

Multimodal AI

Lecture 4.2 – Multimodal Alignment

Paul Liang

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

ppliang@mit.edu

 [@pliang279](https://twitter.com/pliang279)



Assignments for This Coming Week

I want to meet every group at least once regarding their project ideas. Project mentors will be assigned.

Compute credits: 40 x \$50 Kimi credits, 40 x \$40 other credits.

HW2 due next Wednesday 3/4.

Today's lecture

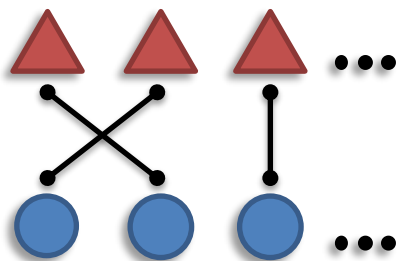
- 1 Multimodal alignment
- 2 Explicit alignment and contrastive learning
- 3 Continuous alignment
- 4 Implicit (emergent) alignment

Challenge 2: Alignment

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

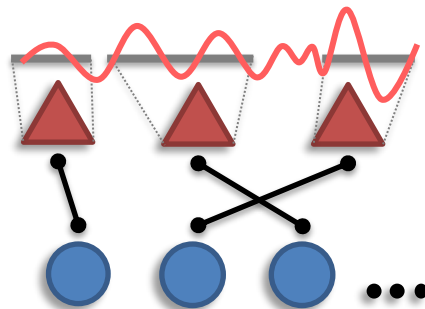
Sub-challenges:

Discrete connections



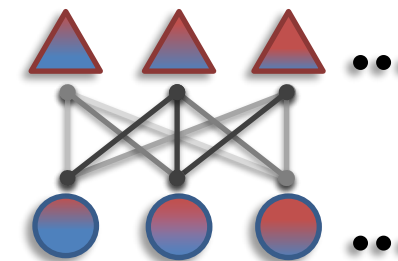
Explicit alignment
(e.g., grounding)

Continuous alignment



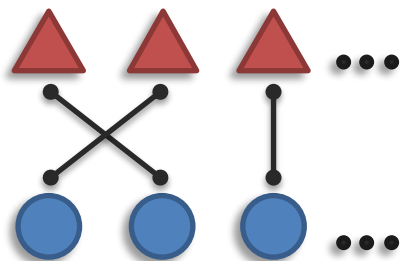
Granularity of individual elements

Contextualized representation



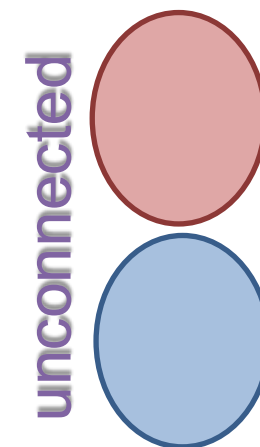
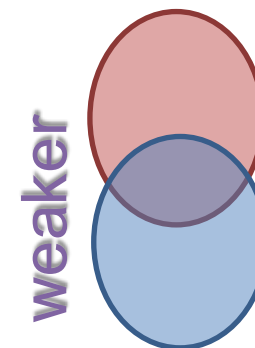
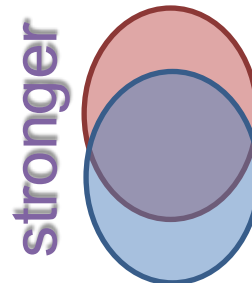
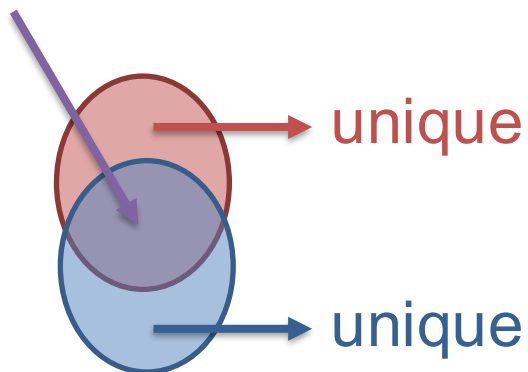
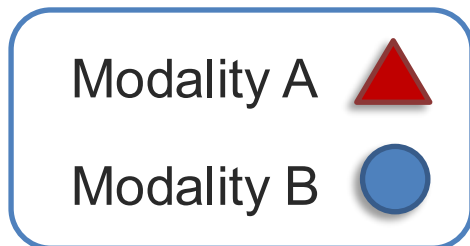
Implicit alignment
+ representation

Challenge 2a: Discrete Alignment

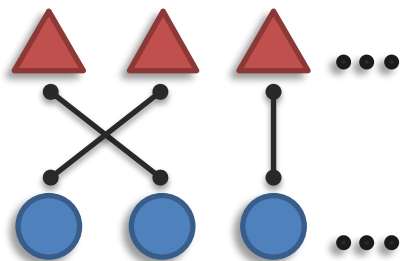


Definition: Identify and model connections between elements of multiple modalities

Shared information that relates modalities



Modality Connections



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

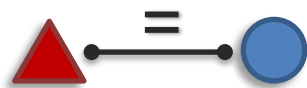


A **woman** reading **newspaper**

Statistical



Association



e.g., correlation,
co-occurrence

Dependency

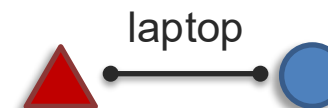


e.g., causal,
temporal

Semantic



Correspondence



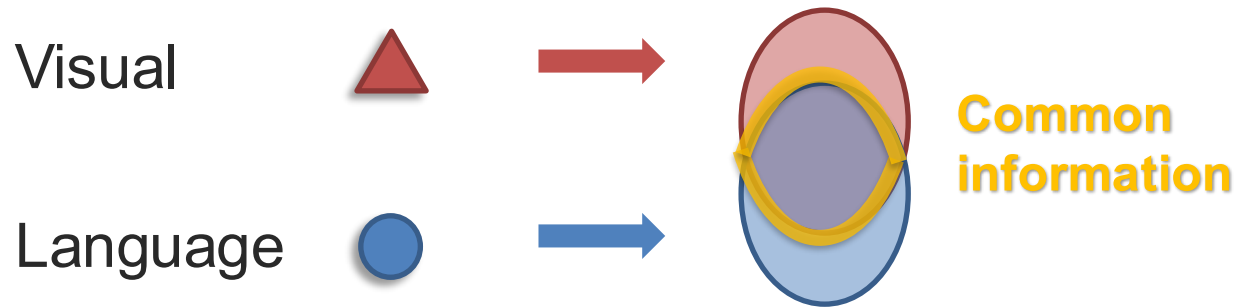
e.g., grounding

Relationship



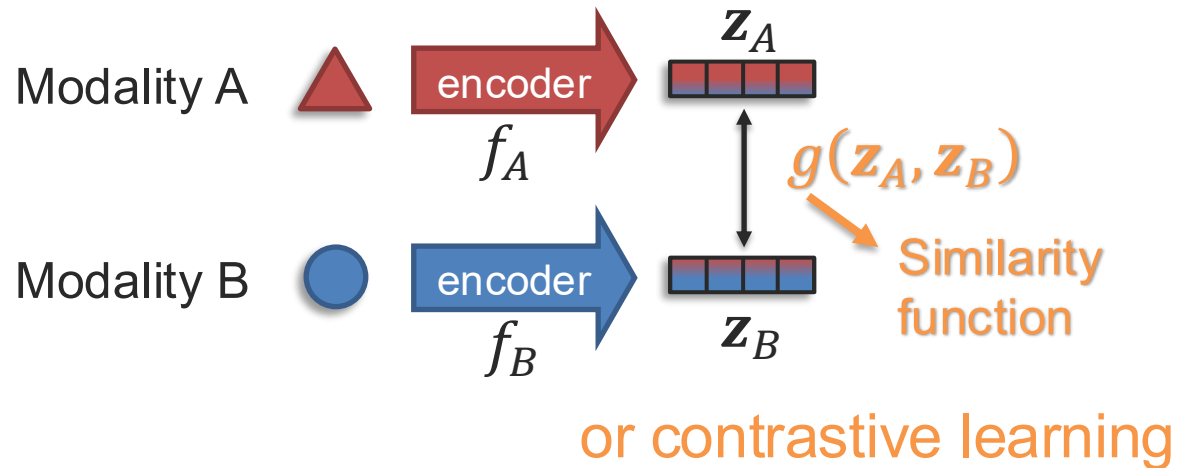
e.g., function

Modality Connections

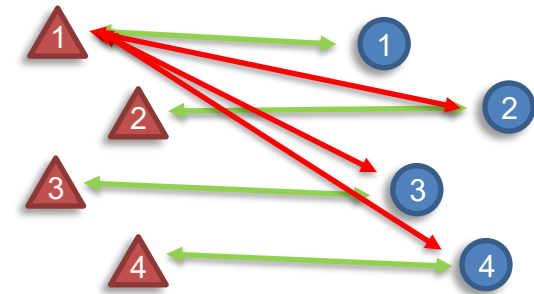


A **woman** reading **newspaper**

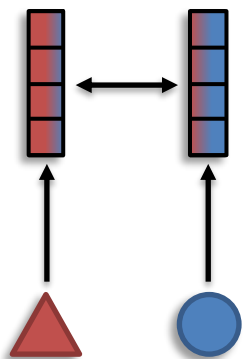
Learning aligned representations:



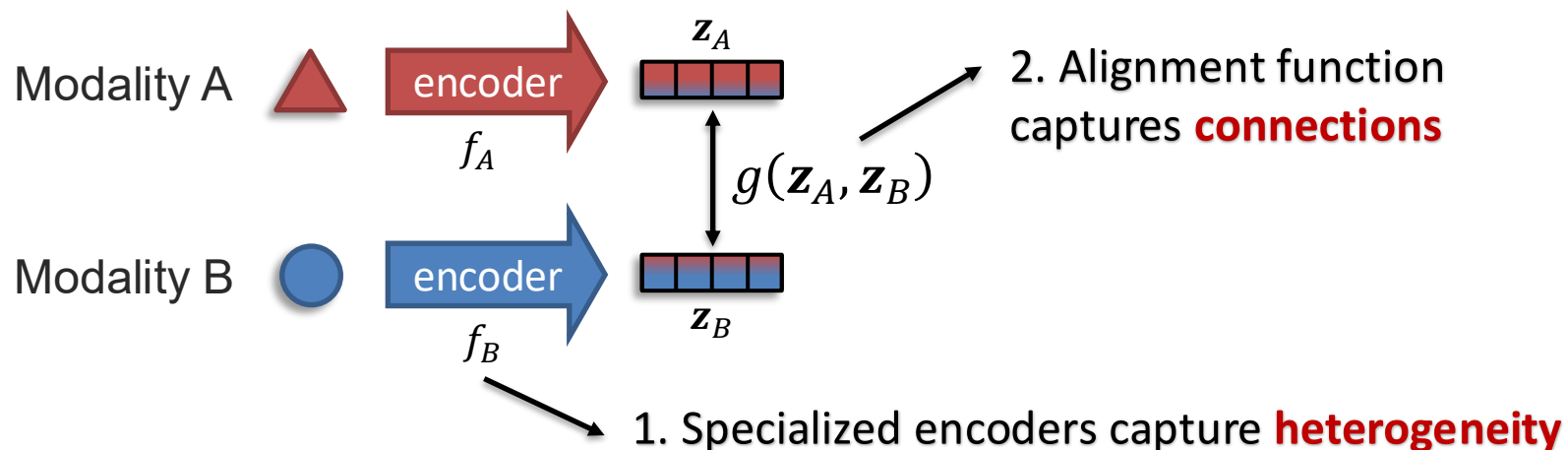
Supervision: Paired data



Aligned Representations



Definition: Learn multimodal representations aligned through their connections.

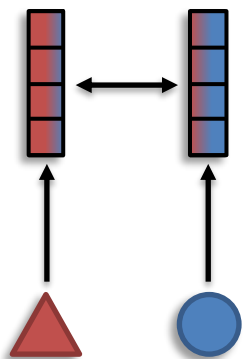


Learning with alignment function:

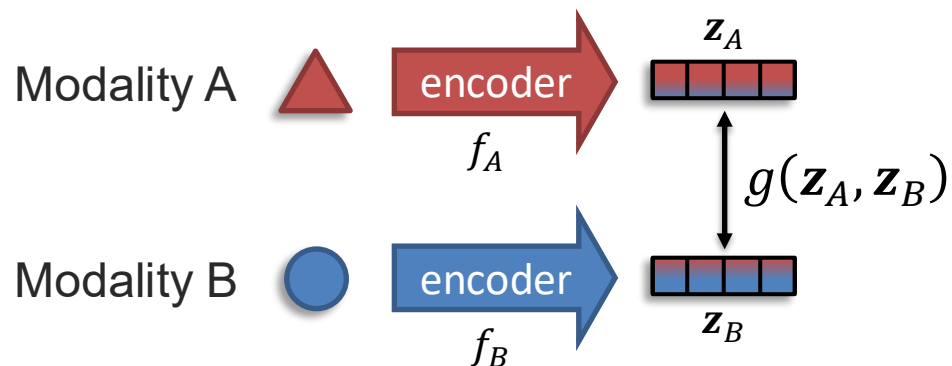
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

Aligned Representations



Definition: Learn multimodal representations aligned through their connections.



Learning with alignment function:

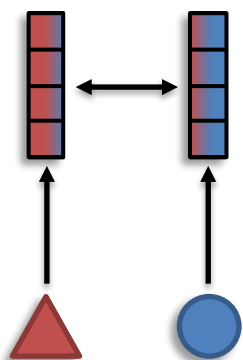
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

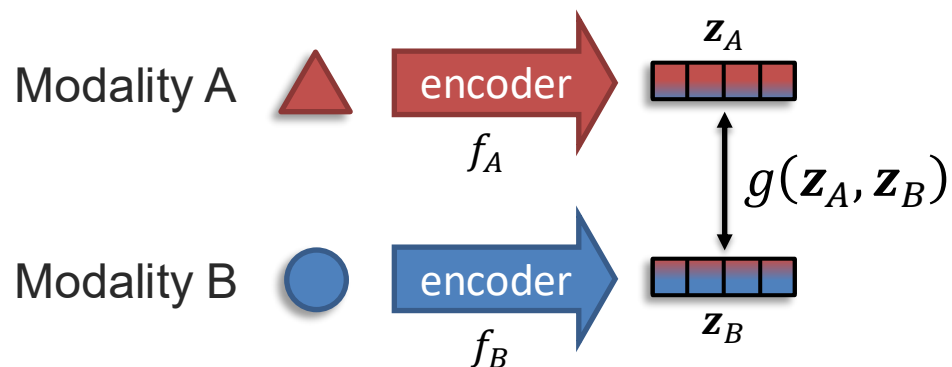
① Cosine similarity:

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

Aligned Representations



Definition: Learn multimodal representations aligned through their connections.



Learning with alignment function:

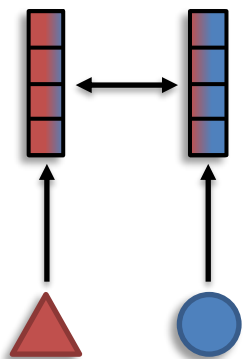
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

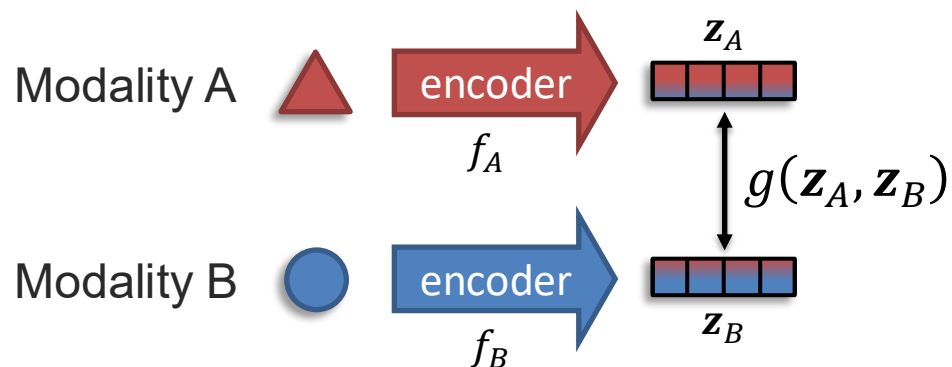
② Kernel similarity functions:

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

Aligned Representations



Definition: Learn multimodal representations aligned through their connections.



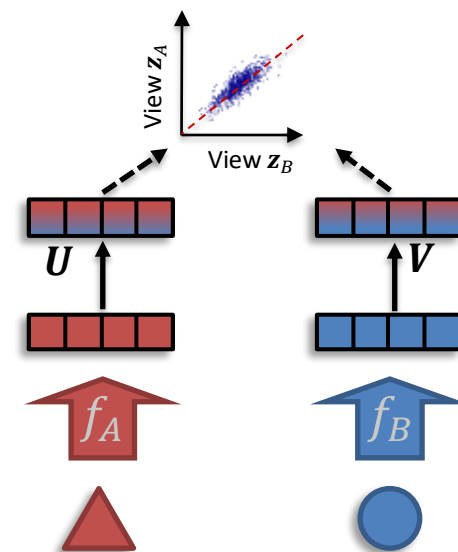
Learning with alignment function:

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

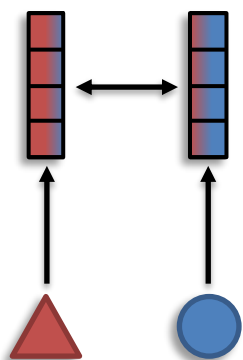
with model parameters θ_g , θ_{f_A} and θ_{f_B}

③ Canonical Correlation Analysis (CCA):

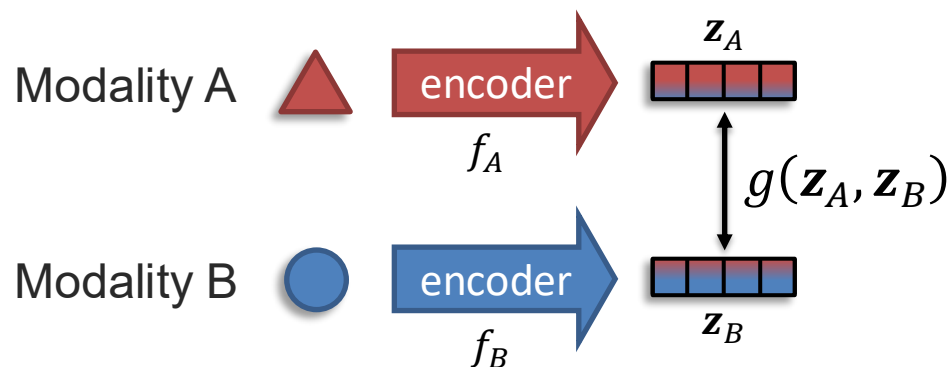
$$\operatorname{argmax}_{V, U, f_A, f_B} \operatorname{corr}(z_A, z_B)$$



Aligned Representations



Definition: Learn multimodal representations aligned through their connections.

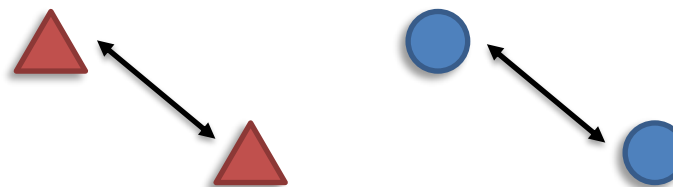


Learning with alignment function:

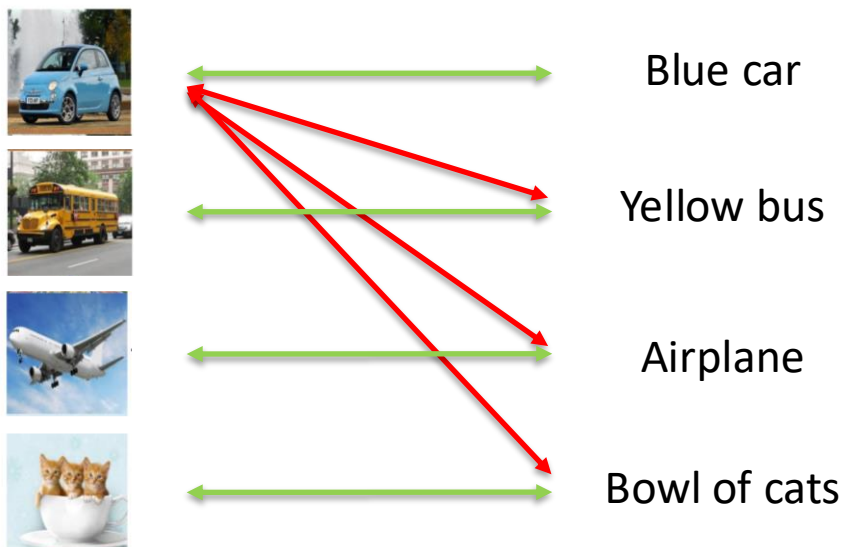
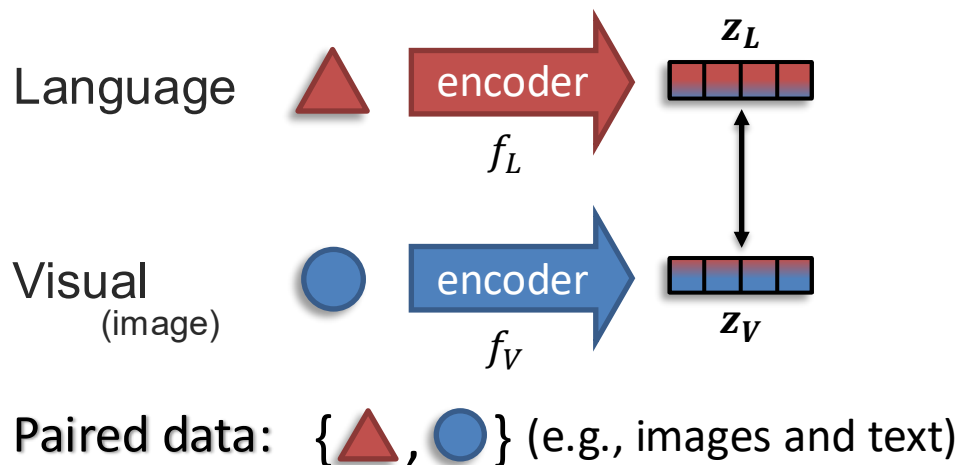
$$\mathcal{L} = g(f_A(\triangle), f_B(\circ))$$

with model parameters θ_g , θ_{f_A} and θ_{f_B}

④ Order, hierarchy, pairwise relationships.



Alignment with Contrastive Learning



Contrastive loss:

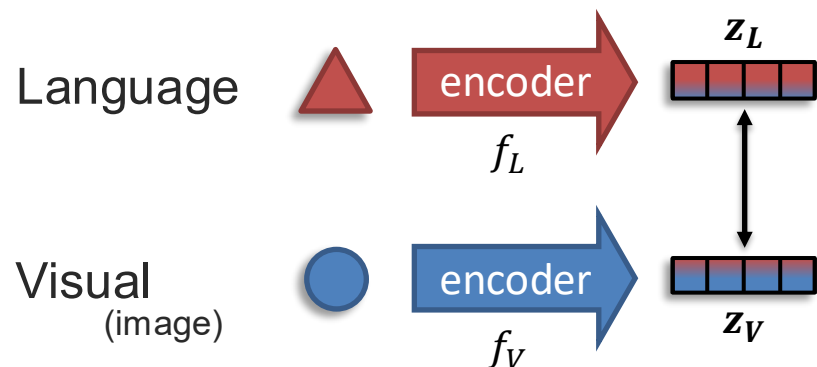
 brings **positive pairs** closer and pushes **negative pairs** apart

Simple contrastive loss:

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

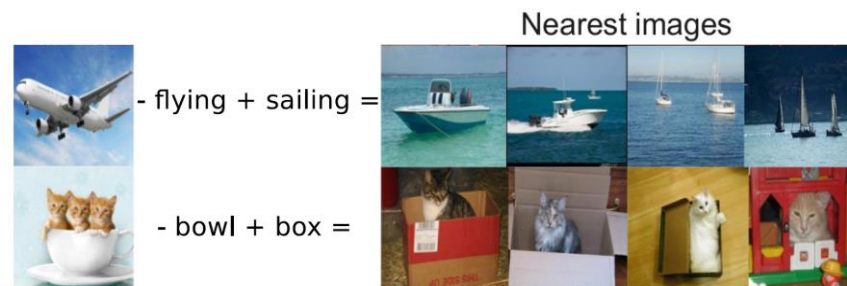
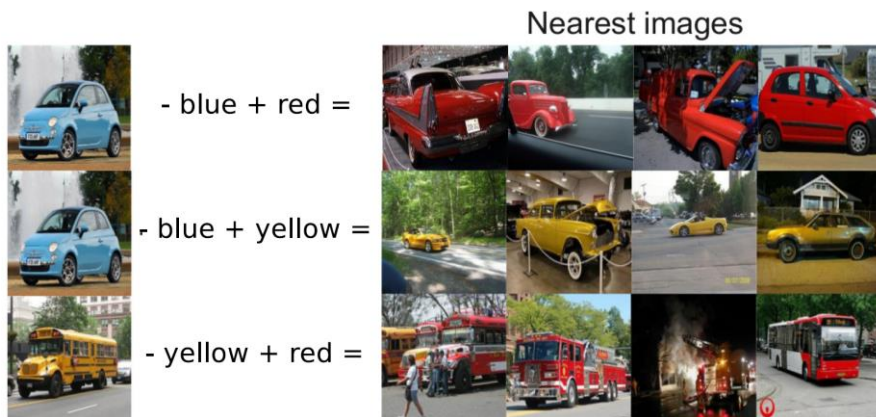
Coordination function (e.g., cosine similarity)

Visual-Semantic Embeddings

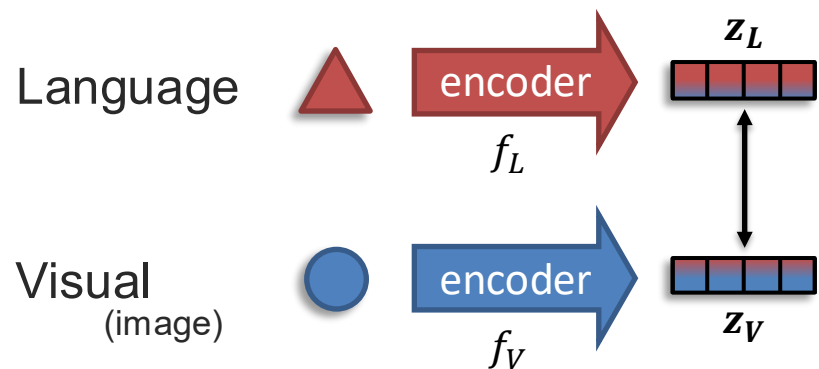


Contrastive loss:

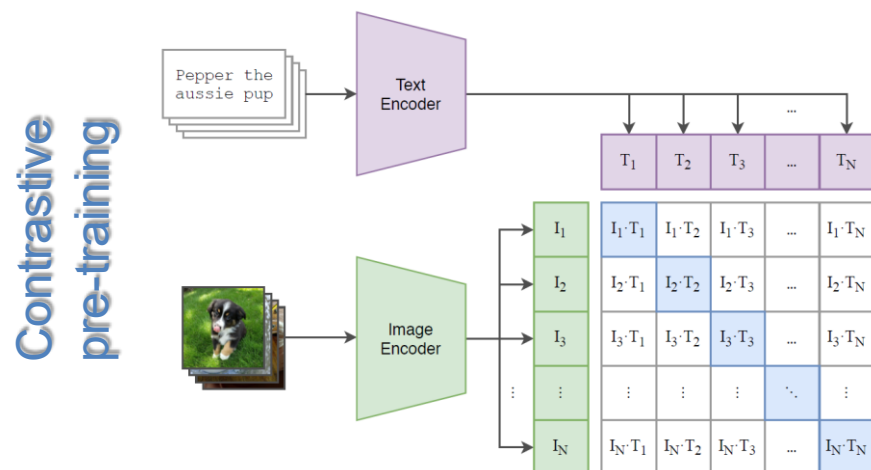
\rightarrow brings **positive pairs** closer and pushes **negative pairs** apart



Contrastive Language Image Pretraining



Positive and negative pairs:



Popular contrastive loss: InfoNCE

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

positive pairs

negative pairs and positive pairs

Similarity function can be cosine similarity

\Rightarrow CLIP encoders (f_L and f_V) are great for language-vision tasks

\Rightarrow z_L and z_V are coordinated but not identical representation spaces

CLIP (Contrastive Language–Image Pre-training)

SUN397

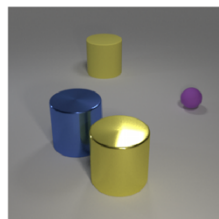
television studio (90.2%) Ranked 1 out of 397



- a photo of a **television studio**.
- a photo of a **podium indoor**.
- a photo of a **conference room**.
- a photo of a **lecture room**.
- a photo of a **control room**.

CLEVR COUNT

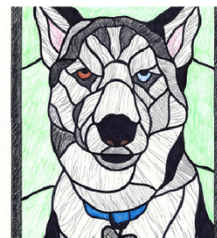
4 (17.1%) Ranked 2 out of 8



- a photo of **3** objects.
- a photo of **4** objects.
- a photo of **5** objects.
- a photo of **6** objects.
- a photo of **10** objects.

IMAGENET-R (RENDITION)

Siberian Husky (76.0%) Ranked 1 out of 200



- a photo of a **siberian husky**.
- a photo of a **german shepherd dog**.
- a photo of a **collie**.
- a photo of a **border collie**.

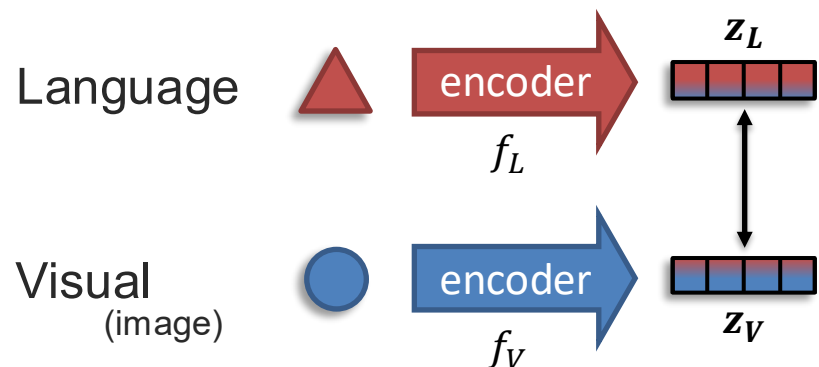
FOOD101

guacamole (90.1%) Ranked 1 out of 101 labels



- a photo of **guacamole**, a type of food.
- a photo of **ceviche**, a type of food.
- a photo of **edamame**, a type of food.
- a photo of **tuna tartare**, a type of food.
- a photo of **hummus**, a type of food.

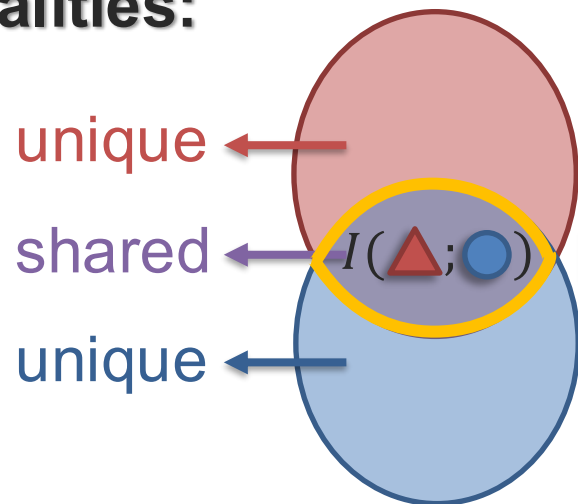
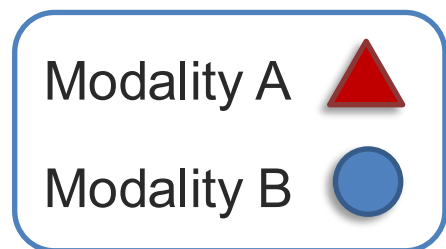
Contrastive Learning and Connected Modalities



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:



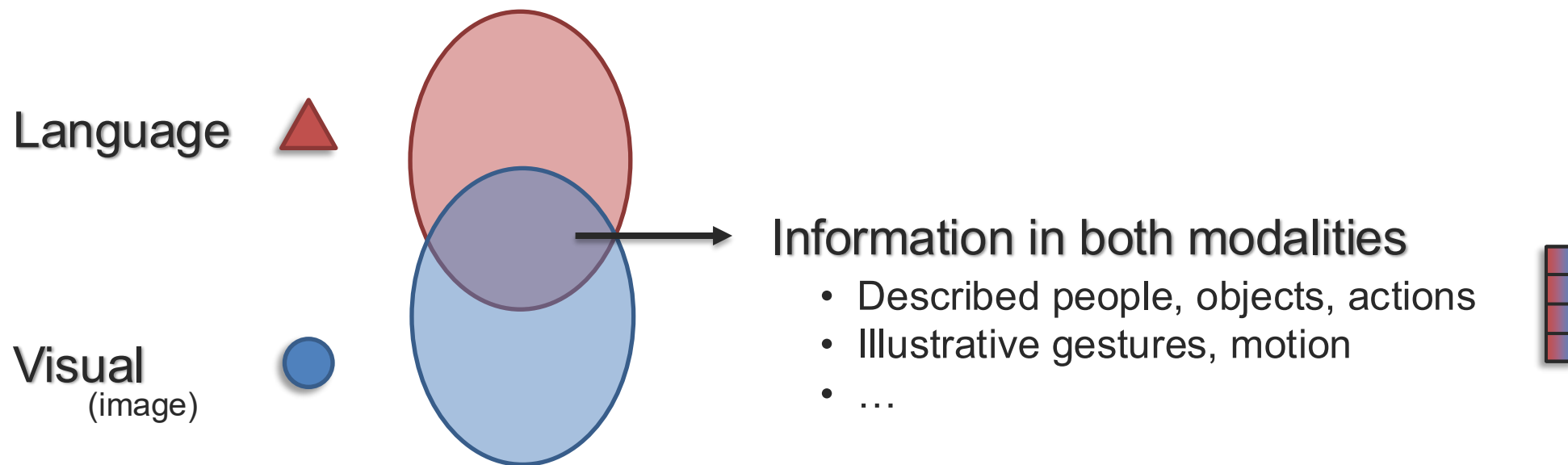
Mutual information $I(X; Y)$

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

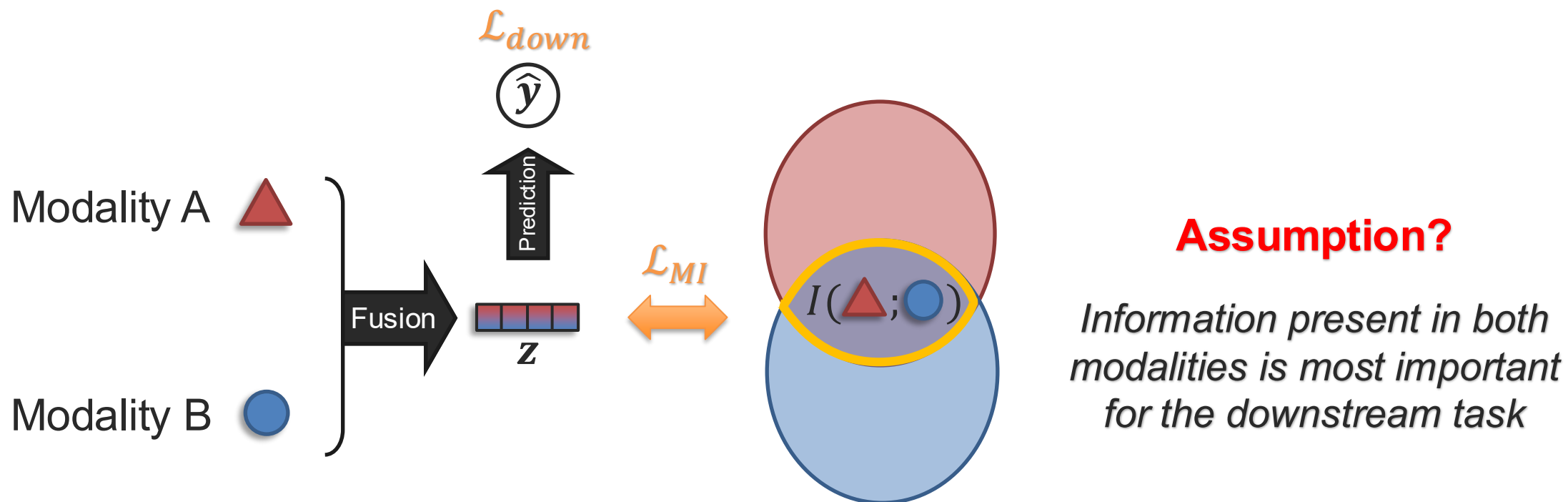


CLIP focuses on shared connections

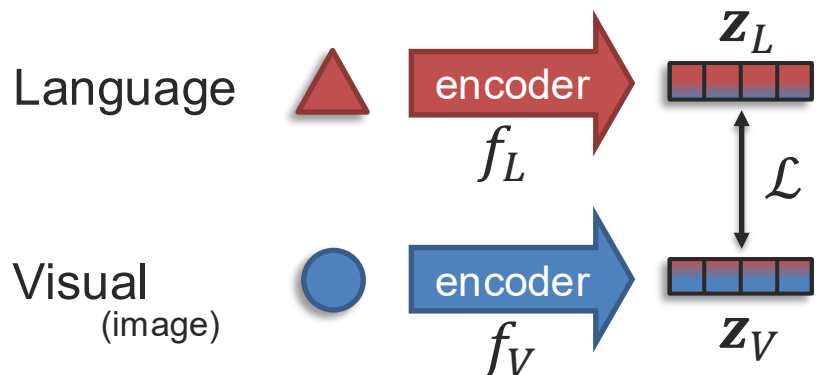
Multimodal Coordination – Information Theory



Multimodal Fusion with Mutual Information



Contrastive Learning and Mutual Information



InfoNCE/CL:

- 'Captures' mutual information
- Optimizes a lower bound on mutual information

In other words:

InfoNCE:

$$\mathcal{L} = -\mathbb{E} \left[\log \frac{f(\mathbf{x}_A^i, \mathbf{x}_B^i)}{\sum_{j=1}^N f(\mathbf{x}_A^i, \mathbf{x}_B^j)} \right]$$

critic function

Critic function f is trained to be a binary classifier distinguishing $\mathbf{x}_A, \mathbf{x}_B \sim p(\mathbf{x}_A, \mathbf{x}_B)$ vs $\mathbf{x}_A, \mathbf{x}_B \sim p(\mathbf{x}_A)p(\mathbf{x}_B)$

At optimal loss, $f^*(\mathbf{x}_A, \mathbf{x}_B) = \frac{p(\mathbf{x}_A, \mathbf{x}_B)}{p(\mathbf{x}_A)p(\mathbf{x}_B)}$.

Plugging f^* back into \mathcal{L} gives:

$$\mathcal{L}^* \geq \mathbb{E} \left[\log \frac{p(\mathbf{x}_A)p(\mathbf{x}_B)}{p(\mathbf{x}_A, \mathbf{x}_B)} N \right] = -I(X_A, X_B) + \log N$$

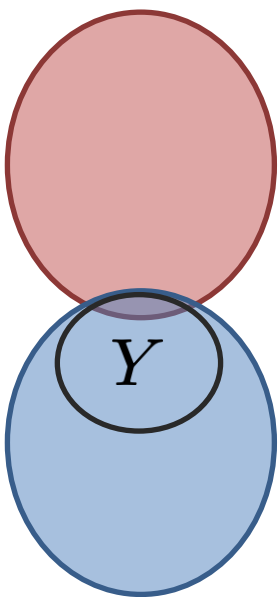
$$I(X_A, X_B) \geq \log N - \mathcal{L}^*$$

Multiview Redundancy and Contrastive Learning

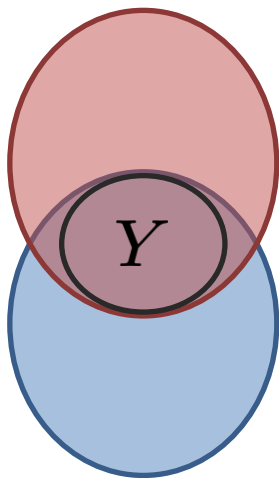
How much information should be shared?

Multi-view redundancy: $I(X_1; X_2) = I(X_1; Y)$

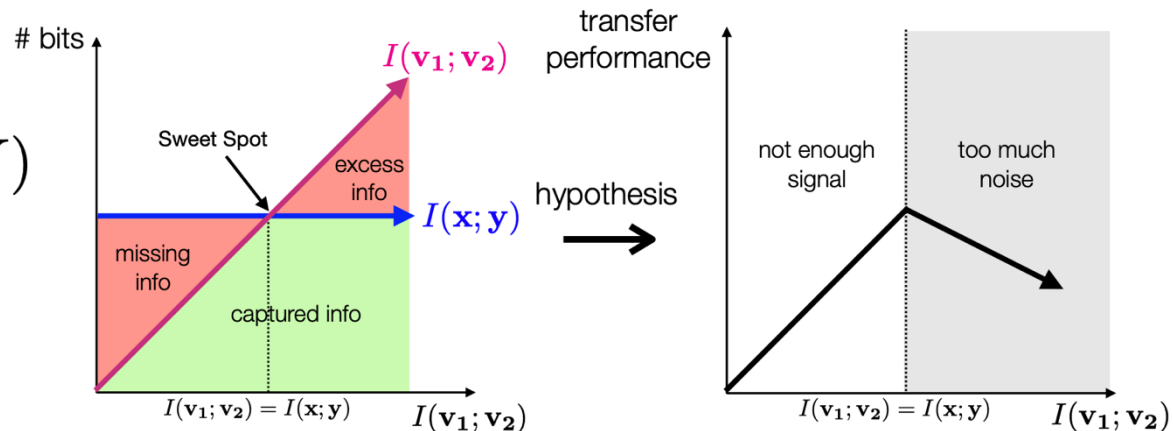
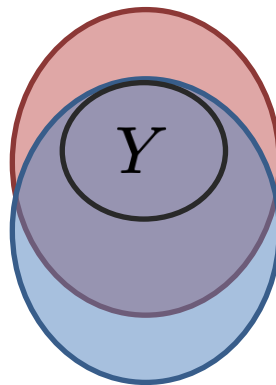
Not enough signal



Just right

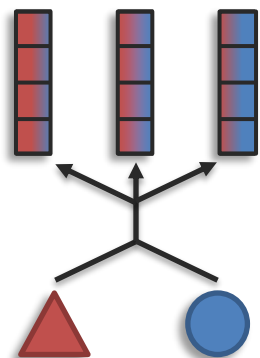


Too much noise



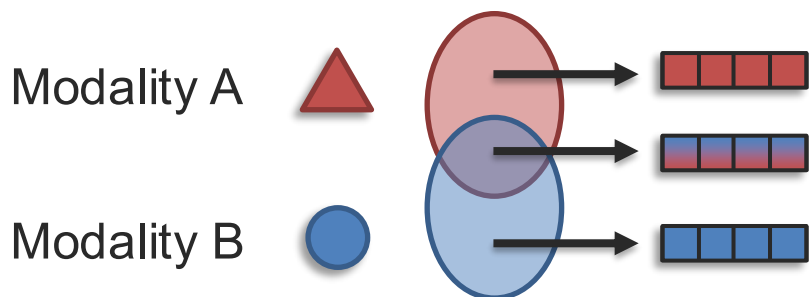
Multi-view redundancy
may not hold for
multimodal problems!

Representation Fission

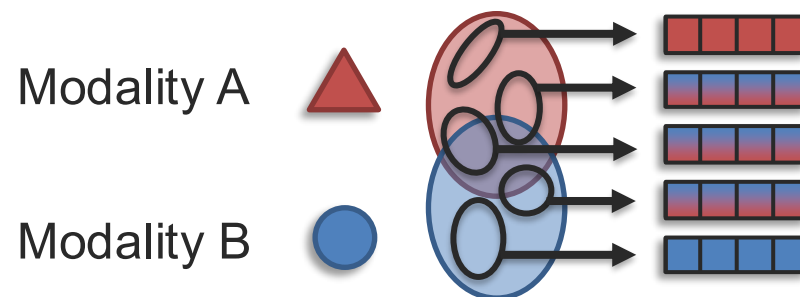


Definition: Learning a new set of representations that reflects multimodal internal structure such as data factorization or clustering

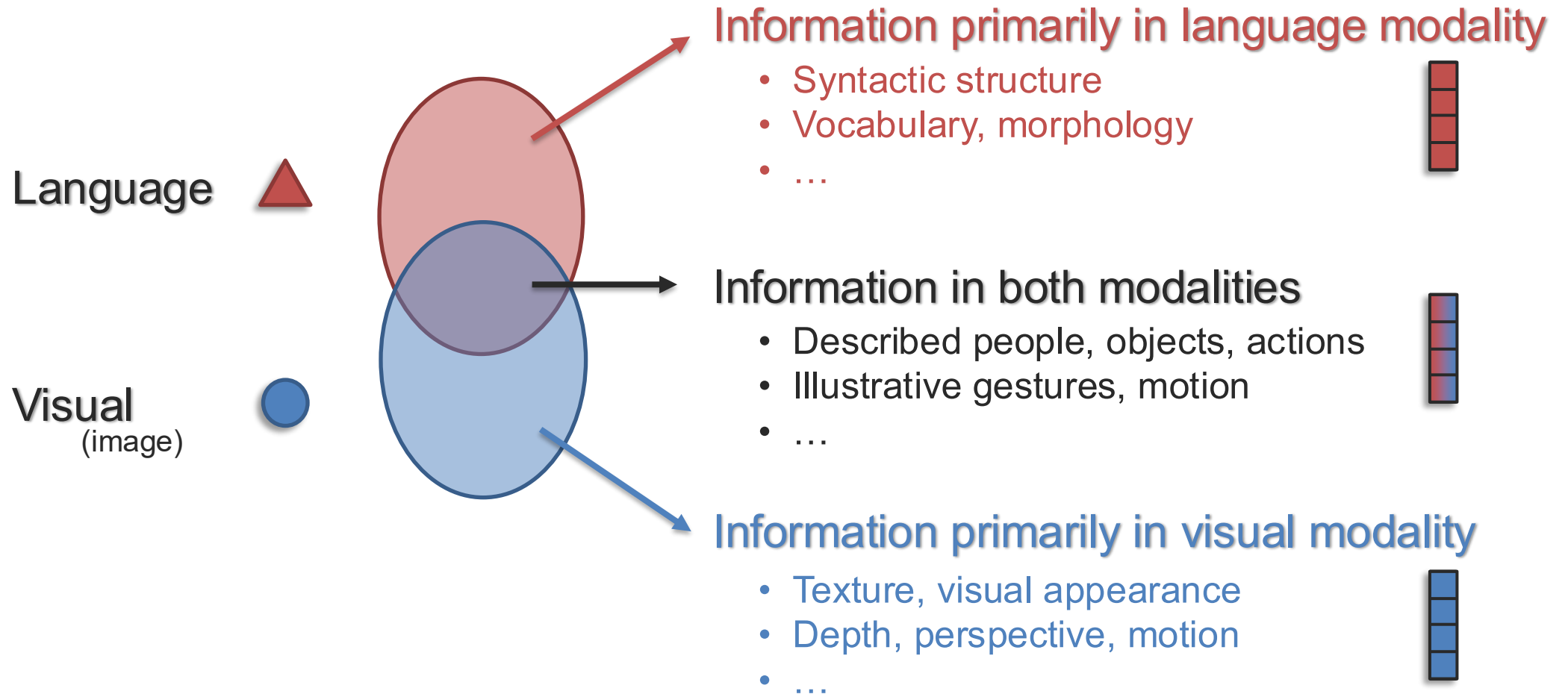
Modality-level fission:



Fine-grained fission:



Modality-level Fission



Factorized Learning of Shared + Unique Information

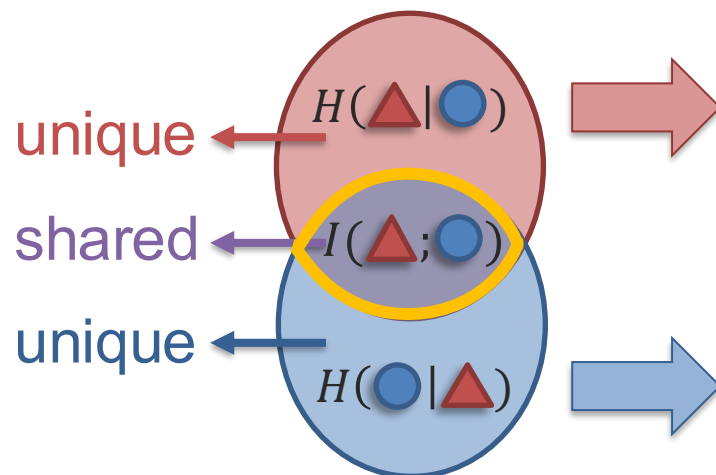


Can you please pass the cow?

Modality A



Modality B



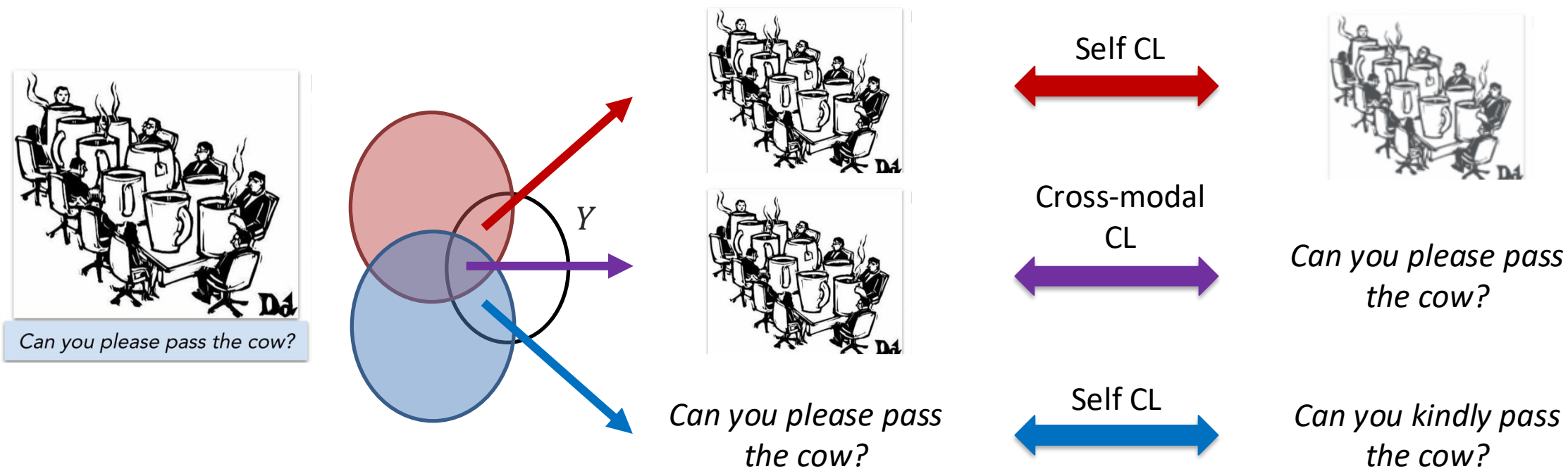
- 1 Maximize the mutual information

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

- 2 Minimize the conditional entropy

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

Factorized Contrastive Learning



Learns both shared and unique information.

Extensions: Global Alignment

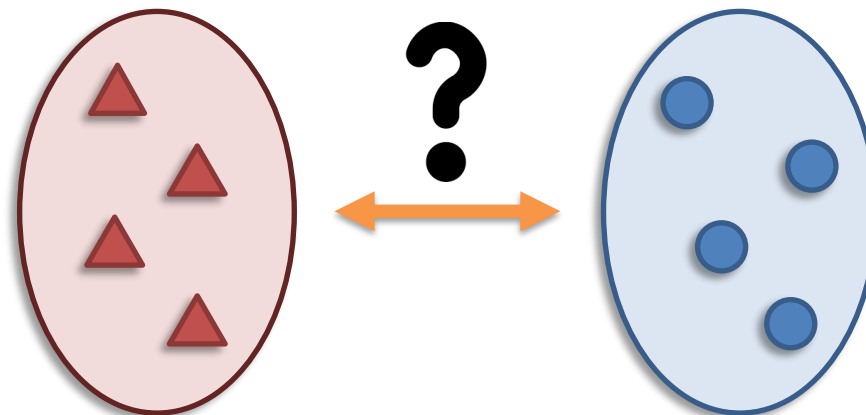
Visual



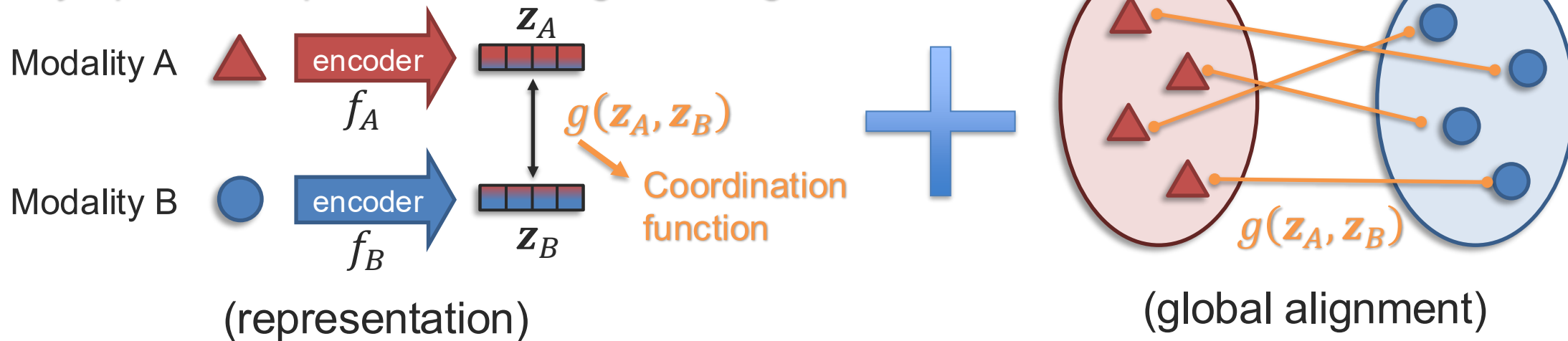
Language



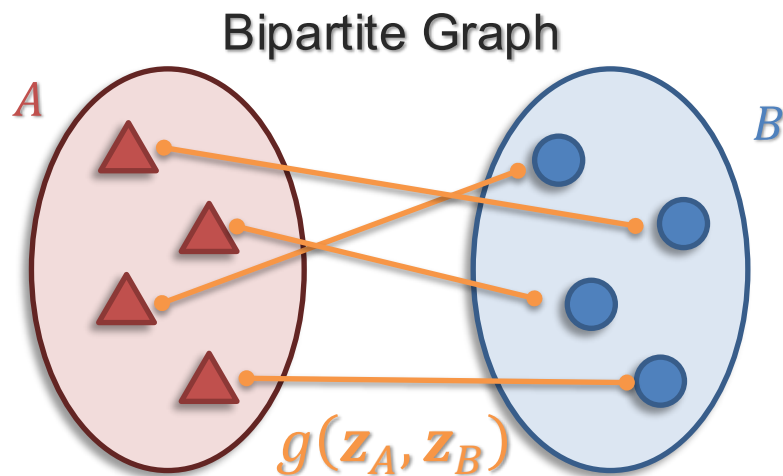
Latent pairing information



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

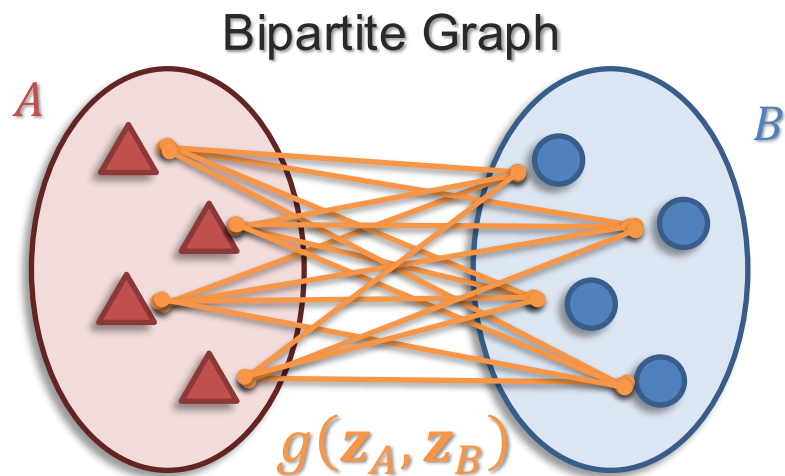
$w_{(i,j)} = g(\mathbf{z}_A^i, \mathbf{z}_B^j)$

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)

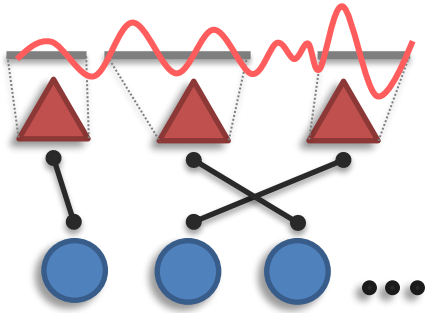
Assignments: $x_{(i,j)}$: soft alignment between z_A^i and z_B^j

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

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

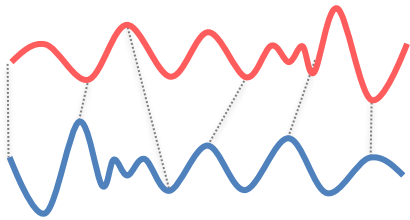
➔ Wasserstein distance give optimal transport

Extensions: Continuous Alignment

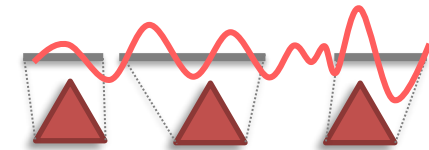


Definition: Model alignment between modalities with continuous signals and no explicit elements

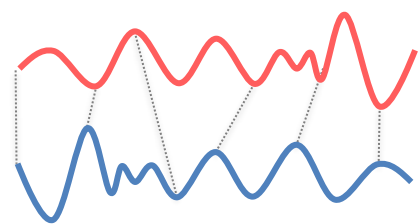
Continuous
warping



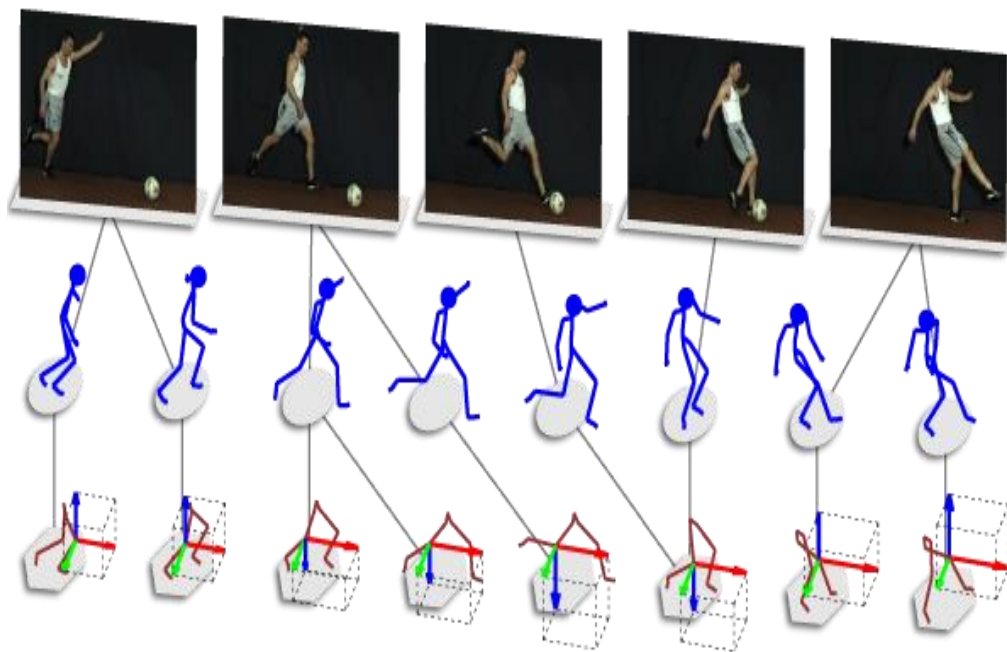
Discretization
(segmentation)



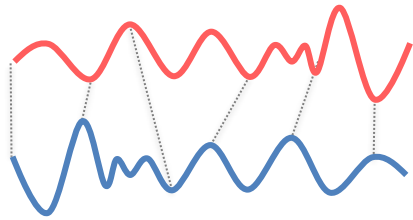
Continuous Alignment



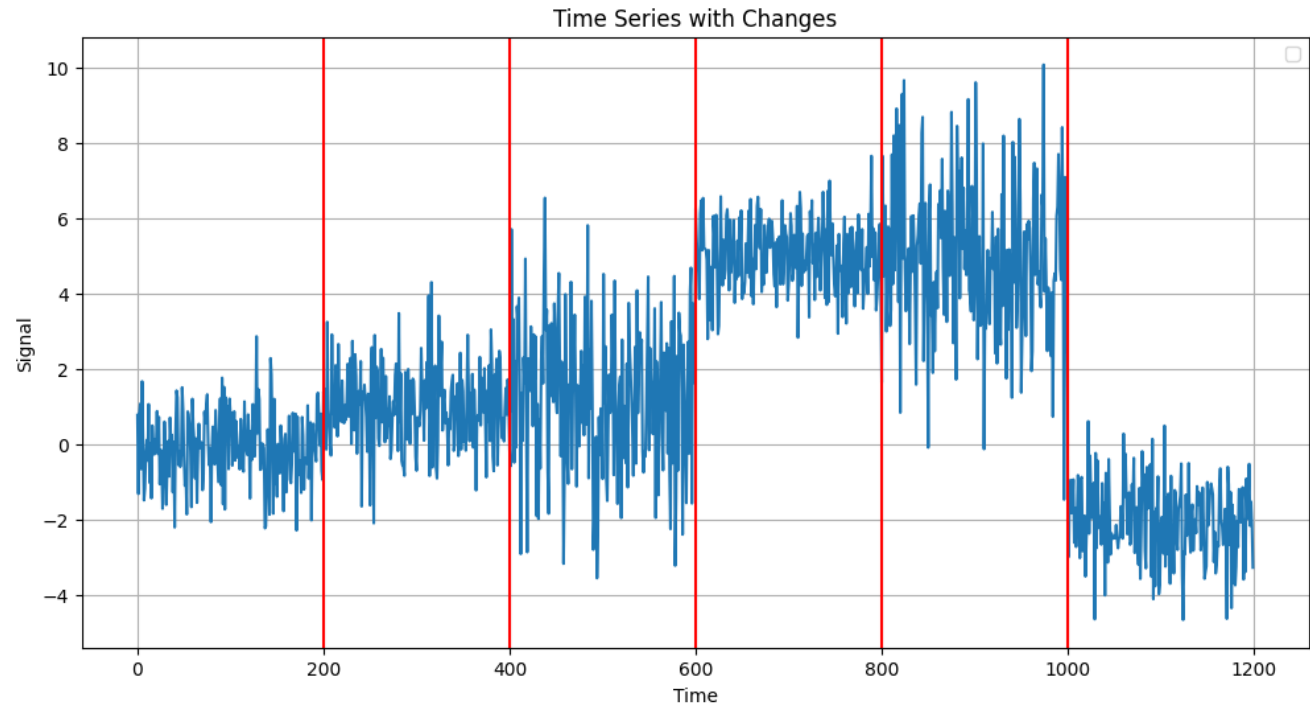
➔ Aligning video sequences



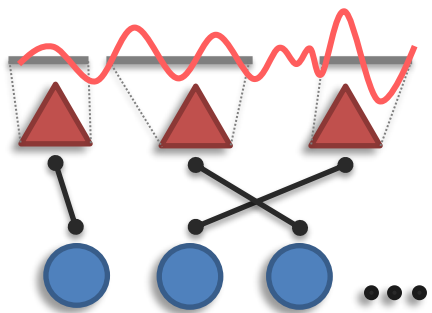
Continuous Alignment



➔ Changepoint detection

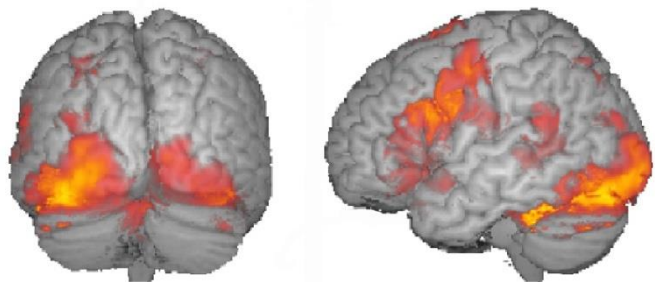


Discretization (aka Segmentation)

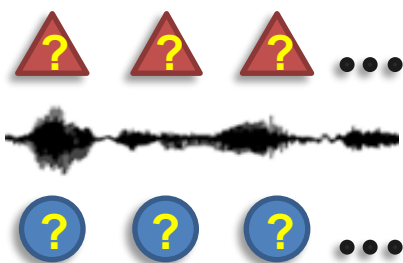


Common assumptions: ① Segmented elements

Examples:



Medical imaging



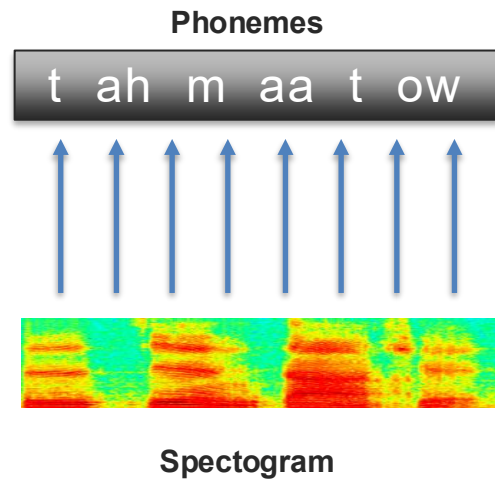
Signals



Images

Discretization – Example

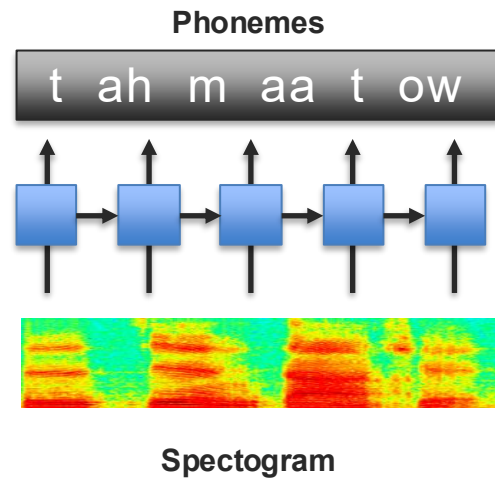
Sequence Labeling and Alignment



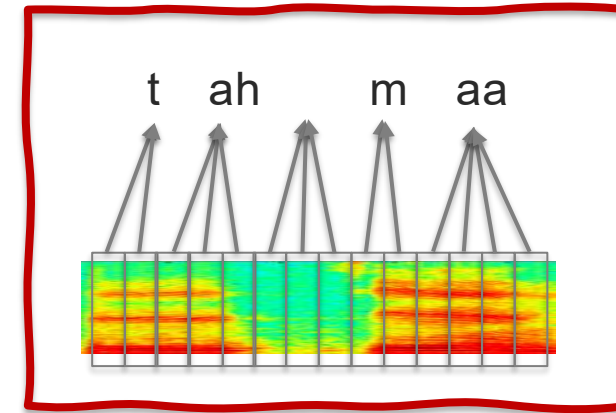
How can we predict the sequence
of phoneme labels?

Discretization – Example

Sequence Labeling and Alignment



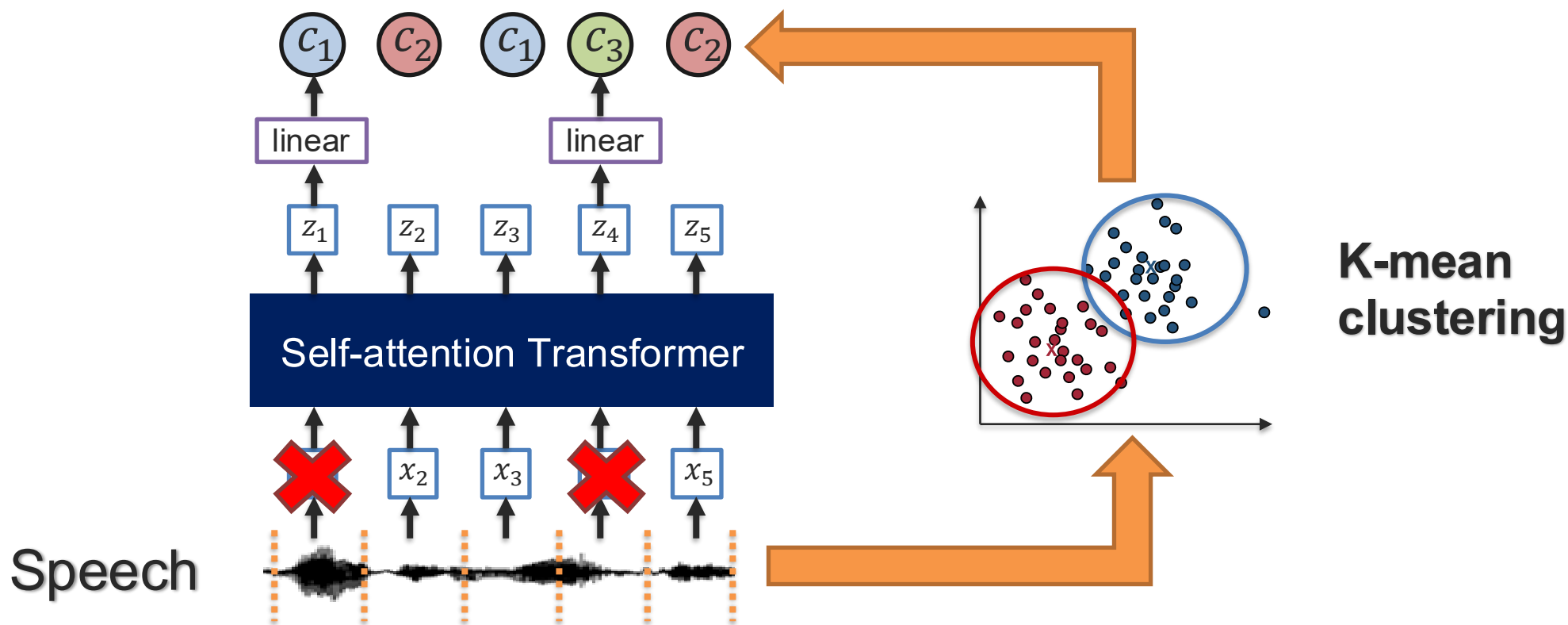
Challenge: many-to-1 alignment



How can we predict the sequence of phoneme labels?

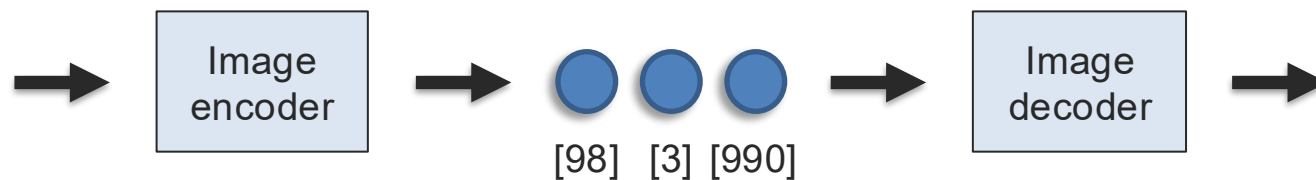
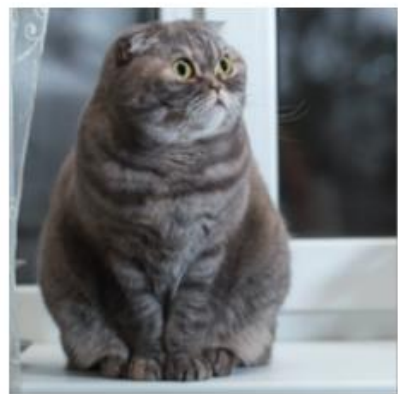
Hidden-Unit Pre-training

HUBERT: Hidden-Unit BERT



VQ-VAE

Using a discrete variational autoencoder to learn discrete visual tokens

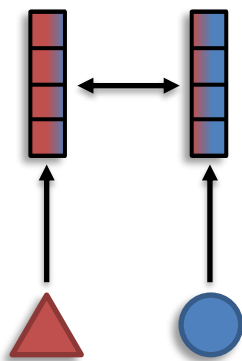


List of digits, [0... 8192]
Each digit is a “visual token”



Extension: Implicit (Emergent) Alignment

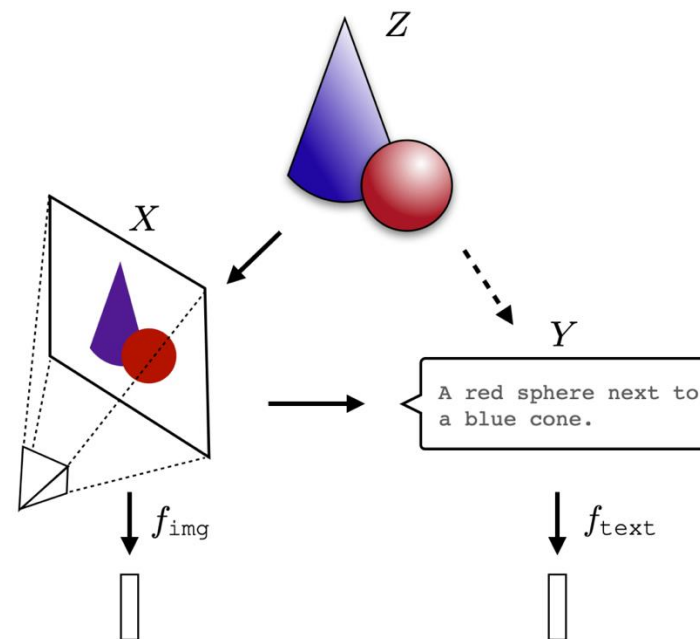
Explicit alignment



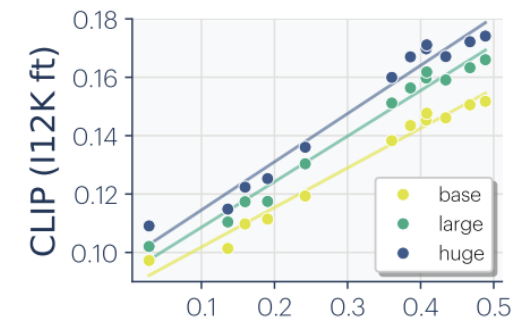
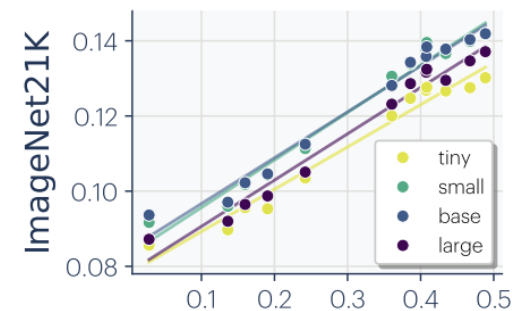
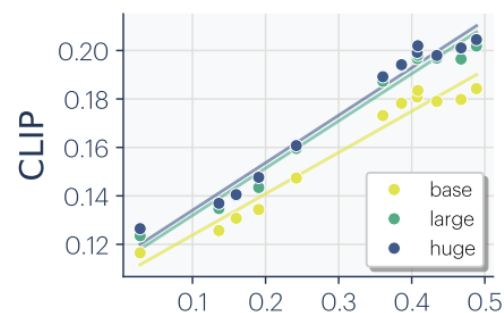
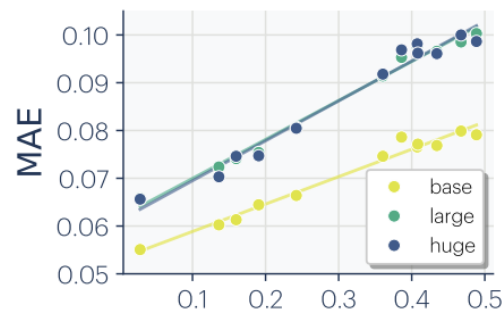
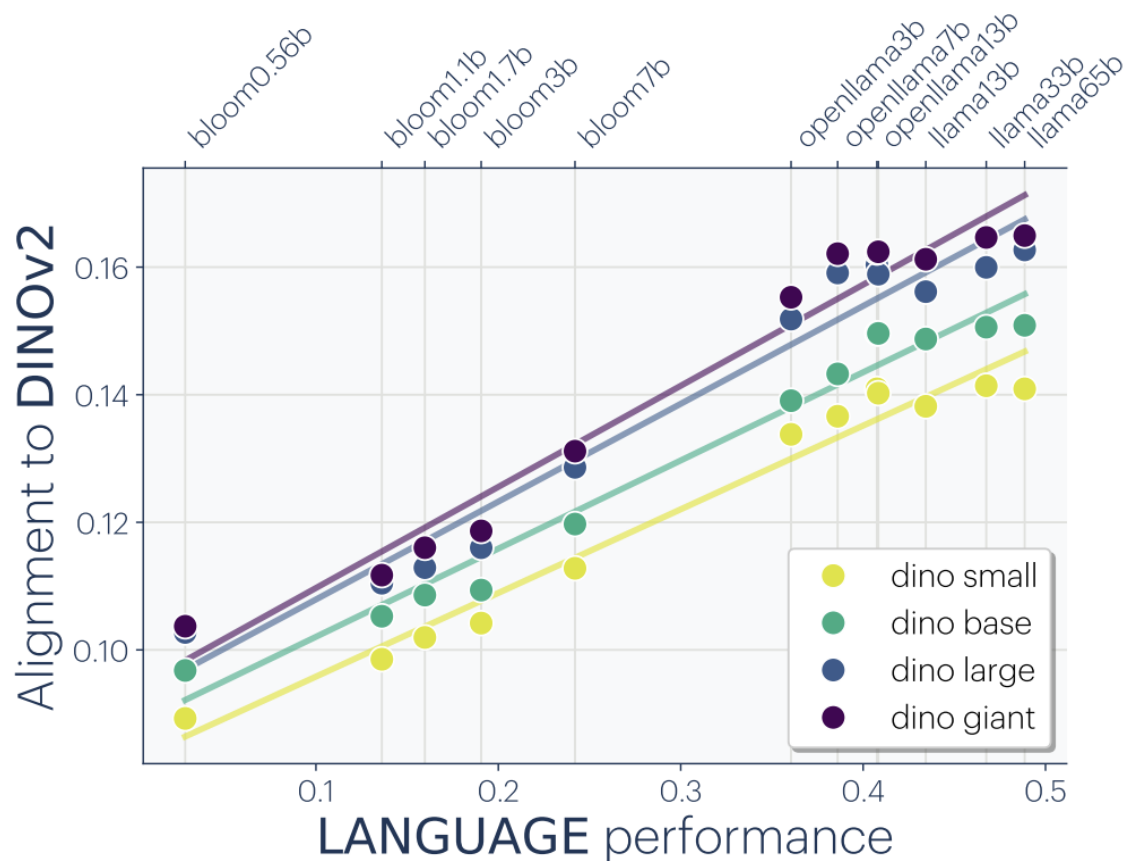
Implicit alignment

The Platonic Representation Hypothesis

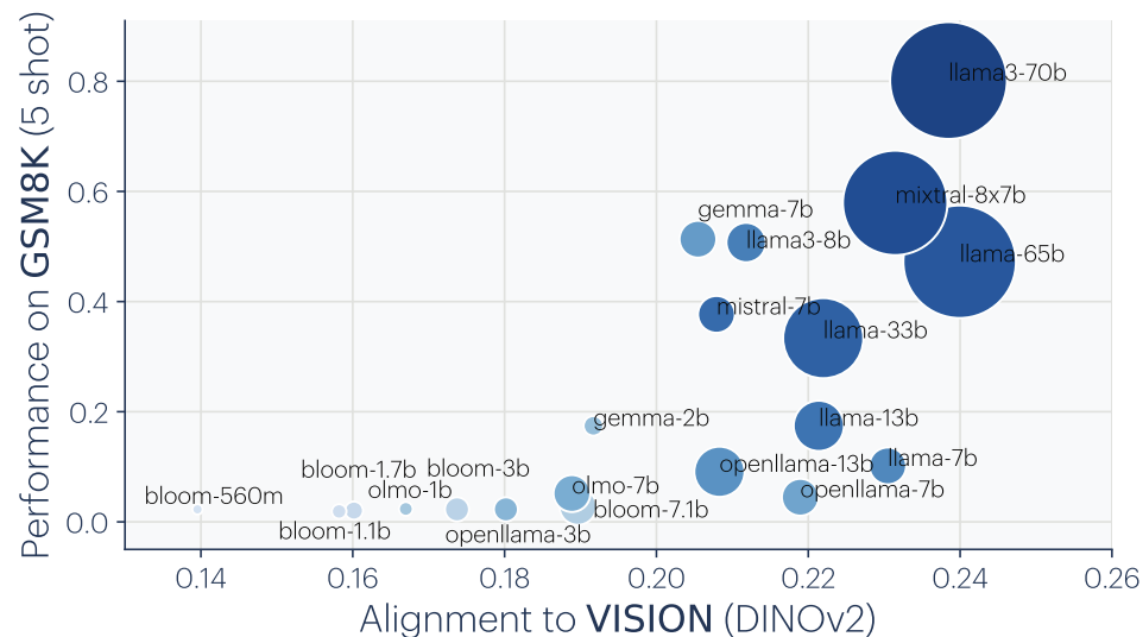
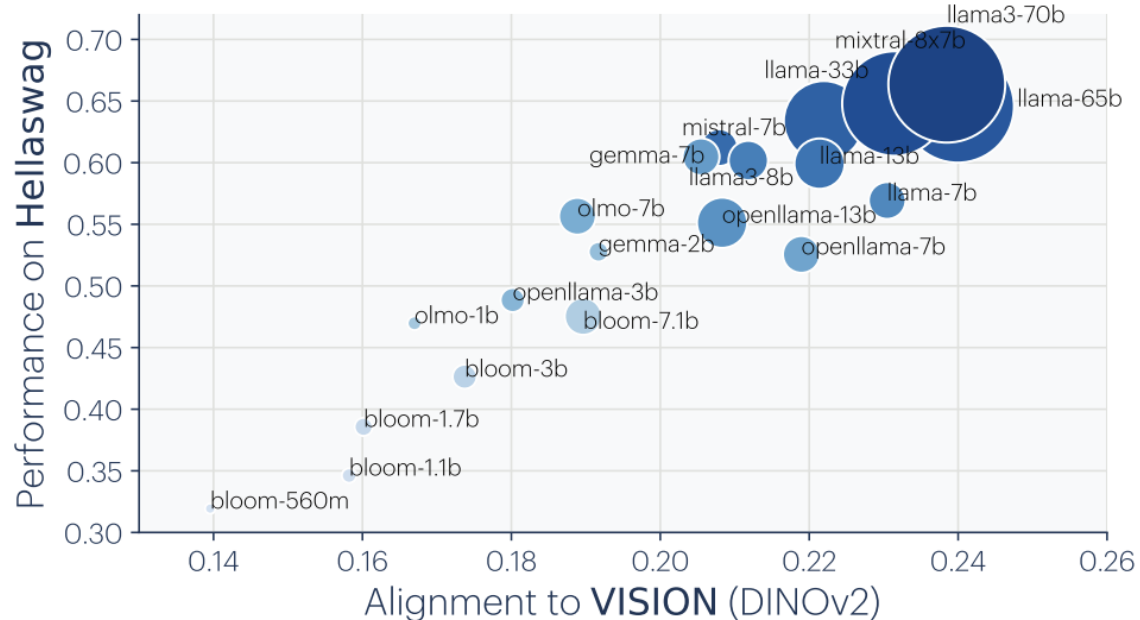
Neural networks, trained with different objectives on different data and modalities, are converging to a shared statistical model of reality in their representation spaces.



Emergence of Alignment

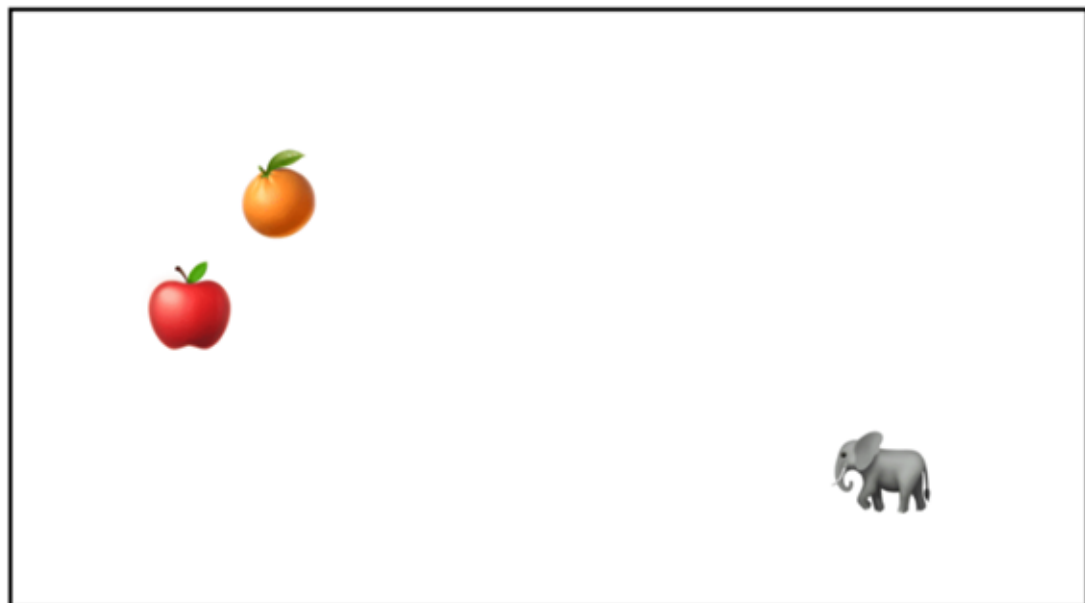


Emergence of Alignment



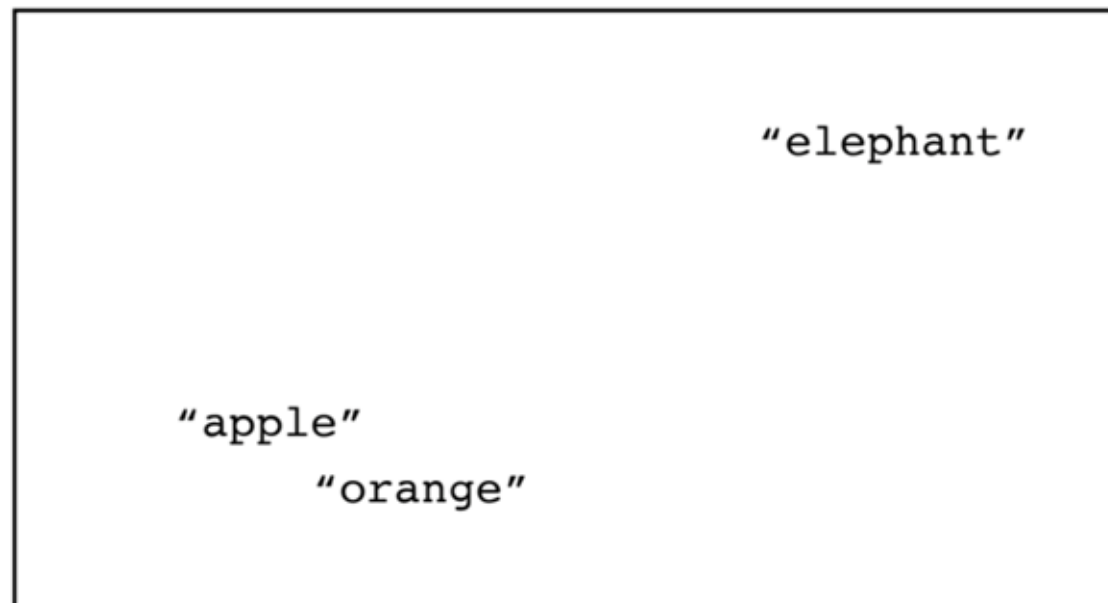
Alignment via Kernel Similarity

Image embeddings



$$K_{\text{img}}(i, j) = \langle E_{\text{img}}(x_i), E_{\text{img}}(x_j) \rangle$$

Text embeddings



$$K_{\text{text}}(i, j) = \langle E_{\text{text}}(y_i), E_{\text{text}}(y_j) \rangle$$

$$K_{\text{img}}(i, j) \approx K_{\text{text}}(i, j) \quad \forall i, j$$

Limitations of Emergence?

Summary

- 1 Multimodal alignment
- 2 Explicit alignment and contrastive learning
- 3 Continuous alignment
- 4 Implicit (emergent) alignment

Assignments for This Coming Week

I want to meet every group at least once regarding their project ideas.

Compute credits: 40 x \$50 Kimi credits, 40 x \$40 other credits.

HW2 due next Wednesday 3/4.