

Multimodal AI

Lecture 7.1 – Midterm Review

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

HW3 due today.

Project proposals will be graded and released this week.

Project midterm instructions will be out this week.

Midterm this Thursday.

Midterm Structure

6 problems, 60 minutes.

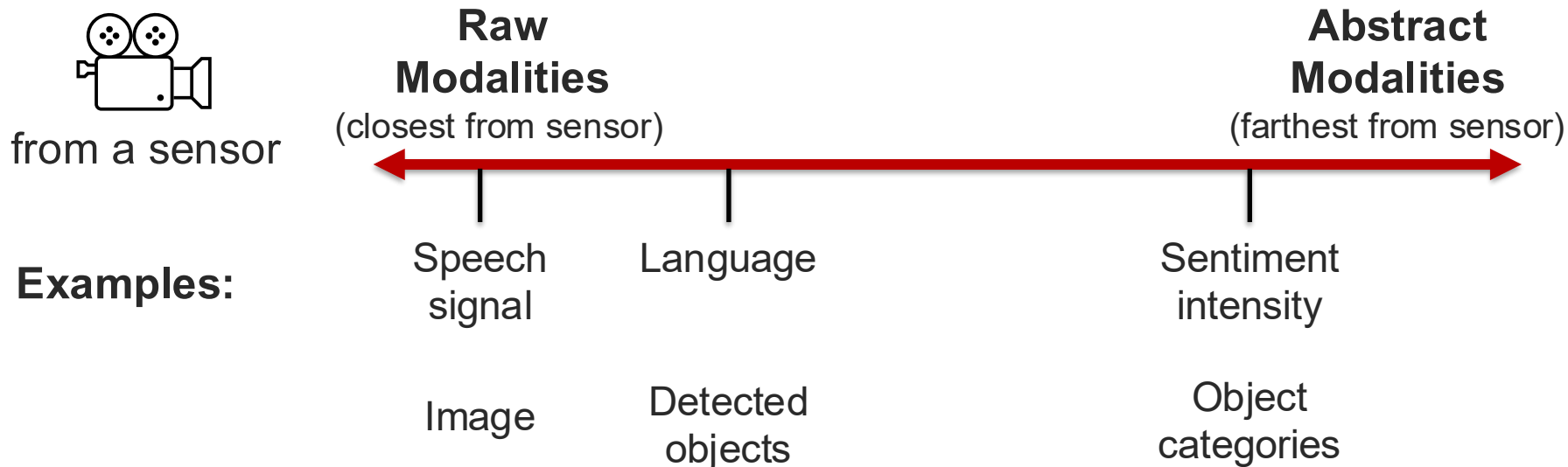
5 problems:

1. MCQ (10, 2 points each = 20 points)
2. Short answers (4, 5 points each = 20 points)
3. Multimodal Fusion (20 points)
4. Multimodal LLMs (20 points)
5. Multimodal Generation (20 points)
6. Bonus open questions (10 points)

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...

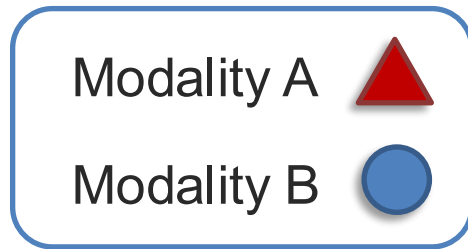
***Multimodal* is the science of**

heterogeneous and interconnected data

Connected + Interacting

Heterogeneous Modalities

Information in different modalities shows diverse qualities, structures, & representations.



Homogeneous Modalities
(with similar qualities)



Images
from 2
cameras



Text from
2 different
languages



Language
and vision



Language
and sensors

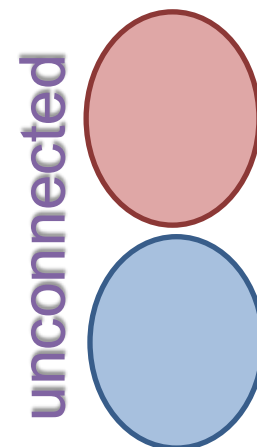
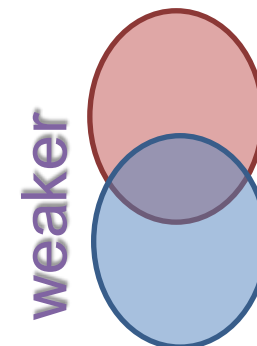
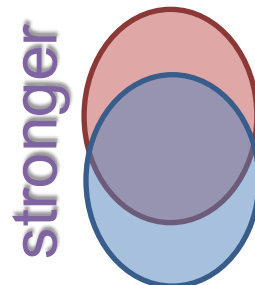
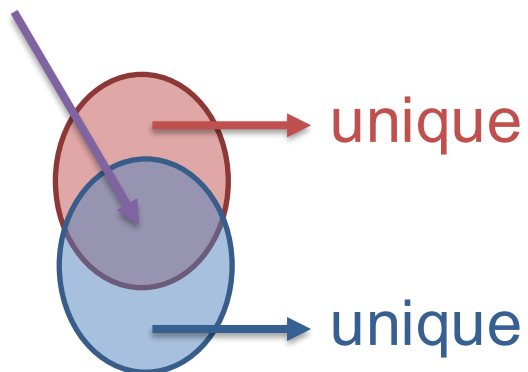
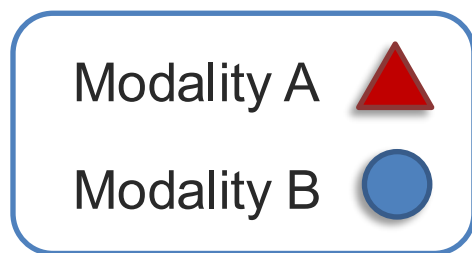
Heterogeneous Modalities
(with diverse qualities)

Examples:

Abstract modalities are more likely to be homogeneous

Connected Modalities

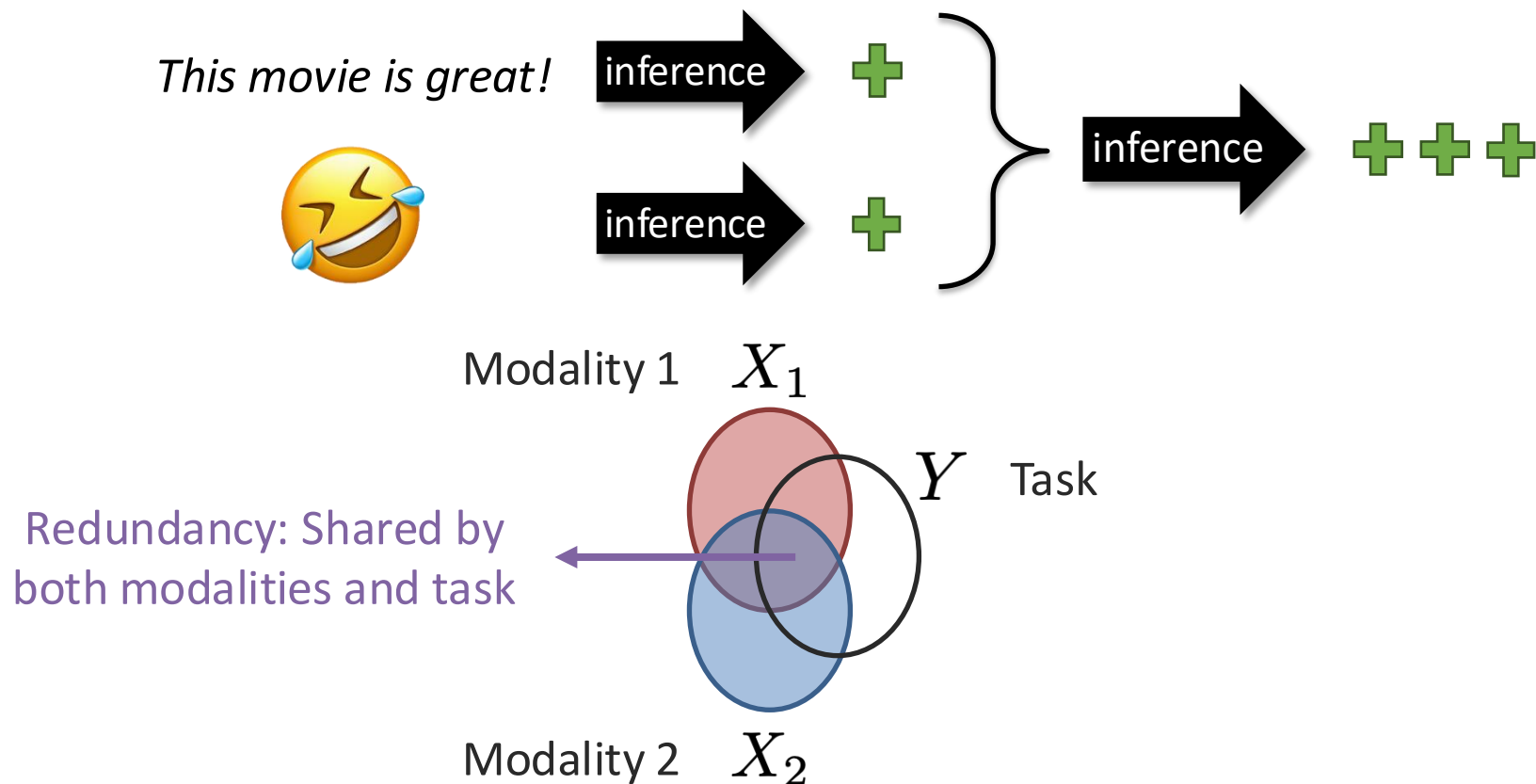
Shared information that relates modalities



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

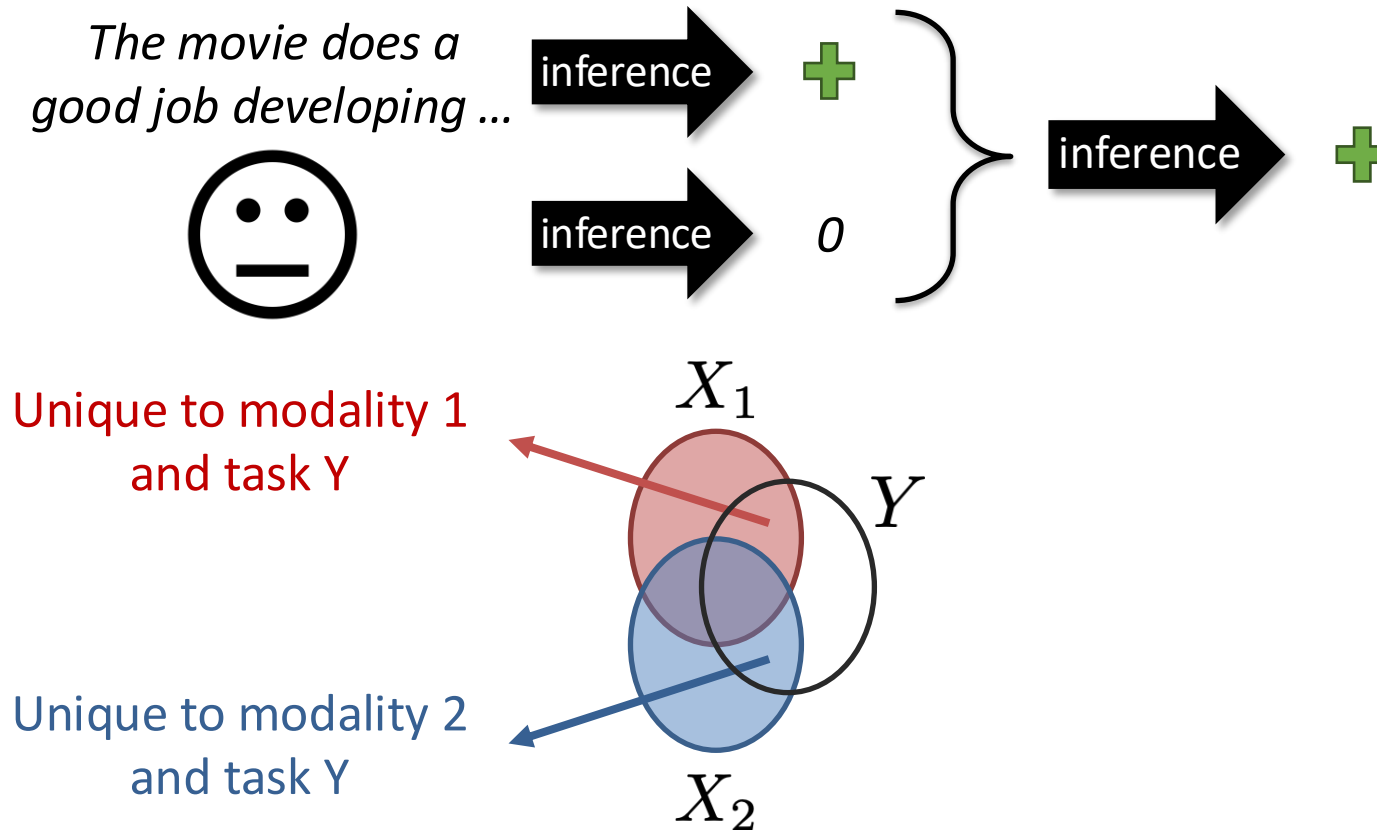
Interacting Modalities

Interactions: How modalities *combine* to provide information for a task.



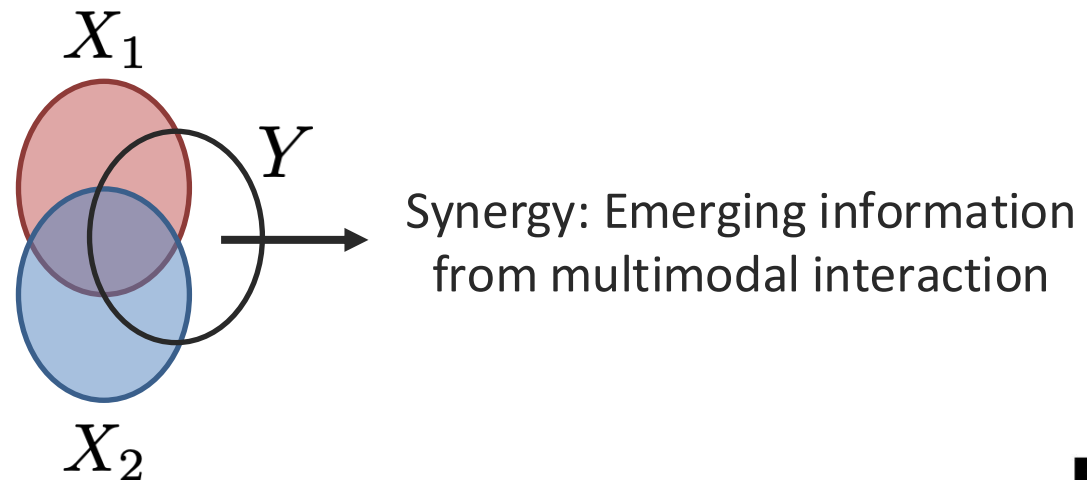
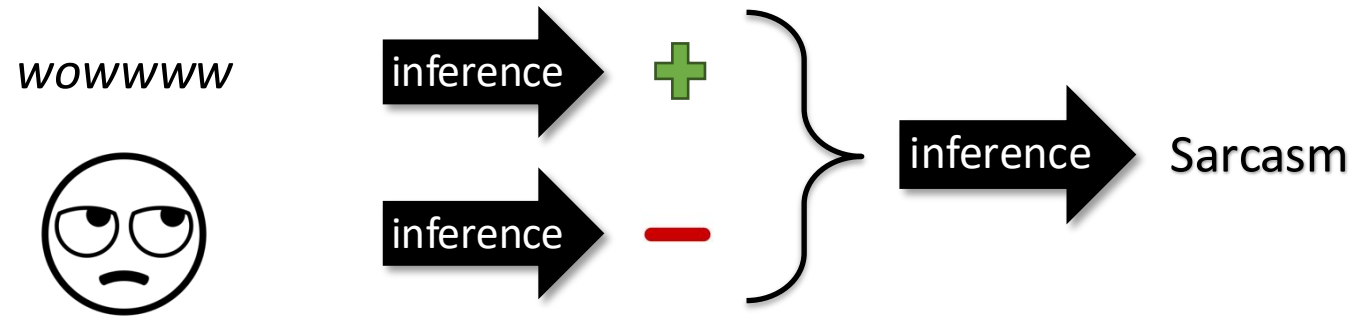
Interacting Modalities

Interactions: How modalities *combine* to provide information for a task.



Interacting Modalities

Interactions: How modalities *combine* to provide information for a task.



*What is
Multimodal?*



Why is it hard?



What is next?

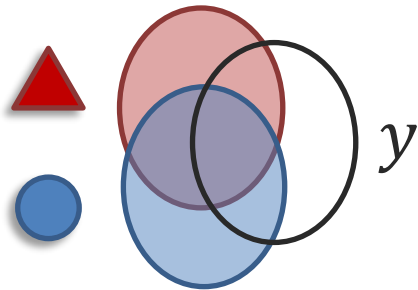
Heterogeneous



Connected

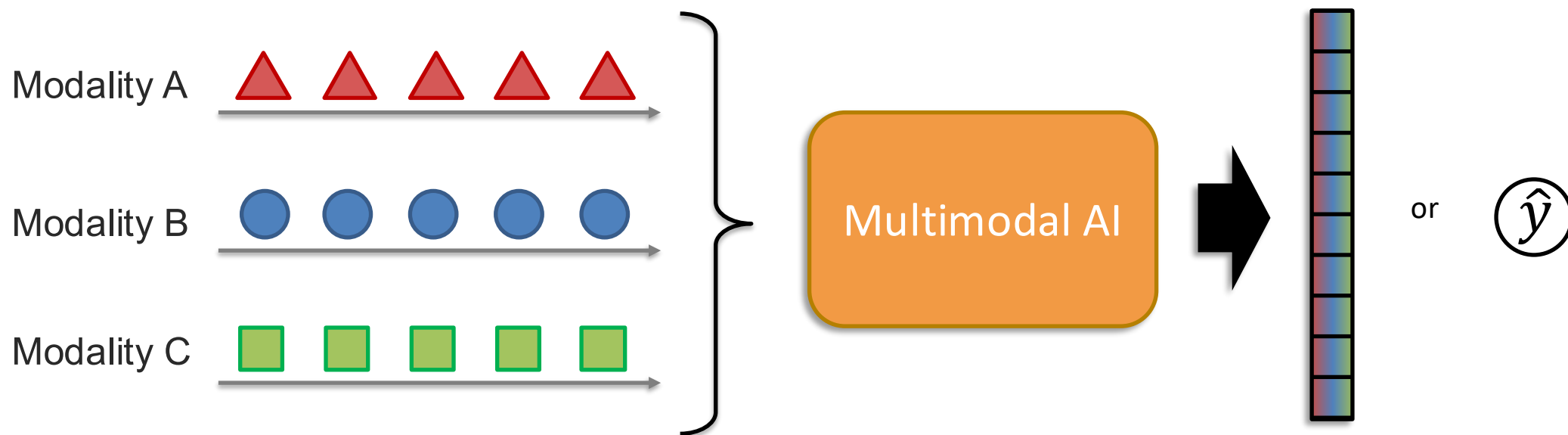


Interacting

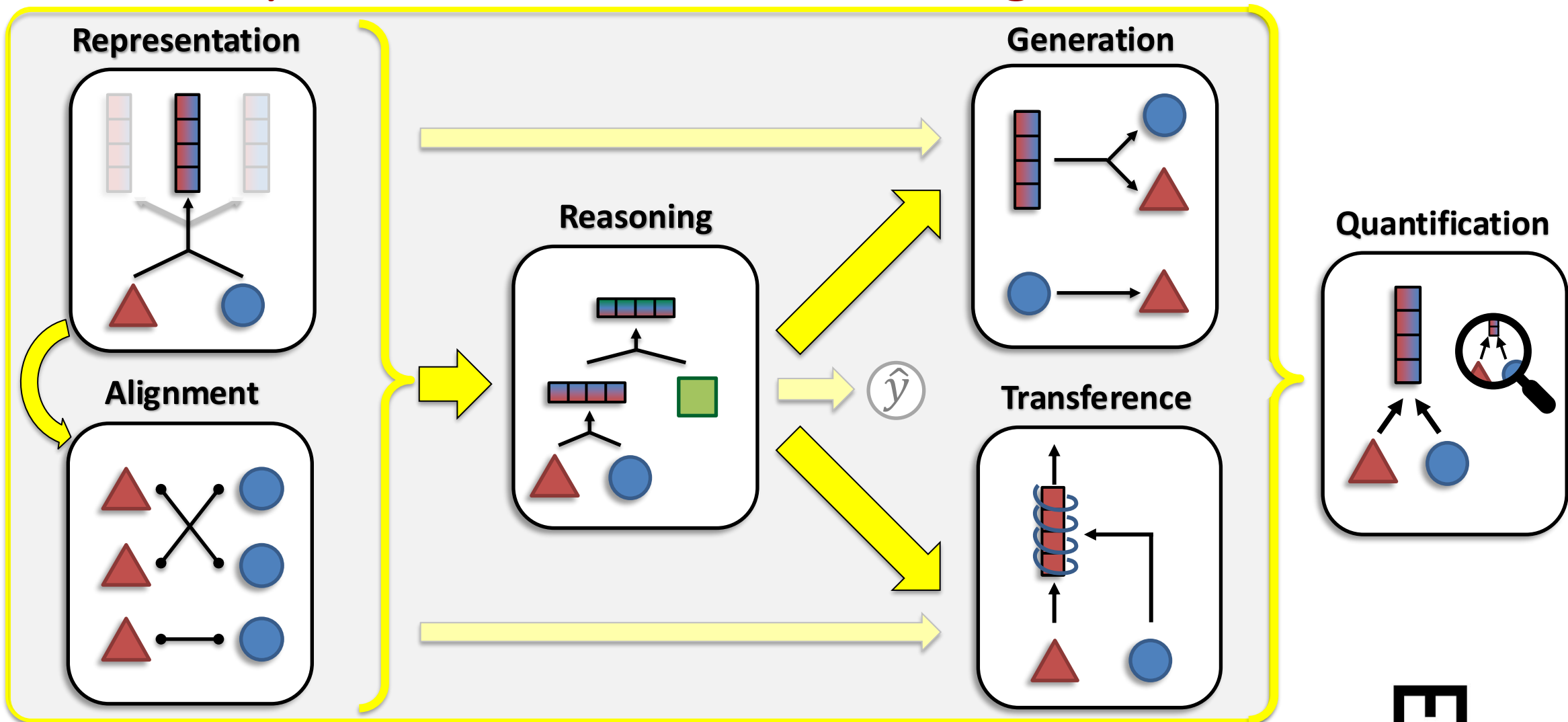


**Multimodal is the scientific
study of heterogeneous and
interconnected data 😊**

Multimodal AI Challenges



Summary of Core Multimodal Challenges

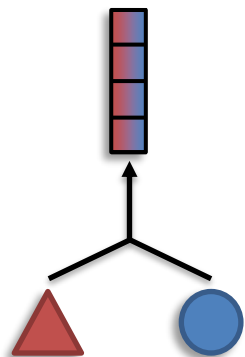


Challenge 1: Representation

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

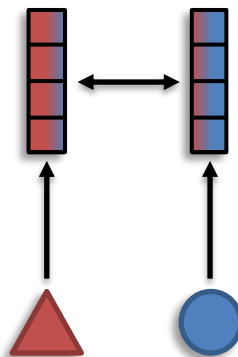
Sub-challenges:

Fusion



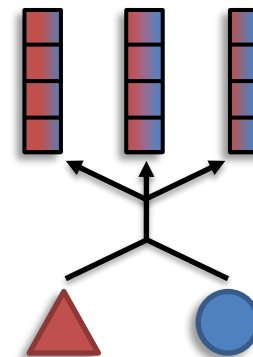
modalities \gt # representations

Coordination



modalities = # representations

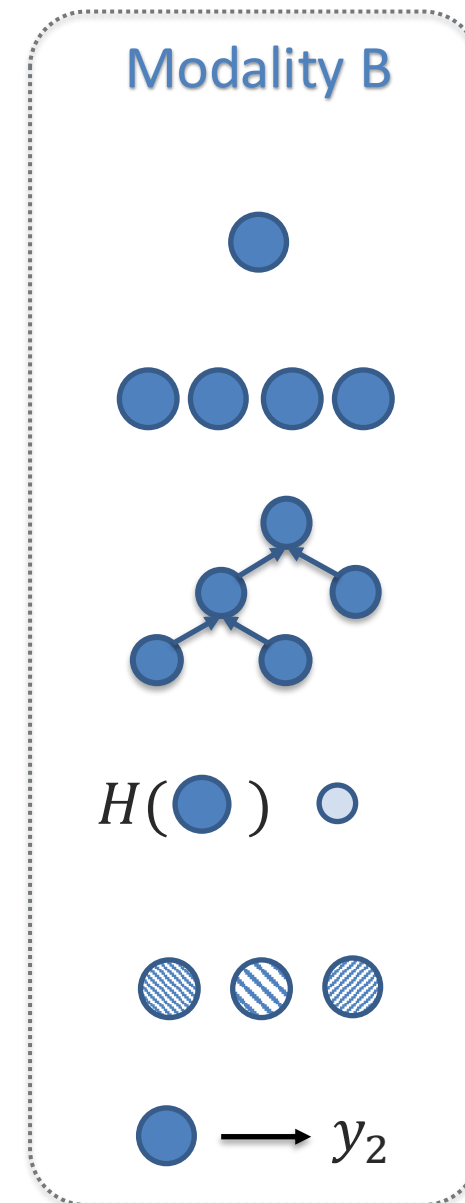
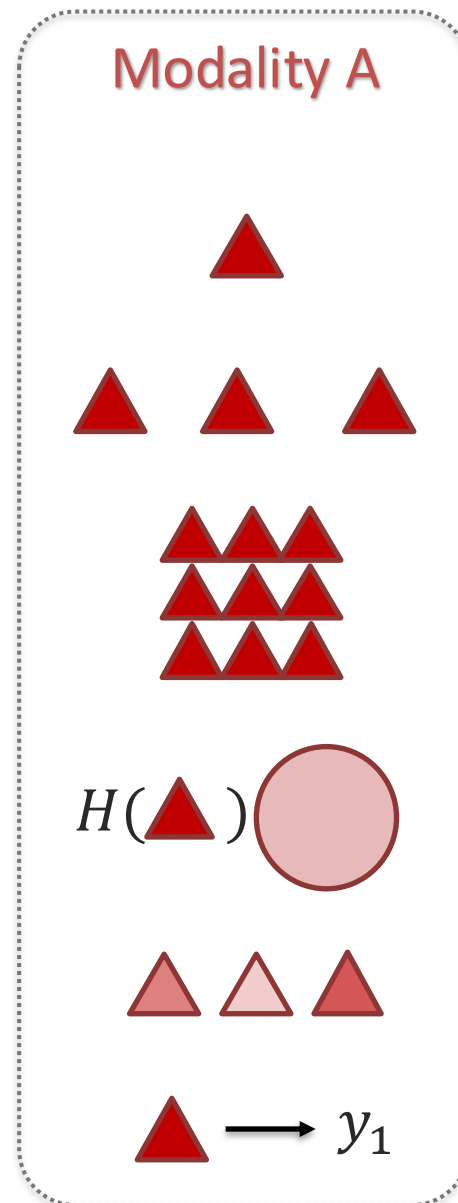
Fission



modalities \lt # representations

Modality Profile

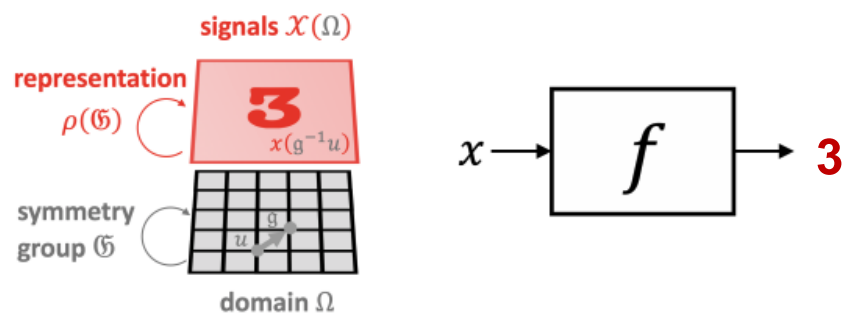
- 1 **Element representations:**
Discrete, continuous, granularity
- 2 **Element distributions:**
Density, frequency
- 3 **Structure:**
Temporal, spatial, latent, explicit
- 4 **Information:**
Abstraction, entropy
- 5 **Noise:**
Uncertainty, noise, missing data
- 6 **Relevance:**
Task, context dependence



Unimodal Structure

Data invariances – example of image classification

A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{Y}$ is \mathfrak{G} -invariant if $f(\rho(\mathfrak{g})x) = f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$ and $x \in \mathcal{X}(\Omega)$, i.e., its output is unaffected by the group action on the input.

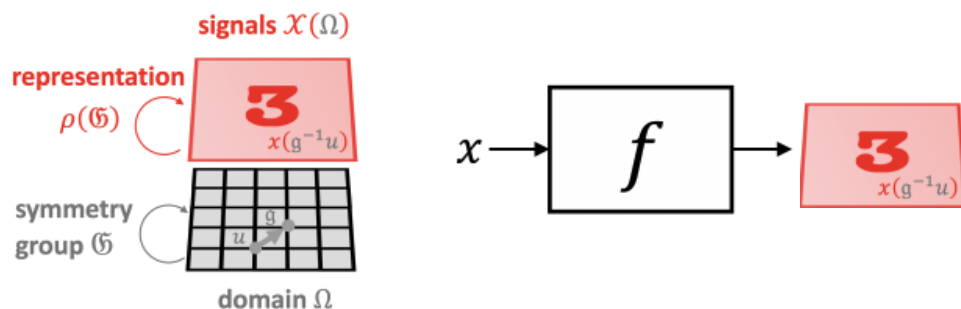


sunset

Unimodal Structure

Data equivariances – example of image segmentation

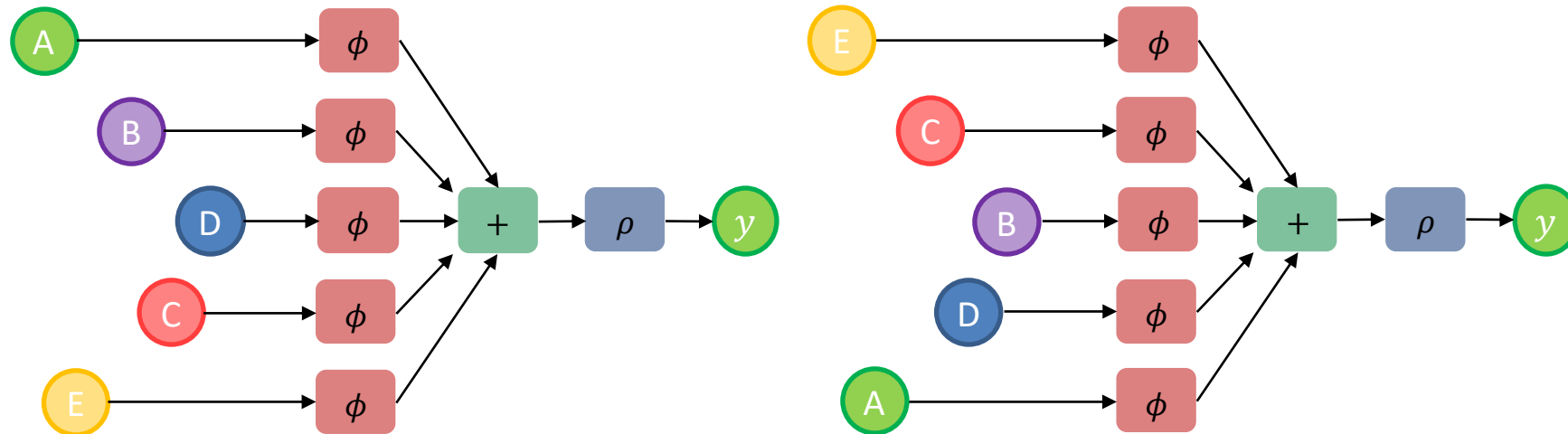
A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{X}(\Omega)$ is \mathfrak{G} -equivariant if $f(\rho(\mathfrak{g})x) = \rho(\mathfrak{g})f(x)$ for all $\mathfrak{g} \in \mathfrak{G}$, i.e., group action on the input affects the output in the same way.



Unimodal Models

Models for set-based data must be invariant to element order.

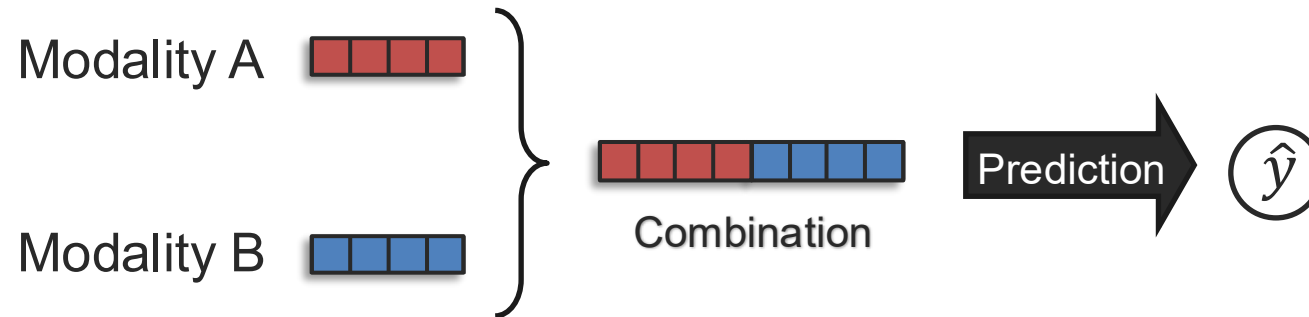
1. Parameter sharing for each set element
2. Permutation invariant aggregation function



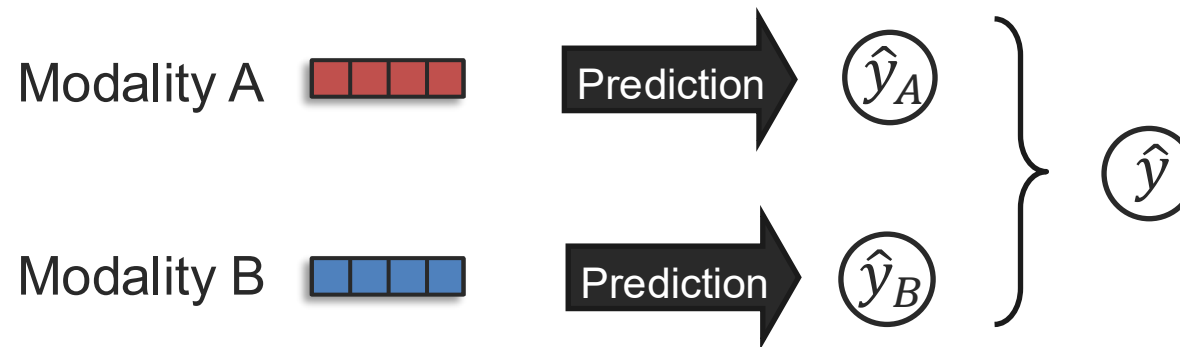
Giving ABDCE also gives ECBDA, BCAED etc...

Early and Late Fusion

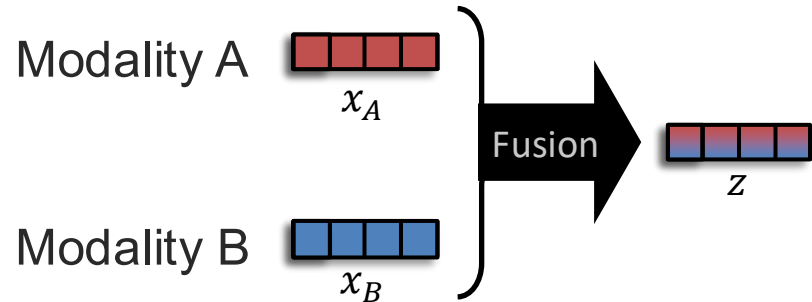
Early fusion:



Late fusion:



Basic Concepts for Representation Fusion



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

→ Let's study the univariate case first

↳ (only 1-dimensional features)

Linear regression:

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

constant Additive terms Multiplicative term error

① Additive interaction:

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

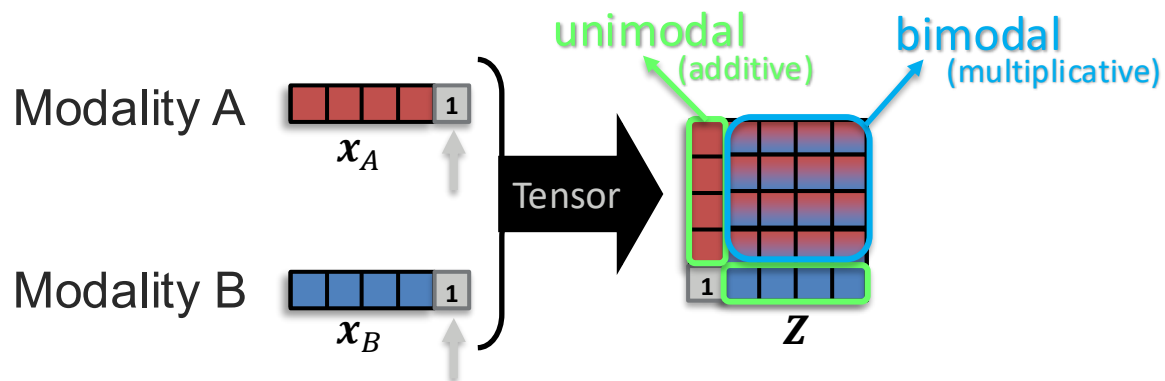
② Multiplicative interaction:

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

③ Additive and multiplicative interactions:

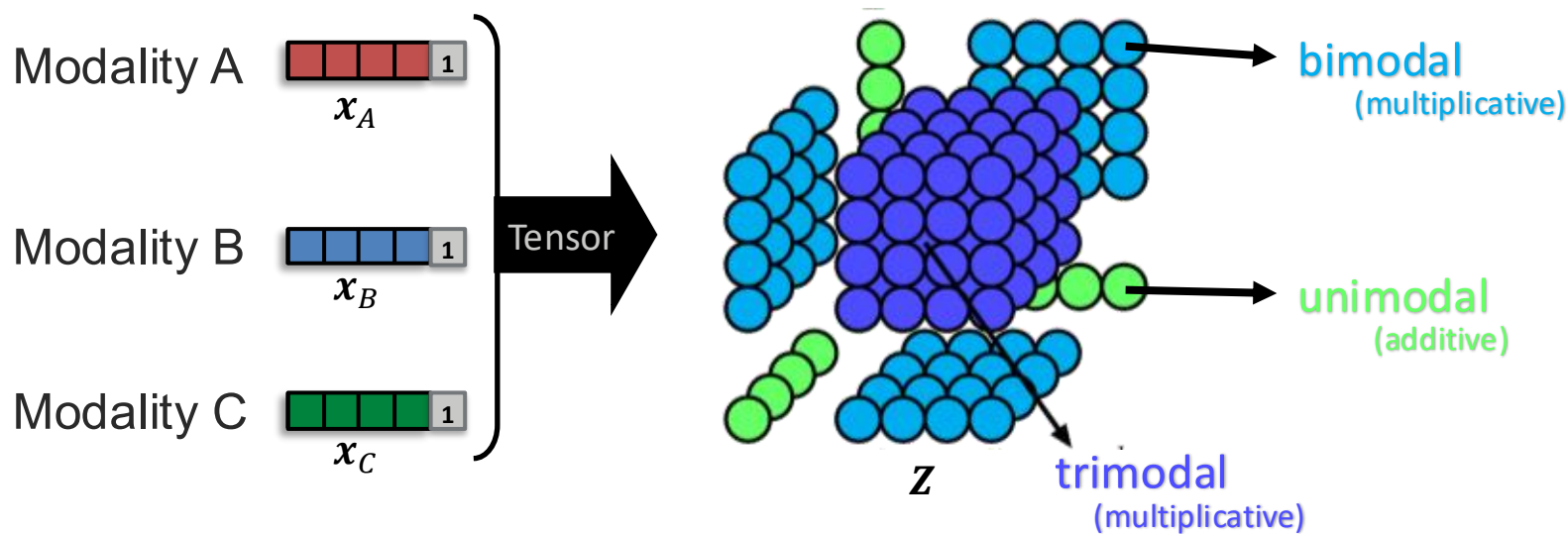
$$z = w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

Tensor Fusion



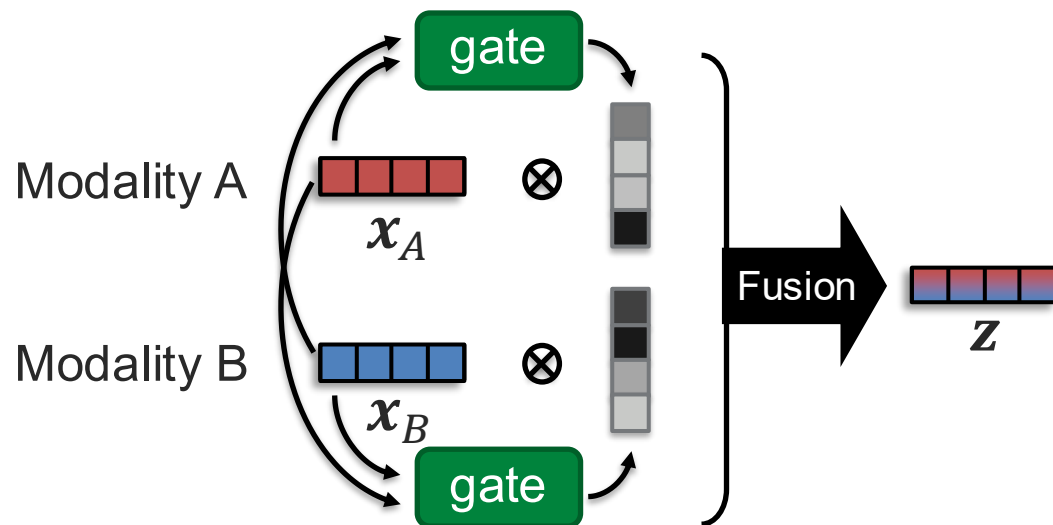
Tensor Fusion (bimodal):

$$Z = w([\mathbf{x}_A \ 1]^T \cdot [\mathbf{x}_B \ 1])$$



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

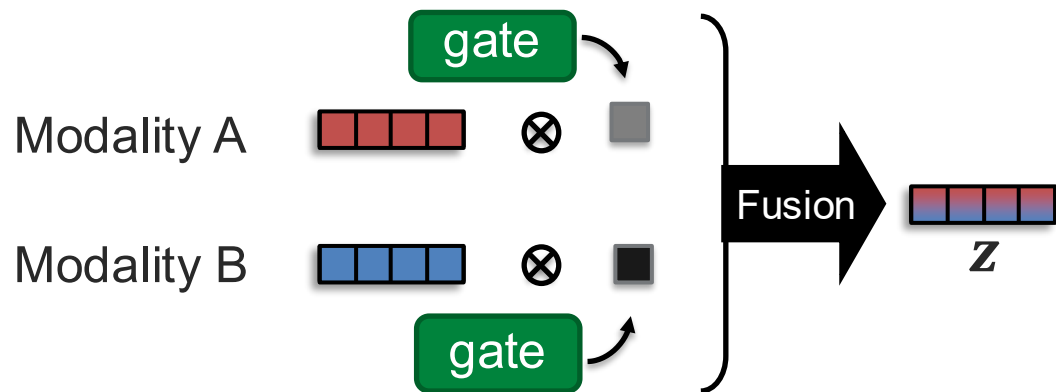
Gated Fusion



Example with additive fusion:

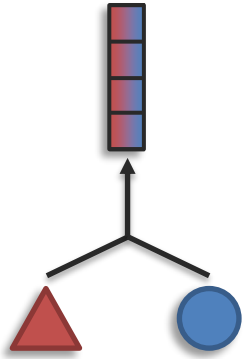
$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

→ g_A and g_B can be seen as attention functions



→ Gating output can be one weight for the whole modality

Summary: How To Multimodal Fusion



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

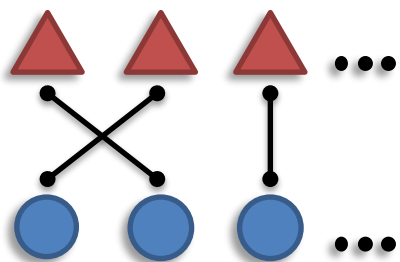


Challenge 2: Alignment

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

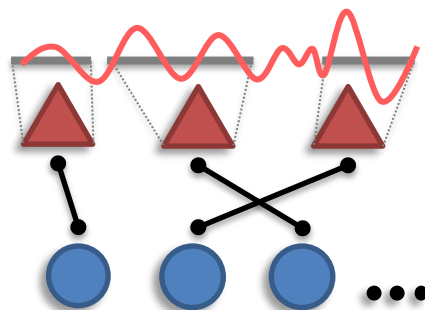
Sub-challenges:

Discrete connections



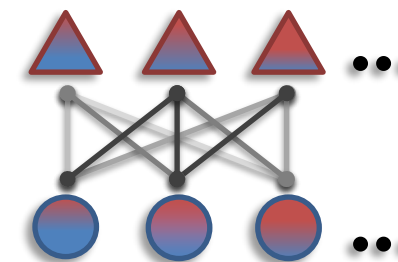
Explicit alignment
(e.g., grounding)

Continuous alignment



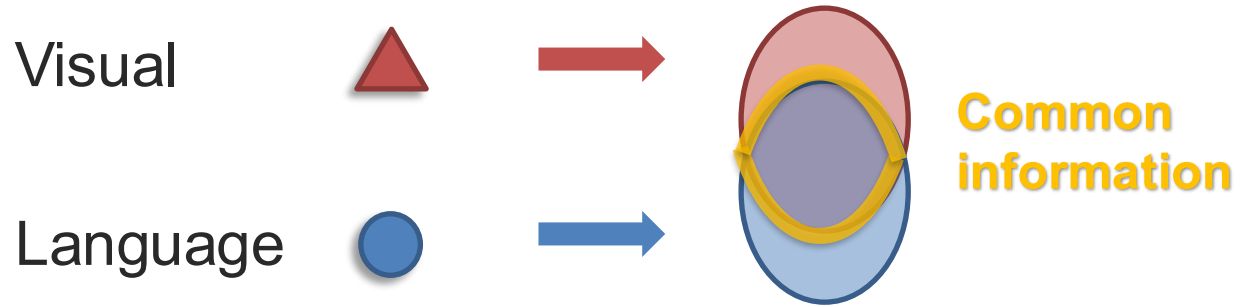
Granularity of
individual elements

Contextualized representation



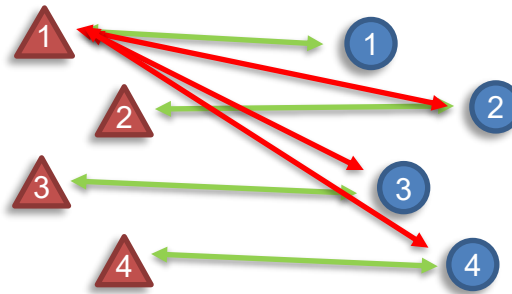
Implicit alignment
+ representation

Modality Connections

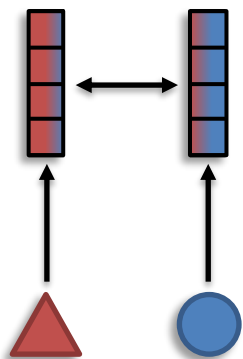


A **woman** reading **newspaper**

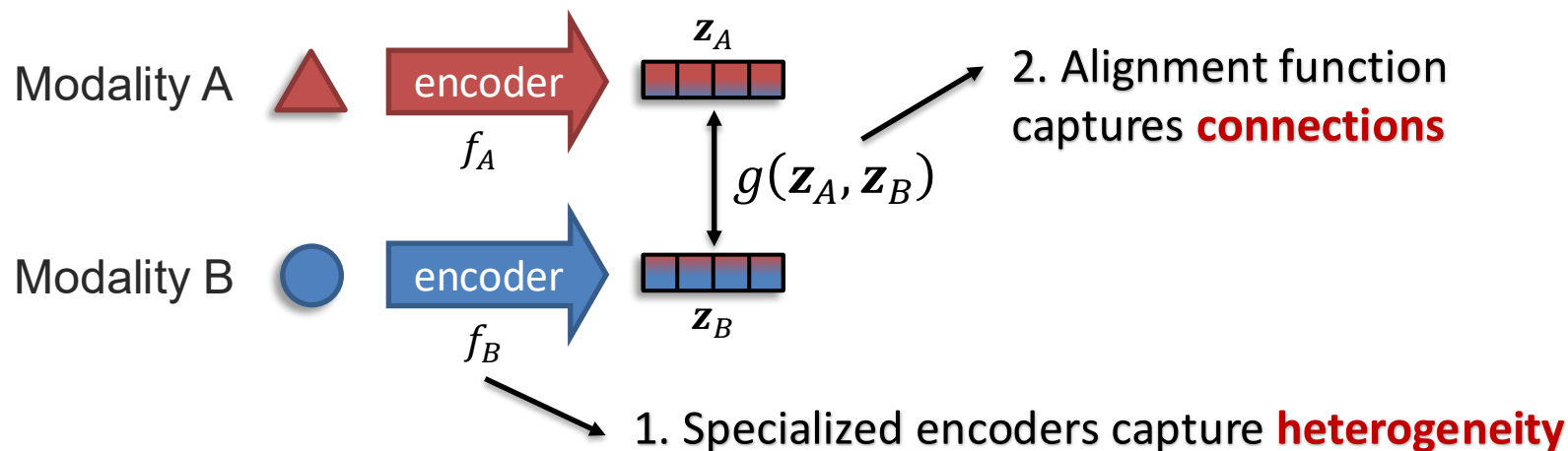
Supervision: Paired data



Aligned Representations



Definition: Learn multimodal representations aligned through their connections.

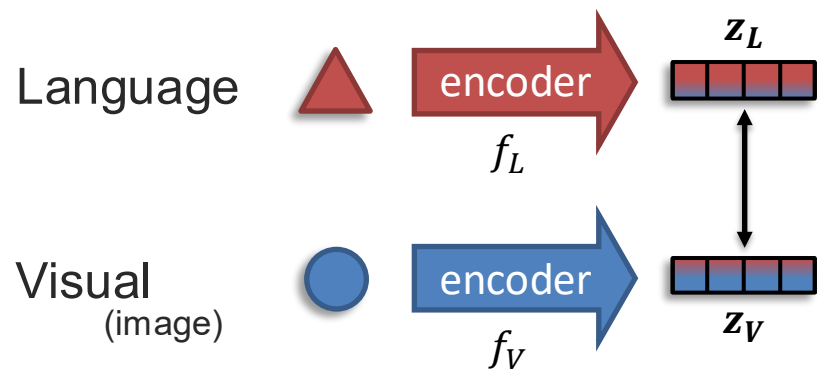


Learning with alignment function:

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

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

Contrastive Language Image Pretraining



Popular contrastive loss: InfoNCE

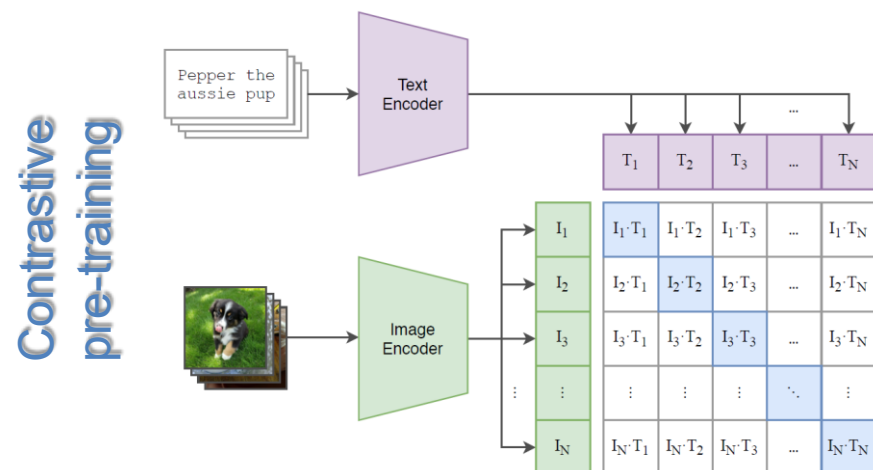
$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\text{sim}(z_A^i, z_B^i)}{\sum_{j=1}^N \text{sim}(z_A^i, z_B^j)}$$

positive pairs

negative pairs and positive pairs

Similarity function can be cosine similarity

Positive and negative pairs:



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

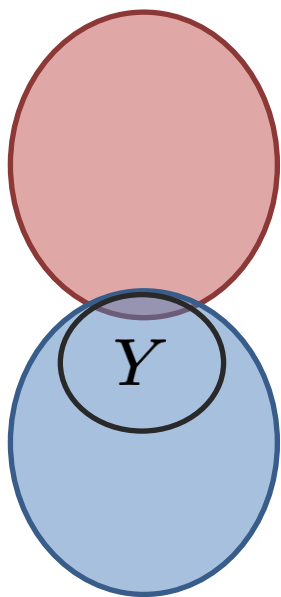
z_L and z_V are coordinated but not identical representation spaces

Multiview Redundancy and Contrastive Learning

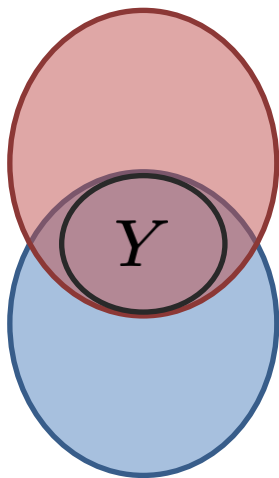
How much information should be shared?

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

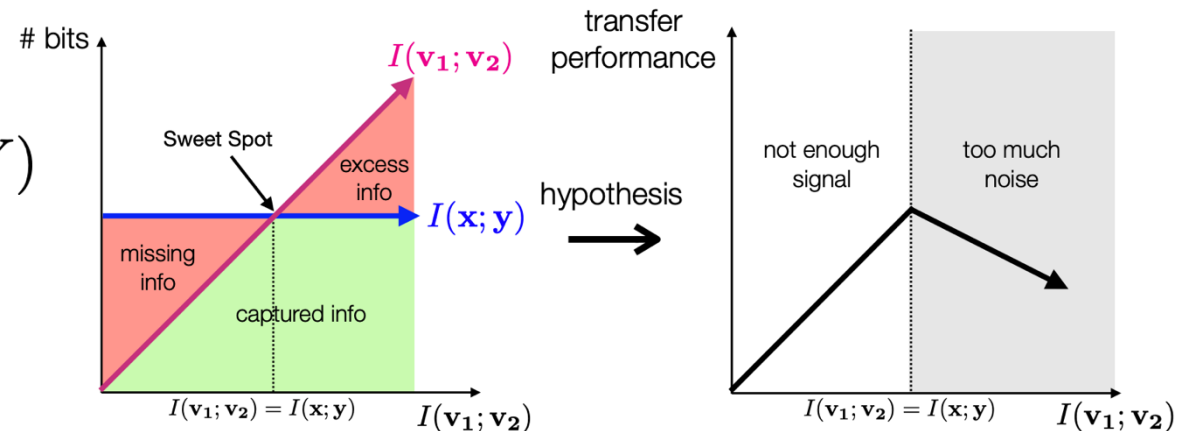
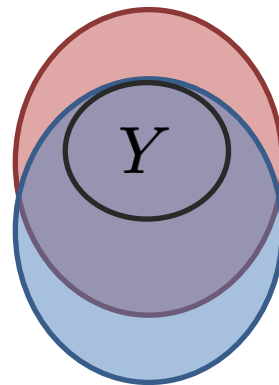
Not enough signal



Just right



Too much noise

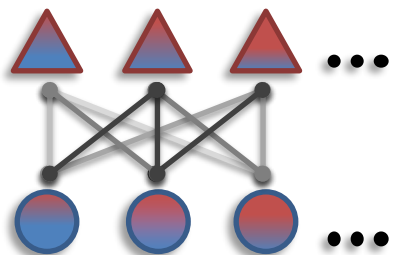


Multi-view redundancy
may not hold for
multimodal problems!

Large Multimodal Models

Part 1: Multimodal foundation model representations of text, video, audio

*It's just a privilege to
watch your mind at work.*

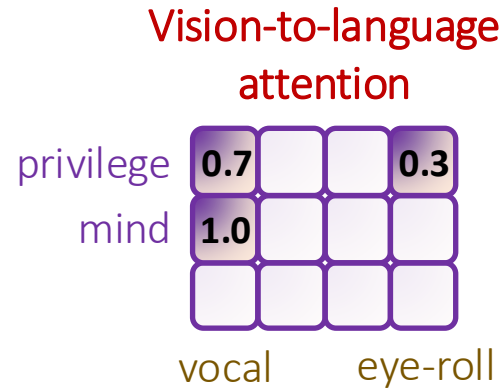
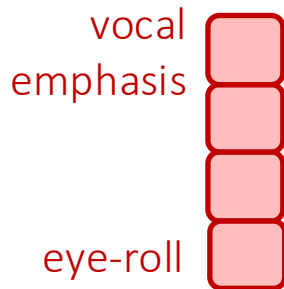
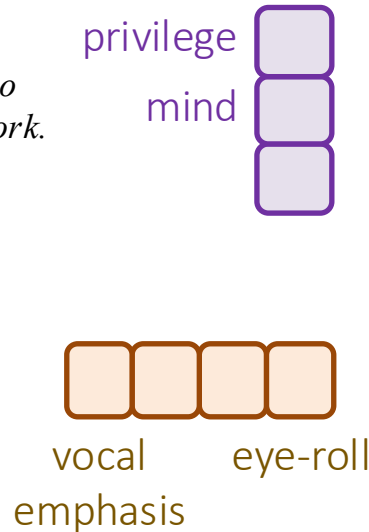


Multimodal
representation



Multimodal Transformers

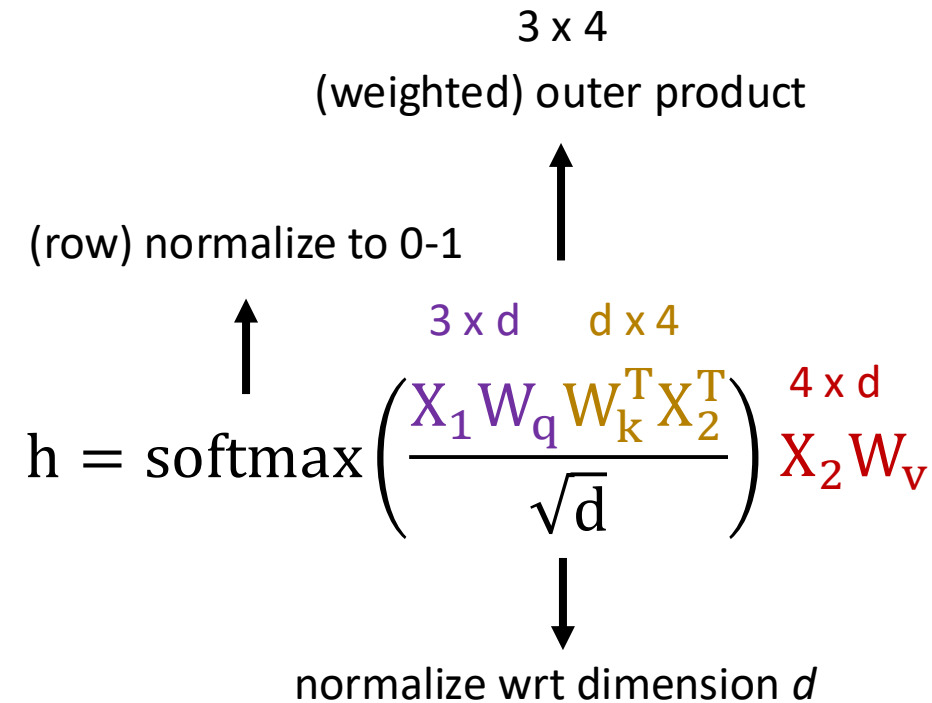
It's just a privilege to watch your mind at work.



privilege
mind



New **language** representation
contextualized with vision



Sarcasm

Visual-and-Language Transformer (ViLT)

Example of alignment between modalities:



a display of **flowers** growing out and over the retaining **wall** in front of **cottages** on a **cloudy** day.



flowers



wall



cottages



cloudy



a room with a **rug**, a **chair**, a **painting**, and a **plant**.



rug



chair



painting

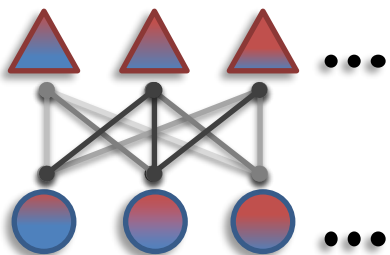


plant

Large Multimodal Models

Part 2: Adapting large language models for multimodal text generation

*It's just a privilege to
watch your mind at work.*



Multimodal
representation



*This person is being sarcastic.
They seem to be close friends.*

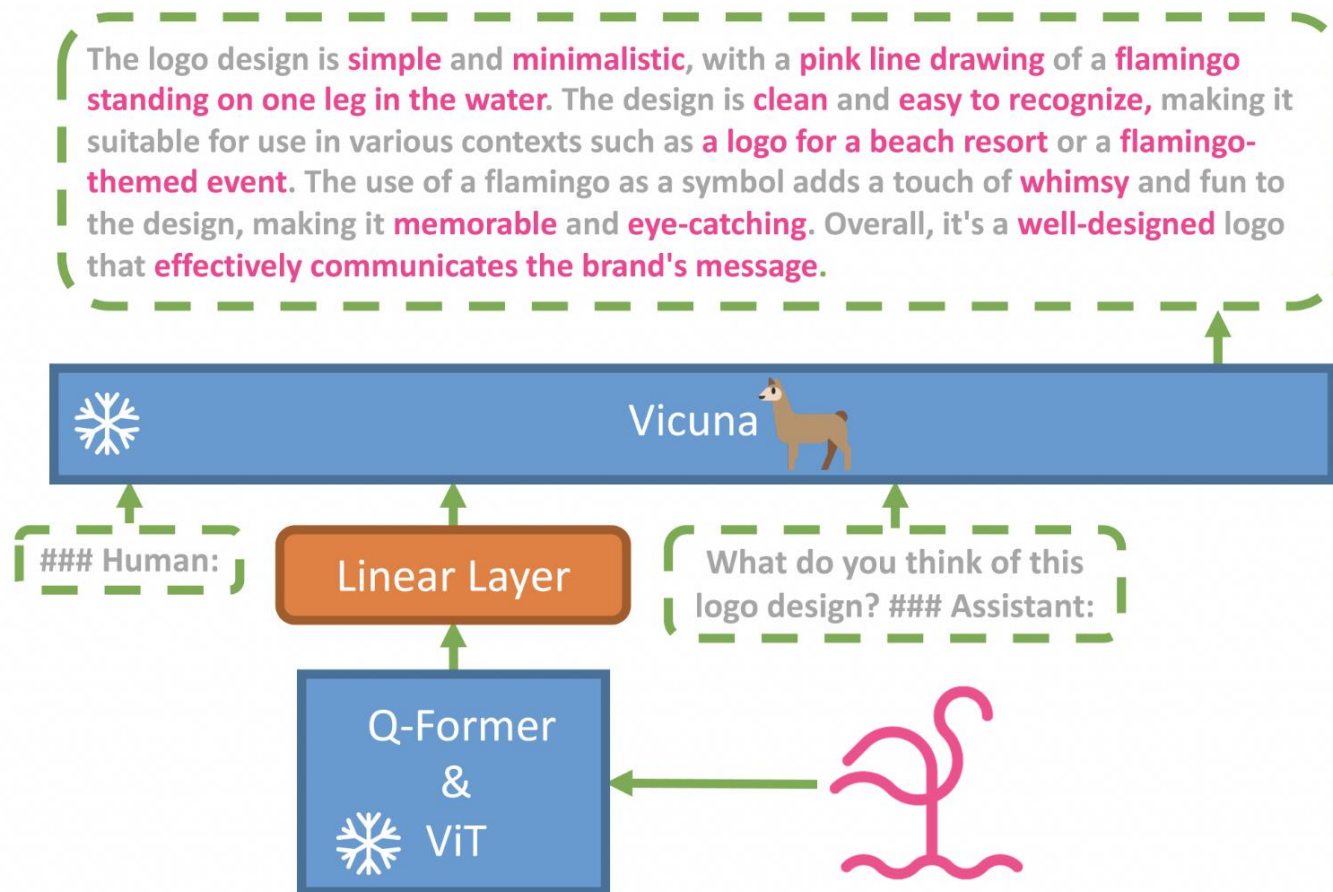


Adapting Large Multimodal Models

Mini-GPT4

Stage 1: **Alignment** using paired image-text data.

Stage 2: **Instruction tuning** using image + text instructions and example completions.



The architecture of MiniGPT-4.

Adapting Large Language Models

Flamingo



PP This is a picture of two teddy bears on the moon.

What are they doing?

PP They are having a conversation.

What object are they using?

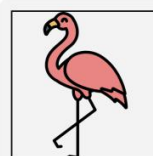
PP It looks like a computer.

Is this surprising?

PP Yes, it is surprising.

Why is this picture surprising to you?

PP I think it is surprising because teddy bears are not usually found on the moon.



What is the common thing about these three images?

PP They are all flamingos.

What is the difference between these three images?

PP The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.



PP This is an apple with a sticker on it.

What does the sticker say?

PP The sticker says "iPod".

Where is the photo taken?

PP It looks like it's taken in a backyard.

Do you think it is printed or handwritten?

PP It looks like it's handwritten.

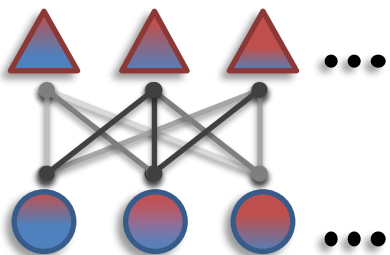
What color is the sticker?

PP It's white.

Large Multimodal Models

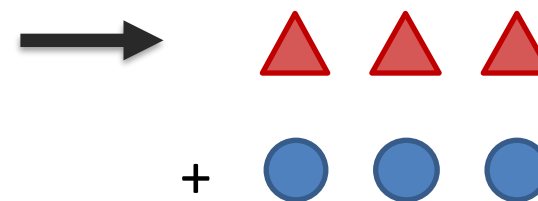
Part 3: Enabling text and image generation

*It's just a privilege to
watch your mind at work.*



Multimodal
representation

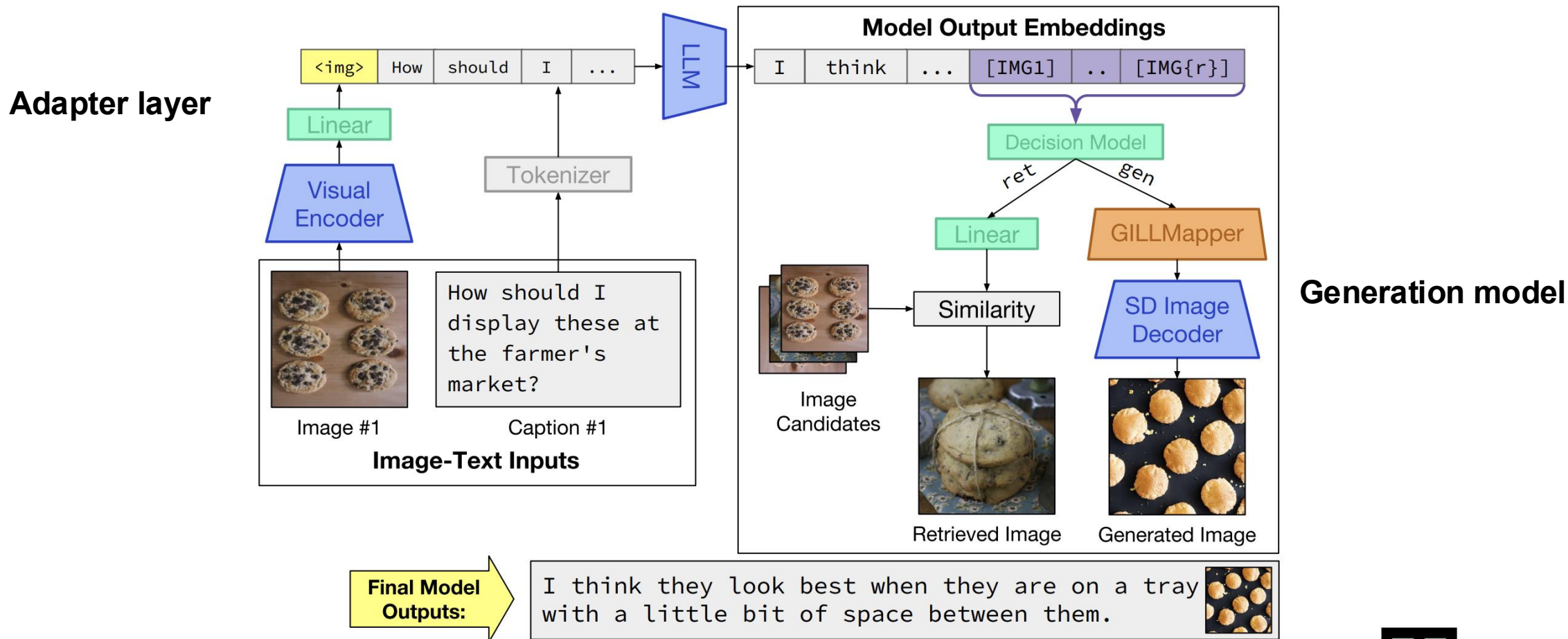
*This person is being sarcastic.
They seem to be close friends.*



*(quote previous episodes)
(highlight multimodal information)*

Grounding LMs for Multimodal Generation

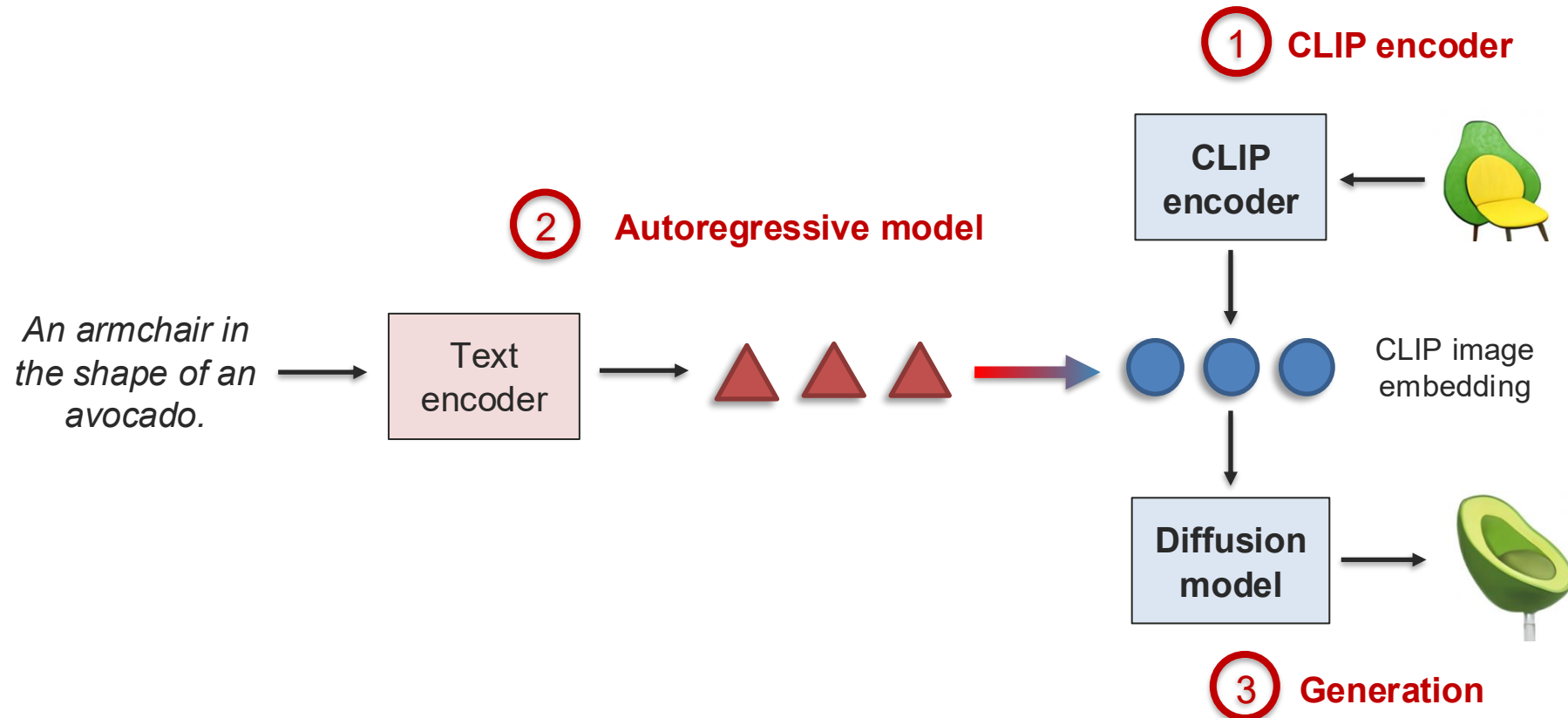
Large multimodal models with image generation



From Text to Multimodal Generation

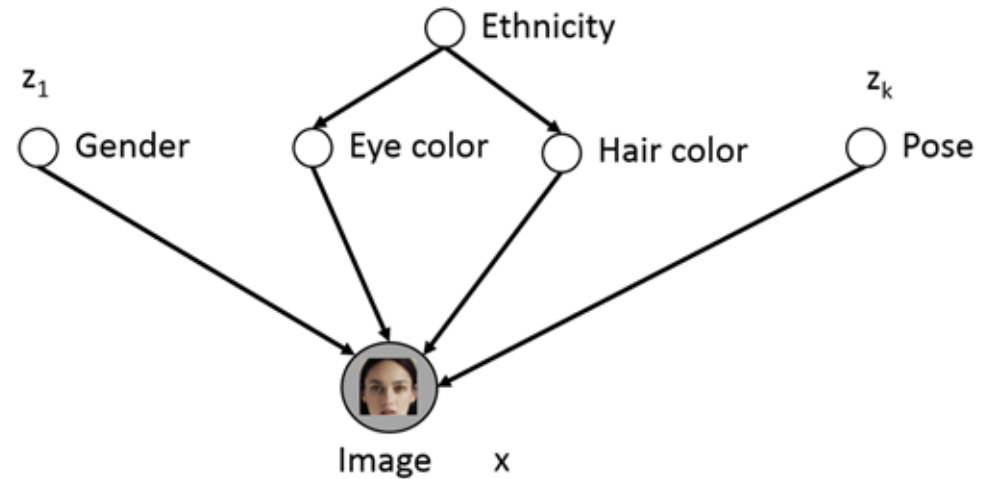
Directly training diffusion models with conditional information

Conditional latent variables are pretrained CLIP embeddings, then diffusion model to generate image.

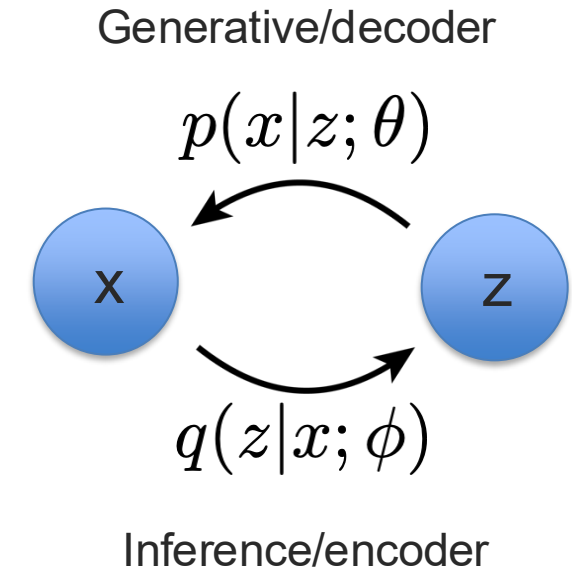
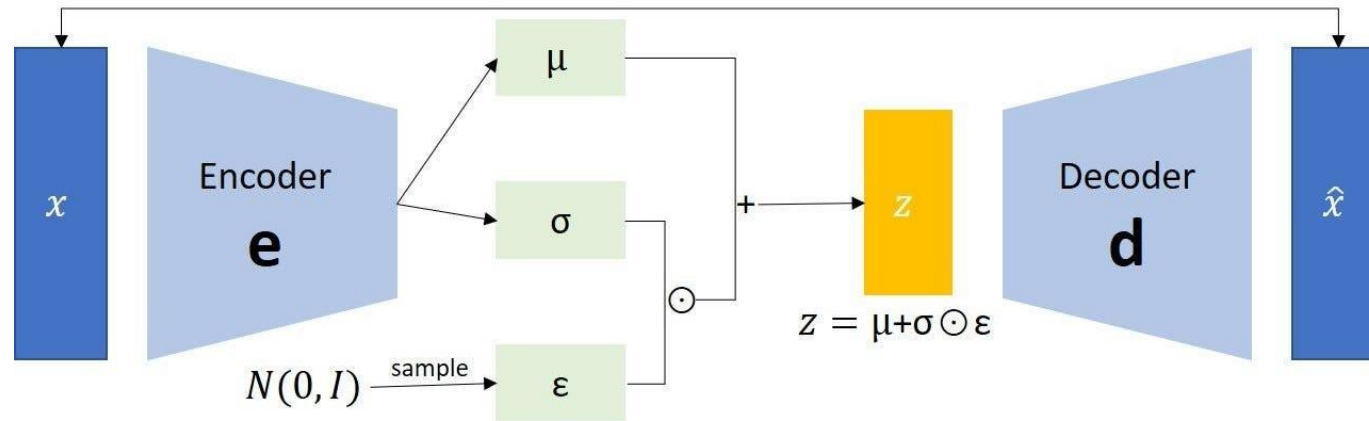


Latent Variable Generative Models

- Lots of variability in images x due to gender, eye color, hair color, pose, etc.
- However, unless images are annotated, these factors of variation are not explicitly available (latent).
- Idea: explicitly model these factors using latent variables z



Learning Parameters of VAEs

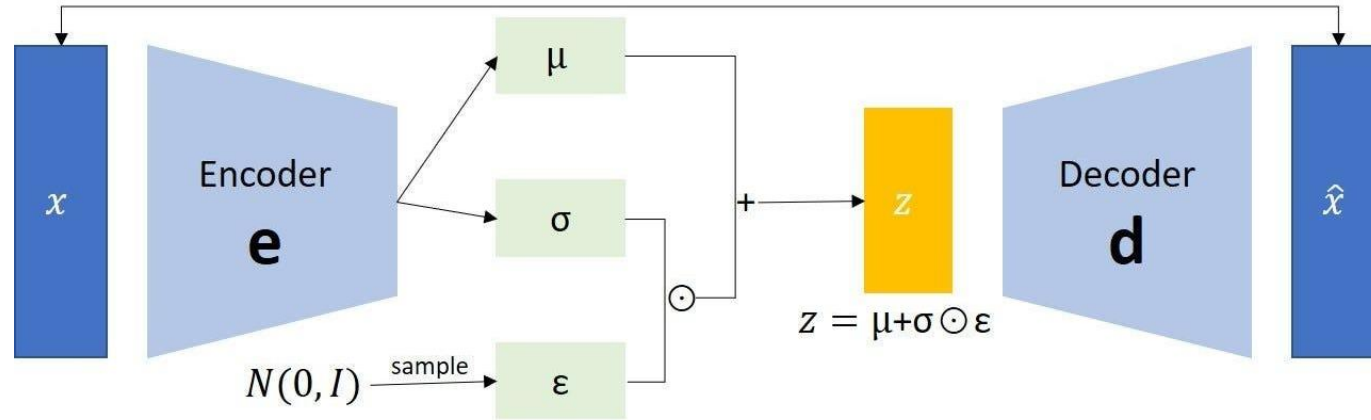


1. Take a datapoint x_i .
2. Map it to μ, σ using $q_\phi(z|x_i)$. **encoder**
3. Sample $\epsilon \sim N(0, I)$ and compute $\hat{z} = \mu + \sigma\epsilon$. **re-parameterize**
4. Reconstruct \hat{x} by sampling from $p(x|\hat{z}; \theta)$. **decoder**

$$\mathcal{L}(\mathbf{x}; \theta, \phi) = \underbrace{E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}; \theta)]}_{\text{reconstruction } \mathbf{x} \rightarrow \mathbf{z} \rightarrow \mathbf{x}} - \underbrace{D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))}_{\text{prior over } \mathbf{z}, \text{ standard Gaussian}}$$

Prior on \mathbf{z} : $\mathbf{z} \sim \mathcal{N}(0, I)$

Learning Parameters of VAEs

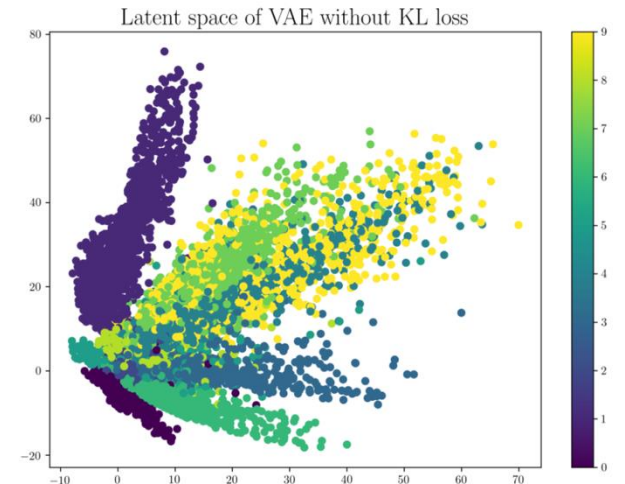
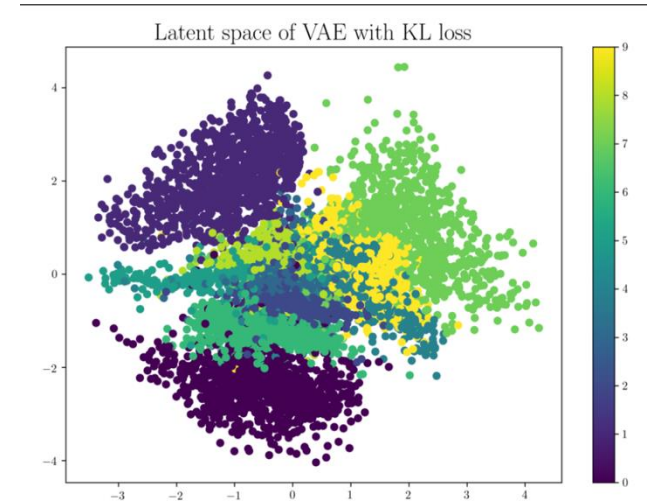


1. Take a datapoint x_i .
2. Map it to μ, σ using $q_\phi(z|x_i)$. **encoder**
3. Sample $\epsilon \sim N(0, I)$ and compute $\hat{z} = \mu + \sigma\epsilon$. **re-parameterize**
4. Reconstruct \hat{x} by sampling from $p(x|\hat{z}; \theta)$. **decoder**

$$\mathcal{L}(\mathbf{x}; \theta, \phi) = \underbrace{E_{q_\phi(z|x)}[\log p(\mathbf{x}|z; \theta)]}_{\text{reconstruction}} - \underbrace{D_{KL}(q_\phi(z|x) \| p(z))}_{\text{prior over } z, \text{ standard Gaussian}}$$

reconstruction
 $x \rightarrow z \rightarrow x$

prior over z ,
standard Gaussian

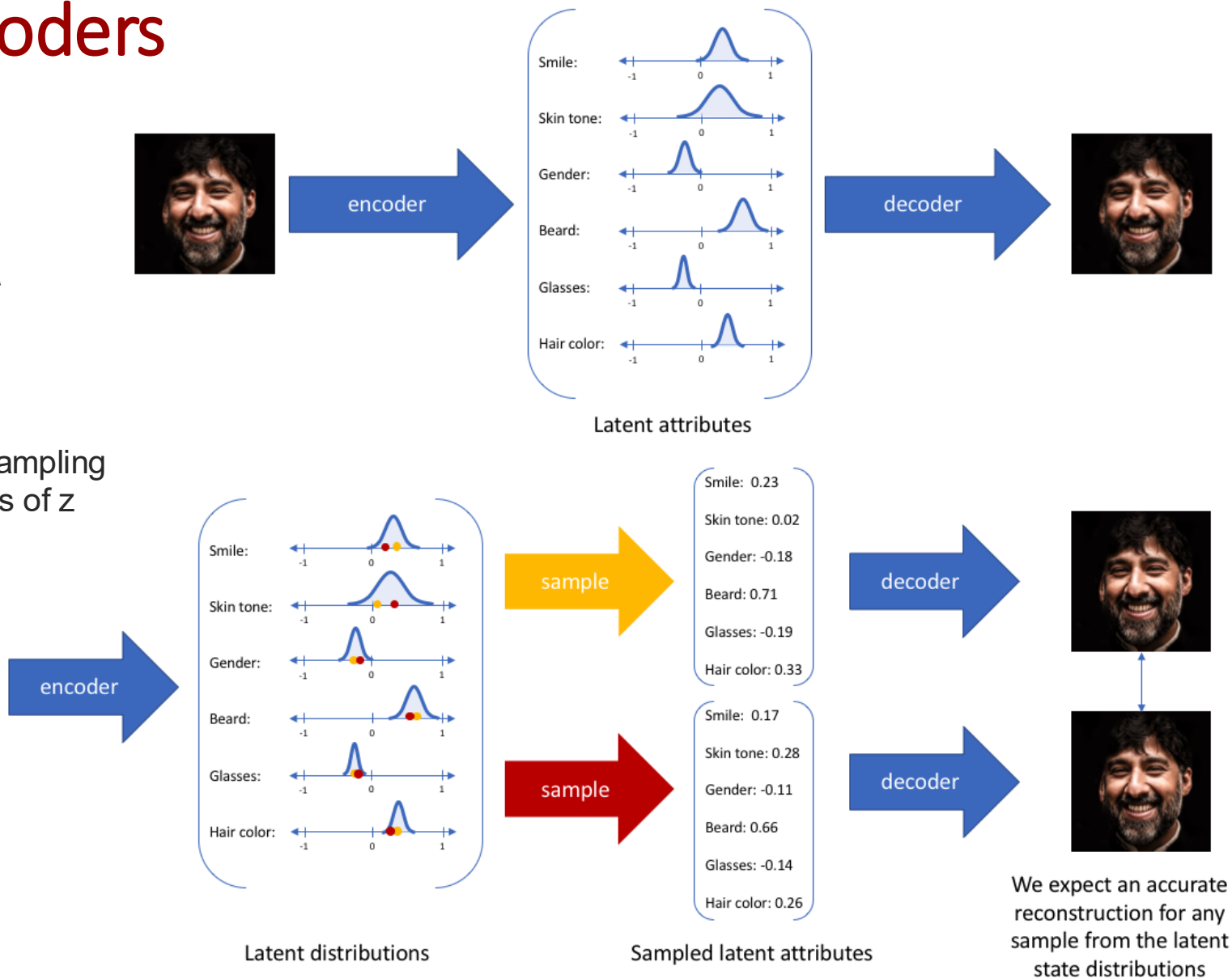


Prior on z : $\mathbf{z} \sim \mathcal{N}(0, I)$

Variational Autoencoders

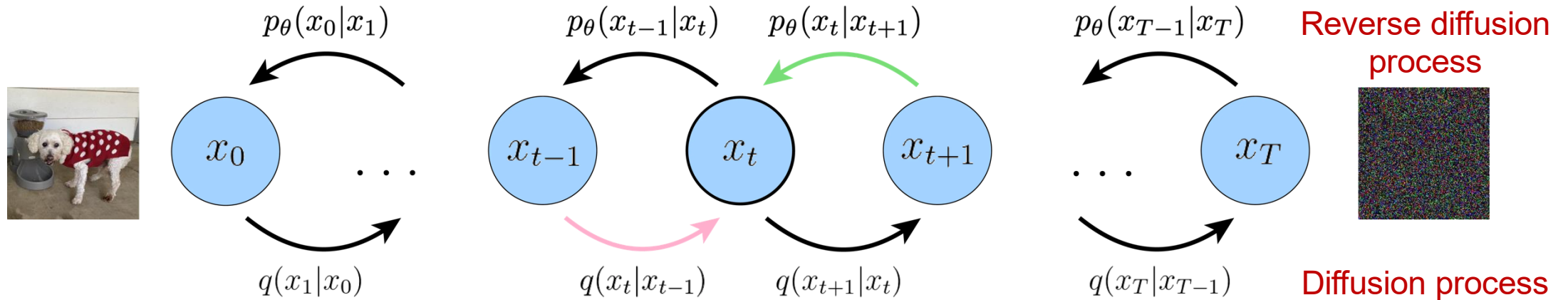
Key ideas:

1. Encoder decoder structure.
2. Simple latent variable z .
3. Complex $p(x|z)$ decoder via neural networks.
4. Reconstruction objective.
5. Prior over latent variable z :
 - smoother latent space permits sampling
 - disentangles different dimensions of z



From VAEs to Diffusion Models

Generative modeling via denoising

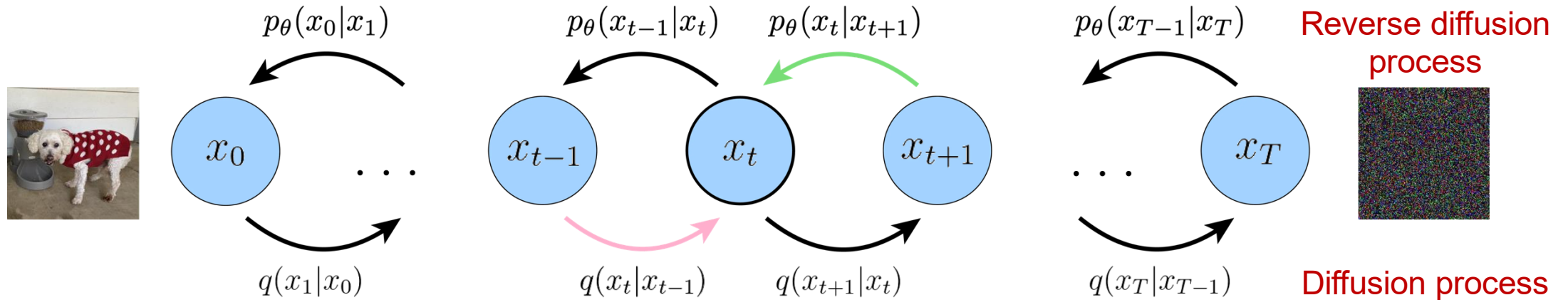


Similar to variational autoencoder, but:

1. The latent dimension is exactly equal to the data dimension.
2. Encoder q is not learned, but pre-defined as a Gaussian distribution centered around the output of previous timestep.
3. Gaussian parameters of latent encoders vary over time such that distribution of final latent is a standard Gaussian.

Learning Diffusion Models

Key idea: use variational inference



$$\begin{aligned}
 \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \\
 &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathcal{D}_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) || p(\mathbf{x}_T))}_{\text{prior matching term}} \\
 &\quad - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [\mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))]}_{\text{denoising matching term}}
 \end{aligned}$$

Multi-level VAE!

Comparisons of Generative Models

Feature	VAE	Diffusion Models	Flow Matching
Core Idea	Encode/decode with latent noise	Add noise and learn to reverse it	Learn a continuous flow from noise
Training Objective	Minimize reconstruction + KL loss	Learn the score (gradient) of data	Match a vector field (ODE-based)
Noise Handling	Noise in latent space	Progressive noise over time	Start from noise, smooth transform
Sampling Speed	Very fast (one pass)	Slow (many denoising steps)	Faster (solving an ODE)
Advantages	Simple, fast, interpretable	Very high-quality outputs	High quality + faster than diffusion
Disadvantages	Blurry samples, limited expressiveness	Expensive, slow sampling	Newer, still developing
Key Examples	VAE (2013), β -VAE	DDPM, Stable Diffusion	Flow Matching (2023), Rectified Flow

Text-to-Image Generation with Latent Diffusion

Semantic Synthesis on Flickr-Landscapes [21]



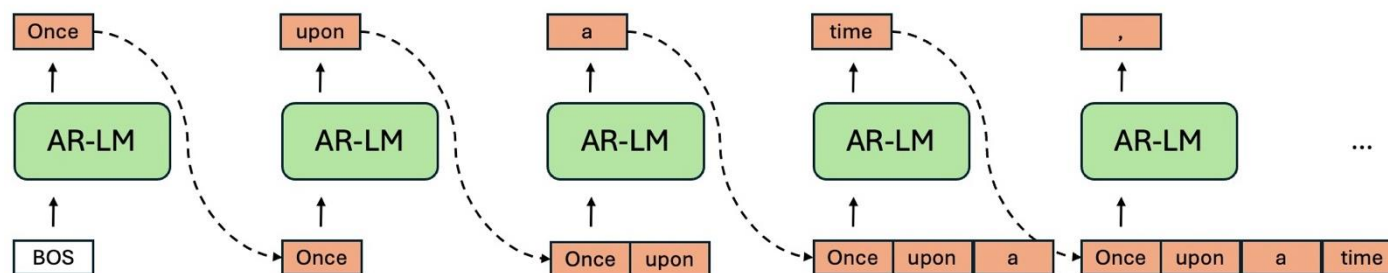
Diffusion Models for Text

Diffusion Language Models: Instead of going token by token, they iteratively refine and predict the *whole* sequence from a noise following a non-autoregressive decoding process.

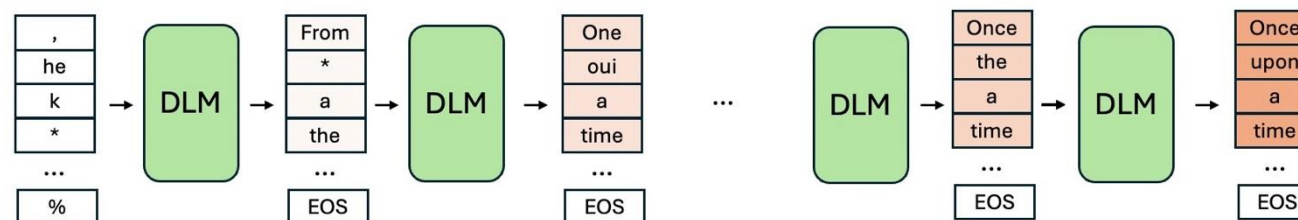
Key advantages:

- Can fix earlier mistakes
- More controllable
- More diversity
- Faster

Autoregressive Language Model



Diffusion Language Model



Midterm Structure

6 problems, 60 minutes.

5 problems:

1. MCQ (10, 2 points each = 20 points)
2. Short answers (4, 5 points each = 20 points)
3. Multimodal Fusion (20 points)
4. Multimodal LLMs (20 points)
5. Multimodal Generation (20 points)
6. Bonus open questions (10 points)

Assignments for This Coming Week

HW3 due today.

Project proposals will be graded and released this week.

Project midterm instructions will be out this week.

Midterm this Thursday.