# Modality Translation through Conditional Encoder-Decoder

Hyunsoo Lee, Yoonsang Lee, Maria Pak, Jinri Kim
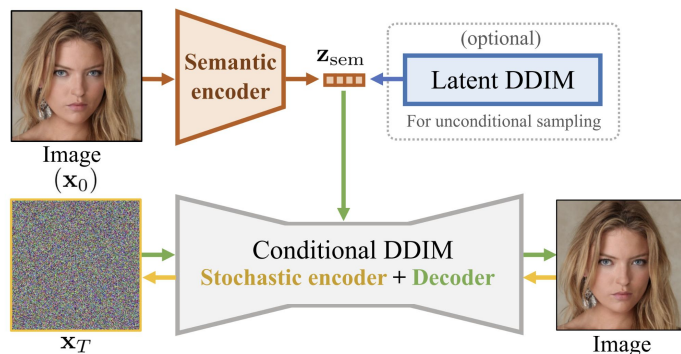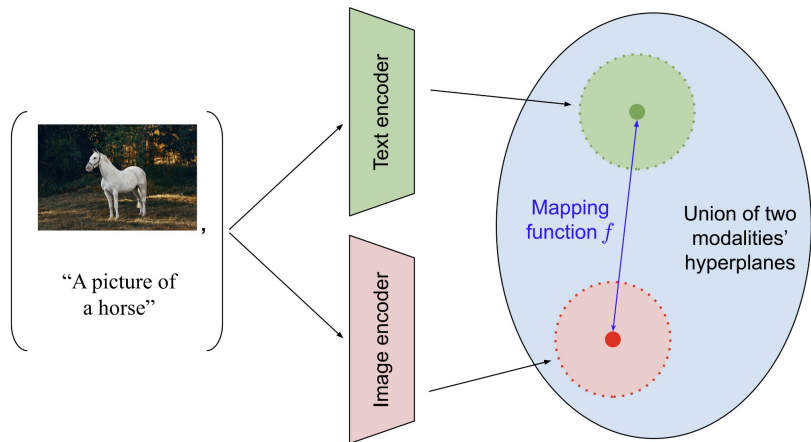
**Seoul National University**
Graduate School of Data Science

# Problem Formulation

# Multimodal Feature Representation

- Recent models are designed and trained specifically for individual tasks
  - Need general-purpose model like CLIP
  - However, cosine similarity between $\mathbf{z}_{\text{txt}}$ and $\mathbf{z}_{\text{img}}$ is only 0.3 😭



💡 Propose a model architecture utilizing **conditional DDIM**
   ↳ Enhance cosine similarity and hope to perform better in general multi-modal downstream tasks

# Conditional Denoising Diffusion Implicit Models

- Distribution modeling using diffusion process
- Forward process
  - Perturbation of data using **gaussian noise** with total $T$ steps
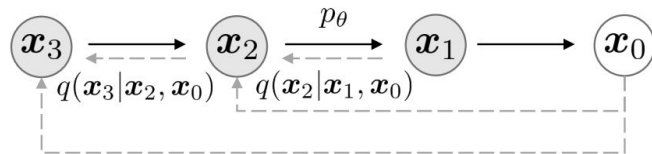  - $\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\epsilon$



- Backward process
  - **Denoising** step : recovering the original data
  - Defined as deterministic process
  - $\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\left(\dfrac{\mathbf{x}_t - \sqrt{1-\alpha_t}\epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})}{\sqrt{\alpha_t}}\right) + \sqrt{1-\alpha_t}\epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})$

- Training objectives
  - Maximizing the log-likelihood of estimated data distribution
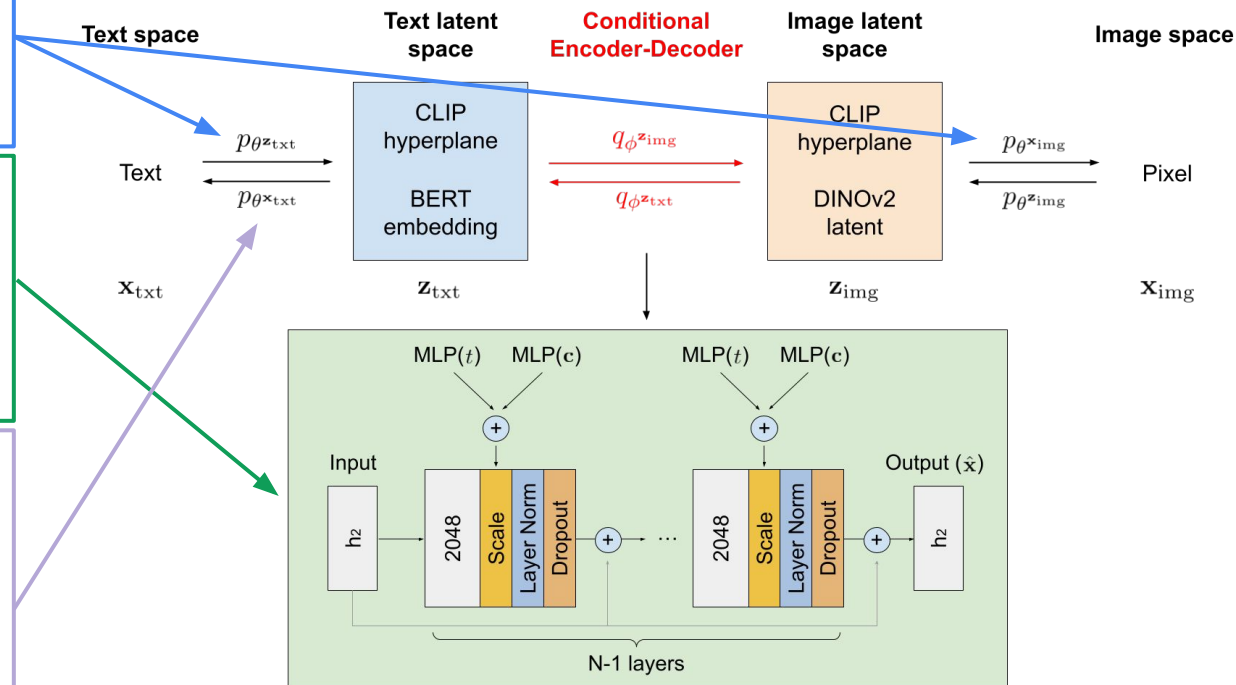  - $\mathbb{E}_{t,\mathbf{x}_0,\epsilon}\left[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})\|^2\right]$

# Conditional Encoder-Decoder Architecture

Step 1: latent encoders $p_{\theta^{\mathbf{x}_{\mathrm{txt}}}}$ and $p_{\theta^{\mathbf{x}_{\mathrm{img}}}}$ extract the text and image embeddings $\mathbf{z}_{\mathrm{txt}}$ and $\mathbf{z}_{\mathrm{img}}$ respectively.

Step 2: conditional encoder-decoder enables bidirectional transformation between two latent vectors (**DDIM**).

Step 3: finally, the transformed embeddings are converted back to the image and text modalities using the latent decoders, $p_{\theta^{\mathbf{x}_{\mathrm{txt}}}}$ and $p_{\theta^{\mathbf{x}_{\mathrm{img}}}}$.
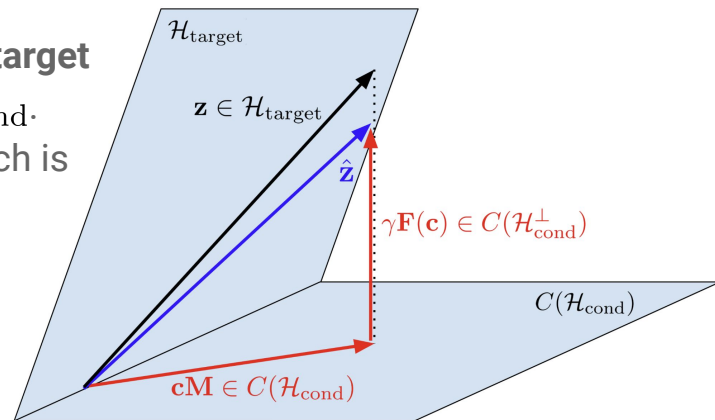
# Sampling Strategy

- Notation
  - Target vector : $\mathbf{z} \in \mathcal{H}_{\text{target}}$ (*e.g.* text embedding)
  - Estimated target vector : $\hat{\mathbf{z}}$
  - Condition vector : $\mathbf{c} \in \mathcal{H}_{\text{cond}}$ (*e.g.* image embedding)

- Novel sampling method

$$\hat{\mathbf{z}} = \mathbf{c}\mathbf{M} + \gamma \mathbf{F}(\mathbf{c}), \quad \mathbf{F}(\mathbf{c}) = \text{Normalize}(\text{CDIM}(\mathbf{c}))$$

- Justification (1)
  - After the optimization, first term predicts a **projection of target** vector ($\mathbf{z}$) onto the **column space** of the hyperplane $\mathcal{H}_{\text{cond}}$.
  - The second term is guided to predict a **residual term** which is perpendicular to the column space.

# Sampling Strategy

- Justification (2)
  - For a given value of hyperparameter $\gamma > 0$, our method ensures that **cosine similarity** between target vector $\mathbf{z}$ and estimated target vector $\hat{\mathbf{z}}$ to be greater or equal than constant $\alpha$ with probability at least

$$P(\textit{Cosine-Sim}(\mathbf{z}, \hat{\mathbf{z}}) \geq \alpha)$$

$$= 1 - \int_{-1}^{\beta} \frac{\Gamma(h_2/2 + 1/2)}{\sqrt{\pi}\Gamma(h_2/2)} (1 - u^2)^{h_2/2 - 1} du$$

  - We empirically choose $\gamma = 1.0$.

# Training Details

- Training objectives
  - Weighted sum of conditional DDIM loss and reconstruction loss
  - Motivated from pix2pix

$$\min_{q_\phi, \mathbf{M}} \lambda_1 \mathbb{E}_{\mathbf{z}, \mathbf{c}}[\|\mathbf{z} - \hat{\mathbf{z}}\|_1] + \lambda_2 \mathbb{E}_{t, \mathbf{z}, \epsilon}\left[\|\epsilon - \epsilon_\theta\left(\mathbf{z}_t, t, \mathbf{c}\right)\|^2\right]$$

- Details
  - Applied classifier-free guidance
  - Uses dropout and residual connections
  - Adam optimizer with learning rate $10^{-4}$, weight decay $10^{-2}$
  - $\lambda_1 = 1.0, \lambda_2 = 2.0$
- Metric
  - Modality translation task
    - Cosine similarity
  - Downstream tasks
    - Widely used metrics for each task

# Results

# Modality Translation Results

**Cosine similarity between $\mathbf{z}$ and $\hat{\mathbf{z}}$ compared with other baselines**

| Dataset | MS-COCO [16] | | CC3M [29] | |
|---|---|---|---|---|
| Method | $\text{Sim}_{txt}$ | $\text{Sim}_{img}$ | $\text{Sim}_{txt}$ | $\text{Sim}_{img}$ |
| LAFITE [34] | 0.0965 | - | 0.0912 | - |
| CLIP-GEN [33] | 0.3042 | - | 0.2896 | - |
| VDLGAN [12] | 0.6104 | 0.7655 | 0.6237 | 0.7105 |
| Ours | **0.8394** | **0.8233** | **0.7389** | **0.7443** |

Table 2. Results of cosine similarity between $\mathbf{z}$ and $\hat{\mathbf{z}}$ from text, image modalities. We used CLIP ViT-B/32 [25] model for both $p_\theta \mathbf{z}_{\text{txt}}$ and $p_\theta \mathbf{z}_{\text{img}}$. Bold number indicates the best performance among the column and '-' indicates unavailability.

**Using various type of latent encoders**

| Dataset | | MS-COCO [16] | |
|---|---|---|---|
| $p_\theta \mathbf{z}_{\text{txt}}$ | $p_\theta \mathbf{z}_{\text{img}}$ | $\text{Sim}_{txt}$ | $\text{Sim}_{img}$ |
| CLIP ViT-L/14 [25] | CLIP ViT-L/14 | 0.7765 | 0.8192 |
| CLIP ViT-L/14 | BERT [3] | 0.9745 | 0.7796 |
| DINOv2 [21] | CLIP-RNx50 | 0.7917 | 0.5207 |

Table 3. Results of cosine similarity between $\mathbf{z}$ and $\hat{\mathbf{z}}$ from text, image modalities using various type of latent encoders.
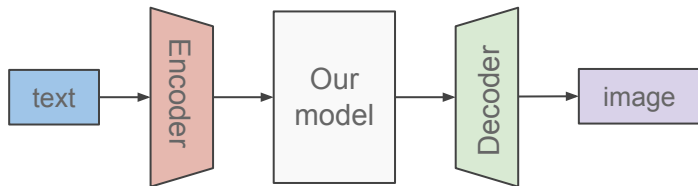
**Cross-domain modality translation result**

| Dataset | MS-COCO [16] → CelebA [17] | | CC3M [29] → MS-COCO | |
|---|---|---|---|---|
| $p_\theta \mathbf{z}_{\text{txt}}, p_\theta \mathbf{z}_{\text{img}}$ | $\text{Sim}_{txt}$ | $\text{Sim}_{img}$ | $\text{Sim}_{txt}$ | $\text{Sim}_{img}$ |
| CLIP ViT-B/32 [25] | 0.8237 | 0.5974 | - | - |
| CLIP ViT-L/14 | 0.6885 | 0.5993 | 0.6817 | 0.7300 |

Table 4. Results of cross-domain experiments. We train our model on bigger dataset, and measure the cosine similarity between ground truth and predict target vector on a relatively small dataset.
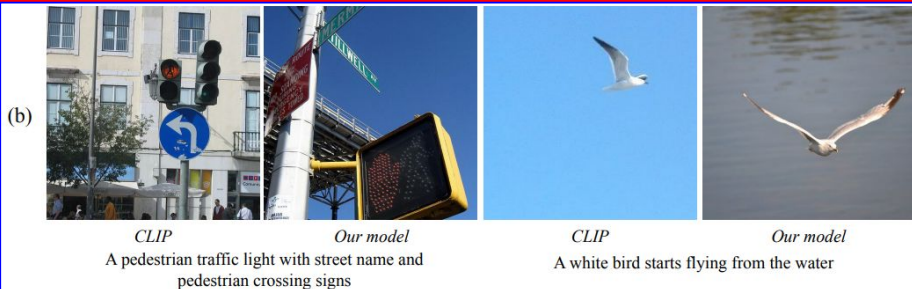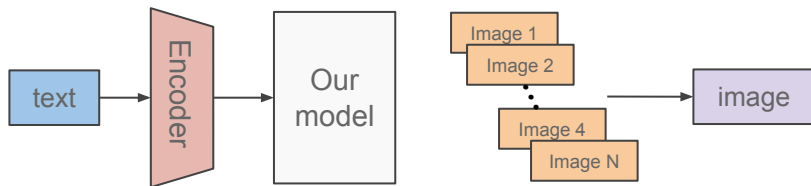
# Downstream Tasks (1)

- **Text-to-image Generation**
  - Encoder : CLIP ViT-L/14
  - Decoder : Karlo

  text → Encoder → Our model → Decoder → image

- **Image Retrieval**
  - Encoder : CLIP ViT-B/32

  text → Encoder → Our model → [Image 1, Image 2, ... Image 4, Image N] → image



(a)

A white boat floating on a lake under mountains

A woman is walking a dog in the city

A room with blue walls and a white sink and door

A large passenger airplane flying through the air

(b)

*CLIP*  *Our model*  *CLIP*  *Our model*

A pedestrian traffic light with street name and pedestrian crossing signs

A white bird starts flying from the water

*CLIP*  *Our model*  *CLIP*  *Our model*

A man that is standing in the grass with a soccer ball

Graffiti covered train stopped at the train platform

# Downstream Tasks (2)

- ## Image Captioning

  - ### Encoder : CLIP RNx50

  - ### Decoder : ClipCap + CapDec

| Method | B@1 | B@4 | R-L |
|--------|-----|-----|-----|
| ClipCap [19] | **74.7** | **33.5** | - |
| CapDec [20] | 69.2 | 26.4 | **51.8** |
| Ours + ClipCap | 65.9 | 23.6 | 47.7 |
| Ours + CapDec | 67.7 | 25.5 | 48.7 |

Table 5. Results for image captioning on MS-COCO dataset.

- ## Image Classification

  - ### Encoder : CLIP ViT-L/14 & B/32

| Method | Accuracy (%) |
|--------|-------------|
| CLIP ViT-L/14 [25] | **85.05** |
| CLIP ViT-B/32 | **69.69** |
| Ours + CLIP ViT-L/14 | 77.11 |
| Ours + CLIP ViT-B/32 | 60.74 |

Table 6. Results for image classification on CIFAR-10 dataset.



(a)

Ours + CapDec — A man holding a tennis racquet on a tennis court
Ours + ClipCap — A man holding a tennis racquet on a tennis court

References
A person holding a tennis racket in the air on a tennis court
A man with a hat and sunglasses playing tennis
A man holding a tennis racquet on a tennis court
A man in sunglasses and a hat is getting ready to hit a tennis ball
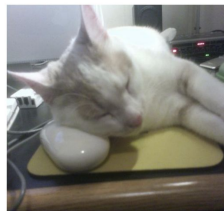A tennis player hits the ball back to his opponent

(b)

Ours + CapDec — A little girl that is holding a toothbrush in her mouth
Ours + ClipCap — A child is brushing her teeth with a toothbrush

References
A small girl with long hair brushing her teeth
A little girl brushing her teeth with an electric toothbrush
a close up of a small child brushing her teeth
A girl in pajamas brushing her teeth with an crayon toothbrush
A little girl brushing her teeth with a tooth brush

(c)

Ours + CapDec — A cat laying on top of a laptop computer
Ours + ClipCap — A cat laying on top of a laptop

References
A cat that is laying with its head down on a mouse
A white cat laying on the computer mouse
A white cat is taking a nap on a mouse
a kitty sleeping on a mouse pad and a mouse
A cat is sleeping on a desk with its head on a computer mouse

(d)

Ours + CapDec — A man is preparing food in a kitchen
Ours + ClipCap — A person in a kitchen baking food in a oven

References
a person with a black oven mit is taking a pan out of the oven
A person reaches into an oven to take out some muffins
A person getting muffins out of an oven
A man in black jacket removing tin of muffins from oven
A muffin tray that is inside of a oven

# Limitations / Future Study

- Extend to other modalities
  - Using audio modality

- Exploit different pretrained models
  - Pretrained Encoder / Decoder

- Attack various downstream tasks
  - Visual Question Answering / Image-Document Retrieval / Text Retrieval

- Reduce time complexity
  - Lightweight models / Efficient sampling strategy

# References

- DDIM (Song, Jiaming et al. "Denoising Diffusion Implicit Models." *ArXiv* abs/2010.02502 (2020): n. pag.)

- Diffusion Autoencoders (Preechakul, Konpat et al. "Diffusion Autoencoders: Toward a Meaningful and Decodable Representation." *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021): 10609-10619.)

- CLIP (Radford, Alec et al. "Learning Transferable Visual Models From Natural Language Supervision." *International Conference on Machine Learning* (2021).)

- VDLGAN (Kang, Minsoo et al. "Variational Distribution Learning for Unsupervised Text-to-Image Generation." *ArXiv* abs/2303.16105 (2023): n. pag.)

- LAFITE (Zhou, Yufan et al. "Towards Language-Free Training for Text-to-Image Generation." *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021): 17886-17896.)

- MS-COCO (Lin, Tsung-Yi et al. "Microsoft COCO: Common Objects in Context." *European Conference on Computer Vision* (2014).)