



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Multimodal Machine Learning

## Lecture 4.2: 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

# Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures
<b>Week 1</b> 8/30 & 9/1	<b>Course introduction</b> <ul style="list-style-type: none"> <li>Multimodal core challenges</li> <li>Course syllabus</li> </ul>	<b>Multimodal applications and datasets</b> <ul style="list-style-type: none"> <li>Research tasks and datasets</li> <li>Team projects</li> </ul>
<b>Week 2</b> 9/6 & 9/8 <i>Read due: 9/9</i>	<b>Basic concepts: neural networks</b> <ul style="list-style-type: none"> <li>Loss functions and neural networks</li> <li>Gradient and optimization</li> </ul>	<b>Unimodal representations</b> <ul style="list-style-type: none"> <li>Dimensions of heterogeneity</li> <li>Visual representations</li> </ul>
<b>Week 3</b> 9/13 & 9/15 <i>Read due: 9/16</i> <i>Proj. Due: 9/14</i>	<b>Unimodal representations</b> <ul style="list-style-type: none"> <li>Language representations</li> <li>Signals, graphs and other modalities</li> </ul>	<b>Multimodal representations</b> <ul style="list-style-type: none"> <li>Cross-modal interactions</li> <li>Multimodal fusion</li> </ul>
<b>Week 4</b> 9/20 & 9/22 <i>Proj. due: 9/25</i>	<b>Multimodal representations</b> <ul style="list-style-type: none"> <li>Coordinated representations</li> <li>Multimodal fission</li> </ul>	<b>Multimodal alignment</b> <ul style="list-style-type: none"> <li>Explicit alignment</li> <li>Multimodal grounding</li> </ul>
<b>Week 5</b> 9/27 & 9/29 <i>Read due: 9/30</i>	<i>Project hours (Research ideas)</i>	<b>Aligned representations</b> <ul style="list-style-type: none"> <li>Self-attention transformer models</li> <li>Masking and self-supervised learning</li> </ul>
<b>Week 6</b> 10/4 & 10/6 <i>Proj. due: 10/9</i>	<b>Multimodal aligned representations</b> <ul style="list-style-type: none"> <li>Multimodal transformers</li> <li>Video and graph representations</li> </ul>	<b>Multimodal Reasoning</b> <ul style="list-style-type: none"> <li>Structured and hierarchical models</li> <li>Memory models</li> </ul>

First assignment due on Sunday 9/25

Second assignment due on Sunday 10/9

# Lecture Schedule

Classes	Tuesday Lectures	Thursday Lectures
<b>Week 7</b> 10/11 & 10/13 Read due: 10/14	<b>Multimodal Reasoning</b> <ul style="list-style-type: none"><li>Reinforcement learning</li><li>Discrete structure learning</li></ul>	<b>Multimodal Reasoning</b> <ul style="list-style-type: none"><li>Logical and causal inference</li><li>External knowledge</li></ul>
<b>Week 8</b> 10/18 & 10/20	<b>Fall Break – No lectures</b>	
<b>Week 9</b> 10/25 & 10/27 Proj. due: 10/30	<b>Generation</b> <ul style="list-style-type: none"><li>Translation, summarization, creation</li><li>Generative models: VAEs</li></ul>	<b>Generation</b> <ul style="list-style-type: none"><li>GANs and diffusion models</li><li>Model evaluation and ethics</li></ul>
<b>Week 10</b> 11/1 & 11/3	<b>Project presentations (midterm)</b>	<b>Project presentations (midterm)</b>
<b>Week 11</b> 11/8 & 11/10 Read due: 11/12	<b>Transference</b> <ul style="list-style-type: none"><li>Modality transfer</li><li>Multimodal co-learning</li></ul>	<b>Quantification</b> <ul style="list-style-type: none"><li>Heterogeneity and interactions</li><li>Biases and fairness</li></ul>
<b>Week 12</b> 11/15 & 11/17 Read due: 11/21	<b>Project hours (Research ideas)</b>	<b>New research directions</b> <ul style="list-style-type: none"><li>Recent approaches in multimodal ML</li></ul>

Midterm assignment due on Sunday 10/30

# Lecture Schedule

---

Classes	Tuesday Lectures	Thursday Lectures
<b>Week 13</b> 11/22 & 11/24	<b>Thanksgiving Week – No Class –</b>	
<b>Week 14</b> 11/30 & 12/2	<b>Language, Vision, and Actions</b> <ul style="list-style-type: none"><li>• Robots, navigation and embodied AI</li><li>• Guest lecturer: Yonatan Bisk</li></ul>	<b>Multimodal Language Grounding</b> <ul style="list-style-type: none"><li>• Grounded semantics and pragmatics</li><li>• Guest lecturer: Daniel Fried</li></ul>
<b>Week 15</b> 12/6 & 12/8 <i>Proj. due: 12/11</i>	<b>Project presentations (final)</b>	<b>Project presentations (final)</b>

Final assignment due on Sunday 12/11

## Team Meetings with Instructor

---

**Sign-up deadline: Sunday 9/25 at 11pm**

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



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Multimodal Machine Learning

## Lecture 4.2: 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

---

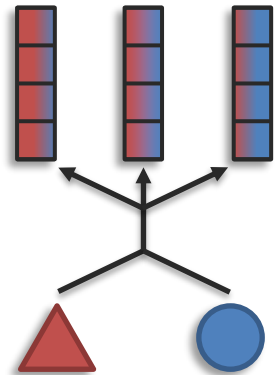
- Fine-grained fission
  - Cluster-based approach
- Discrete alignment
  - Local alignment
    - Coordinated representations; hard and soft attention
  - Global alignment
    - Assignment problem and optimal transport
- Continuous alignment
  - Continuous warping
    - Dynamic time warping
  - Discretization and segmentation



# Fine-Grained Fission

---

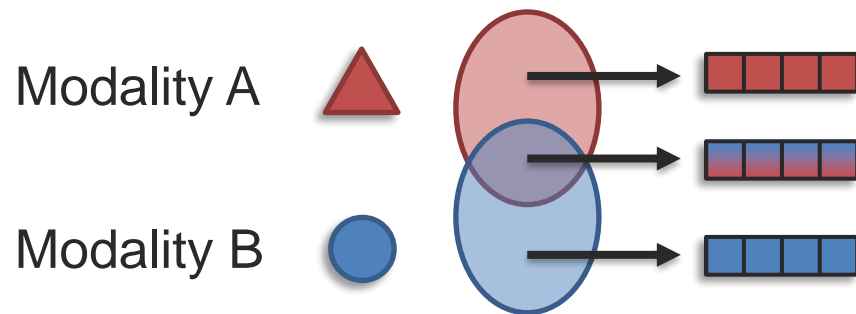
# Sub-Challenge 1c: Representation Fission



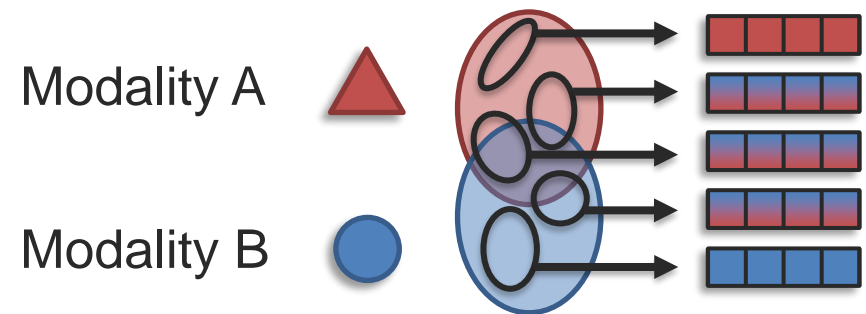
**Definition:** learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering

How to automatically discover these internal clusters, factors?

Modality-level fission:

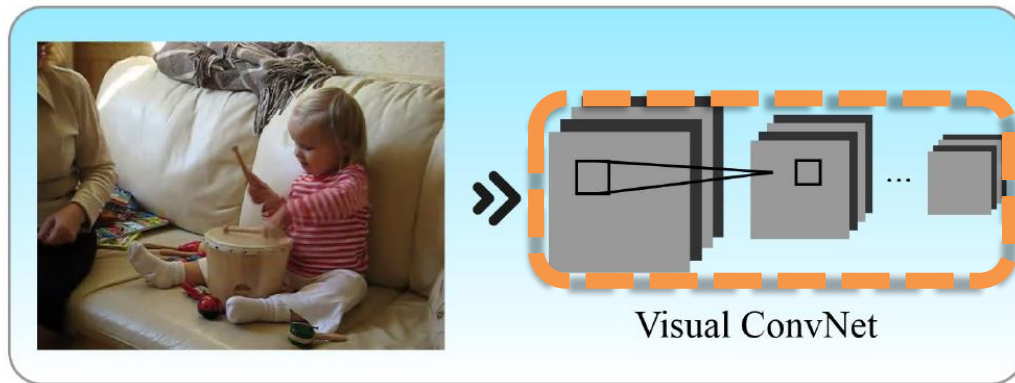


Fine-grained fission:

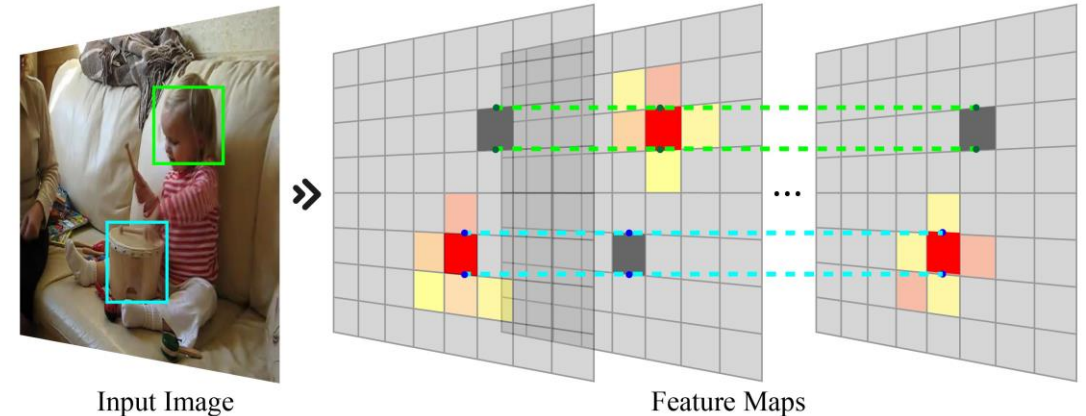
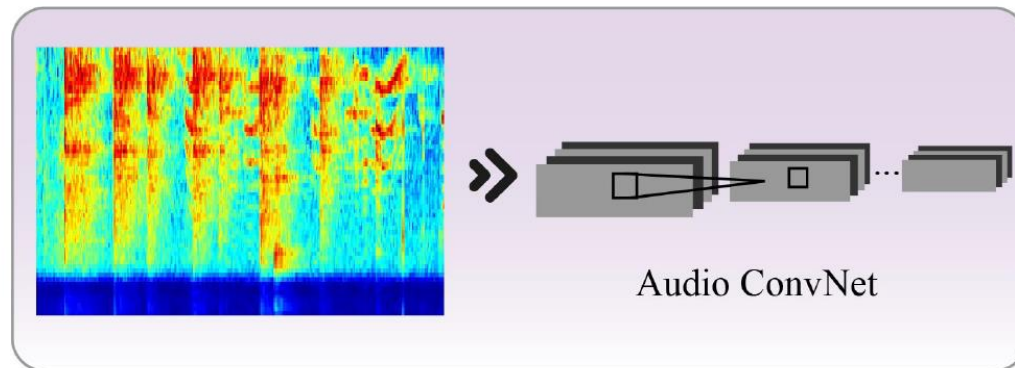


# Fine-Grained Fission – A Clustering Approach

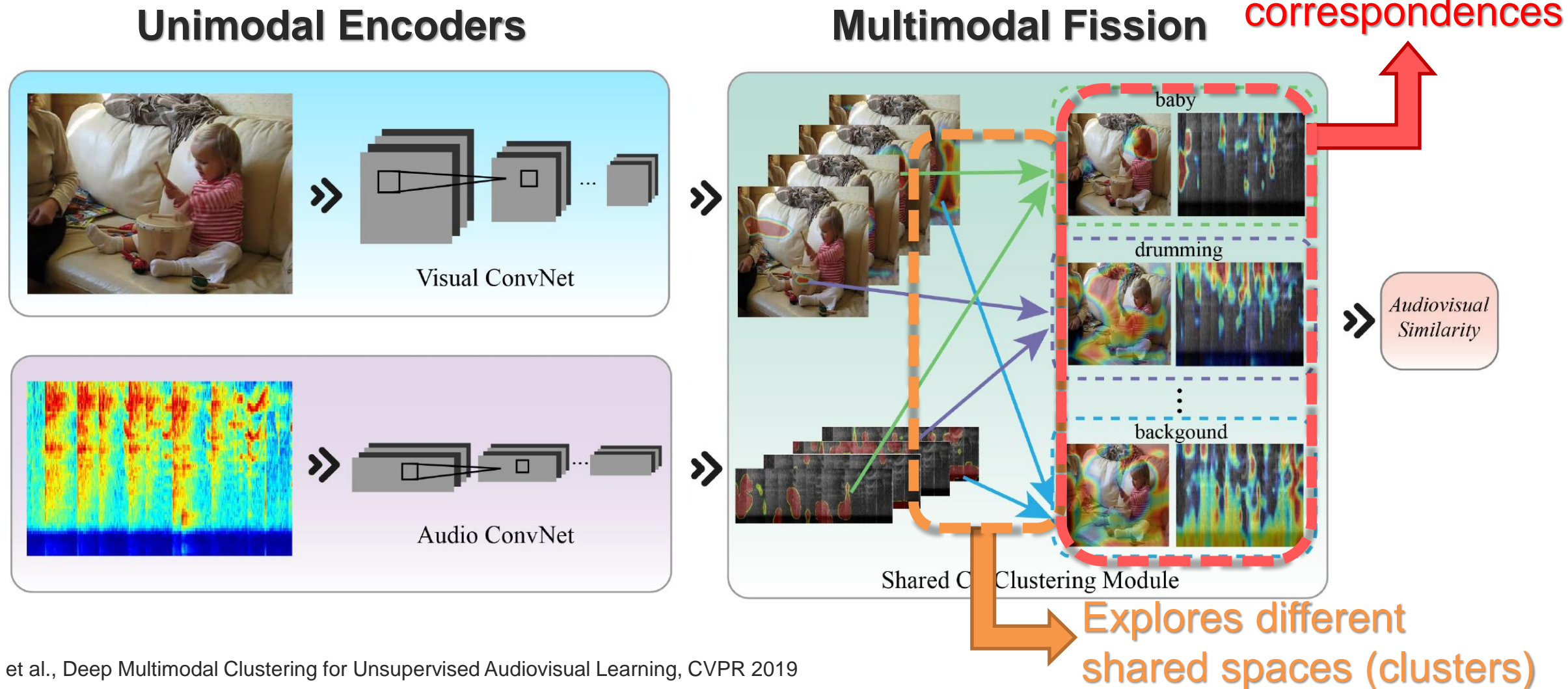
## Unimodal Encoders



Localized activations for different objects



# Fine-Grained Fission – A Clustering Approach



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

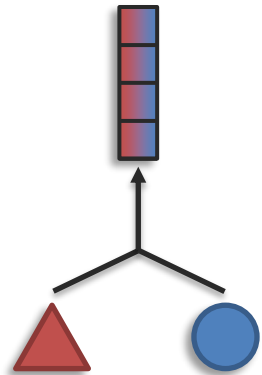
# Challenge 1: Representation

---

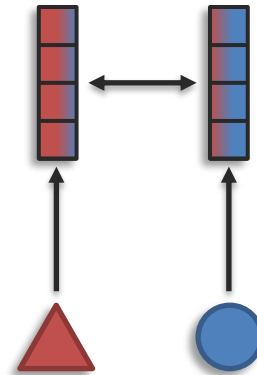
**Definition:** Learning representations that reflect cross-modal interactions between individual elements, across different modalities

## Sub-challenges:

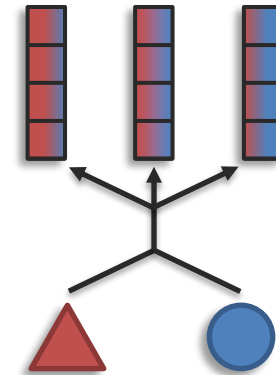
### Fusion



### Coordination



### Fission



# Challenge 2: Alignment

---

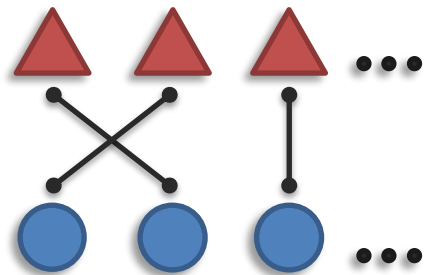
## Challenge 2: Alignment

---

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

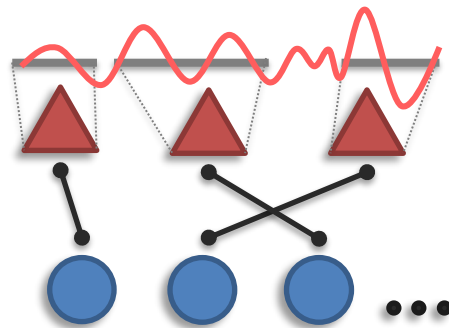
### Sub-challenges:

#### Discrete Alignment



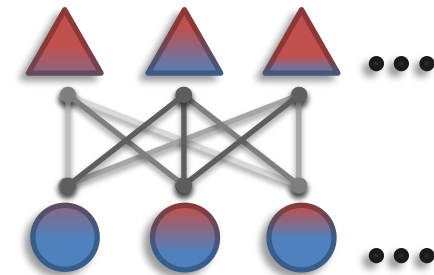
Discrete elements and connections

#### Continuous Alignment



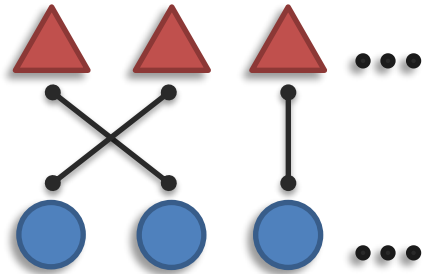
Segmentation and continuous warping

#### Contextualized Representation

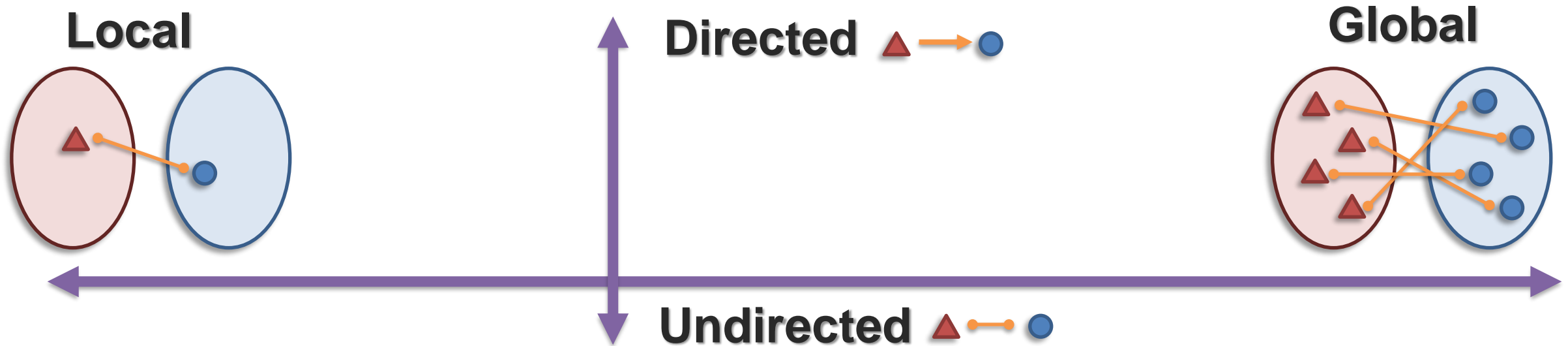


Alignment + representation

## Sub-Challenge 2a: Discrete Alignment

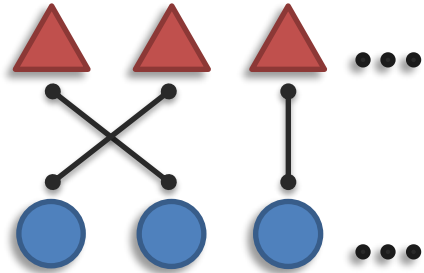


**Definition:** Identify and model connections between elements of multiple modalities





# Connections



Why should 2 elements be connected?

## Statistical



Association

Dependency



e.g., correlation,  
co-occurrence



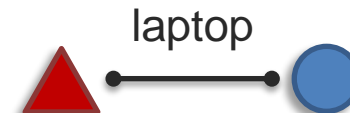
e.g., causal,  
temporal

## Semantic



Correspondence

Relationship

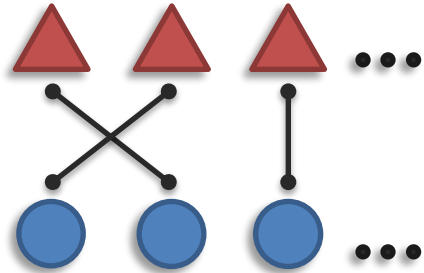


e.g., grounding



e.g., function

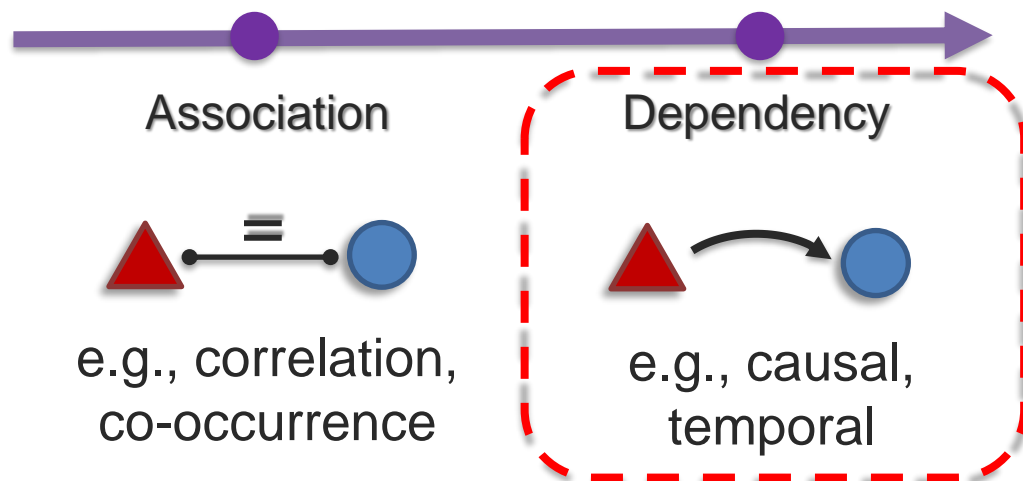
# Connections



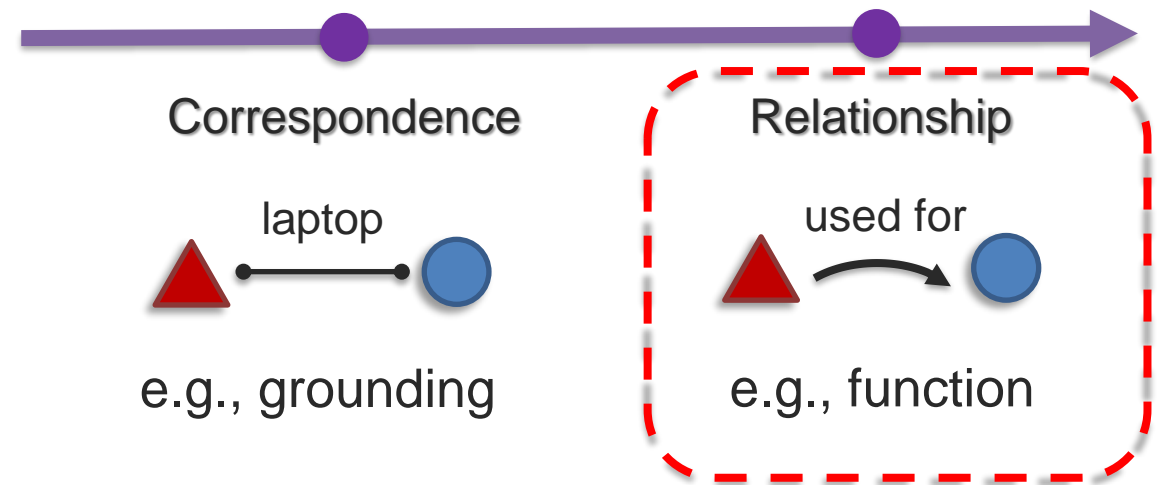
Why should 2 elements be connected?

Relationships and Dependencies will be discussed in more details in **Reasoning** challenge

## Statistical



## Semantic



# Language Grounding

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



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

## Statistical



Association

Dependency



e.g., correlation,  
co-occurrence



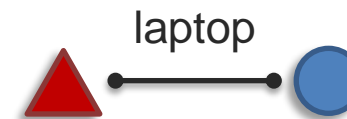
e.g., causal,  
temporal

## Semantic

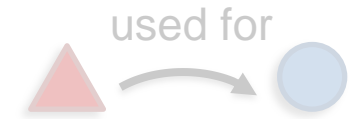


Correspondence

Relationship

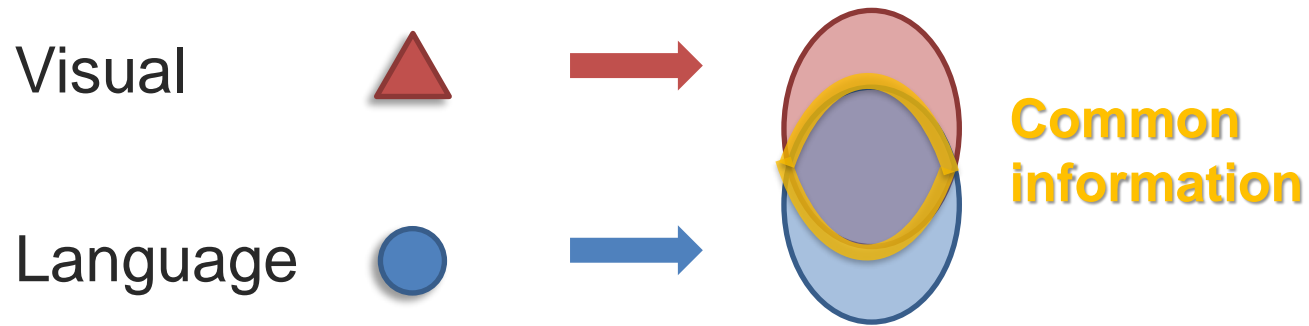


e.g., grounding



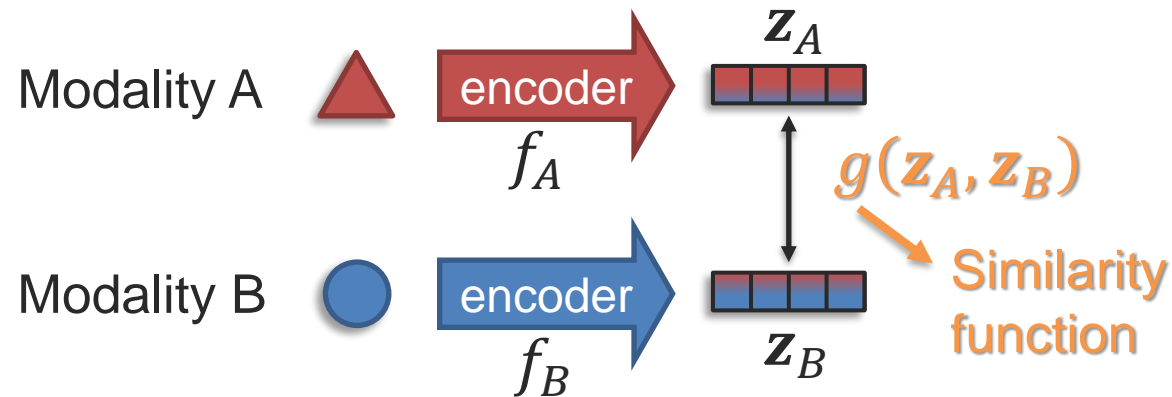
e.g., function

# Local Alignment – Coordinated Representations



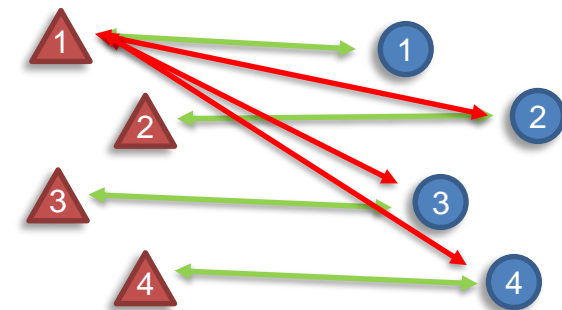
A **woman** reading **newspaper**

Learning coordinated representations:

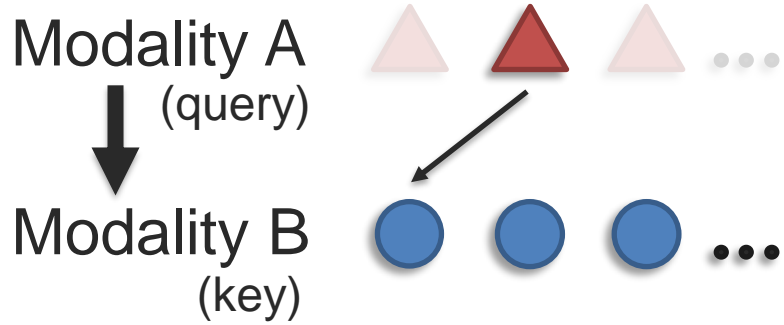


or contrastive learning

Supervision: Paired data



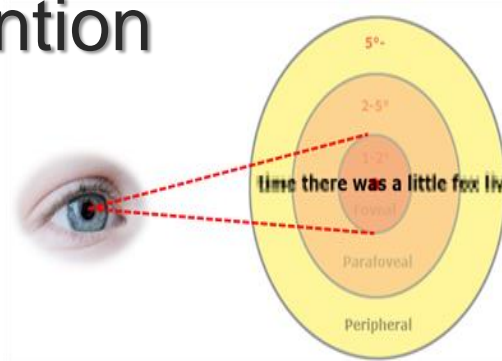
# Directed Alignment



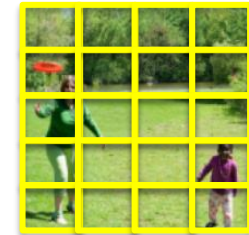
A woman is throwing a frisbee

Which object?

## Attention



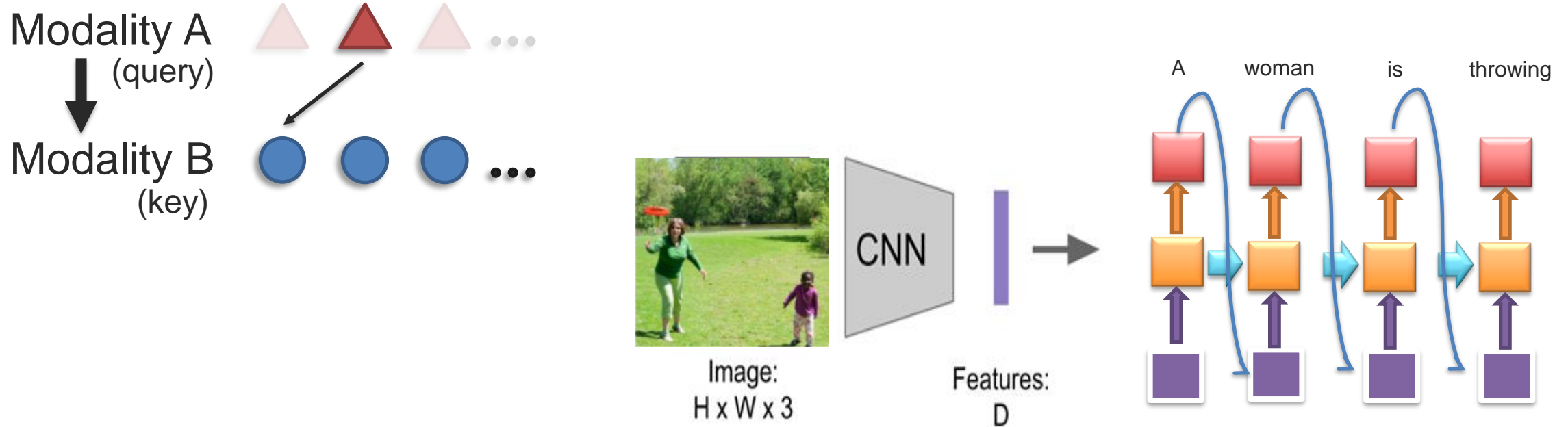
1 Soft attention



2 Hard attention

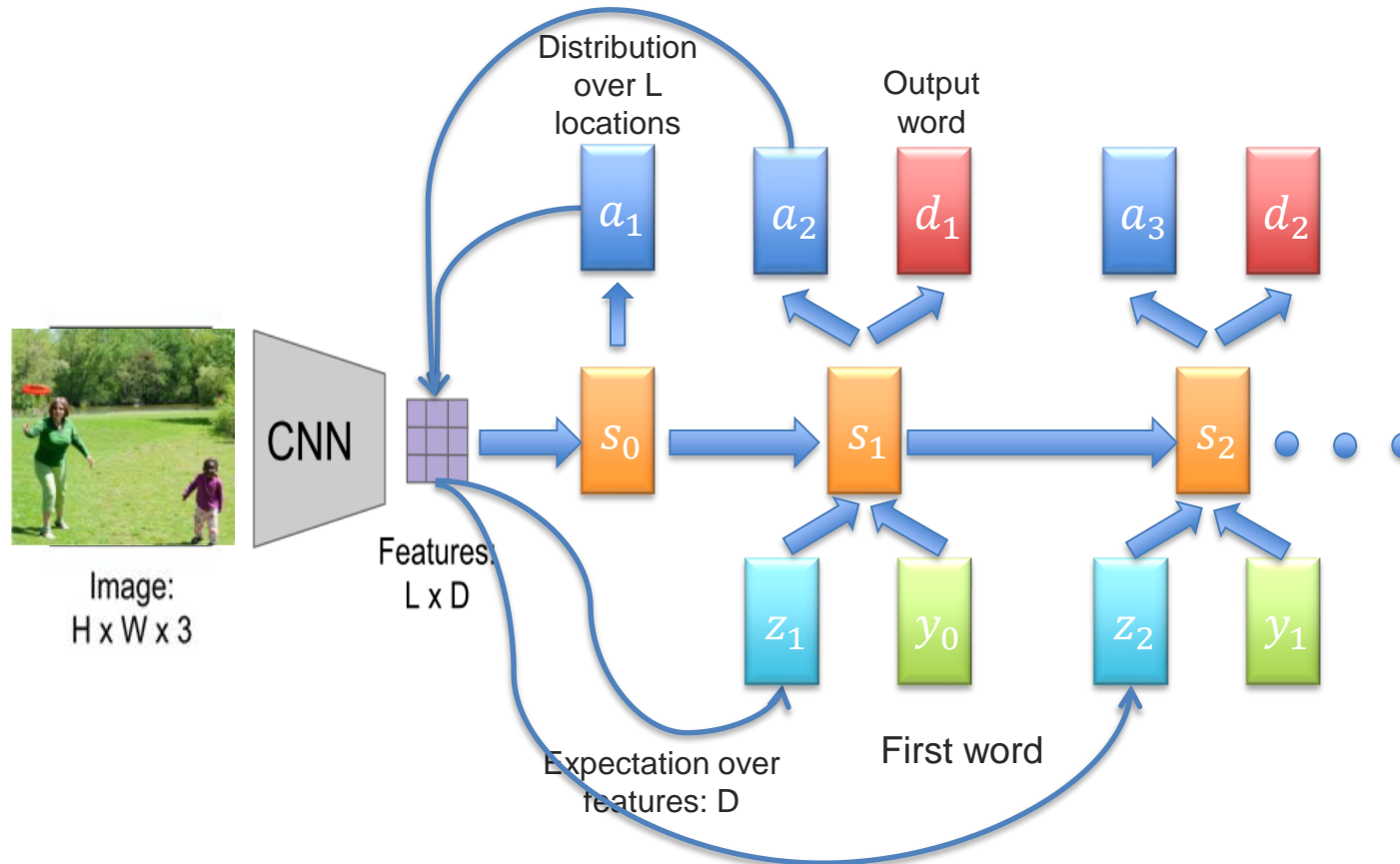


# Directed Alignment – Image Captioning



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)

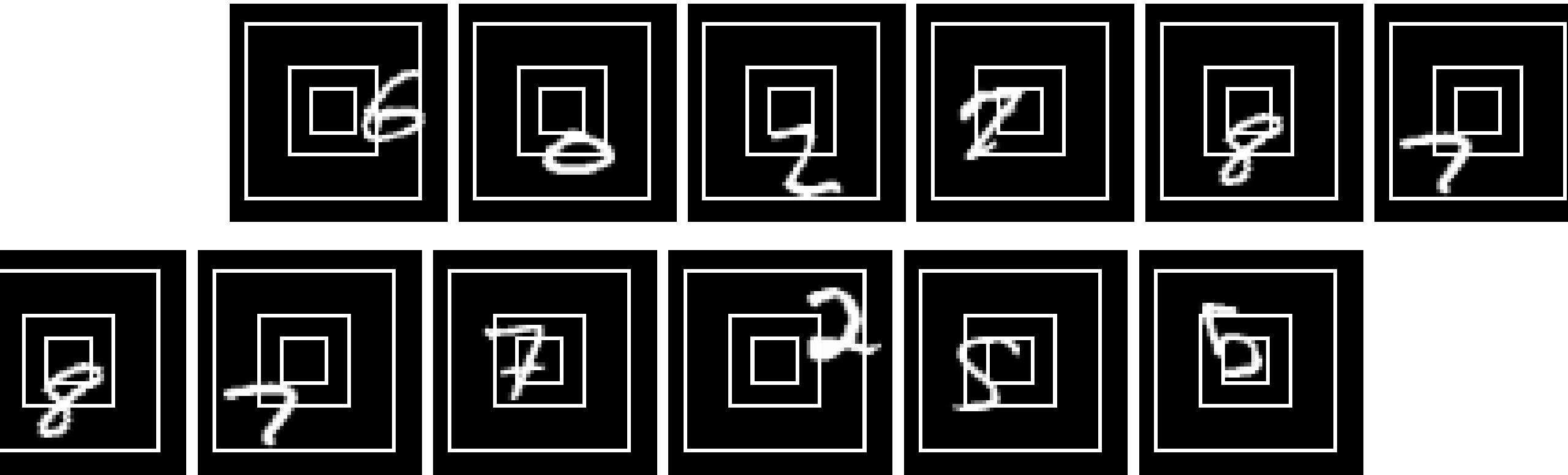
# Example – Image Captioning



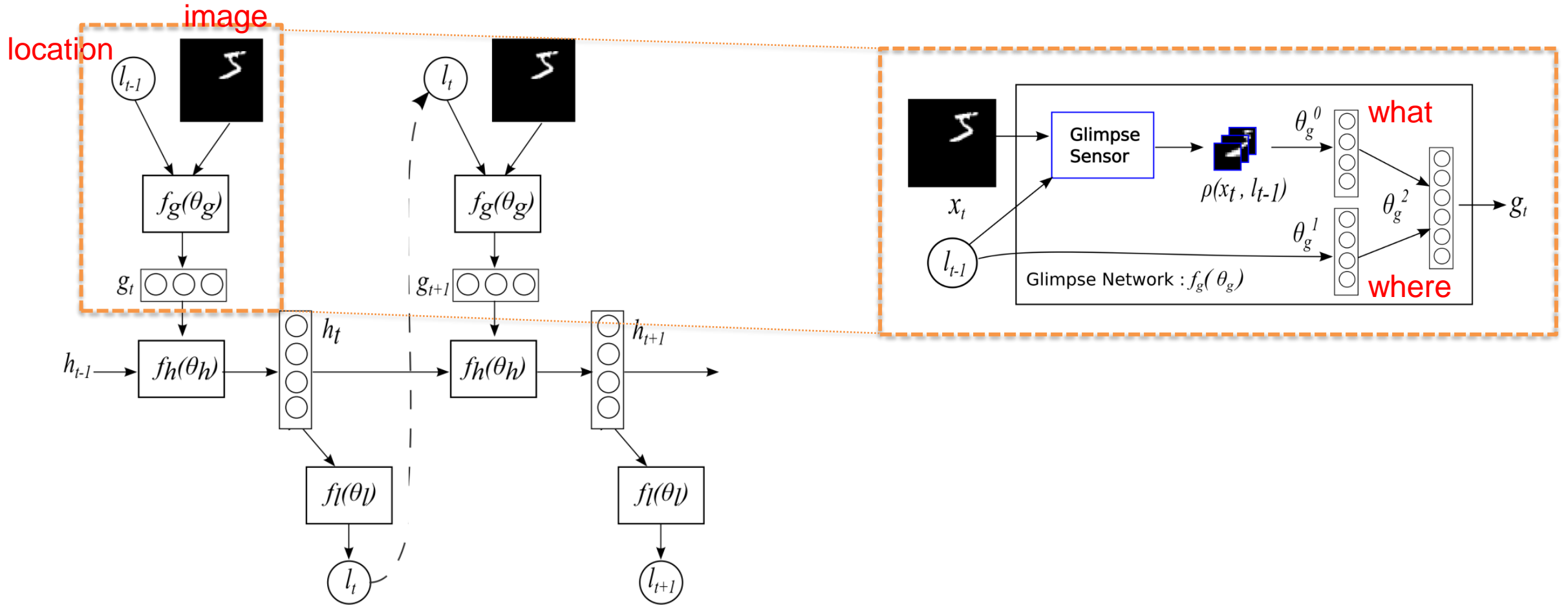
[Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Xu et al., 2015]

# Hard attention - Example

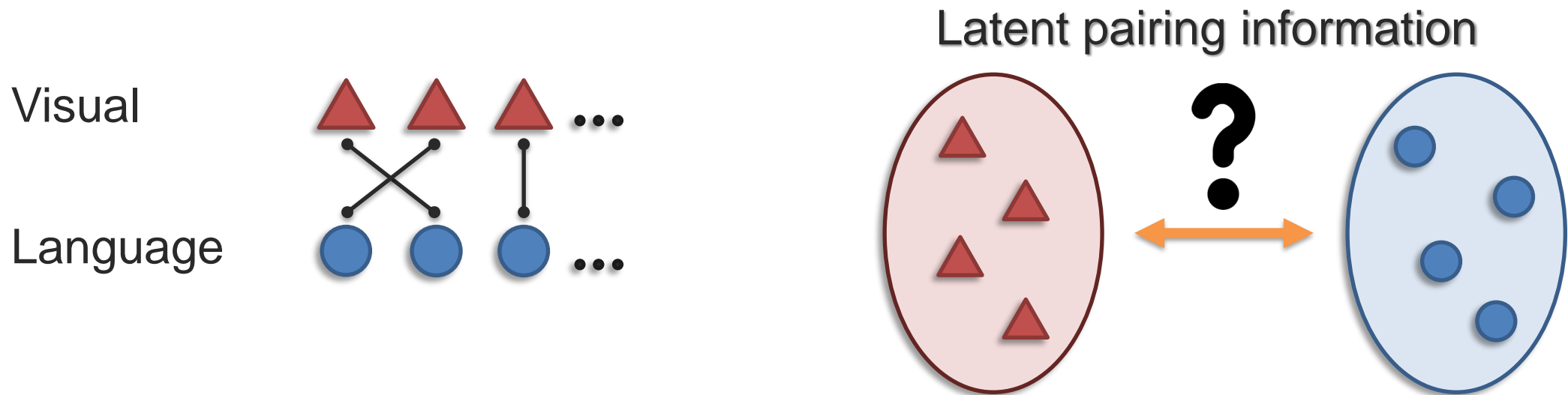
---



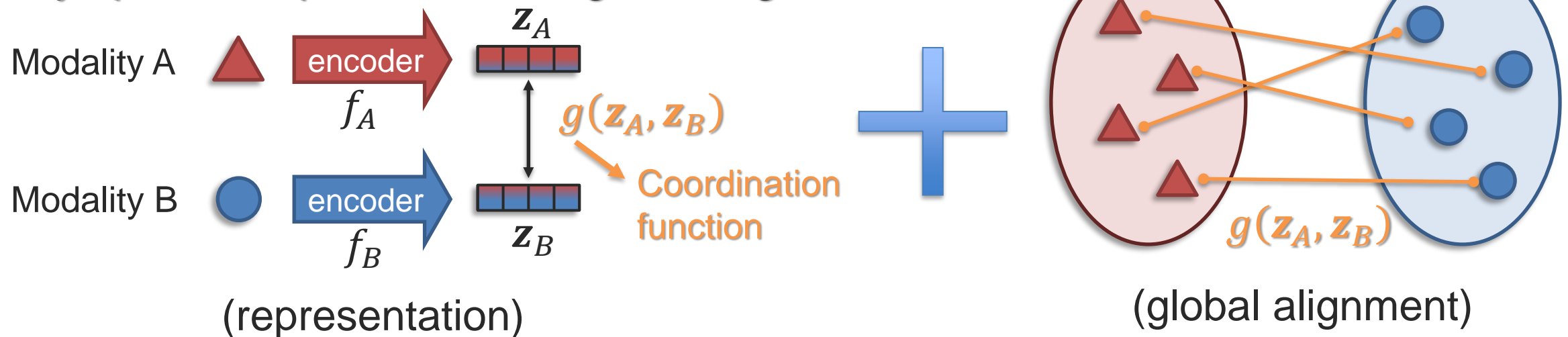
# Hard Attention – Recurrent Model of Visual Attention



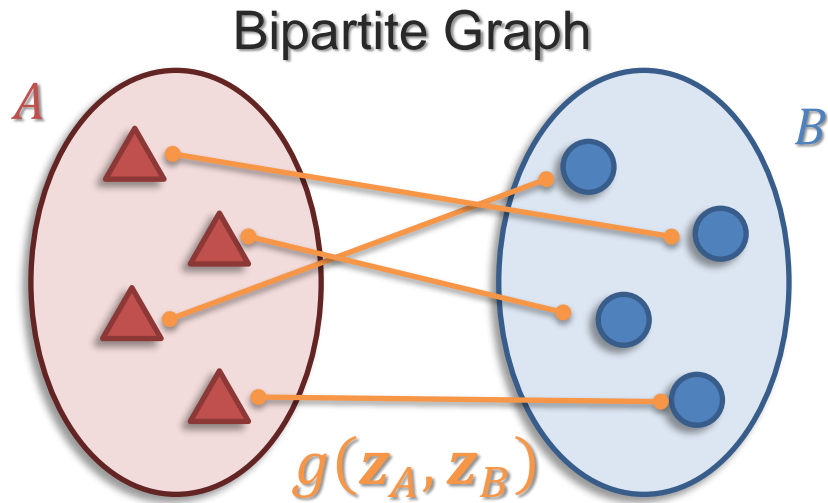
# Global Alignment



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:  ~~$w_{(i,f(i))} = g(\mathbf{z}_A^i, \mathbf{z}_B^{f(i)})$~~

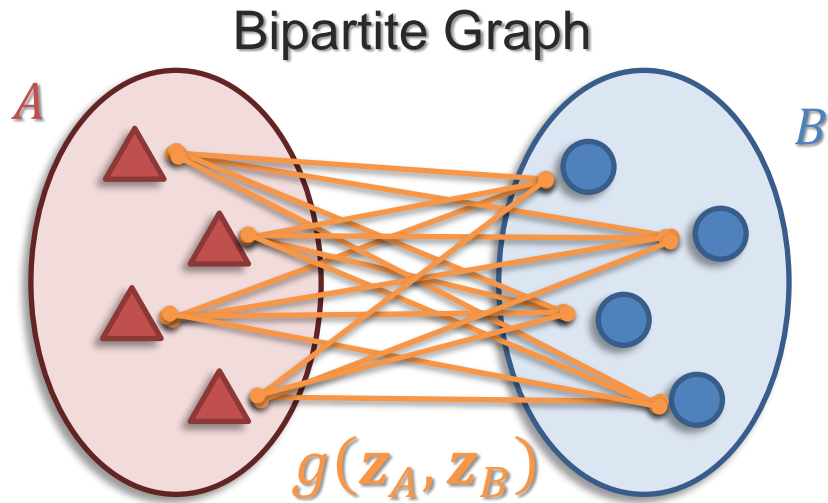
$w_{(i,j)} = g(\mathbf{z}_A^i, \mathbf{z}_B^j)$

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 give optimal transport

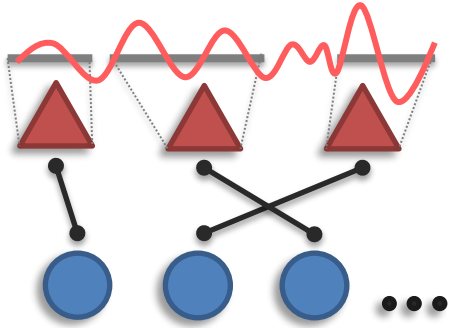
# Continuous Alignment





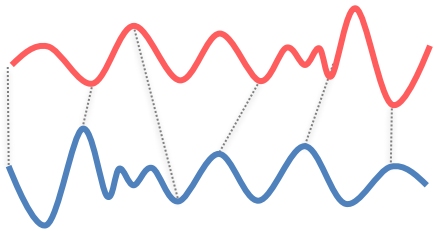
## Challenge 2b: Continuous Alignment

---

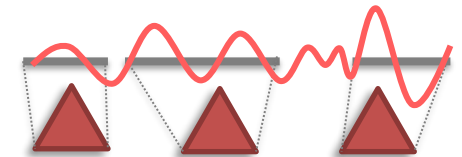


**Definition:** Model alignment between modalities with continuous signals and no explicit elements

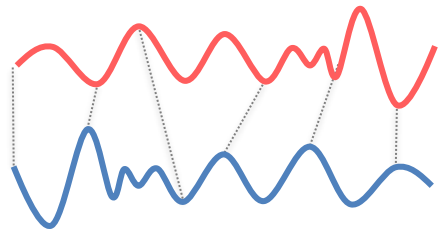
Continuous  
warping



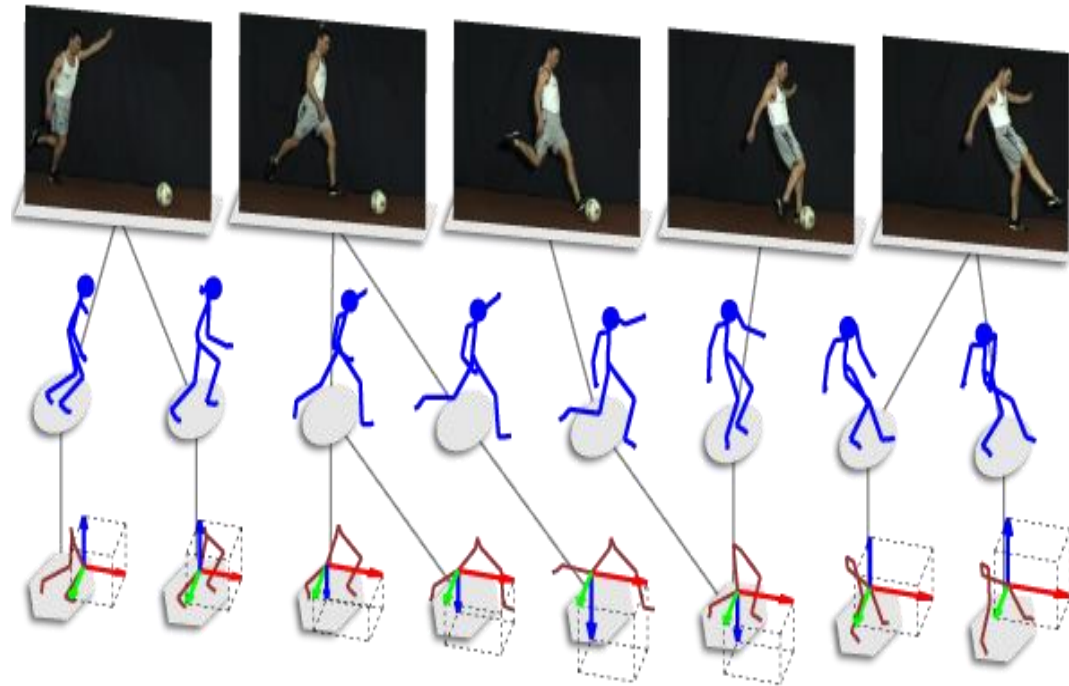
Discretization  
(segmentation)



# Continuous Warping – Example



➔ Aligning video sequences



# Dynamic Time Warping (DTW)

We have two unaligned temporal unimodal signals

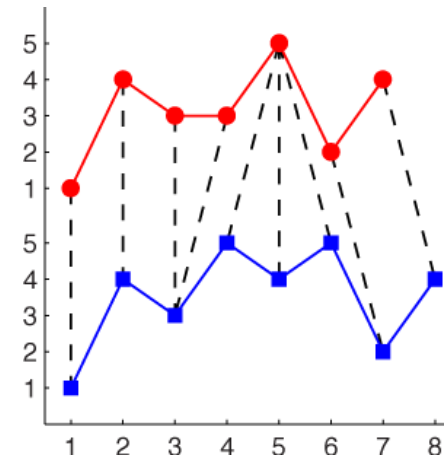
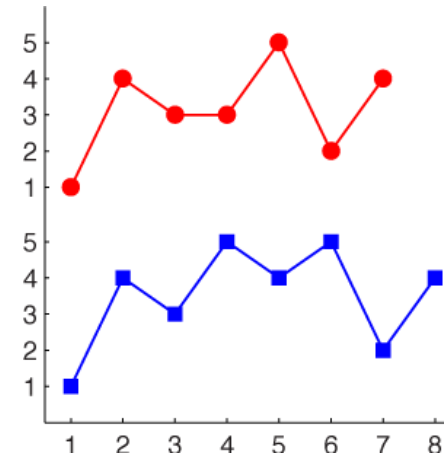
- $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_x}] \in \mathbb{R}^{d \times n_x}$
- $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{n_y}] \in \mathbb{R}^{d \times n_y}$

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!

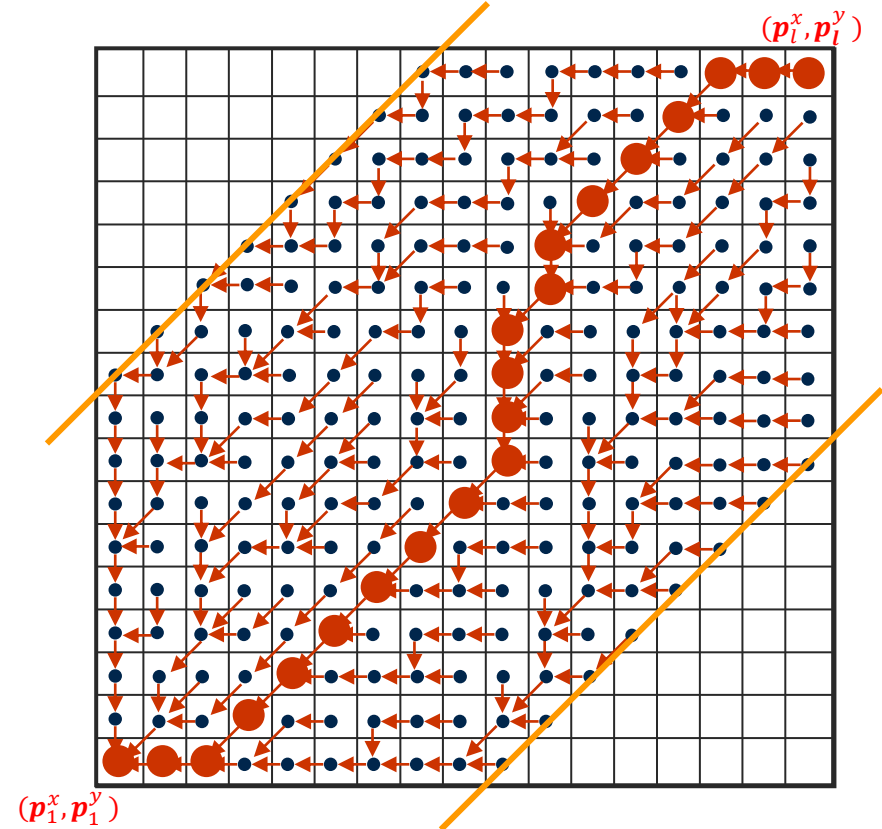


# Dynamic Time Warping (DTW)

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

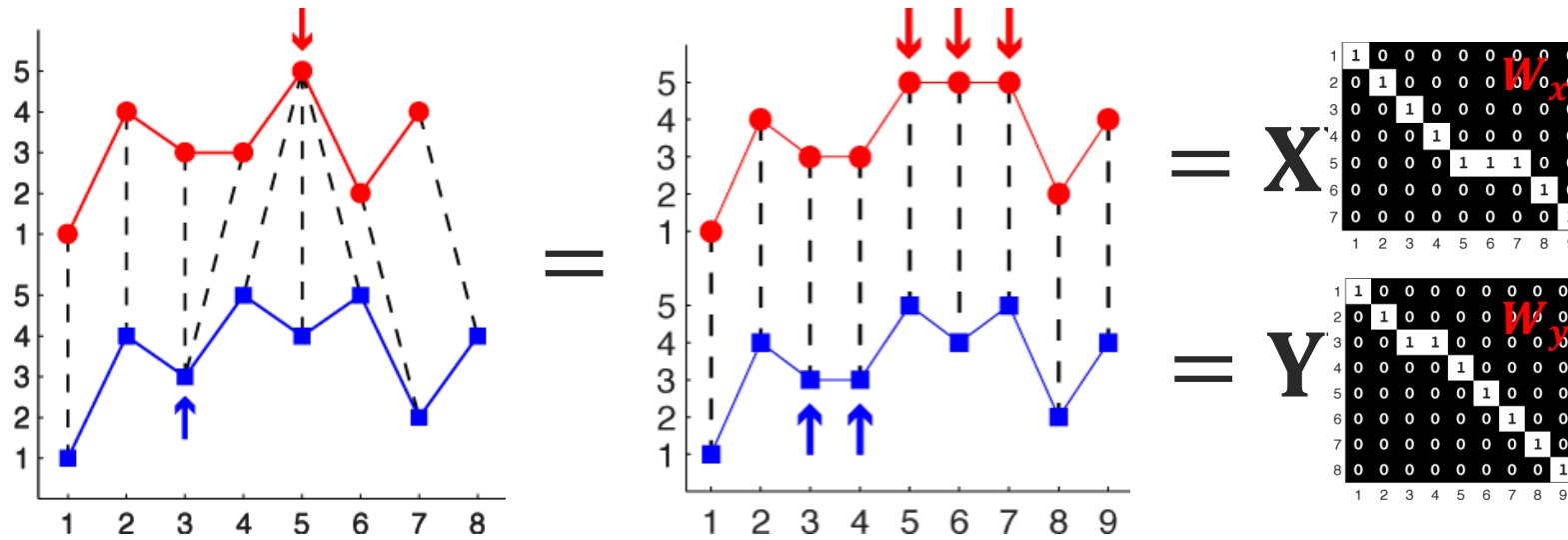
Solved using dynamic programming while respecting the restrictions



# 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:

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

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

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

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

A differentiable version of DTW also exists...

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

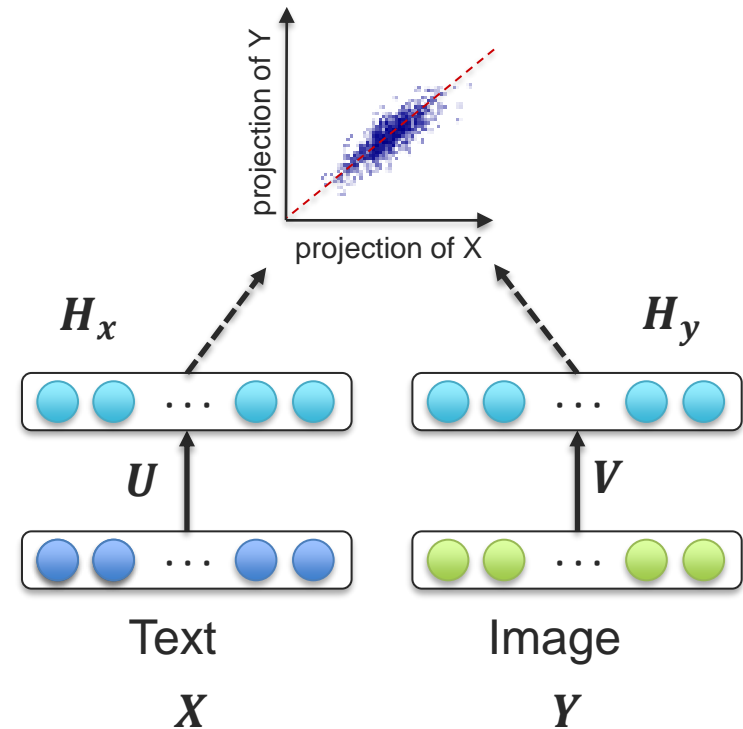
# Canonical Correlation Analysis – Reminder

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

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

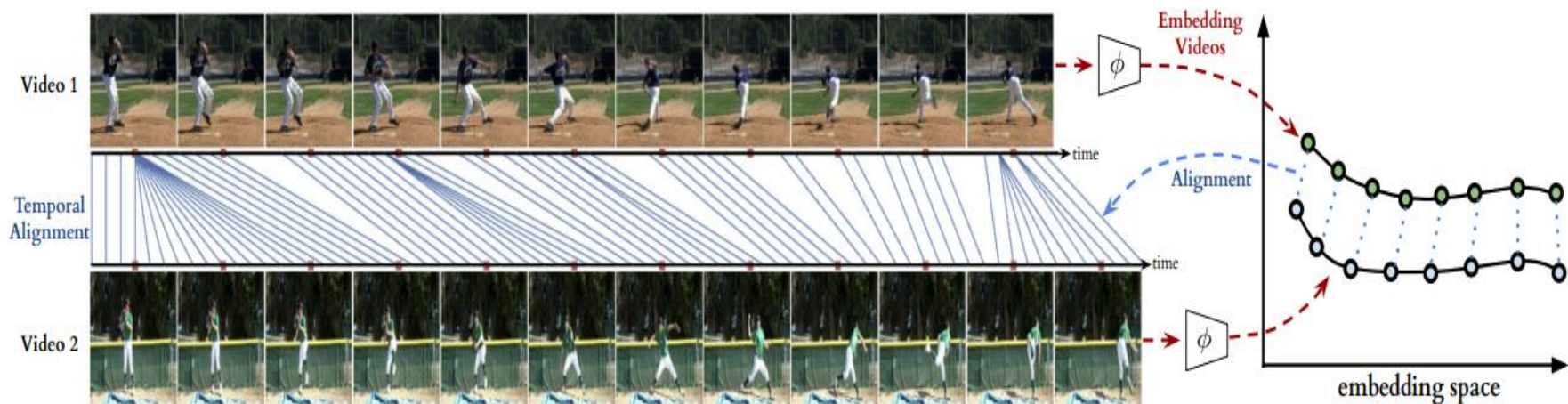
Allows to align multi-modal or multi-view (same modality but from a different point of view)

- $\mathbf{W}_x, \mathbf{W}_y$  – 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

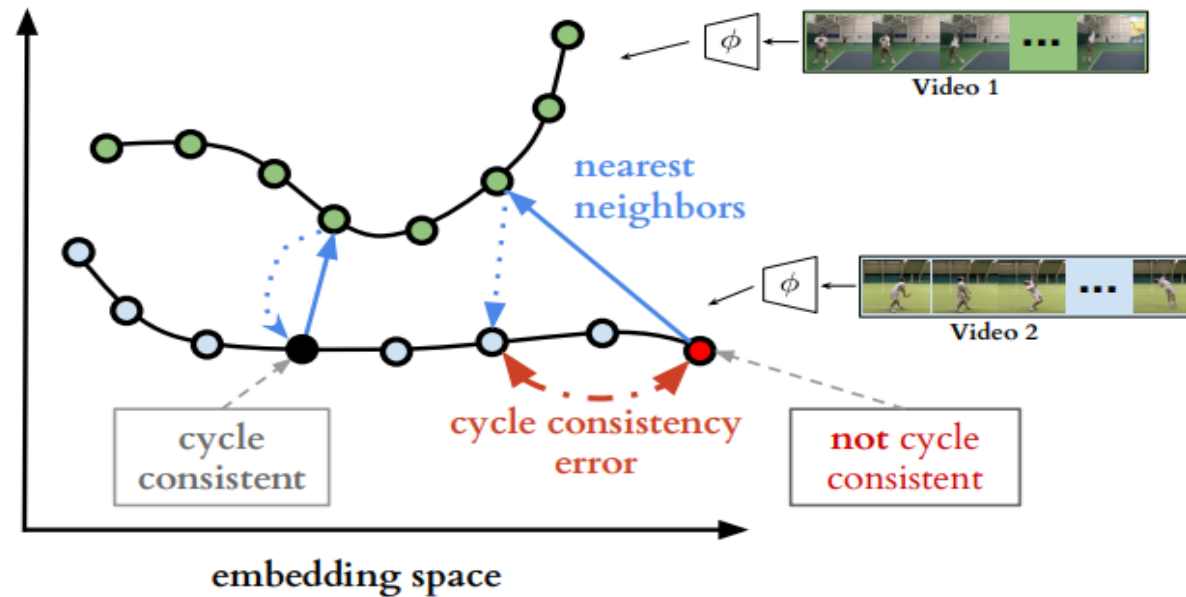


How can we define a loss function to enforce the alignment between sequences while at the same time learning good representations?



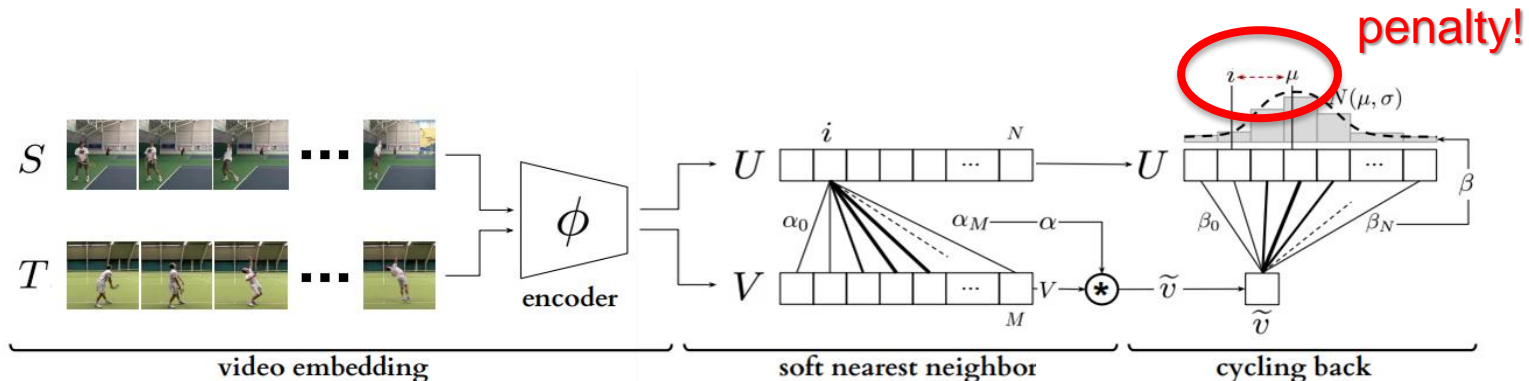
# Temporal Cycle-Consistency Learning

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:

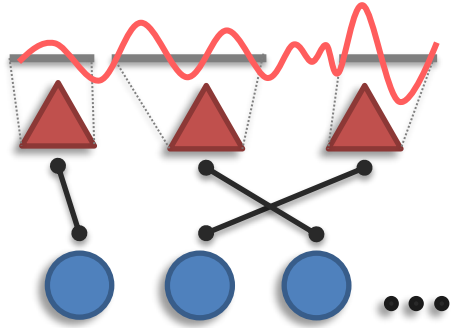
distances:  $\alpha_j = \frac{e^{-\|u_i - v_j\|^2}}{\sum_k^M e^{-\|u_i - v_k\|^2}}$       Soft nearest neighbor:  $\tilde{v} = \sum_j^M \alpha_j v_j$

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

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

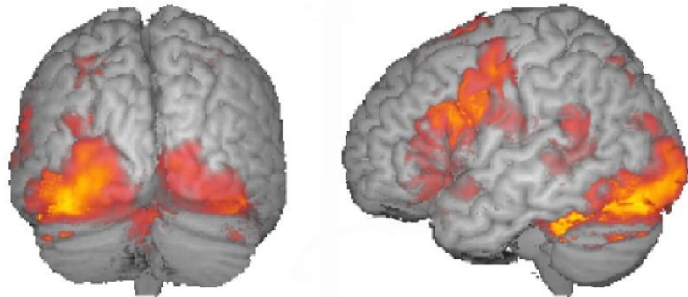
# Discretization (aka Segmentation)

---

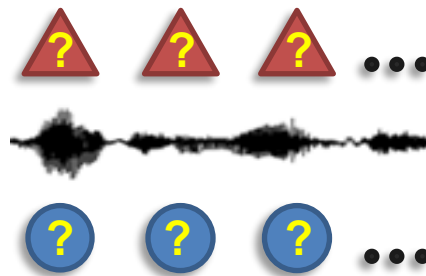


Common assumptions: ① Segmented elements

Examples:



Medical imaging



Signals



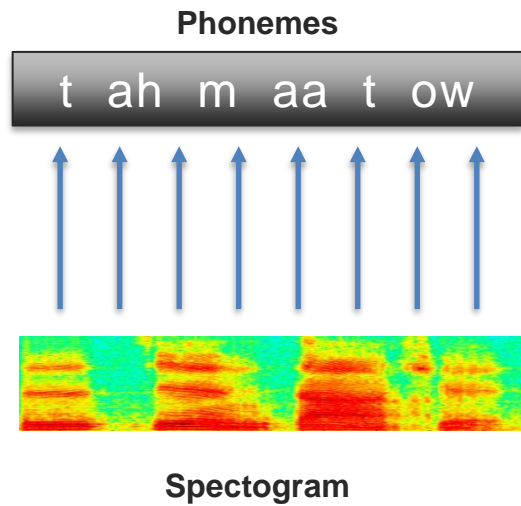
Images

objects

# Discretization – Example

---

## Sequence Labeling and Alignment

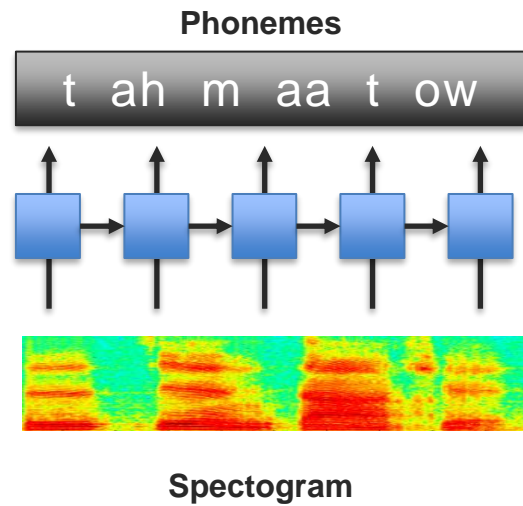


How can we predict the sequence  
of phoneme labels?

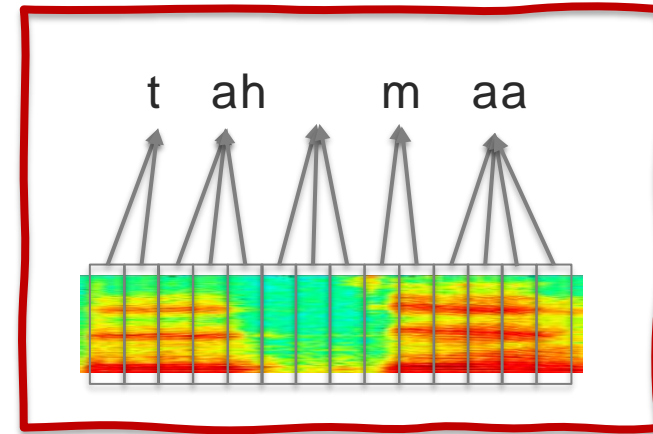
# Discretization – Example

---

## Sequence Labeling and Alignment



Challenge: many-to-1 alignment

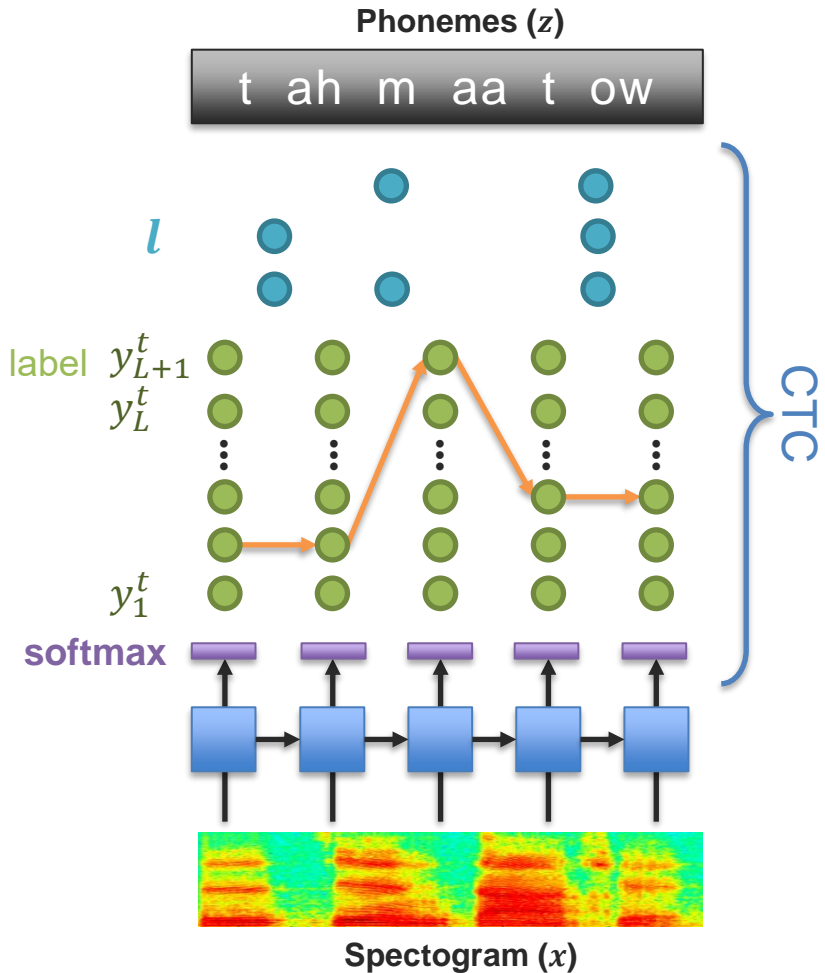


How can we predict the sequence of phoneme labels?

# Discretization – A Classification Approach

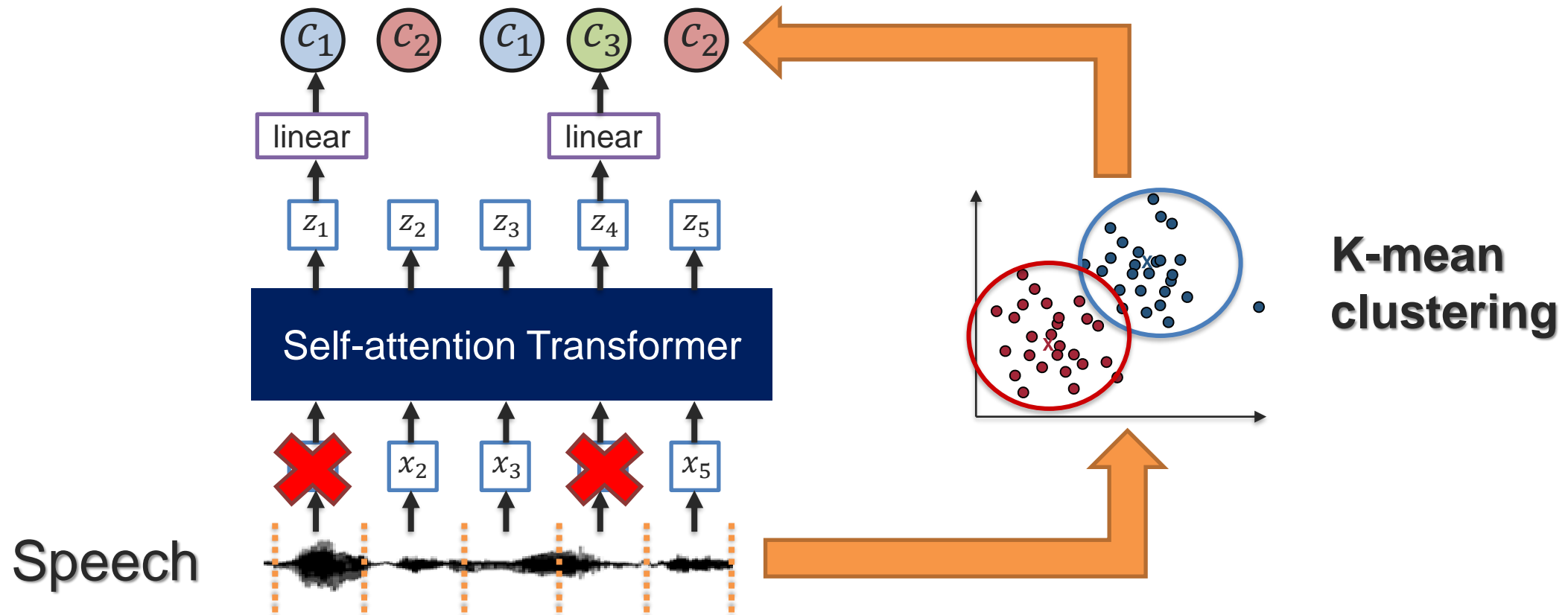
## Connectionist Temporal Classification

- 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



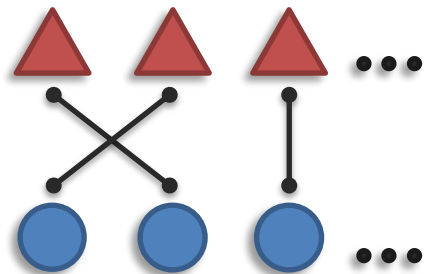
## Challenge 2: Alignment

---

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

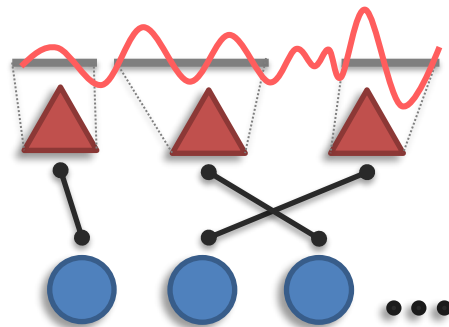
### Sub-challenges:

#### Discrete Alignment



Discrete elements  
and connections

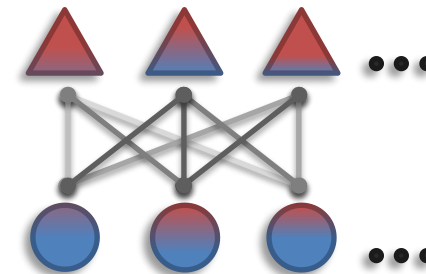
#### Continuous Alignment



Segmentation and  
continuous warping

Next week!

#### Contextualized Representation



Implicit alignment  
+ representation