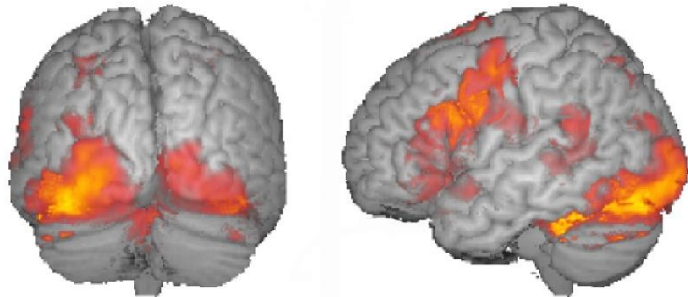
# Discretization (aka Segmentation)

---

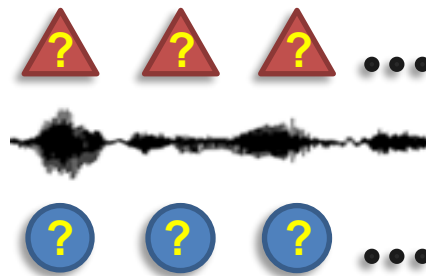


Common assumptions: ① Segmented elements

Examples:



Medical imaging



Signals



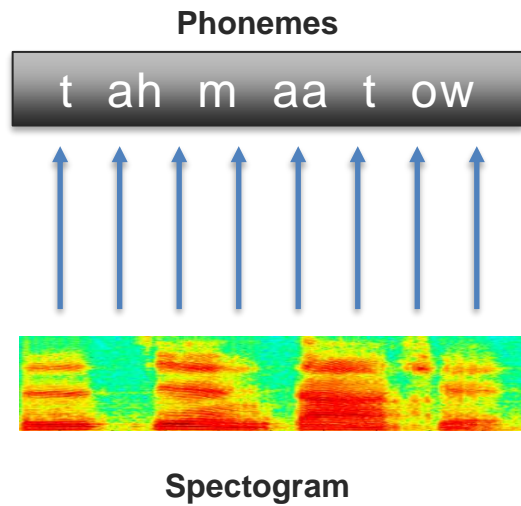
Images

objects

# Discretization – Example

---

## Sequence Labeling and Alignment

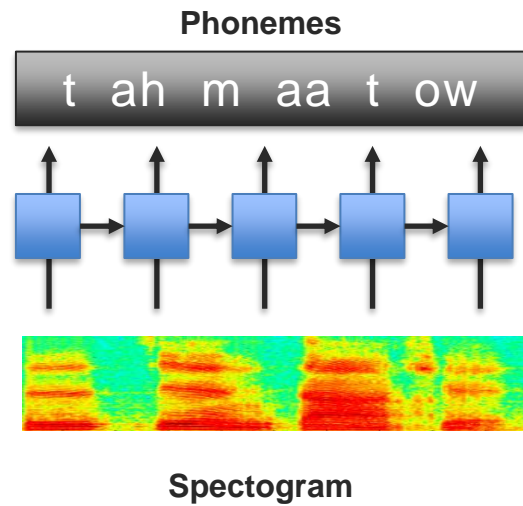


How can we predict the sequence  
of phoneme labels?

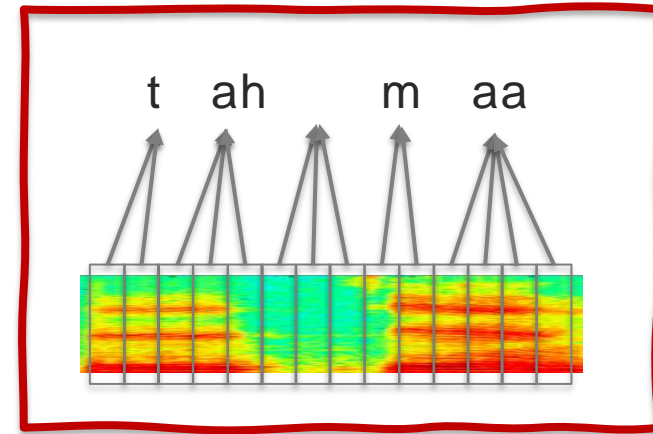
# Discretization – Example

---

## Sequence Labeling and Alignment



Challenge: many-to-1 alignment



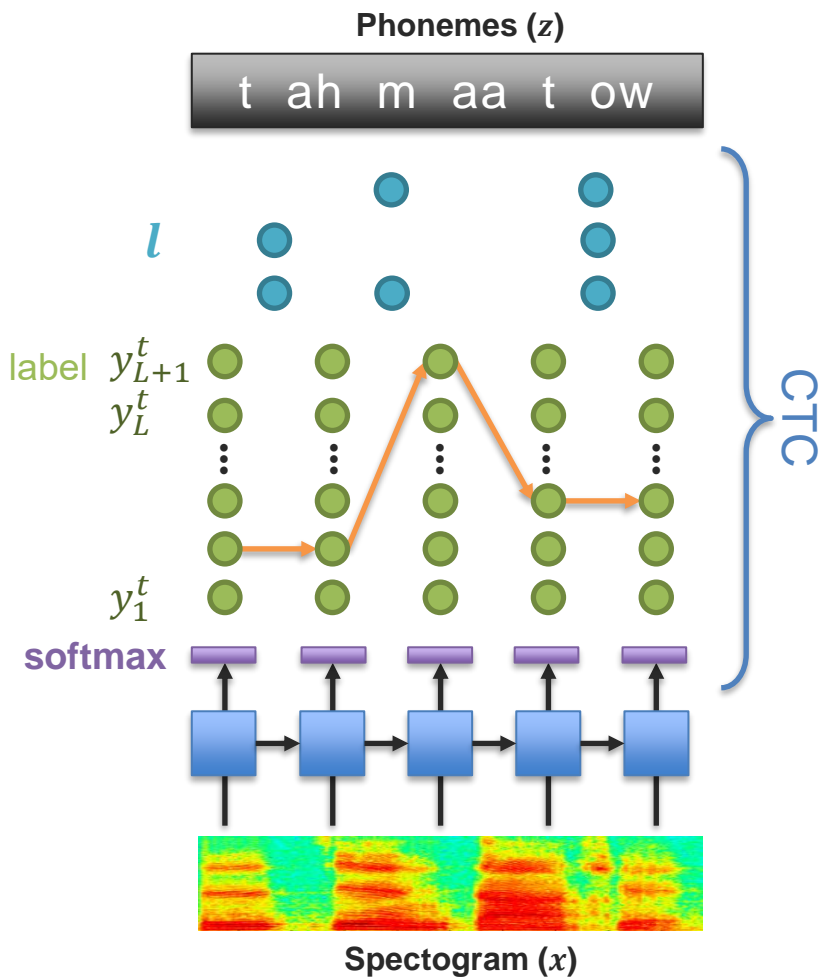
How can we predict the sequence of phoneme labels?



# Discretization – A Classification Approach

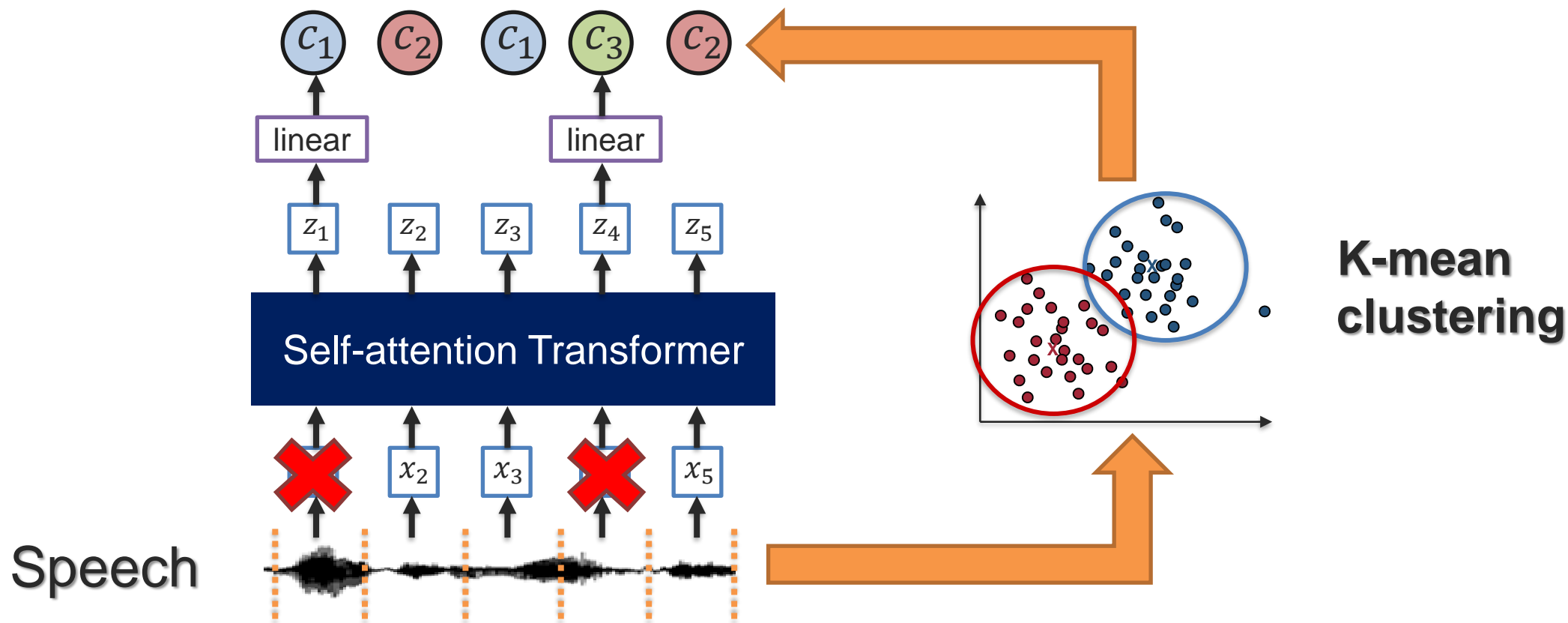
## Connectionist Temporal Classification

- ④ Most probable sequence labels
- ③ Predicted labels  $l$
- ② Path  $\pi$  over the activations:
- ① Output activations (distribution):



# Discretization and Representation – Cluster-based Approaches

## HUBERT: Hidden-Unit BERT

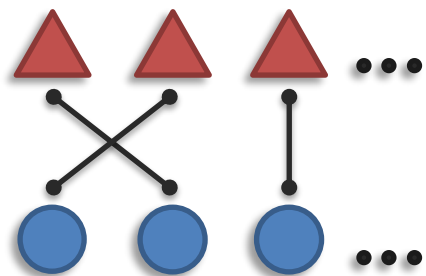


## Challenge 2: Alignment

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

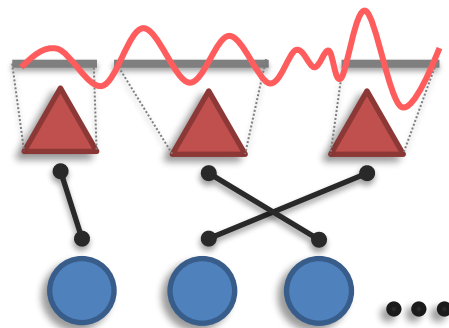
### Sub-challenges:

#### Discrete Alignment



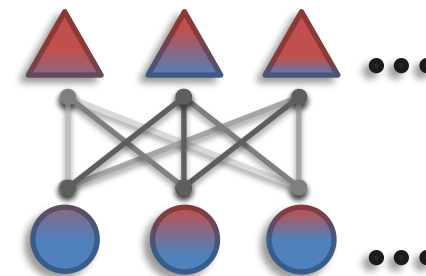
Discrete elements  
and connections

#### Continuous Alignment



Segmentation and  
continuous warping

#### Contextualized Representation



Alignment  
+ representation

Thursday!

