



HYPER-TRANSFORMING LATENT DIFFUSION MODELS

Ignacio Peis

Technical University of Denmark
ipeaz@dtu.dk



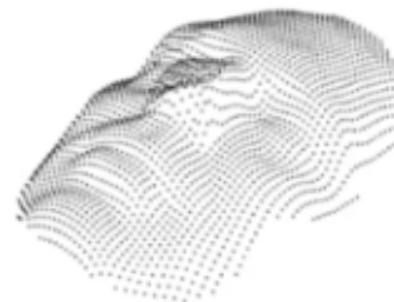
Motivation

- We typically discretized data that are continuous in nature.

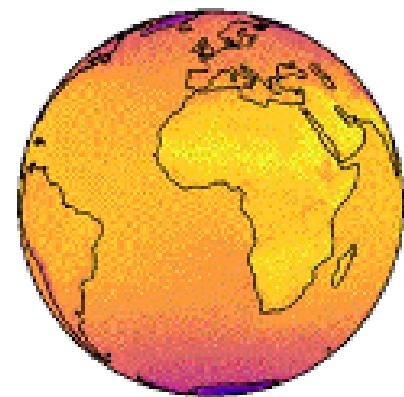
2D Images



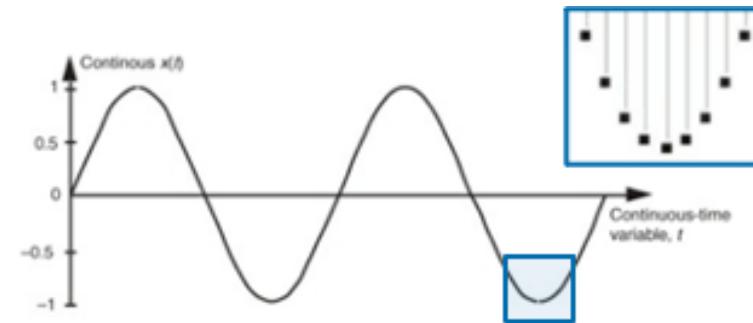
3D Images



Polar data



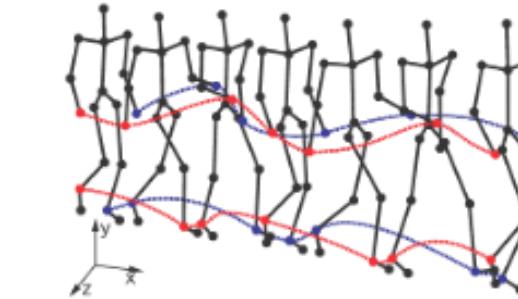
Time series



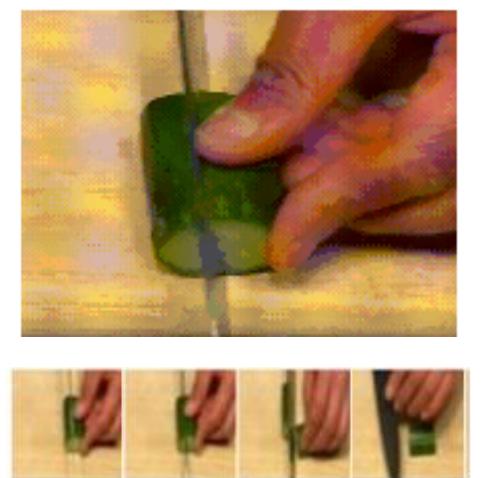
Audio



Motion sequences



Video



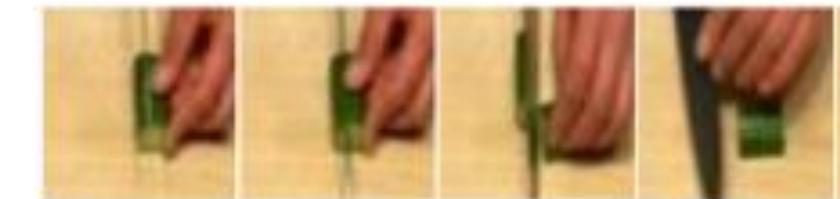
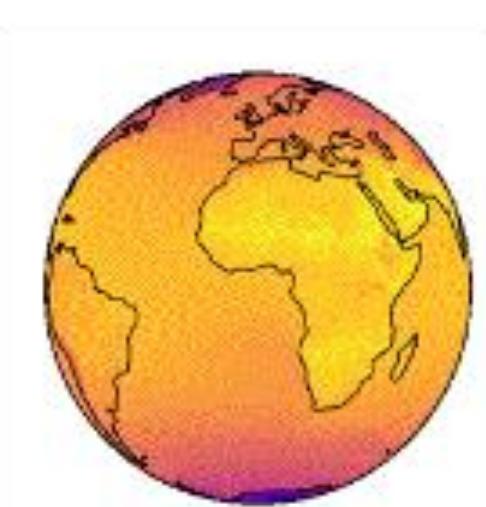
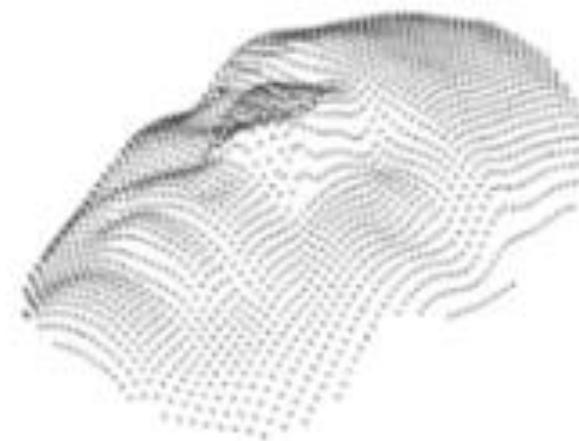
Spatial

Temporal

Spatio-temporal

Motivation

- Real data can be expressed as a function over continuous coordinate systems.



$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^3, f(x_1, x_2) = (r, g, b)$$

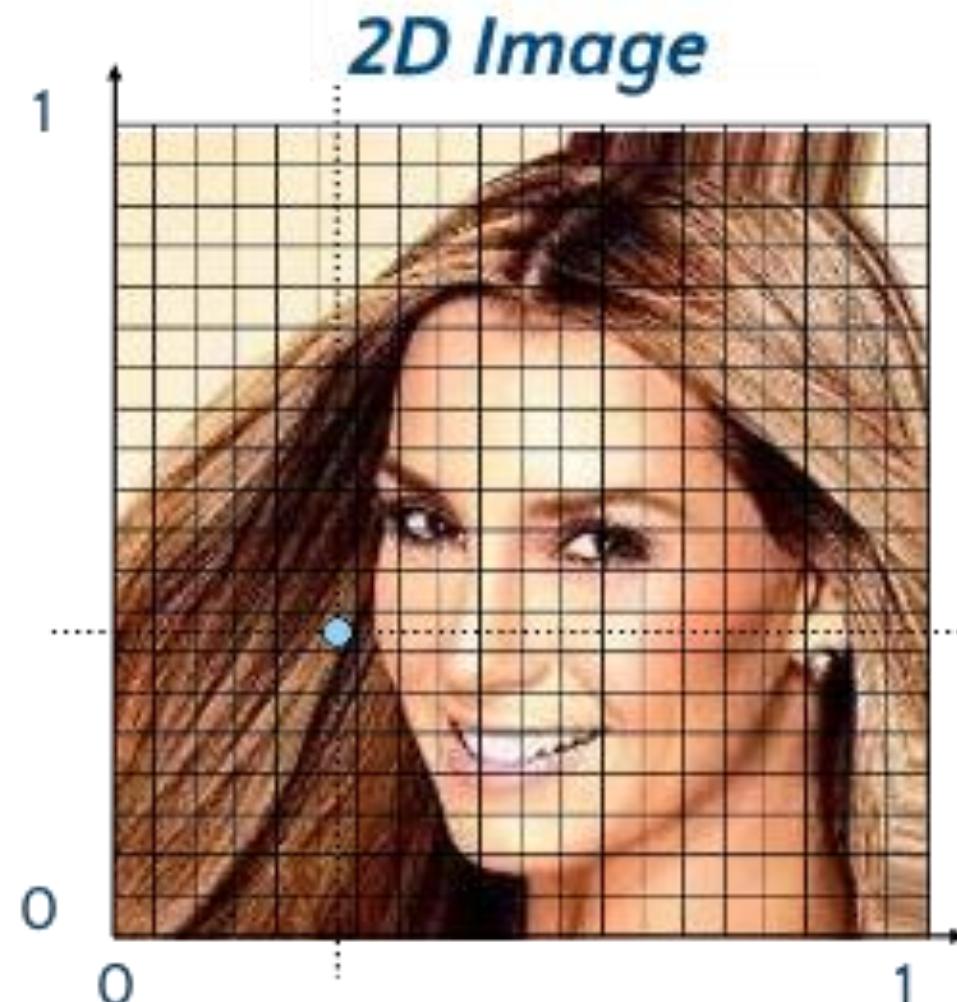
$$f : \mathbb{R}^3 \rightarrow \{0, 1\}, f(x_1, x_2, x_3) = p$$

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(\varphi, \lambda) = T$$

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}^3, f(x_1, x_2, t) = (r, g, b)$$

Motivation

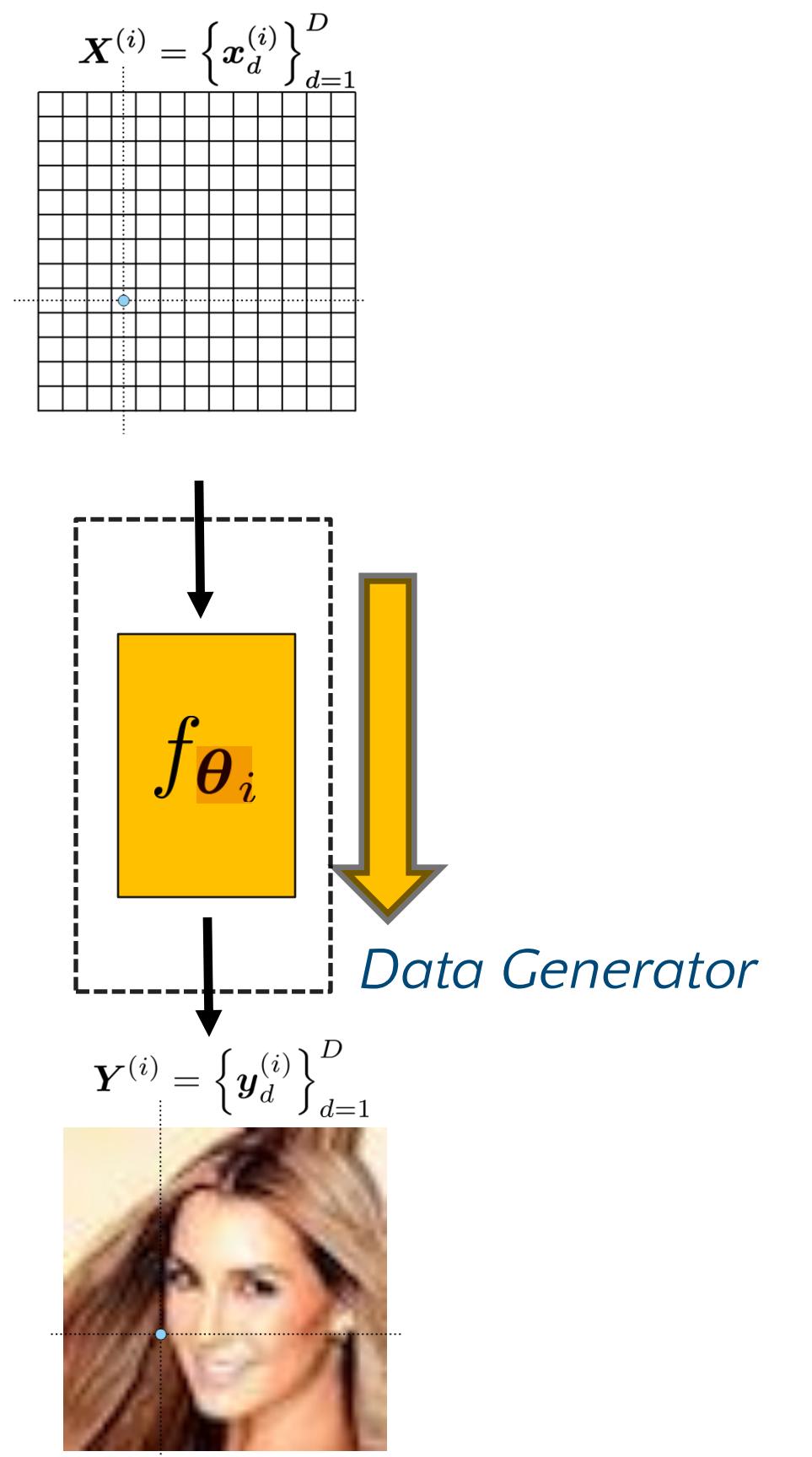
- Focusing on images:



- Generator function $f : \mathbf{X} \rightarrow \mathbf{Y}$ creates this specific image with the mapping $f(\mathbf{x}_d) = \mathbf{y}_d, d \in [1, \dots, D]$
- Each pixel is now a pair $\{\mathbf{x}_d, \mathbf{y}_d\}$ where $x_d \in \mathbb{R}^2, y_d \in \mathbb{R}^3$
- Full image is a pair of sets $\mathbf{X} = \{\mathbf{x}_d\}_{d=1}^D, \mathbf{Y}_d = \{\mathbf{y}_d\}_{d=1}^D$

Implicit Neural Representations

INRs [20-22]



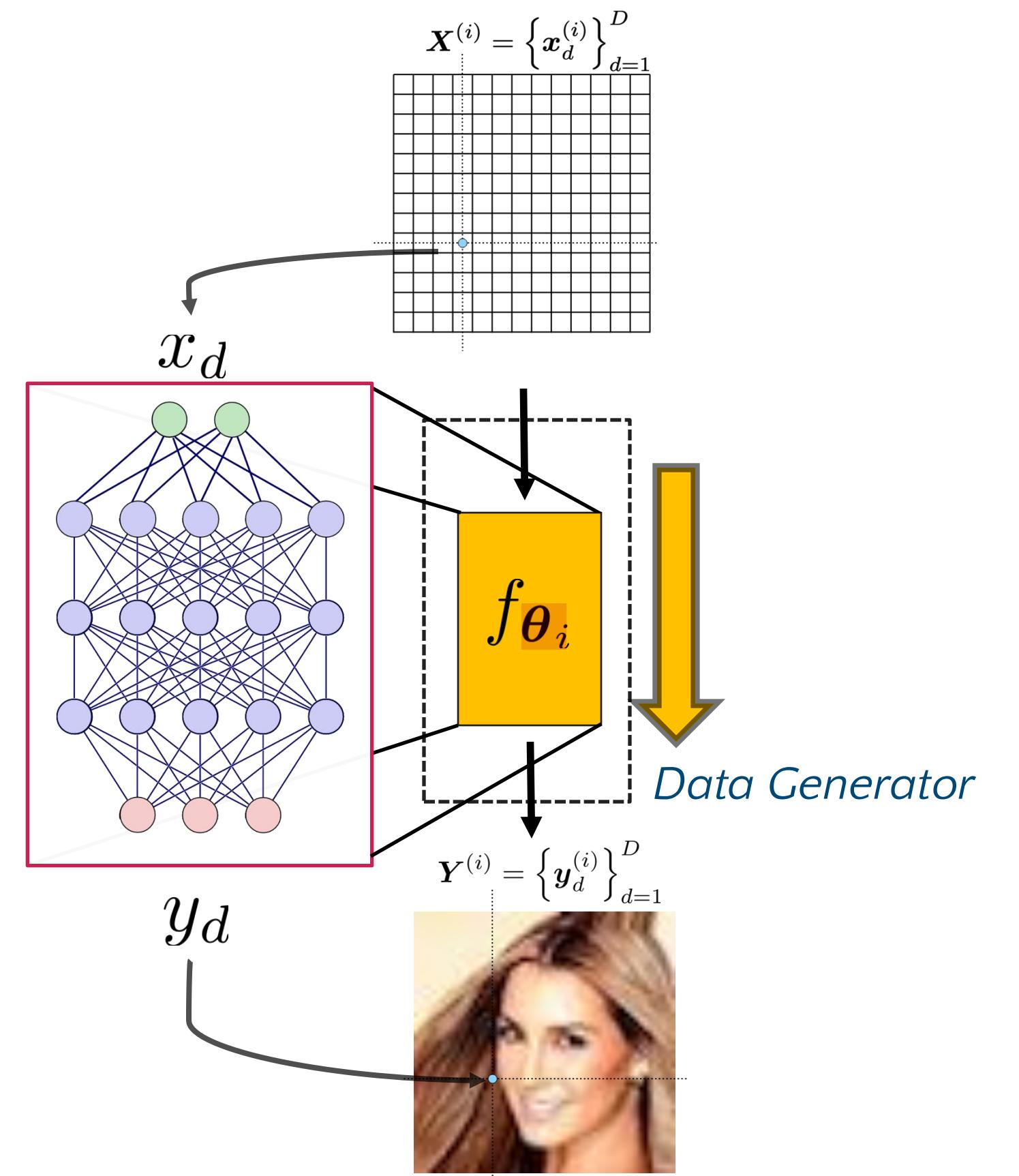
[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

[20] Sitzmann et al., 2019

Implicit Neural Representations

INRs [20-22]



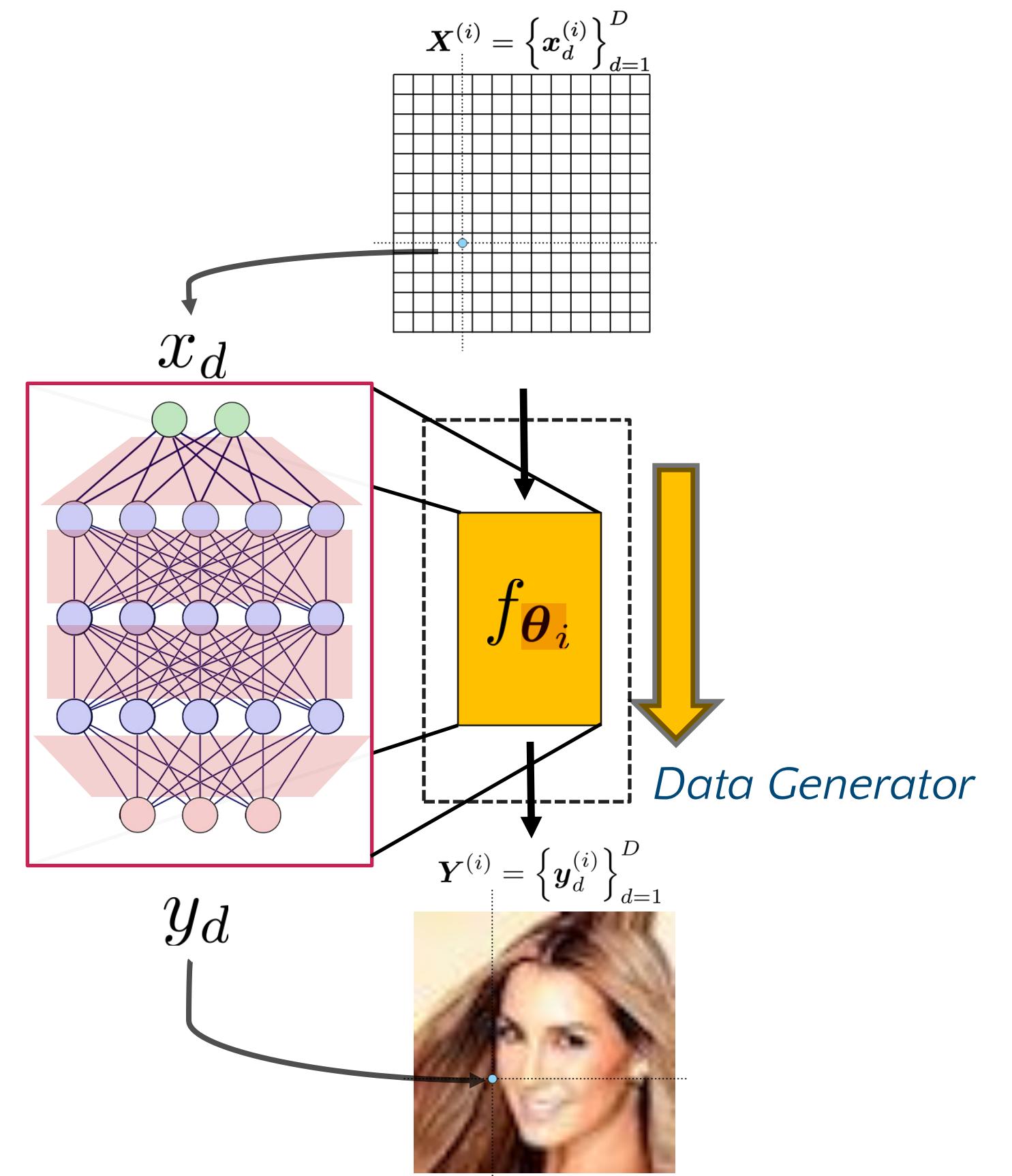
[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

[20] Sitzmann et al., 2019

Implicit Neural Representations

INRs [20-22]



[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

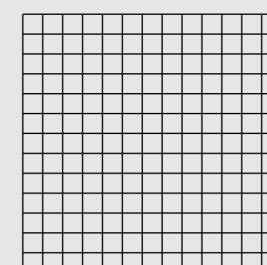
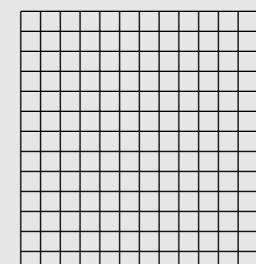
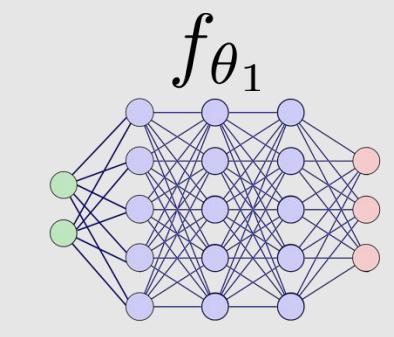
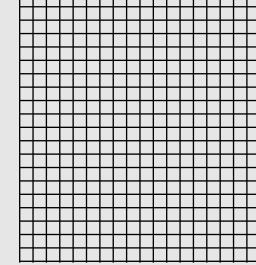
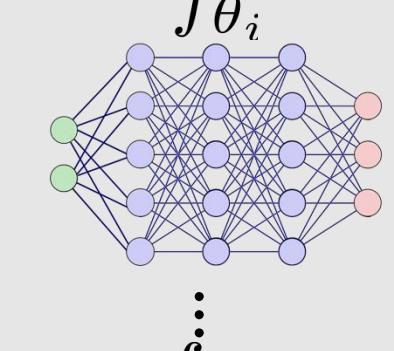
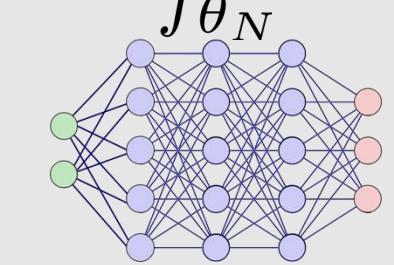
[20] Sitzmann et al., 2019

Implicit Neural Representations

INRs [20-22]

Data generator f_{θ_i} is unique to each image

$$\mathbf{X}^{(i)} = \left\{ \mathbf{x}_d^{(i)} \right\}_{d=1}^D$$


 \vdots

 \vdots

 \vdots

 \vdots


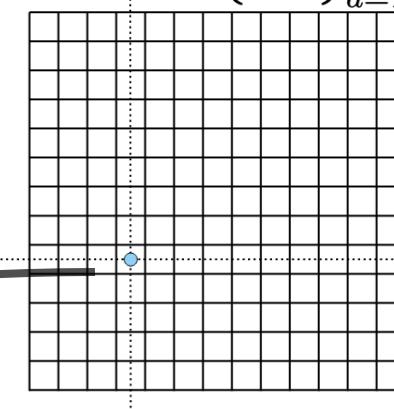
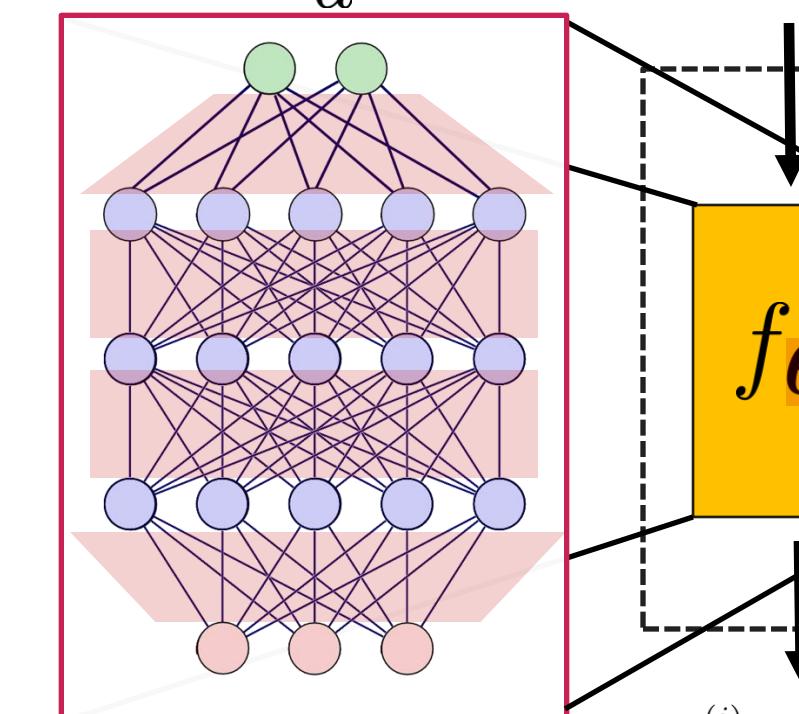
$$\mathbf{Y}^{(i)} = \left\{ \mathbf{y}_d^{(i)} \right\}_{d=1}^D$$


 $i=1$

 \vdots

 $i=N$

$$\mathbf{X}^{(i)} = \left\{ \mathbf{x}_d^{(i)} \right\}_{d=1}^D$$


 x_d


Data Generator

$$\mathbf{Y}^{(i)} = \left\{ \mathbf{y}_d^{(i)} \right\}_{d=1}^D$$

 y_d


[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

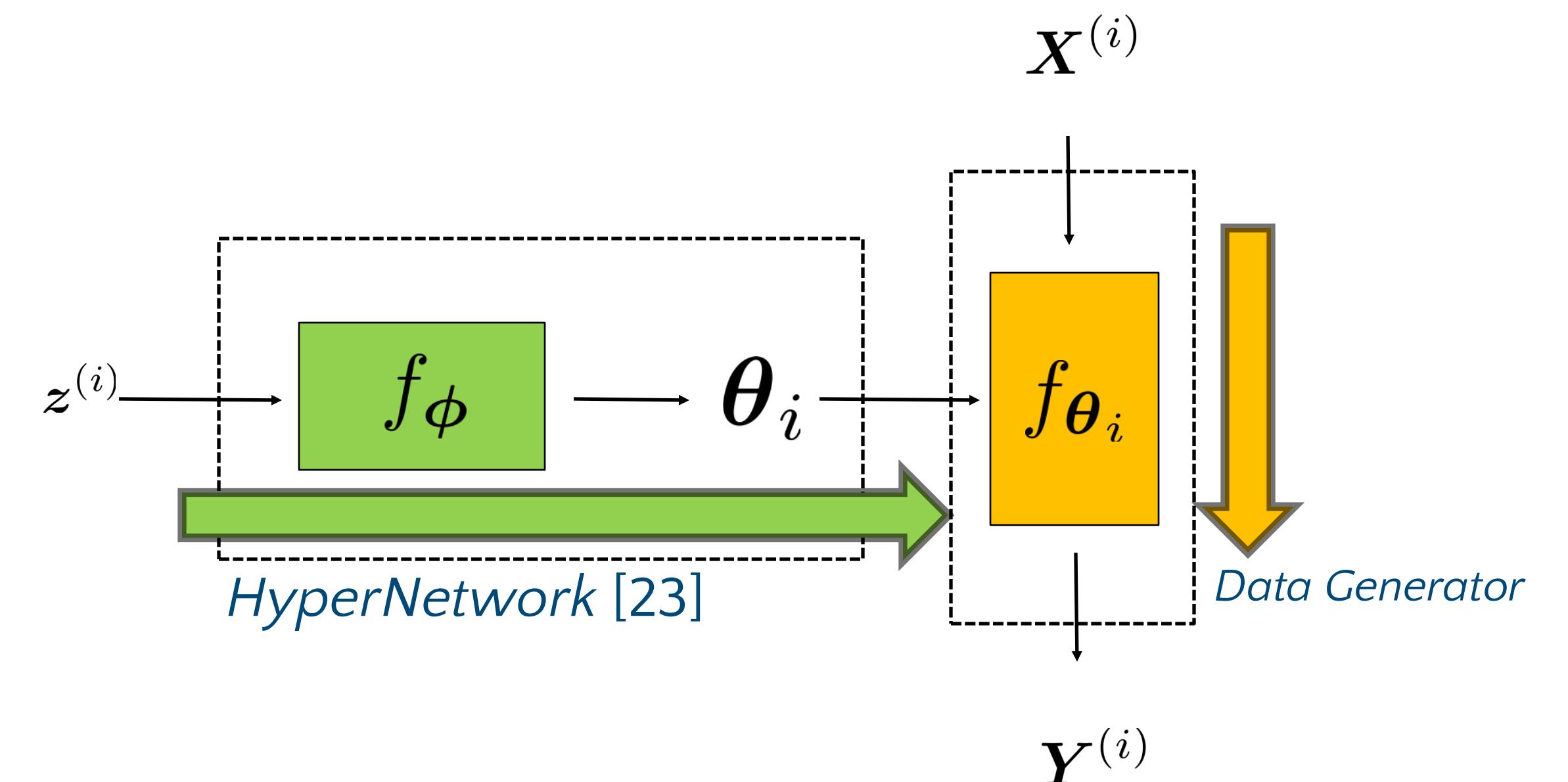
[20] Sitzmann et al., 2019

Deep Generative Models of INRs

How to scale to large datasets?

How to map a latent representation to an INR?

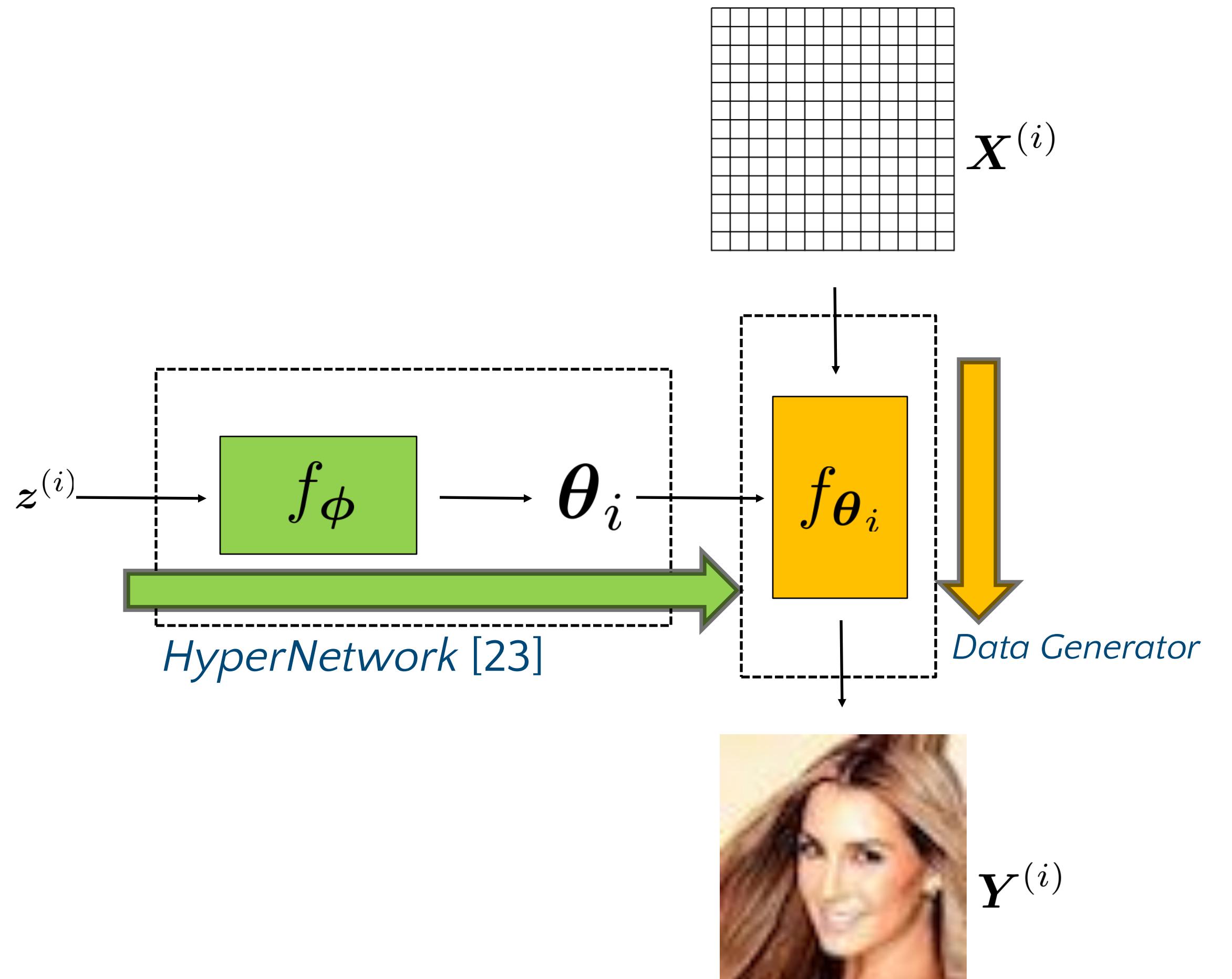
Deep Generative Models of INRs



[23] Ha et al., 2017

Deep Generative Models of INRs

Have $z^{(i)}$, a summary representation of image.

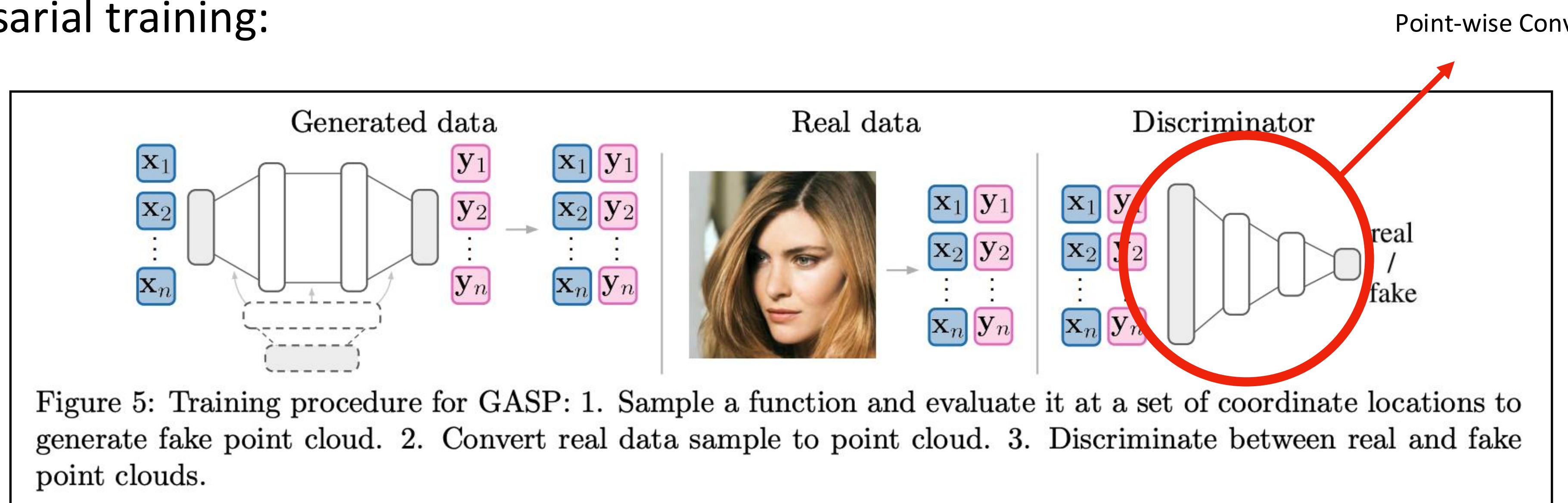


[23] Ha et al., 2017

Previous work

GASP[5]

- Adversarial training:



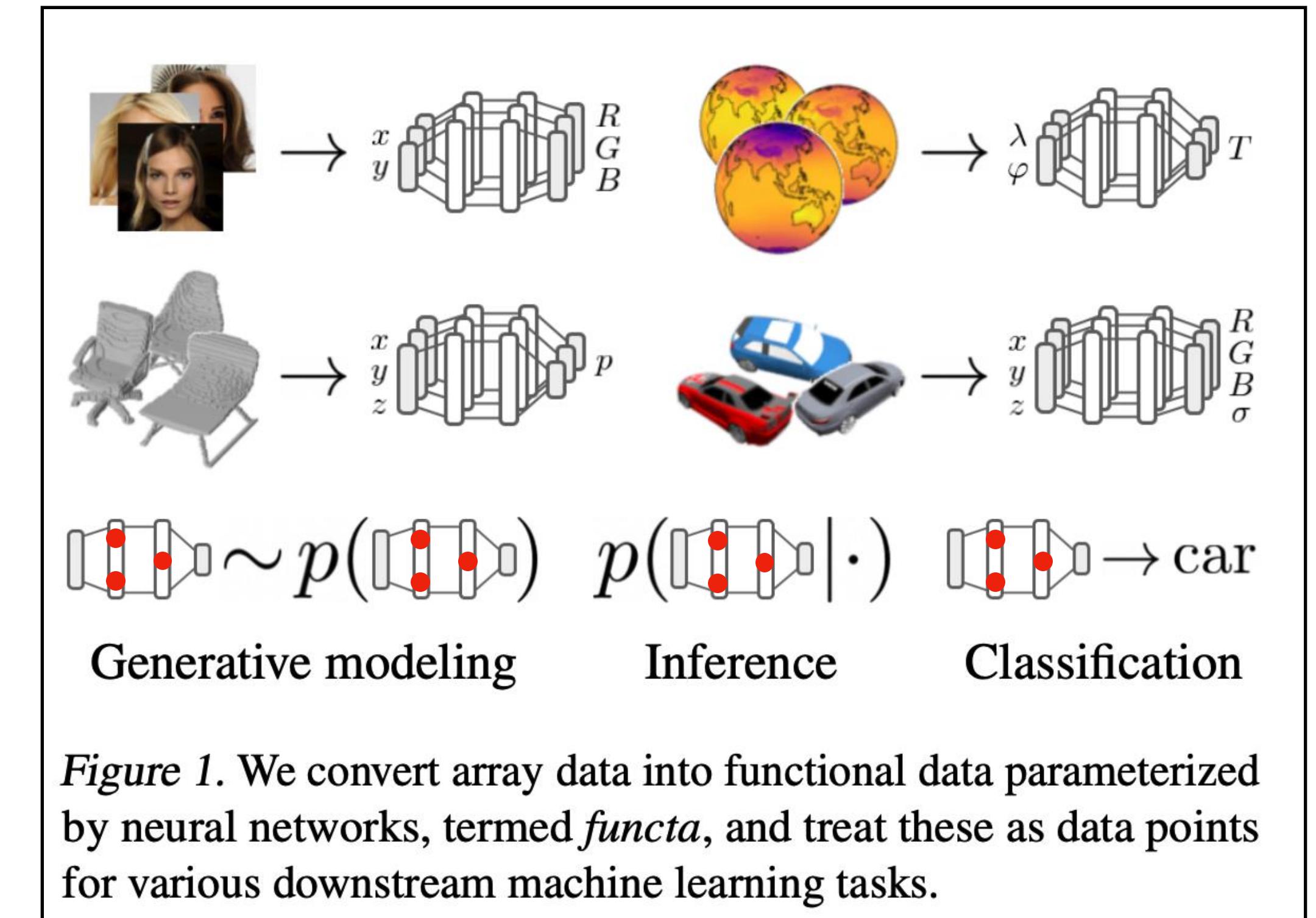
- ✗ Can't tackle inference related tasks.

[5] Dupont et al., 2020

Previous work

Functa^[6]

- Decoupled training:
 1. Fit an INR per datapoint using SIREN^[20] and **modulation vectors**, named **functas**.
 2. Train any generative model on the functa dataset of vectors.
- ✖ Computationally expensive inference.



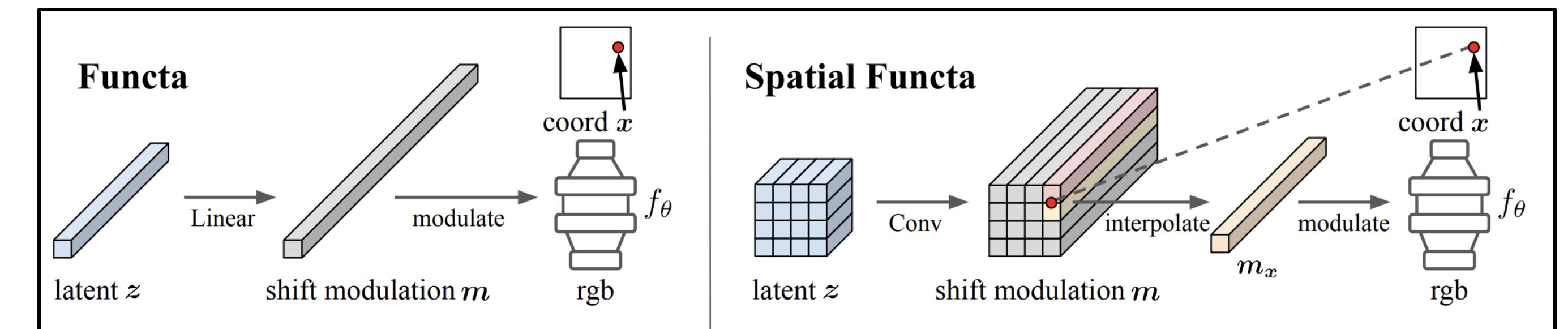
^[6] Dupont et al., 2022

^[20] Sitzmann et al., 2020

Previous work

Spatial Functa^[26]

- Decoupled training:
 1. Fit an INR per datapoint using SIREN^[20] and **modulation tensor**.
 2. Train any generative model on the functa dataset of tensors.
- ✖ Computationally expensive inference.



[26] Bauer et al., 2023

[20] Sitzmann et al., 2020

Deep Generative Models of INRs

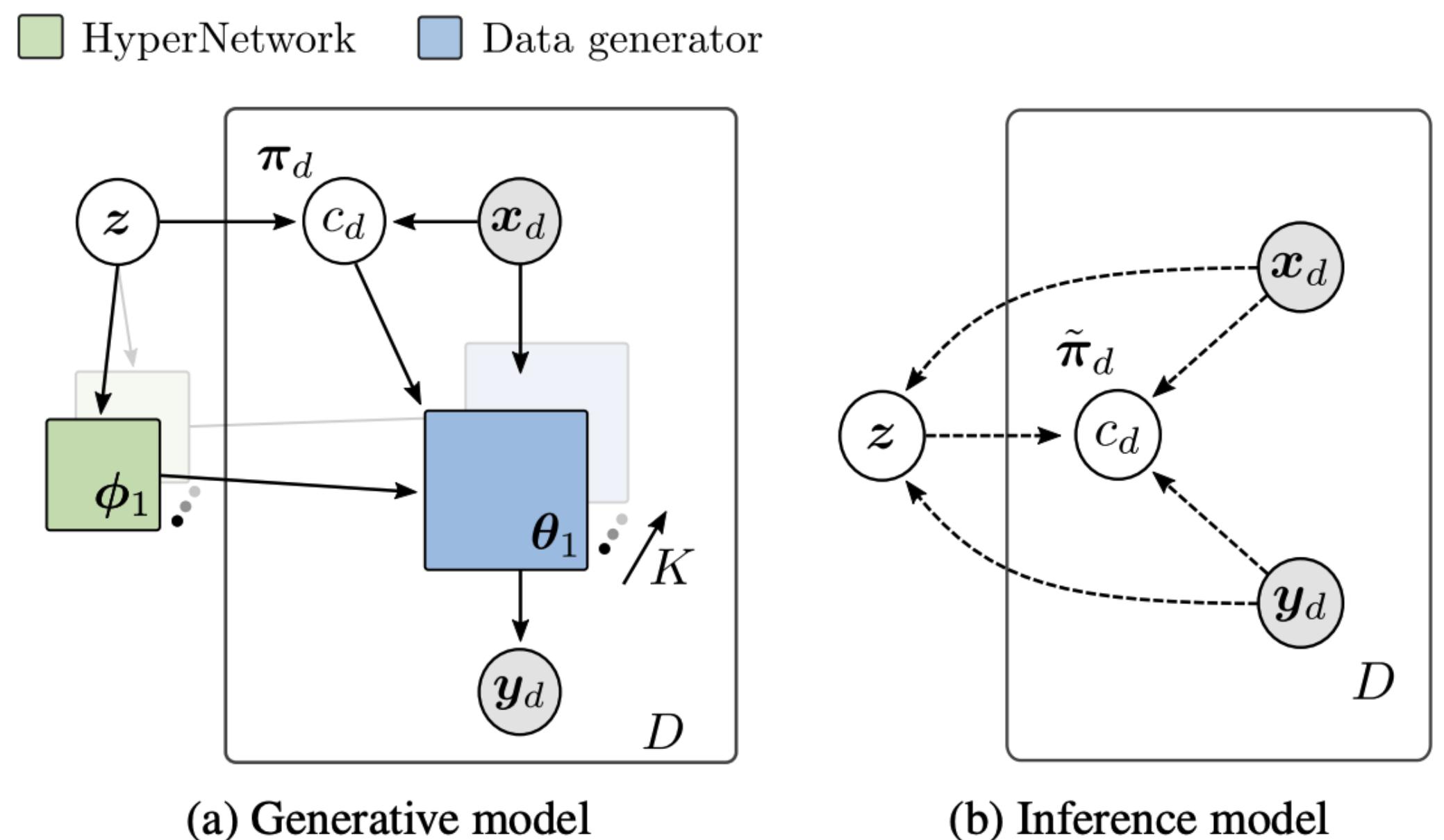
How to infer the latent representation \mathbf{z} ?

$$q_{\gamma}(\mathbf{z}|\mathbf{Y}, \mathbf{X}) \quad p_{\psi}(\mathbf{z})$$

Proposed methods (1)

VAMoH

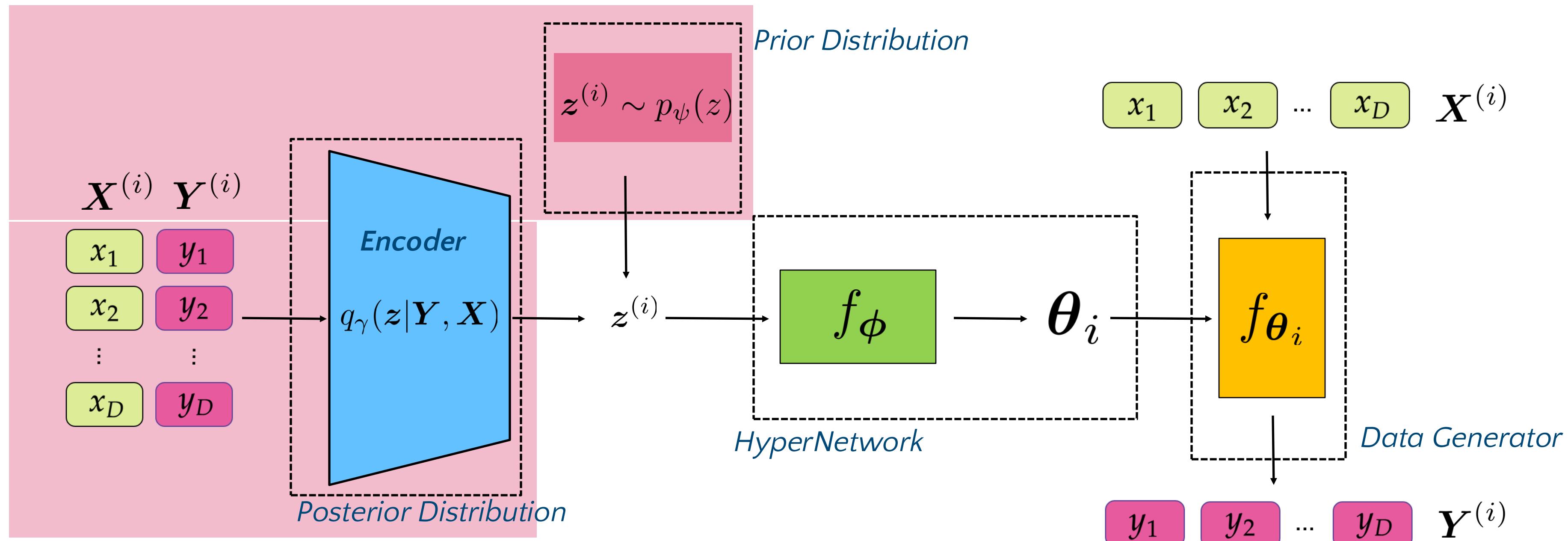
Variational Mixture of HyperGenerators [25]



[25] Koyuncu et al., 2023

VAMoH

Encoder

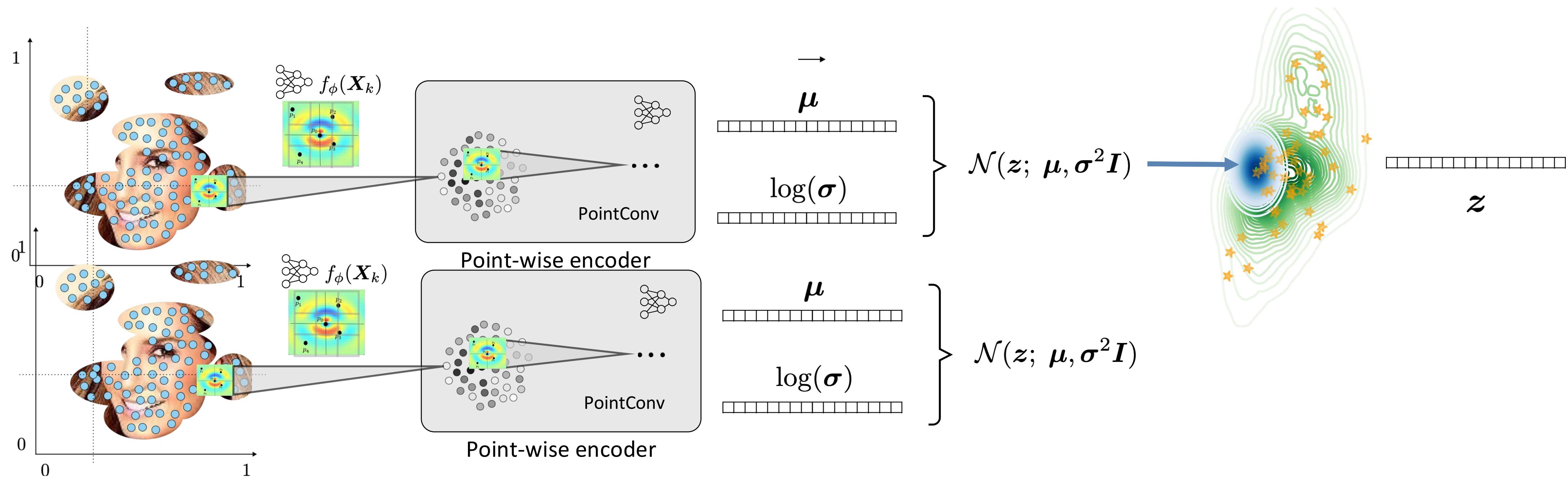


$z^{(i)}$: Latent Variable

VAMoH

Encoder

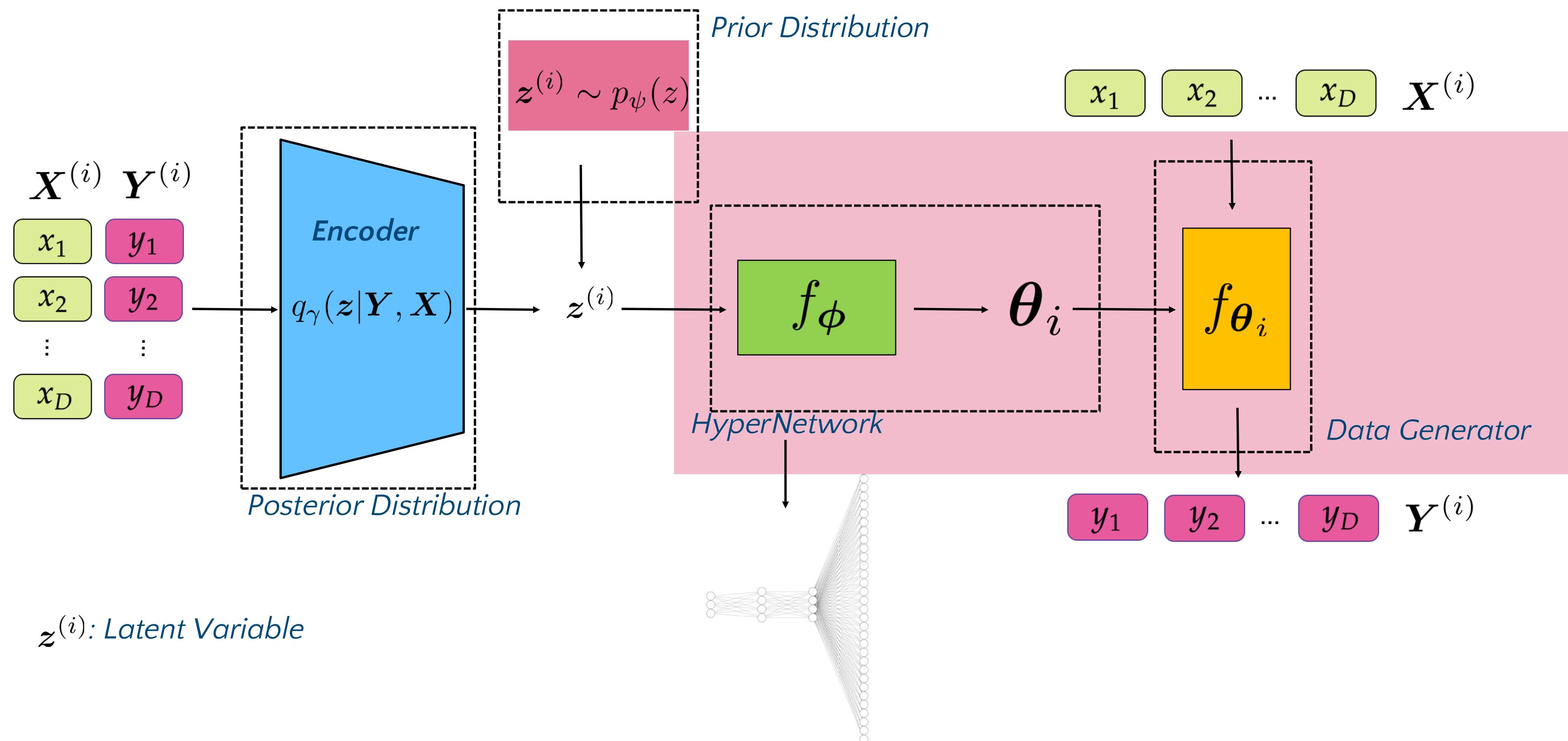
- PointConv^[21] encoder for point clouds.



^[21] Wu et al., 2019

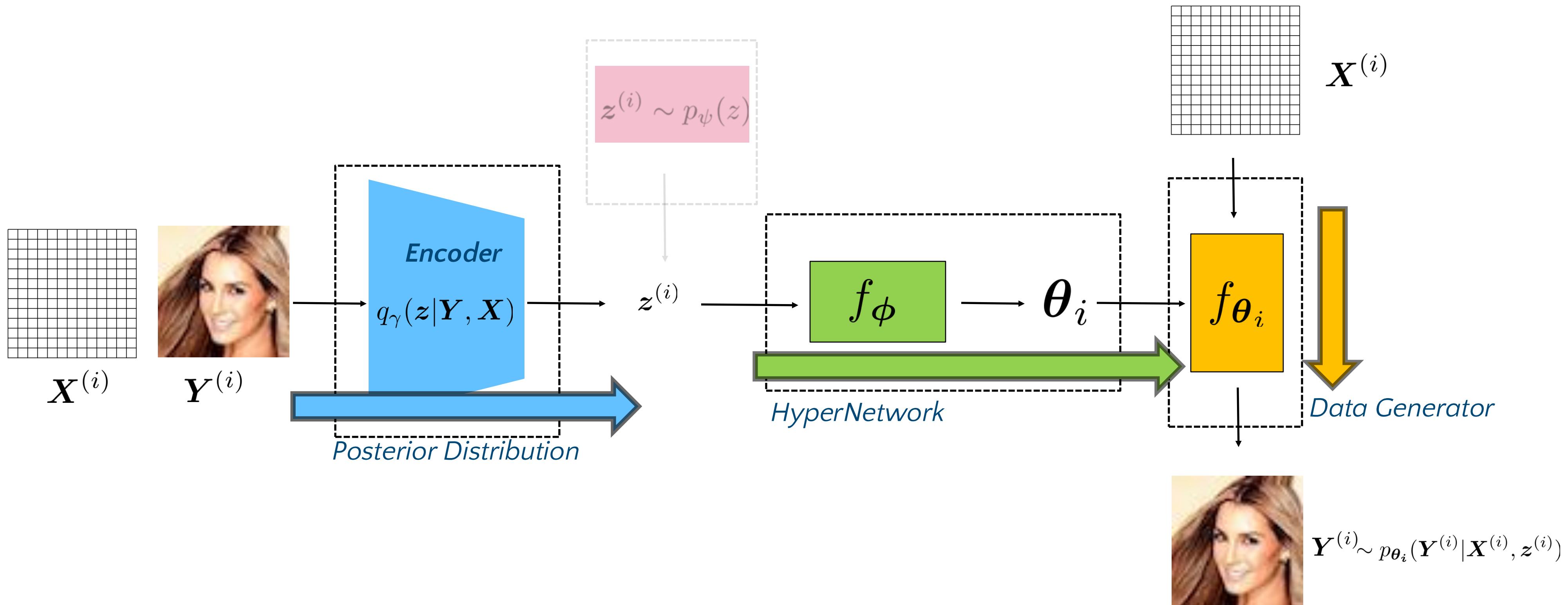
VAMoH

Decoder



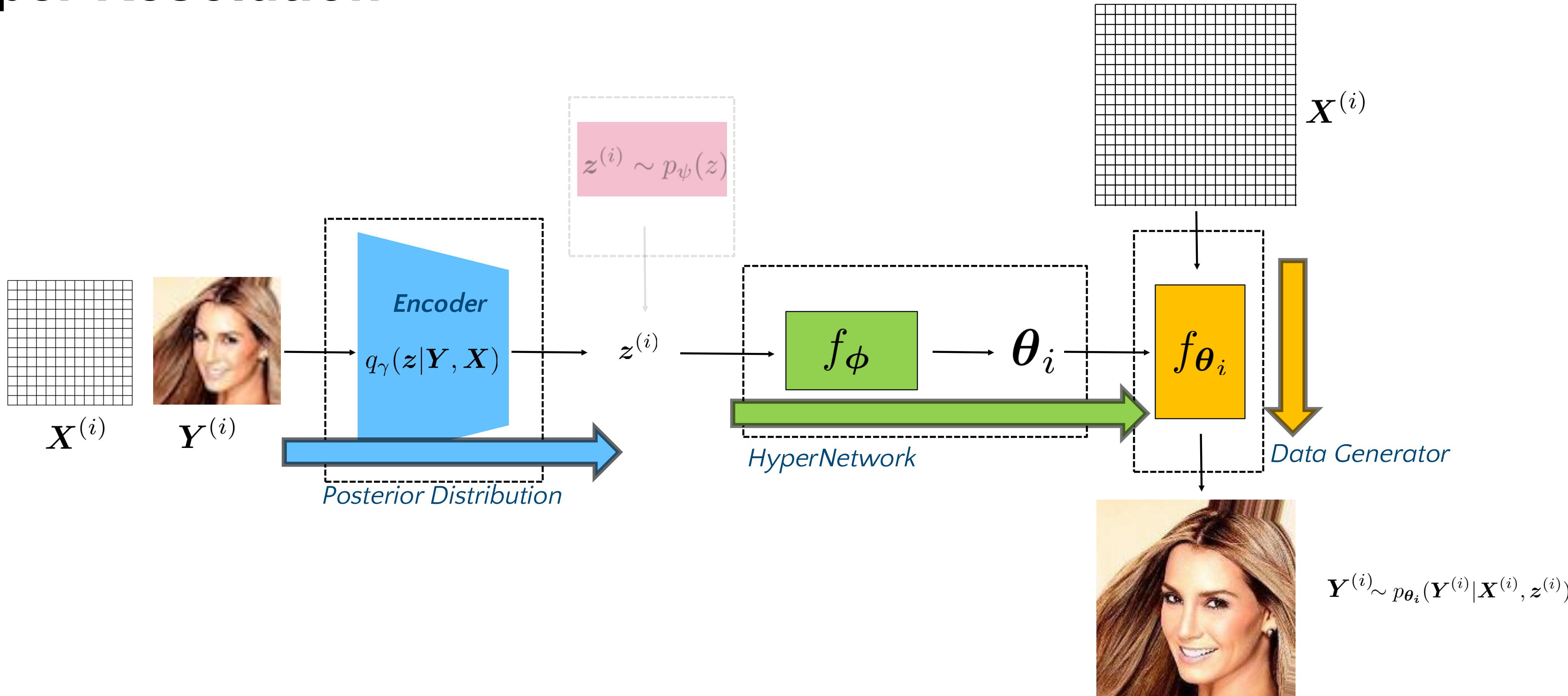
VAMoH

Reconstruction



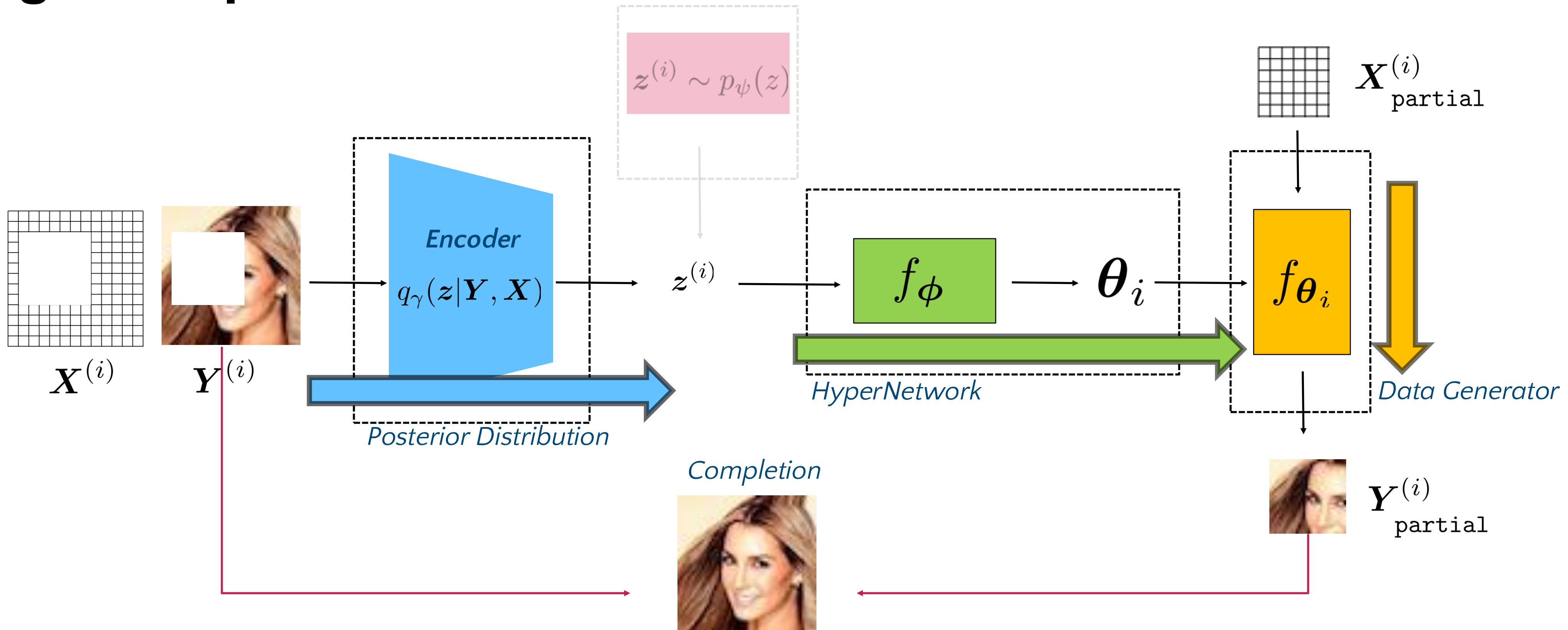
VAMoH

Super Resolution



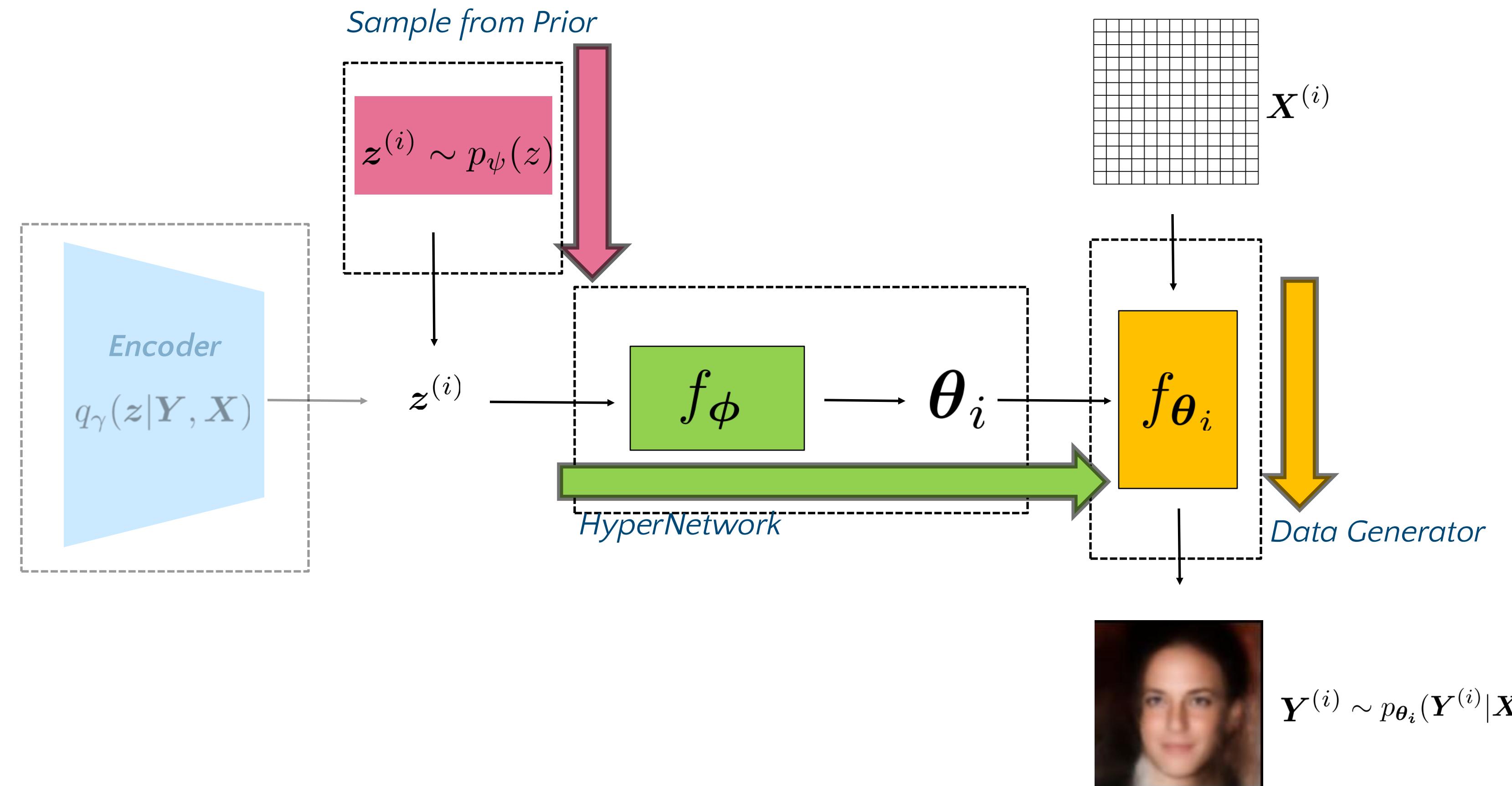
VAMoH

Image Completion



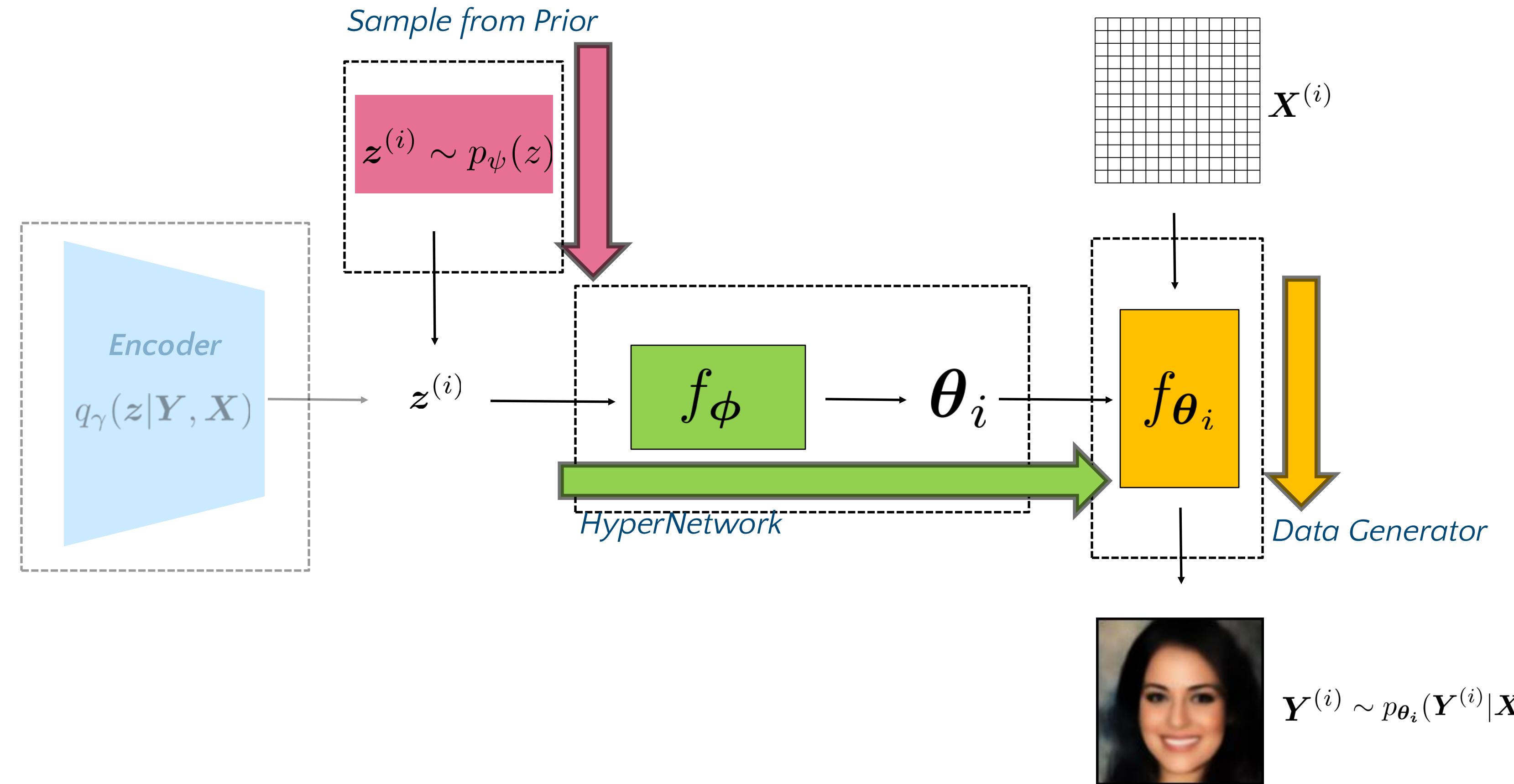
VAMoH

Image Generation



VAMoH

Image Generation



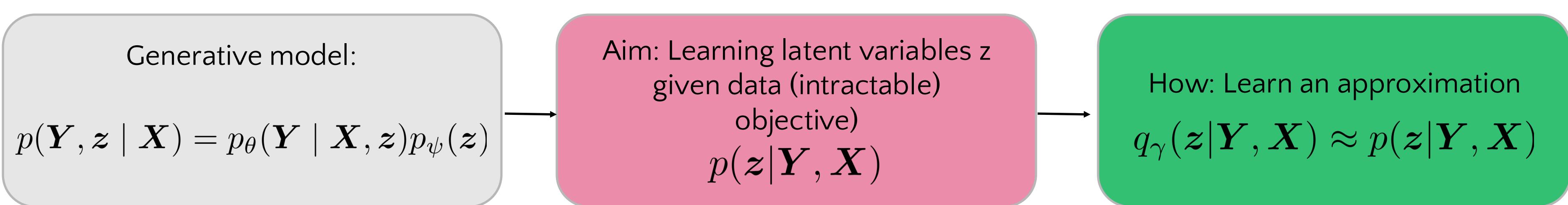
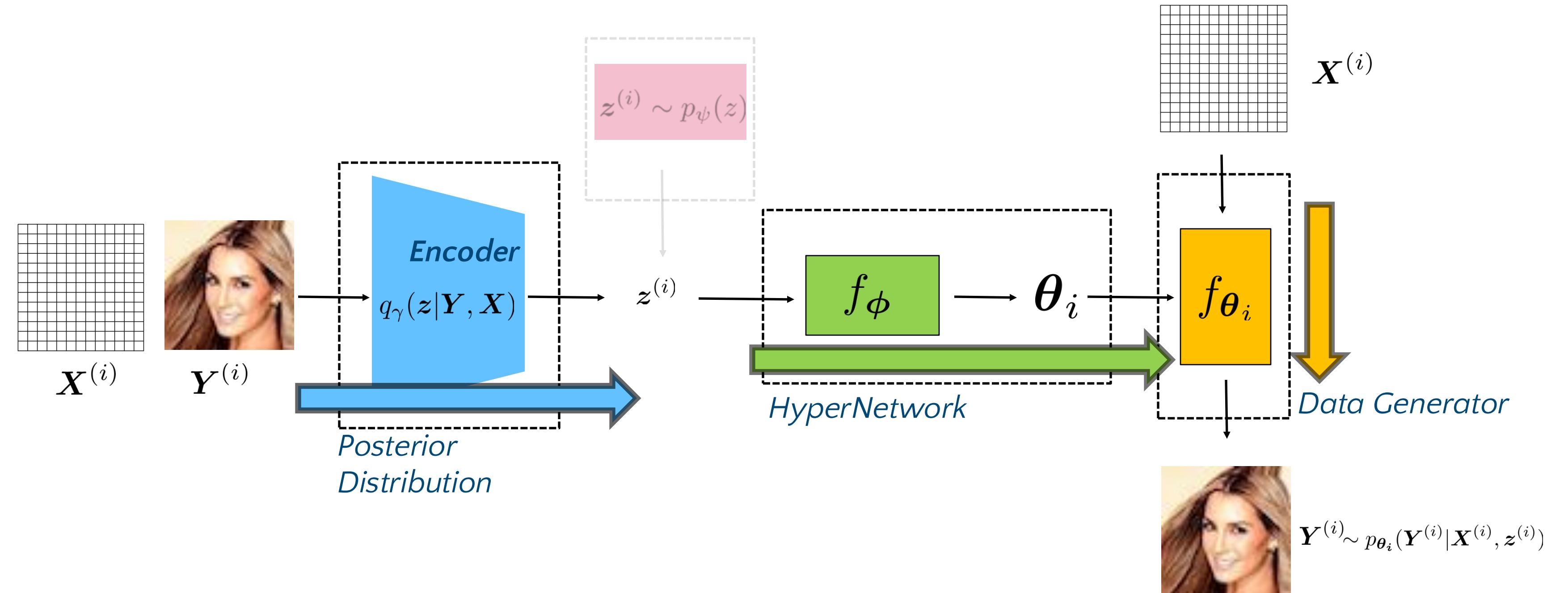
VAMoH

Optimization

How to learn all these steps end-to-end from data?

VAMoH

Optimization



VAMoH

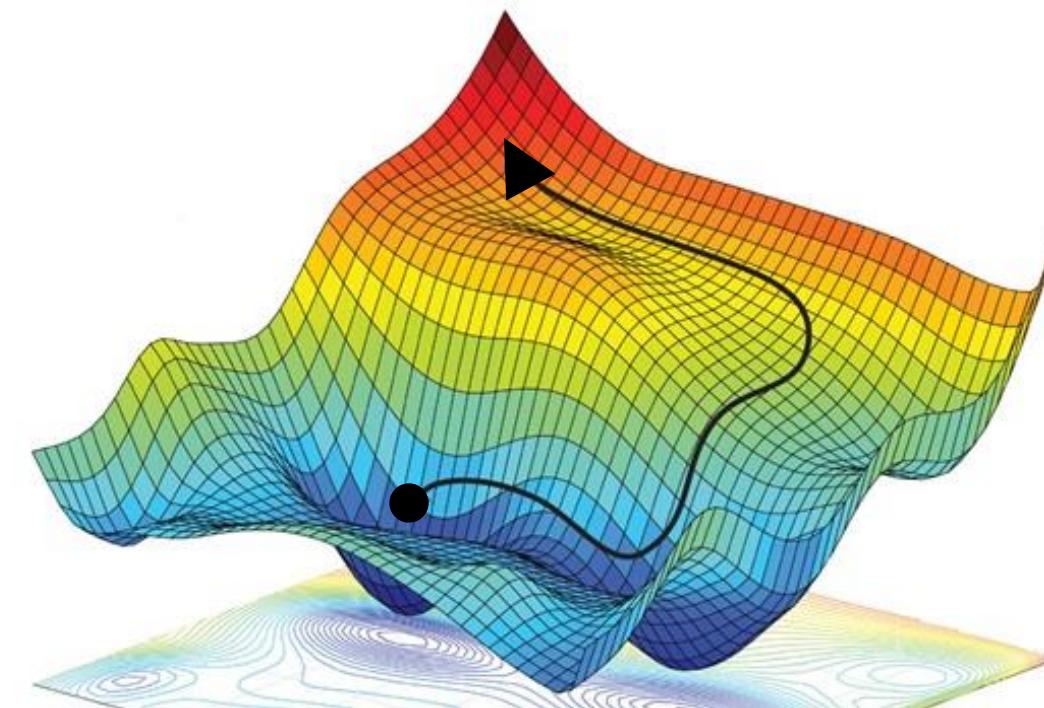
Optimization

- For a single data sample (\mathbf{X}, \mathbf{Y})

$$\max_{\phi, \gamma} \mathcal{L}(\phi, \gamma; \mathbf{Y}, \mathbf{X}) = \underbrace{\max_{\phi, \gamma} \mathbb{E}_{q_\gamma(z|\mathbf{Y}, \mathbf{X})} [\log p_\theta(\mathbf{Y} | \mathbf{X}, z)]}_{\text{Reconstruction}} - \underbrace{D_{KL}(q_\gamma(z|\mathbf{Y}, \mathbf{X}) || p_\psi(z))}_{\text{Regularization}}$$

- For all samples in our dataset $\left(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)}\right), i \in [N]$

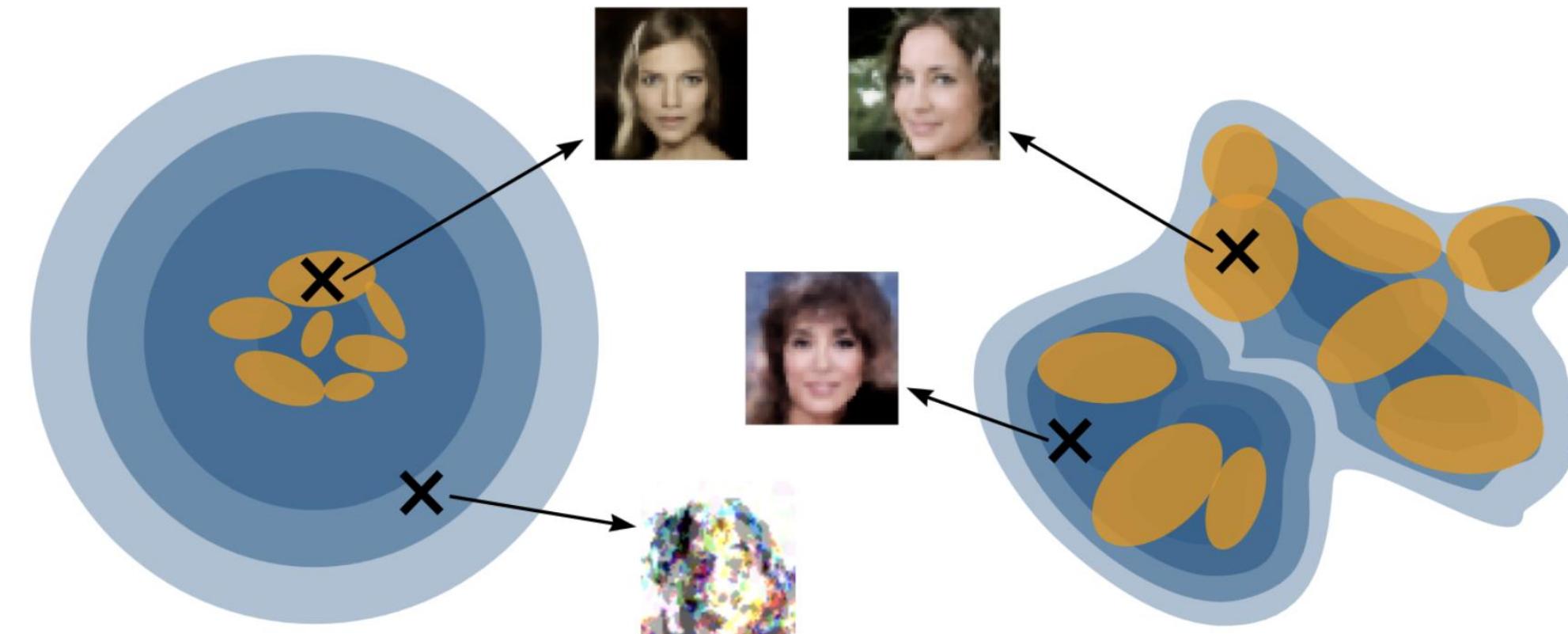
$$\max_{\phi, \gamma} \sum_{i=1}^N \mathcal{L}(\phi, \gamma; \mathbf{Y}^{(i)}, \mathbf{X}^{(i)})$$



VAMoH

'Holes' problem

█ $p(\mathbf{z})$
 █ $q(\mathbf{z}|\mathbf{X}_i, \mathbf{Y}_i)$



Regularization Term:

$$\min_{\gamma} D_{KL}(q_{\gamma}(\mathbf{z} | \mathbf{Y}, \mathbf{X}) \| p_{\psi}(\mathbf{z}))$$

We need to align the approximate posterior with the prior.

$$p_{\psi}(\mathbf{z}) \quad q_{\gamma}(\mathbf{z})$$

Problem:

If the prior is too simple, it hinders generation quality.

Solution:

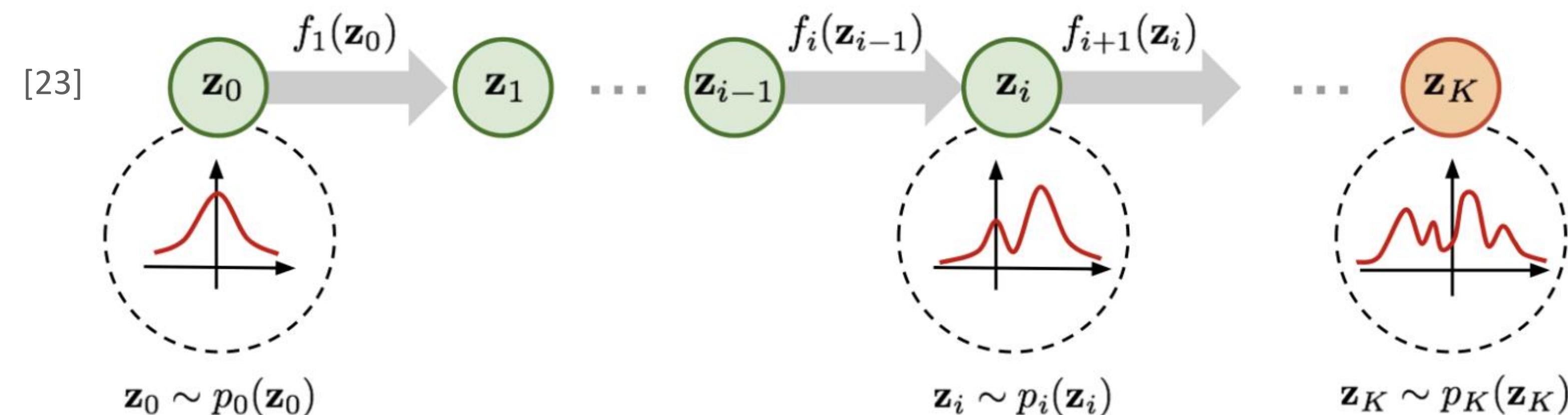
Learn a more complex $p_{\psi}(\mathbf{z})$ with another NN.

$$\min_{\gamma, \psi} D_{KL}(q_{\gamma}(\mathbf{z} | \mathbf{Y}, \mathbf{X}) \| p_{\psi}(\mathbf{z}))$$

VAMoH

Flow-based prior

- More expressive prior using RealNVP (Real-valued, Non-Volume Preserving) Flow.

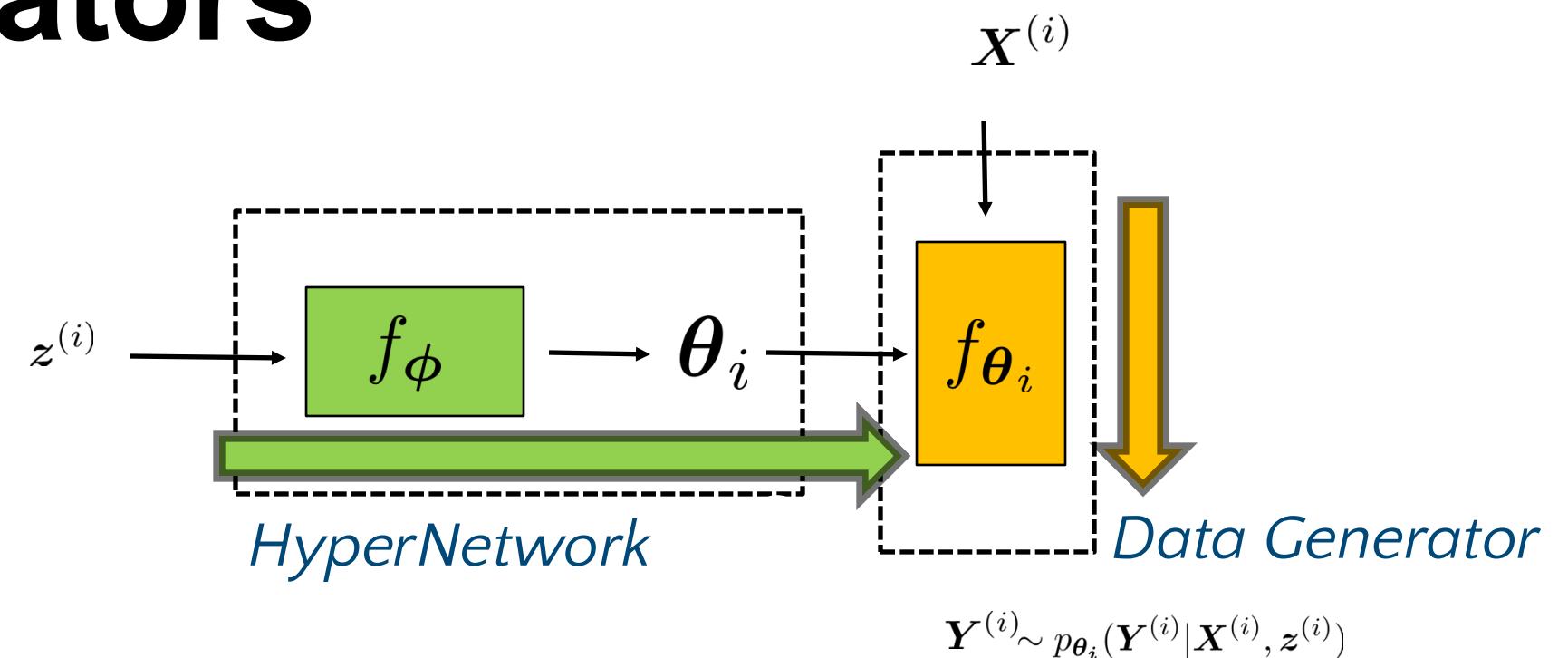


$$\mathbf{z}^{(i)} \sim p_\psi(\mathbf{z})$$

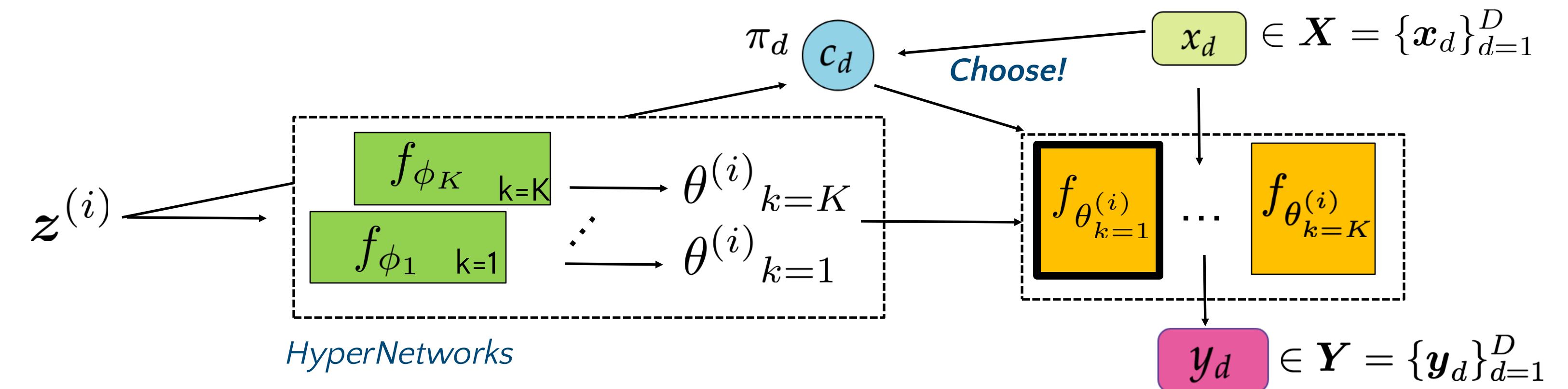
VAMoH

Mixture of HyperGenerators

Single HyperGenerator



Mixture of HyperGenerators



VAMoH

Mixture of HyperGenerators

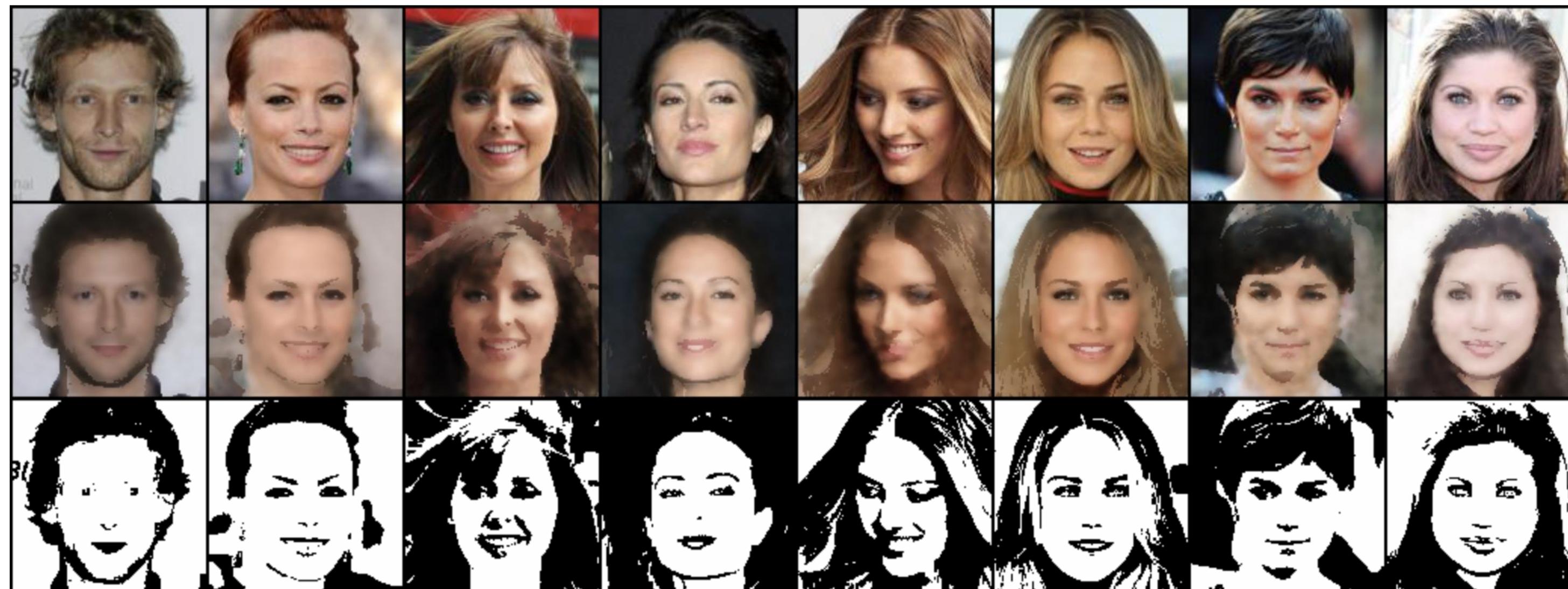


Image Reconstruction with Mixture of HyperGenerators

VAMoH

- For a single data sample

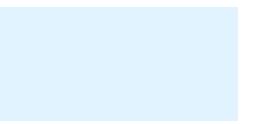
$$(\mathbf{X}, \mathbf{Y})$$

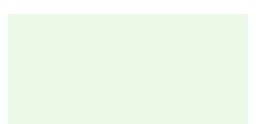
$$\mathcal{L}(\mathbf{Y}, \mathbf{X}; \psi, \phi, \gamma) = \sum_{d=1}^D \mathbb{E}_{q_{\gamma_z}(\mathbf{z} | \mathbf{Y}, \mathbf{X})} \left[\sum_{k=1}^K \log p_{\theta_k} (\mathbf{y}_d | \mathbf{x}_d) \cdot \pi_{dk} \right] - D_{KL}(q_{\gamma_z}(\mathbf{z} | \mathbf{X}, \mathbf{Y}) \| p_{\psi_z}(\mathbf{z}))$$

$$- D_{KL}(q_{\gamma_c}(\mathbf{C} | \mathbf{z}, \mathbf{X}, \mathbf{Y}) \| p_{\psi_c}(\mathbf{C} | \mathbf{z}, \mathbf{X}))$$

- For all samples in our dataset

$$(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)}), i \in [N]$$

 Reconstruction

 KL of the continuous latent variable

 KL of the discrete latent variable

$$\max_{\phi, \gamma, \psi} \sum_{i=1}^N \mathcal{L}(\phi, \gamma, \psi; \mathbf{Y}^{(i)}, \mathbf{X}^{(i)})$$

Experiments

Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	✗

Experiments

Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample 

Experiments

Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Completion
GASP (2021) [5]	GAN	Minimax	Forward Pass	$\min_{\phi} - \log p(\phi) + \lambda \sum_{i \in \mathcal{I}} \ f_{\phi}(\mathbf{x}_i) - \mathbf{f}_i\ _2^2$
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample 

Experiments

Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample 
VaMoH (ours)	VAE-based	Single optimization	Forward Pass	Forward pass 

VAMoH provides a probabilistic generative model that is efficient, robust, and expressive for modeling distribution over functions.

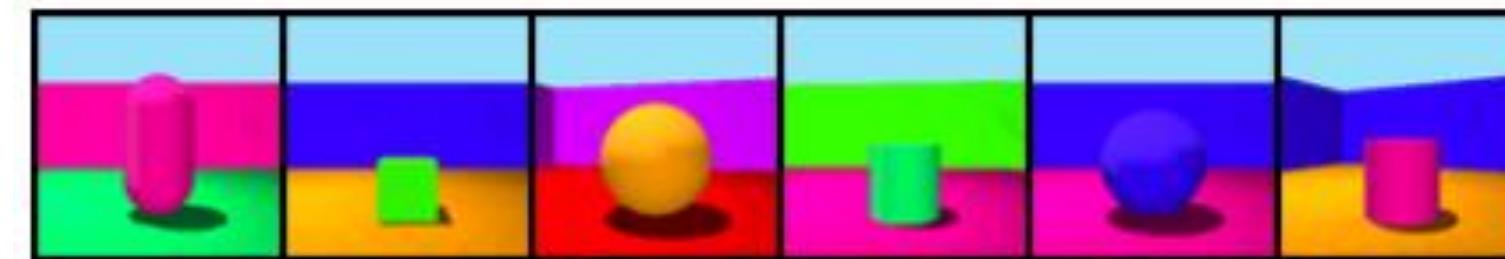
Experiments

Datasets

PolyMNIST (28x28)



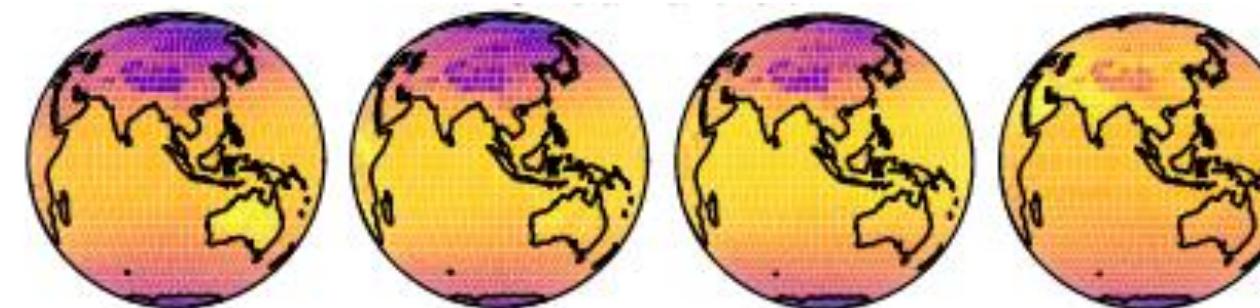
Shapes3D (64x64)



CelebA-HQ (64x64)



ERA5 (Polar)

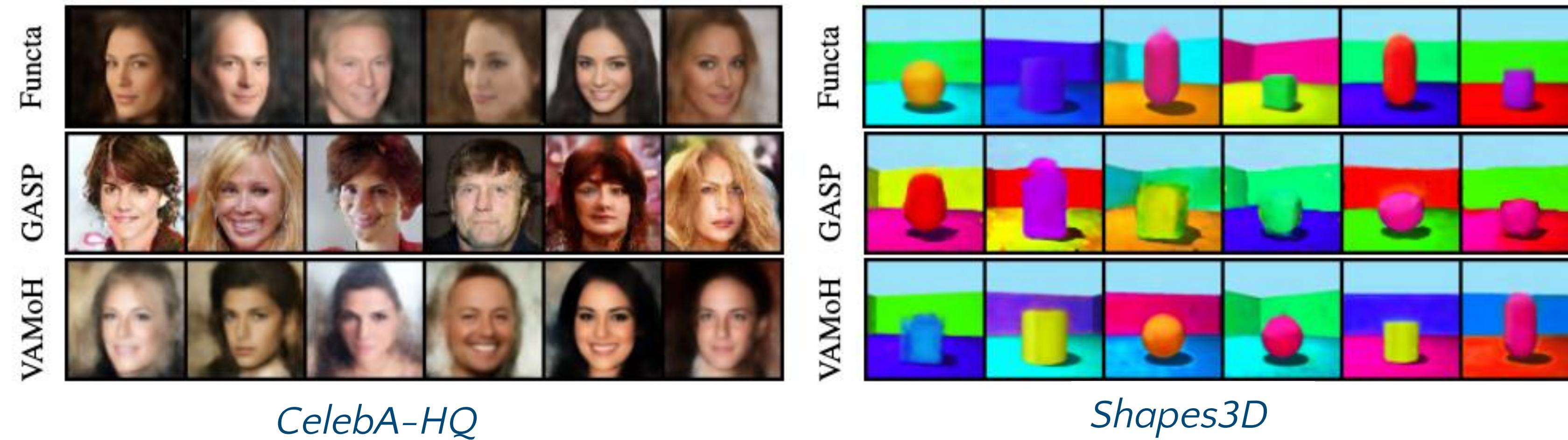


ShapeNET (Voxels)



Experiments

Generation

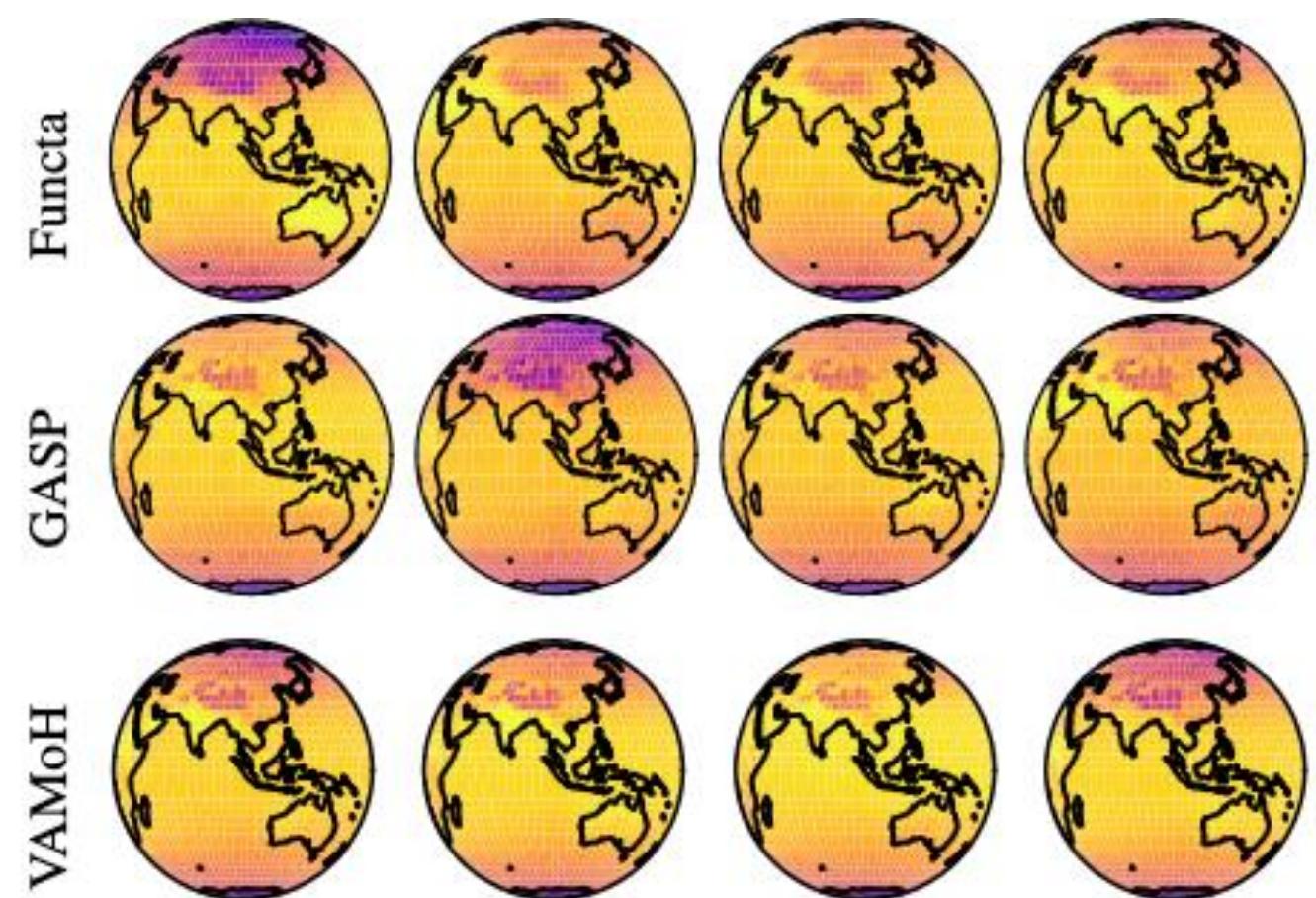


Experiments

Generation



PolyMNIST



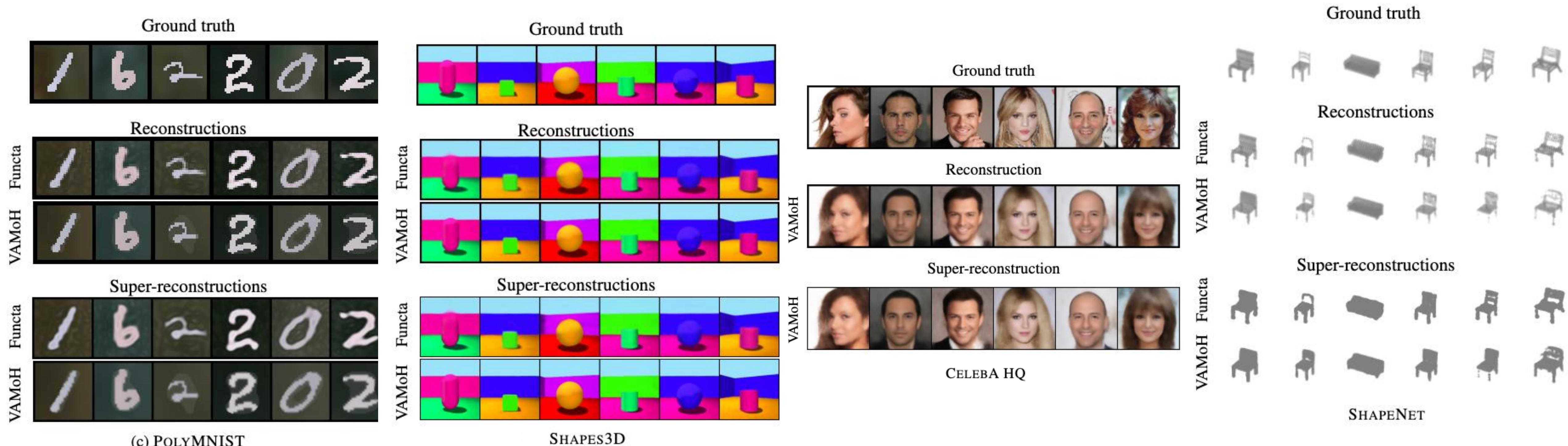
ERA5



ShapeNET

Experiments

Reconstructions



Experiments

Inference times

Table 2: Comparison of inference time (seconds) for reconstruction task of VaMoH and Functa. On the right-most two columns, we show the speed improvement of VaMoH compared to Functa (3) which is trained with 3 gradient steps as suggested in the original paper [Dupont et al., 2022b] and Functa (10) which is trained with 10 gradient step to obtain the results of Functa depicted in Figures 16,17. Please note that these experiments are run on the same GPU device.

Dataset	Model Inference Time (secs)			Speed Improvement	
	VaMoH	Functa (3)	Functa (10)	vs. Functa (3)	vs. Functa (10)
POLYMNIST	0.00453	0.01648	0.05108	x 3.64	x 11.28
SHAPES3D	0.00536	0.01759	0.05480	x 3.28	x 10.22
CELEBA HQ	0.00757	0.01733	0.05381	x 2.29	x 7.11
ERA5	0.00745	0.01899	0.05932	x 2.55	x 7.96
SHAPENET	0.00689	0.02095	0.06576	x 3.04	x 9.54

Reconstruction

Dataset	Model Inference Time (secs)			Speed Improvement	
	VaMoH	Functa (3)	Functa (10)	vs. Functa (3)	vs. Functa (10)
POLYMNIST	0.00455	0.01649	0.05109	x 3.62	x 11.23
SHAPES3D	0.00544	0.01768	0.05489	x 3.25	x 10.09
CELEBA HQ	0.00833	0.01729	0.05377	x 2.08	x 6.46
ERA5	0.00790	0.01997	0.06030	x 2.53	x 7.63
SHAPENET	0.01440	0.02089	0.06569	x 1.45	x 4.56

Super-reconstruction

Experiments

Image completion

The figure displays a 6x6 grid of images illustrating reconstruction results. The columns are labeled "Recons." and the rows are labeled "In".

- Row 1:** Shows six faces with black squares covering the eyes.
- Row 2:** Shows the same six faces with the black squares removed, revealing the eyes.
- Row 3:** Shows six 3D shapes (cylinders and spheres) on colored platforms. The first three images have black squares over their centers; the last three do not.
- Row 4:** Shows the same 3D shapes as Row 3, but with black squares over their bases.
- Row 5:** Shows six dark images with white, stylized, handwritten-like patterns. The first three patterns are slanted, while the last three are more vertical.
- Row 6:** Shows the same six dark images as Row 5, but with the patterns blurred or less distinct.

Missing a patch (in-painting)

The figure displays a 6x6 grid of images illustrating reconstruction results for different datasets. The columns are labeled "Recons." and "In Recons." vertically along the left side. The rows show reconstructions for three types of data:

- Row 1 (Faces):** Shows six original face images at the top, followed by their reconstructed versions below.
- Row 2 (Faces):** Shows six original face images, followed by their reconstructed versions.
- Row 3 (3D Shapes):** Shows six original 3D shape models, followed by their reconstructed versions.
- Row 4 (3D Shapes):** Shows six original 3D shape models, followed by their reconstructed versions.
- Row 5 (Digits):** Shows six original digit images, followed by their reconstructed versions.
- Row 6 (Digits):** Shows six original digit images, followed by their reconstructed versions.

Missing half of the image

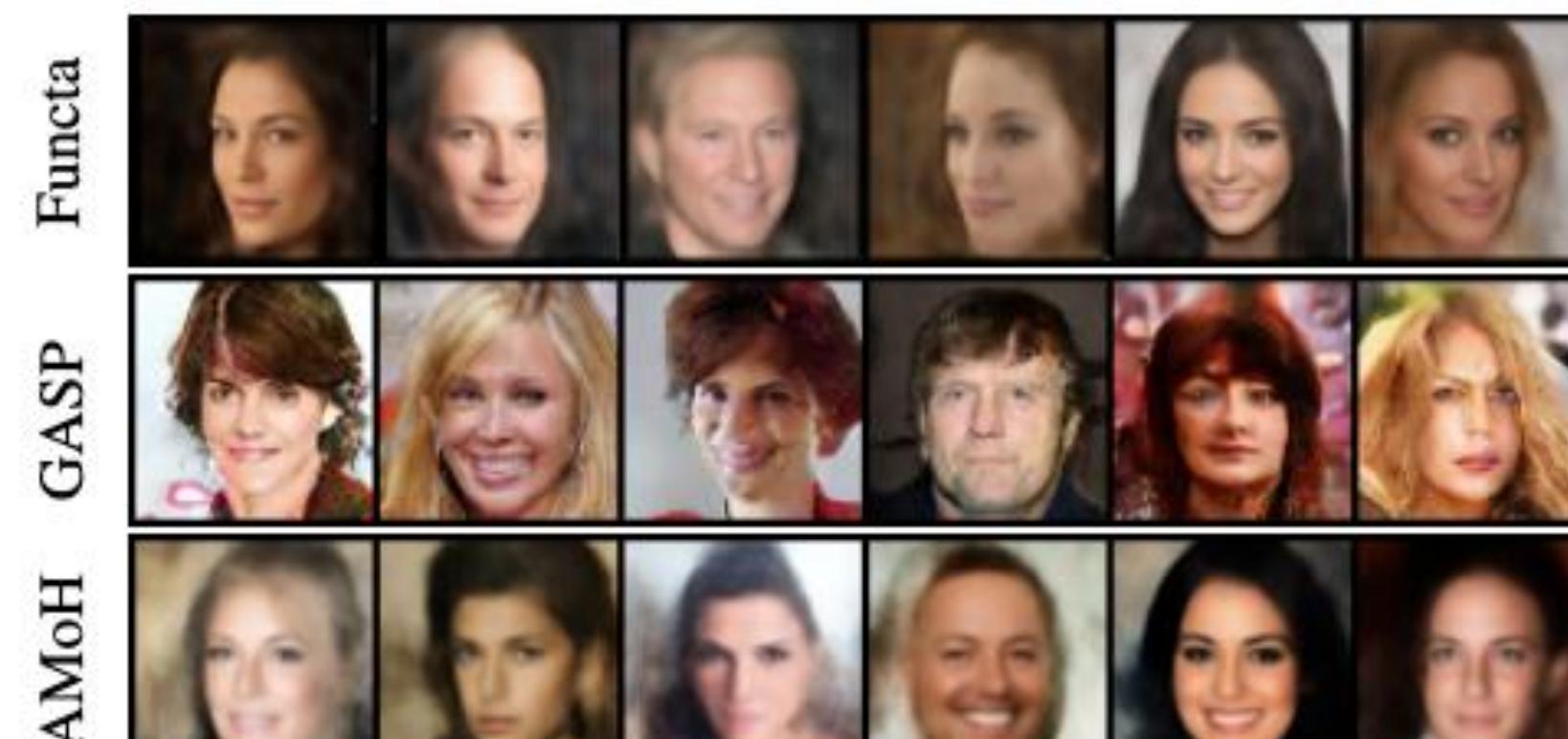
Image out-painting

Proposed method (2)

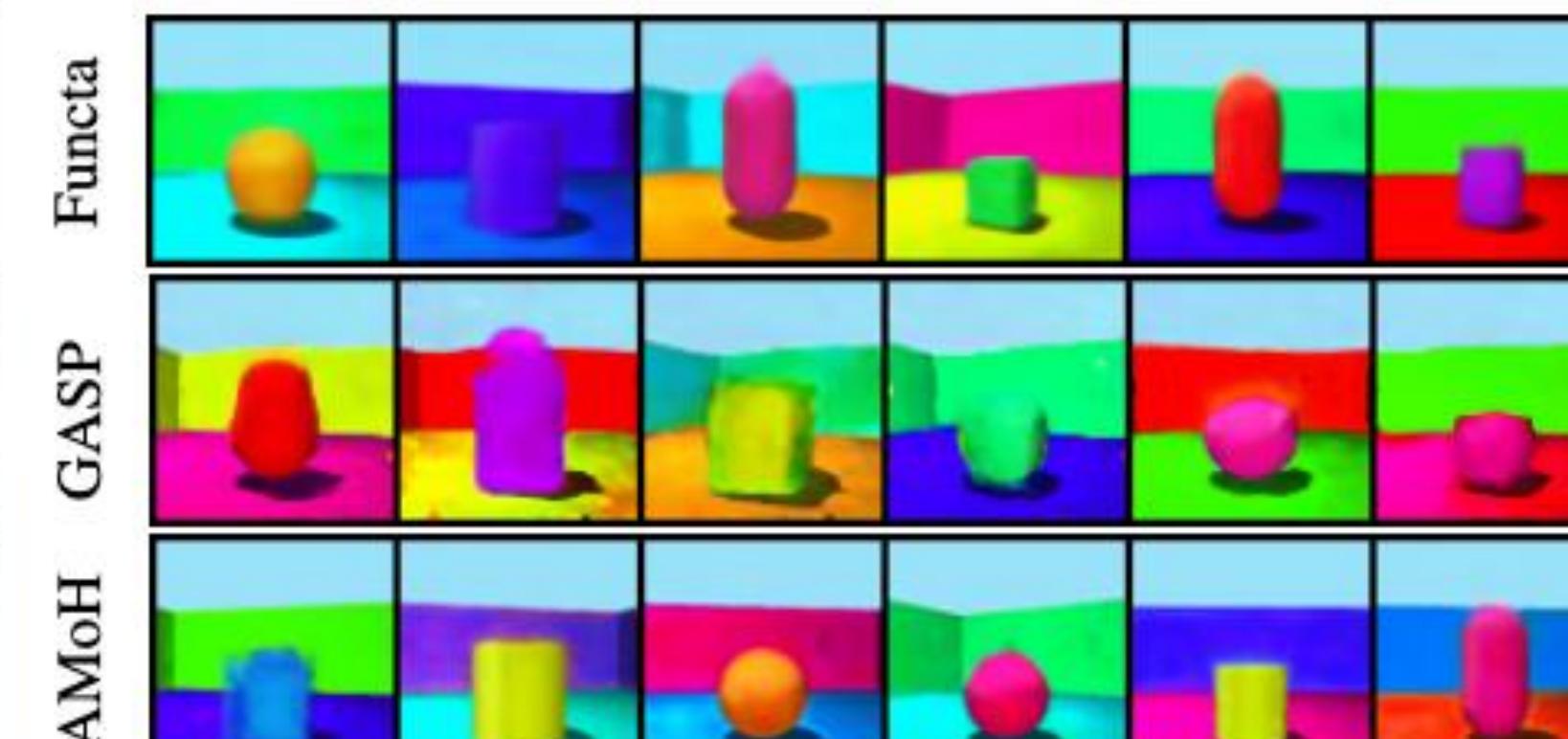
Limitations of previous work

Flexibility of the latent space in [5, 6, 25]

- This makes generation quality poor.



(a) CELEBA HQ



(b) SHAPES3D

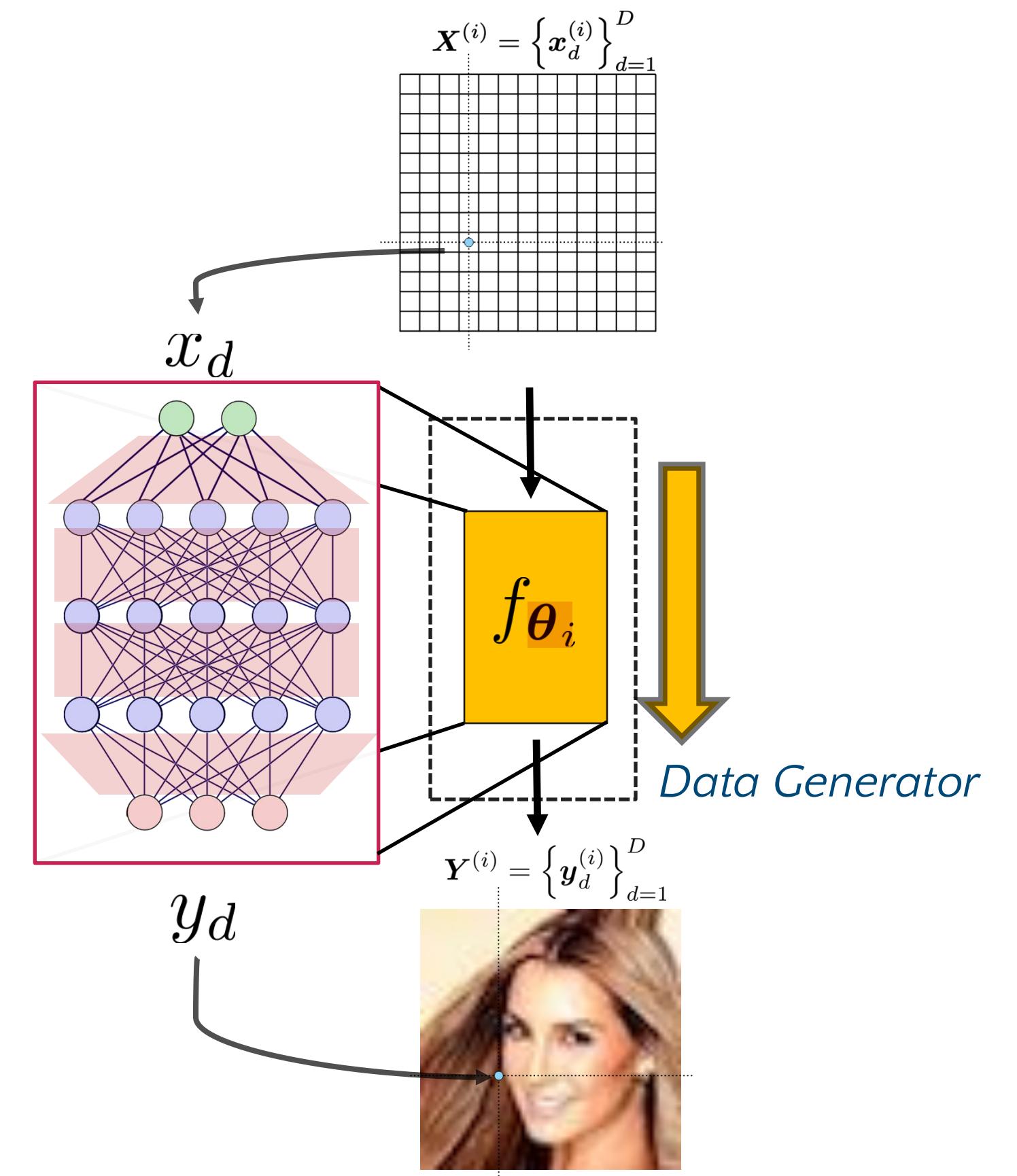
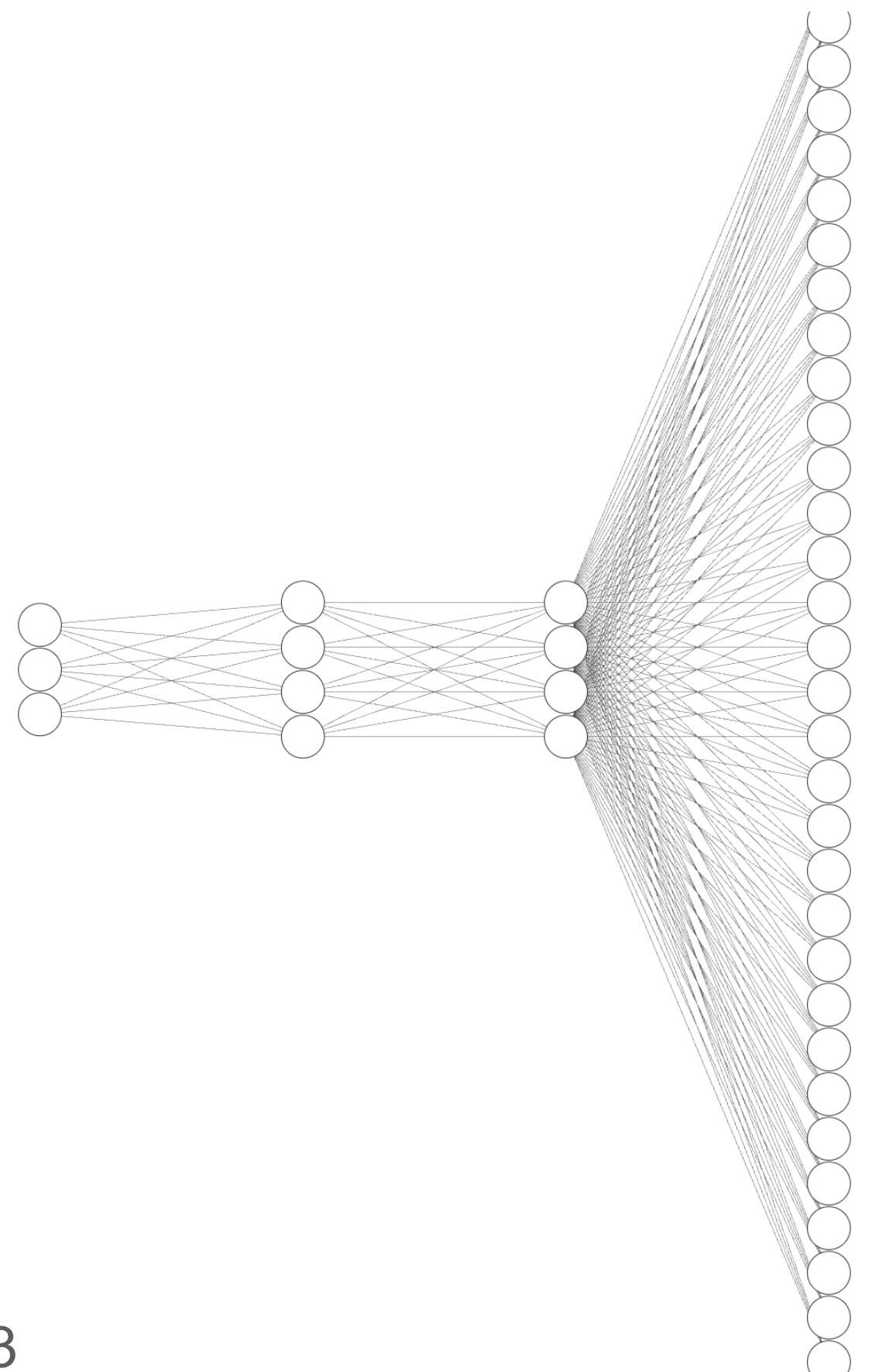
[5] Dupont et al., 2020

[6] Dupont et al., 2022

[25] Koyuncu et al., 2023

Limitations of previous work

Hypernet bottleneck in [5, 25]



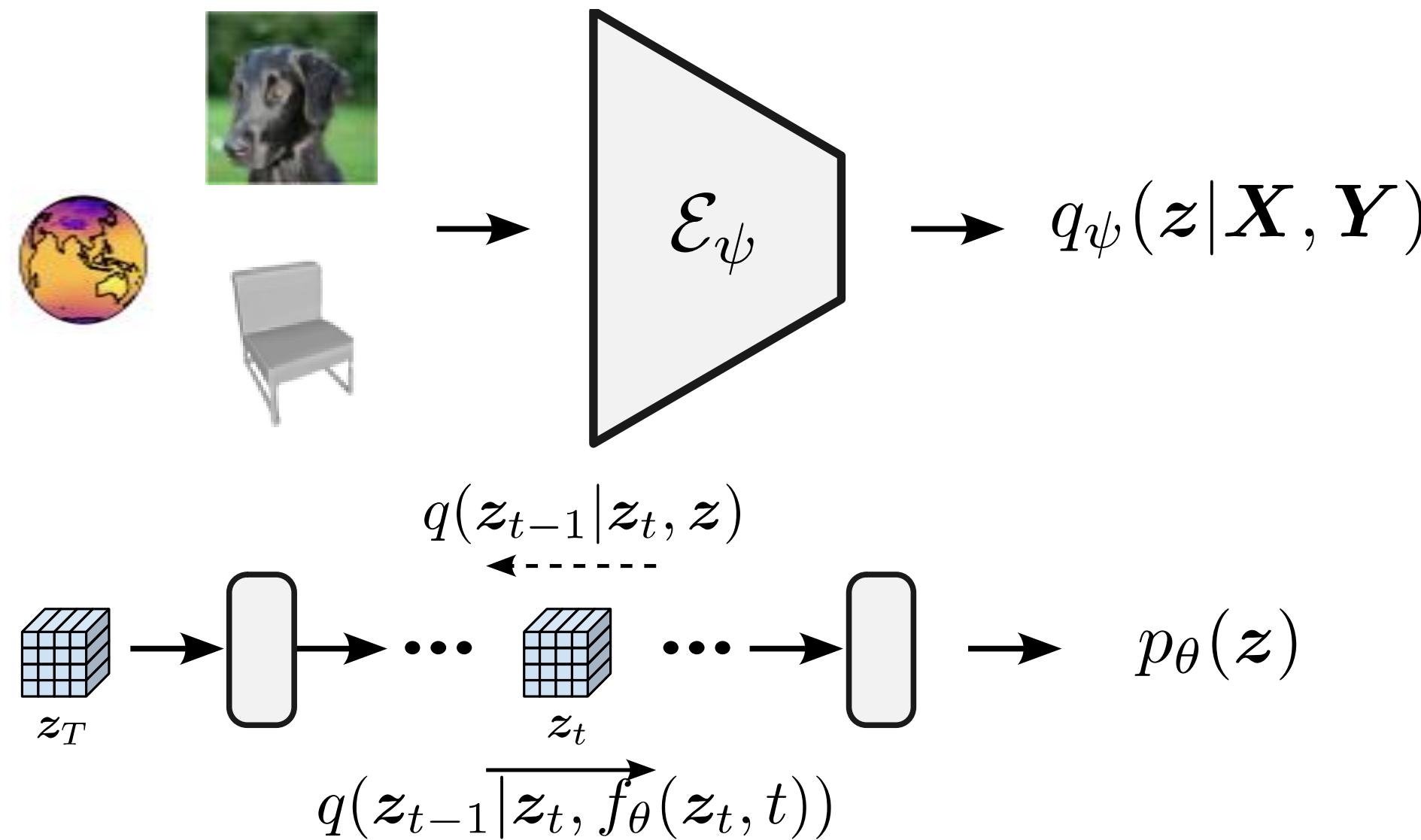
[5] Dupont et al., 2020

[25] Koyuncu et al., 2023

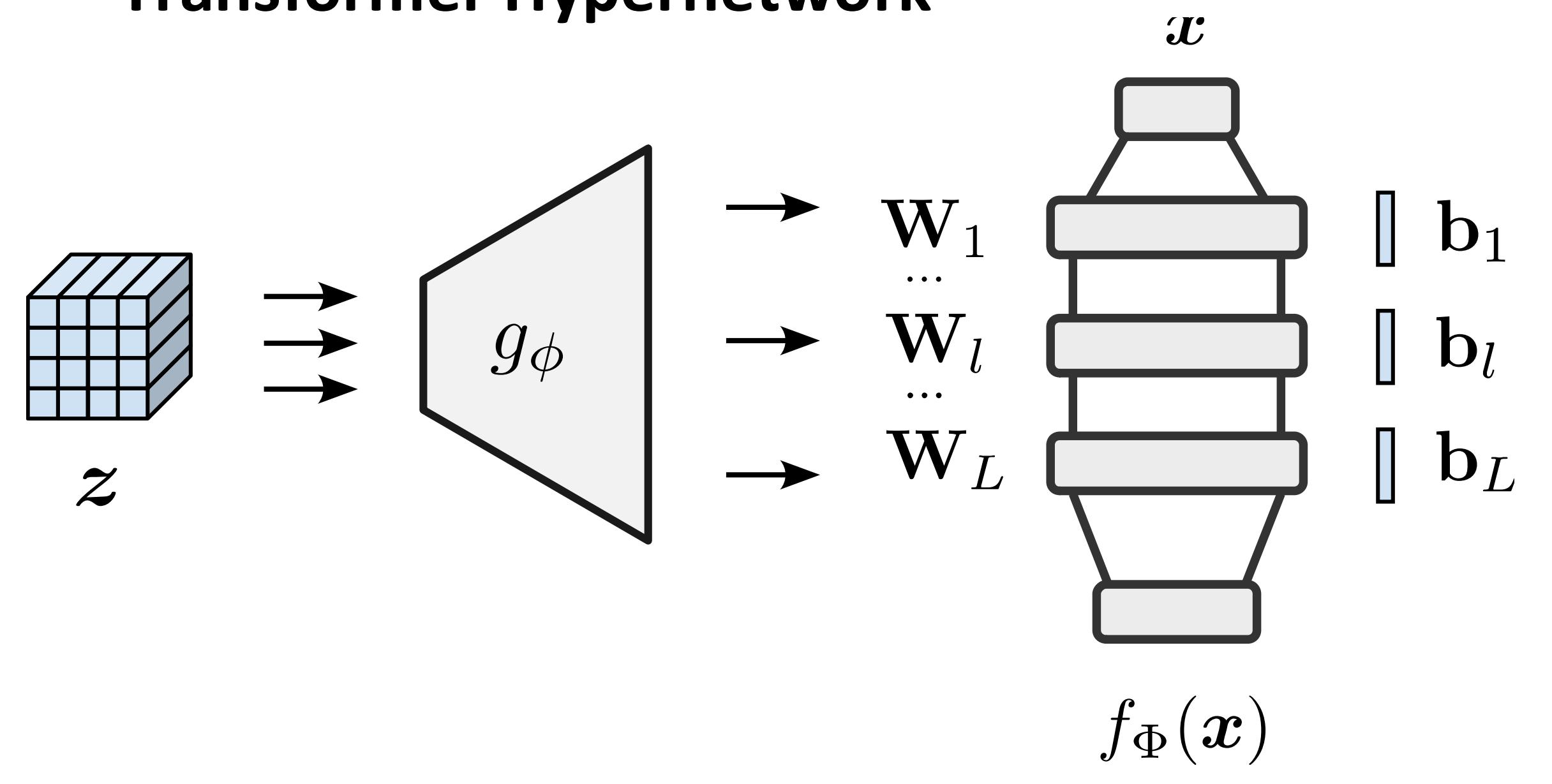
Proposed methods (2)

Hyper-Transforming Latent Variable Models [27] (LDMI)

Latent Diffusion [28]



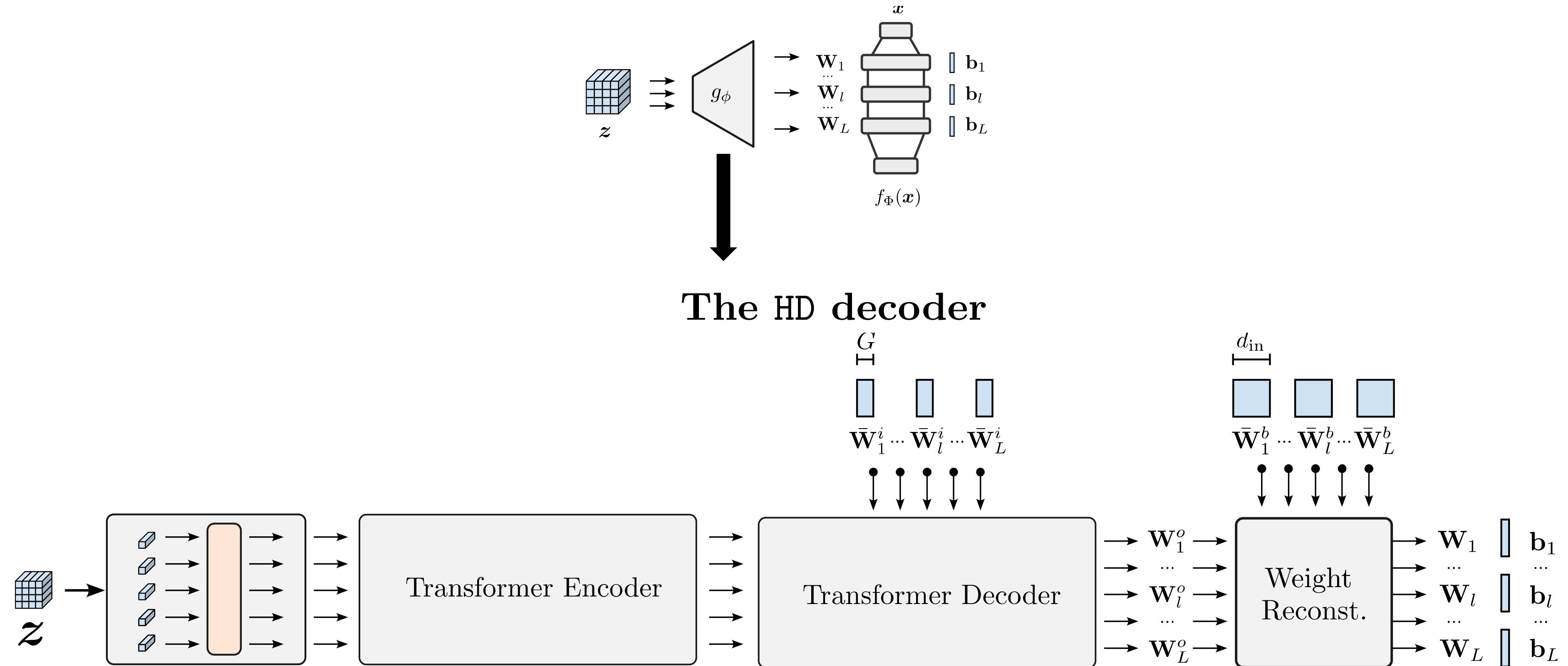
Transformer Hypernetwork



[27] Peis et al., 2025

Proposed methods (2)

The HD decoder

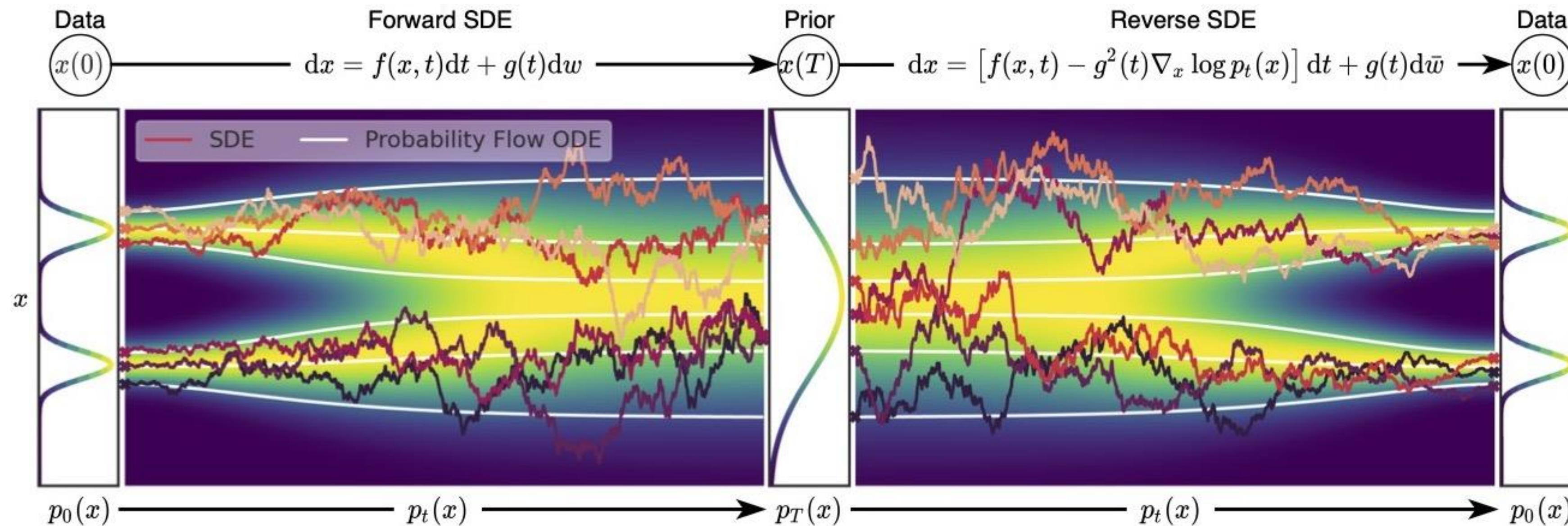


[27] Peis et al., 2025

Diffusion Models [29]

Denoising Score Matching

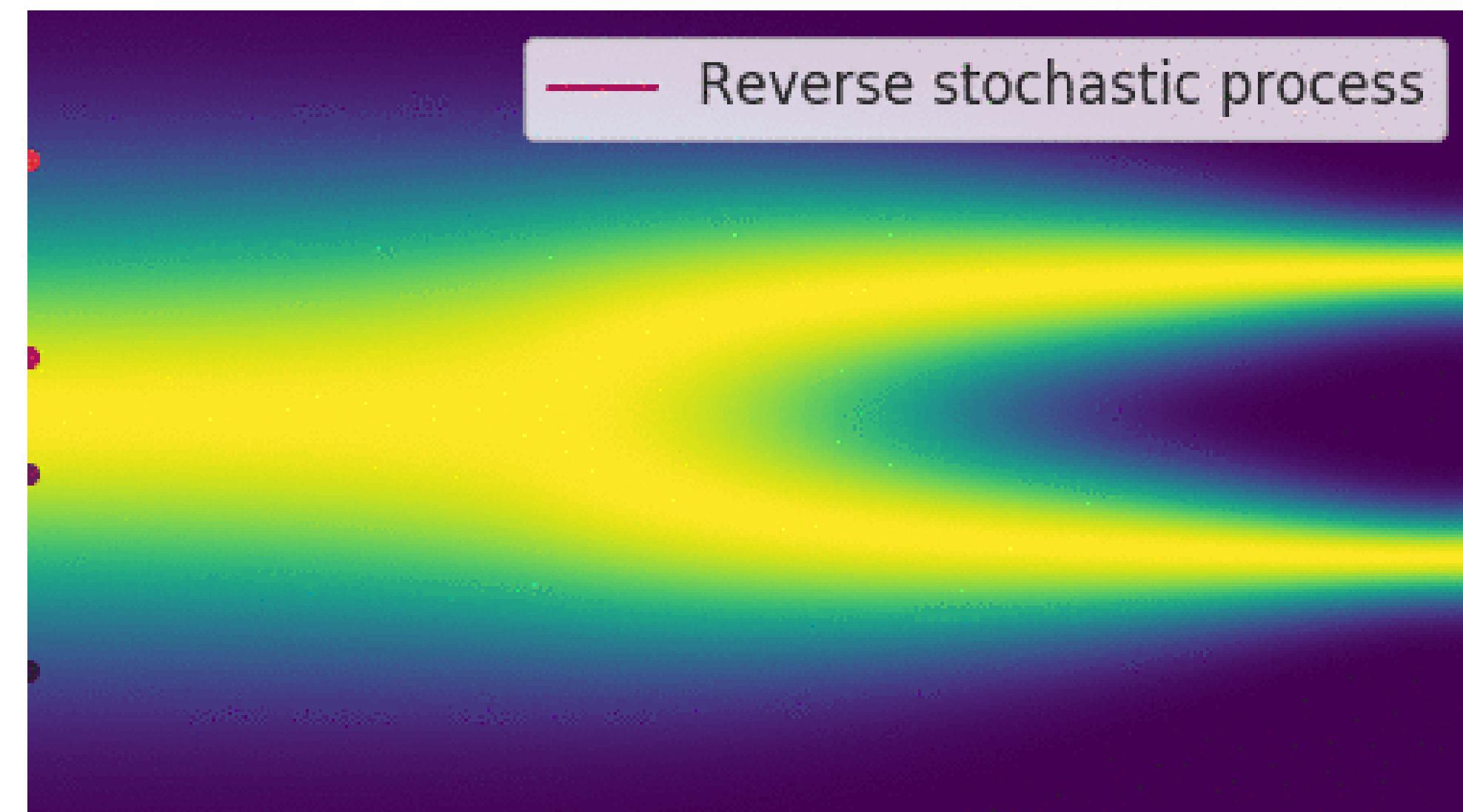
$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t) | \mathbf{x}(0)} \left[\left\| \mathbf{s}_{\theta}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \right\|_2^2 \right] \right\}$$



[29] Song et al., 2020

Diffusion Models [29]

$$s_{\theta}(\mathbf{x}_t, t)$$



[29] Song et al., 2020

DDPM [30]

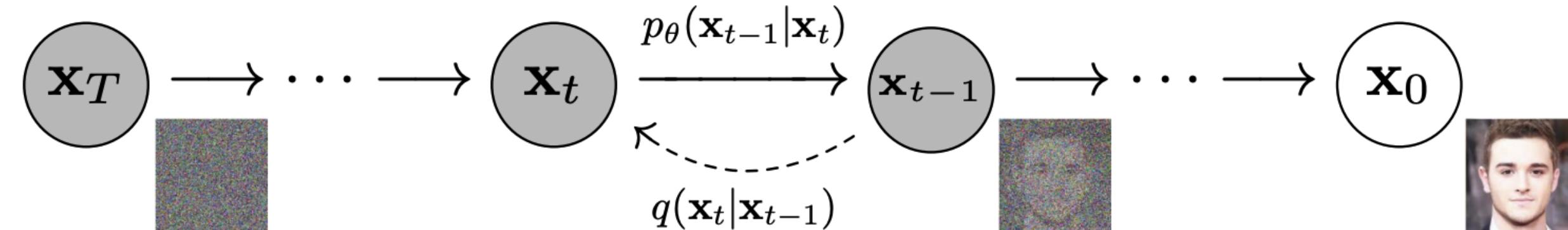


Figure 2: The directed graphical model considered in this work.

$$\begin{aligned}
 p_\theta(\mathbf{x}_{0:T}) &:= p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t), \quad p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \\
 q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) &:= \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1}), \quad q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) := \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}\right) \\
 q(\mathbf{x}_t \mid \mathbf{x}_0) &= \mathcal{N}\left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right) \quad \alpha_t := 1 - \beta_t \quad \bar{\alpha}_t := \prod_{s=1}^t \alpha_s
 \end{aligned}$$

$$\underbrace{\mathbb{E}_q[D_{\text{KL}}(q(\mathbf{x}_T \mid \mathbf{x}_0) \| p(\mathbf{x}_T))]}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1)}_{L_0}$$

[30] Ho et al., 2020

DDPM [30]

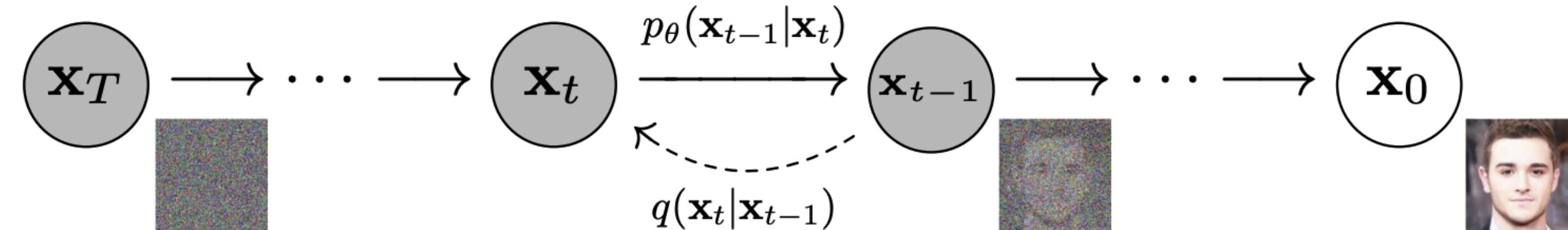


Figure 2: The directed graphical model considered in this work.

$$\mathbb{E}_q \underbrace{[D_{\text{KL}}(q(\mathbf{x}_T \mid \mathbf{x}_0) \| p(\mathbf{x}_T))]}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t))}_{L_{t-1}} \underbrace{- \log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1)}_{L_0}$$



$$\mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t \right) \right\|^2 \right]$$

DDPM [30]

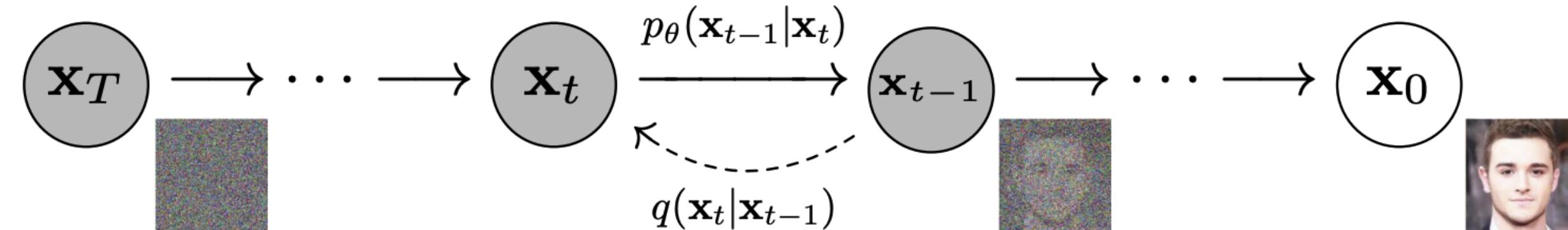


Figure 2: The directed graphical model considered in this work.

$$\mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t \right) \right\|^2 \right]$$

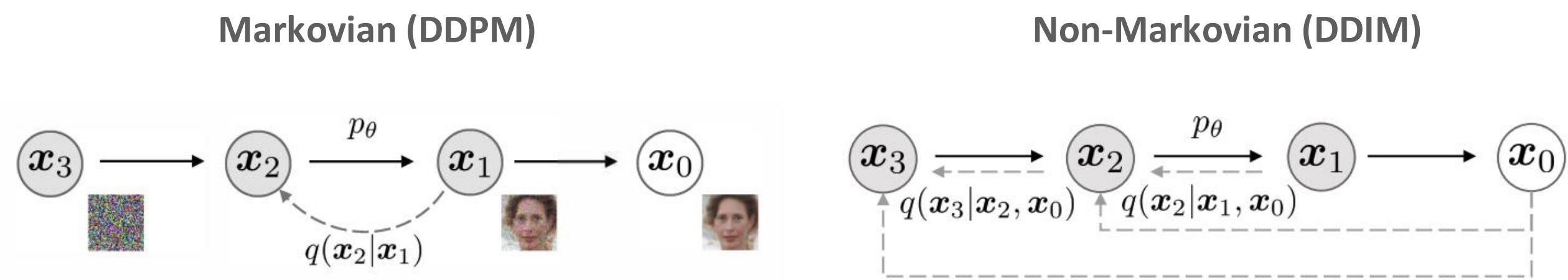


$$L_{\text{simple}} (\theta) := \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \left[\left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t \right) \right\|^2 \right]$$

LDMI

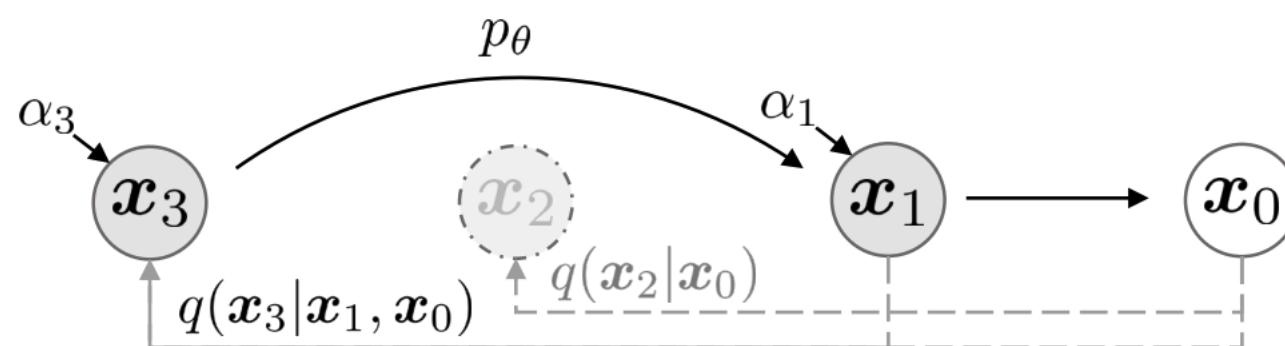
DDIM [31]

- Define a Non-Markovian Inference Model.
- The objective is the same!



$$L_{\text{simple}} (\theta) := \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \left[\left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t \right) \right\|^2 \right]$$

- Using the same model, you can sample in fewer steps!



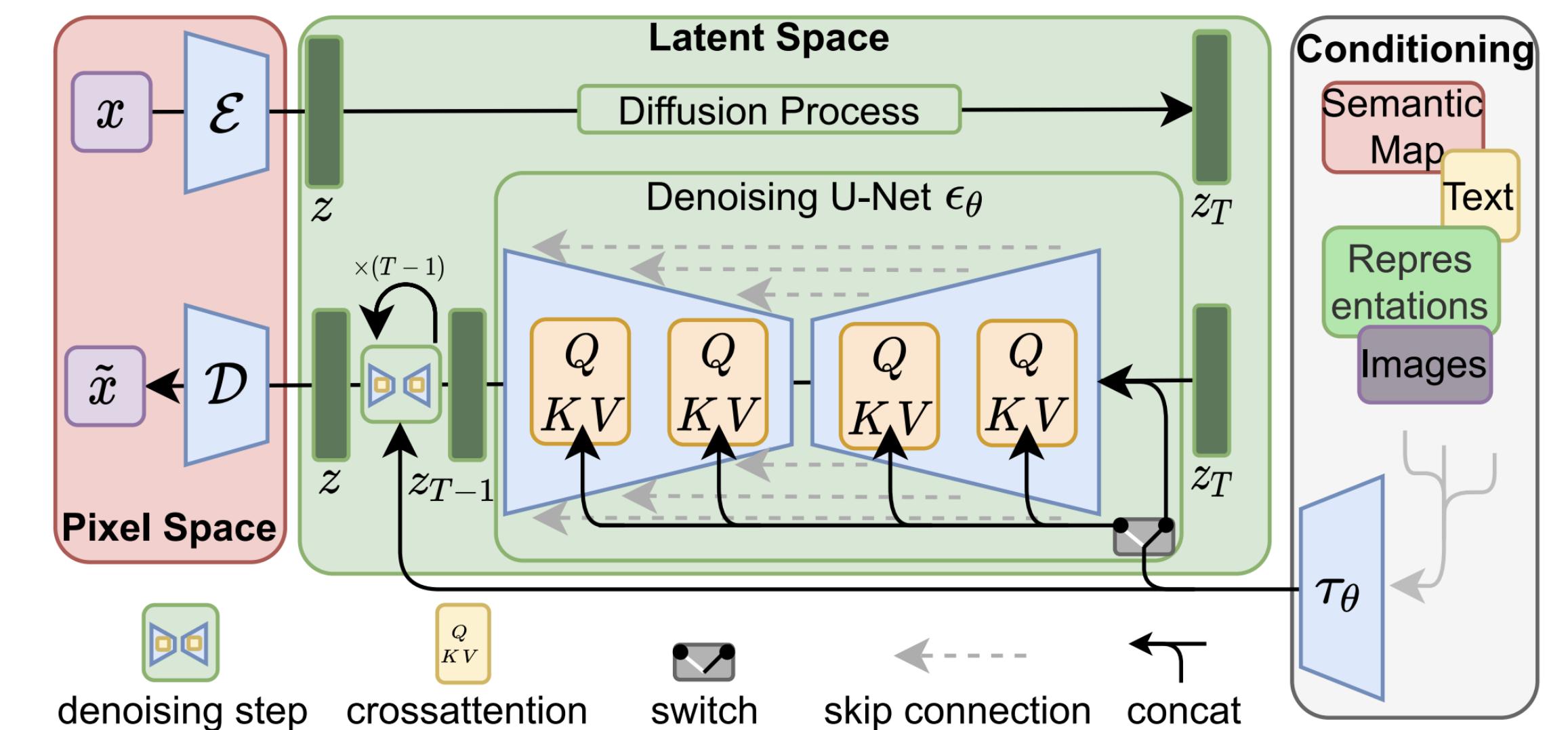
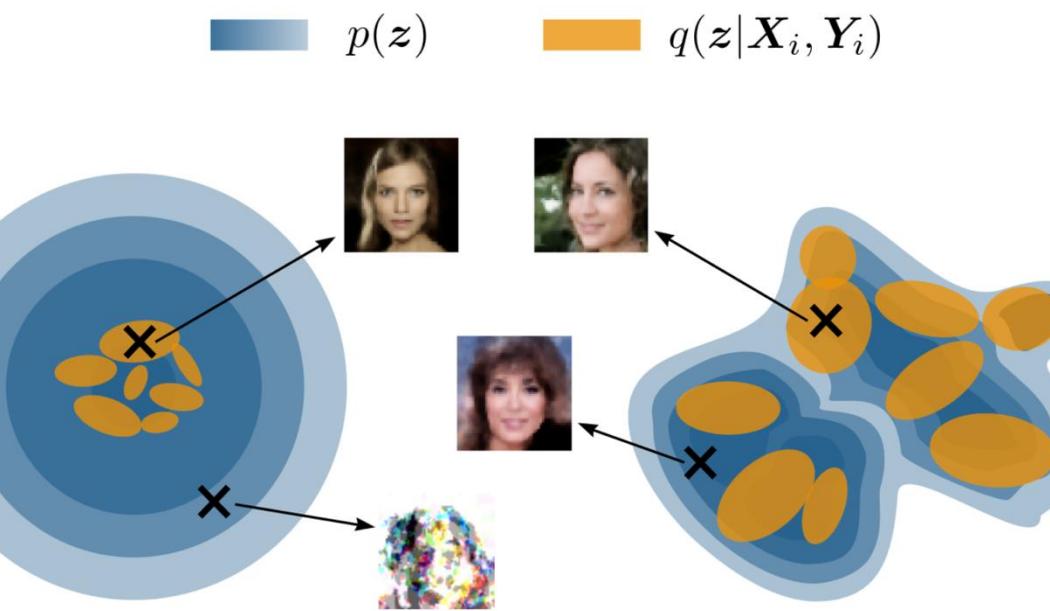
Latent Diffusion Models [28]

- First stage:

$$\begin{aligned}\mathcal{L}_{\text{VAE}}(\phi, \psi) = & \mathbb{E}_{q_\psi(z|X)} [\log p_\Phi(X)] \\ & - \beta \cdot D_{\text{KL}}(q_\psi(z|X) \| p(z)), \\ & + \mathcal{L}_{\text{perceptual}} + \mathcal{L}_{\text{GAN}}\end{aligned}$$

- Second stage:

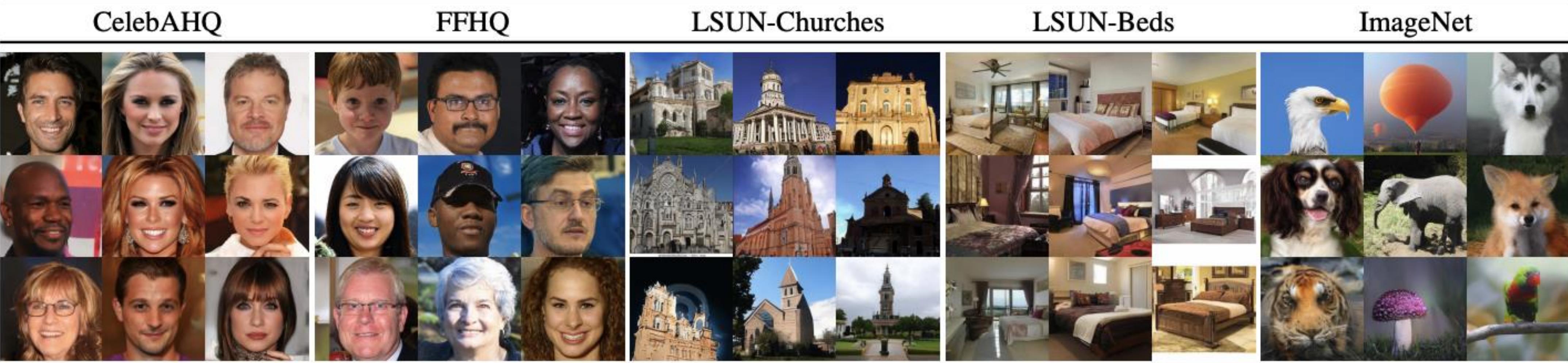
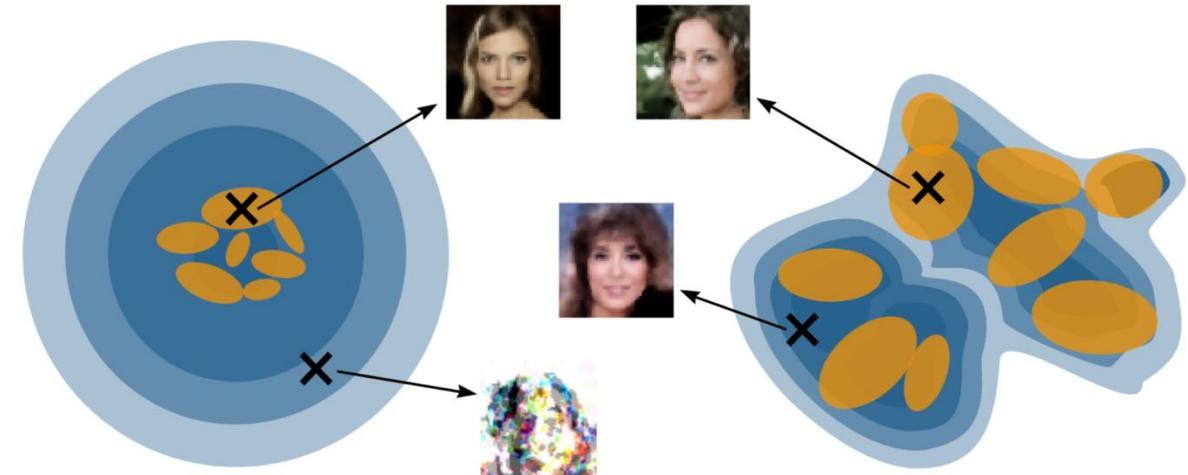
$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{X, z, \epsilon, t} \left[\lambda(t) \|\epsilon - \epsilon_\theta(z_t, t)\|^2 \right],$$



LDMI

Latent Diffusion Models

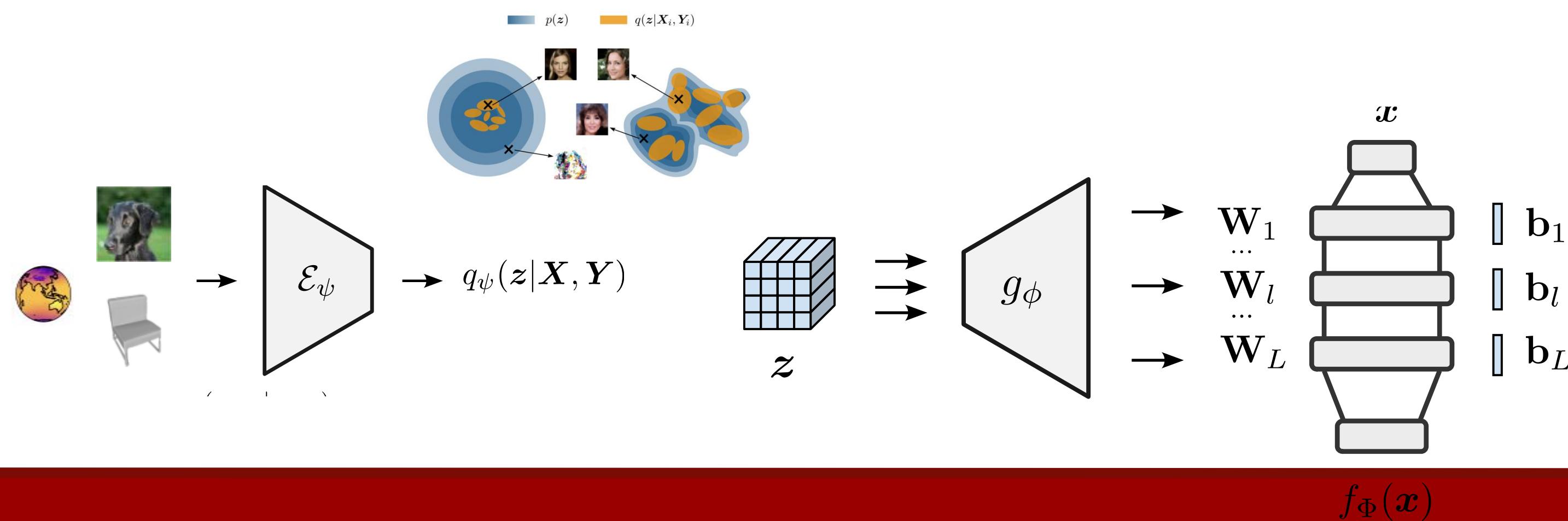
$$\textcolor{blue}{p(\mathbf{z})} \quad \textcolor{orange}{q(\mathbf{z}|\mathbf{X}_i, \mathbf{Y}_i)}$$



Latent Diffusion Models for Implicit Neural Representations

- We will train an “*under-regularized*” autoencoder (VAE or VQ-VAE) to accurately represent data in a (tensor-shaped) latent space.
 - The latents are mapped into INRs using our **transformer–based hypernetwork decoder**.

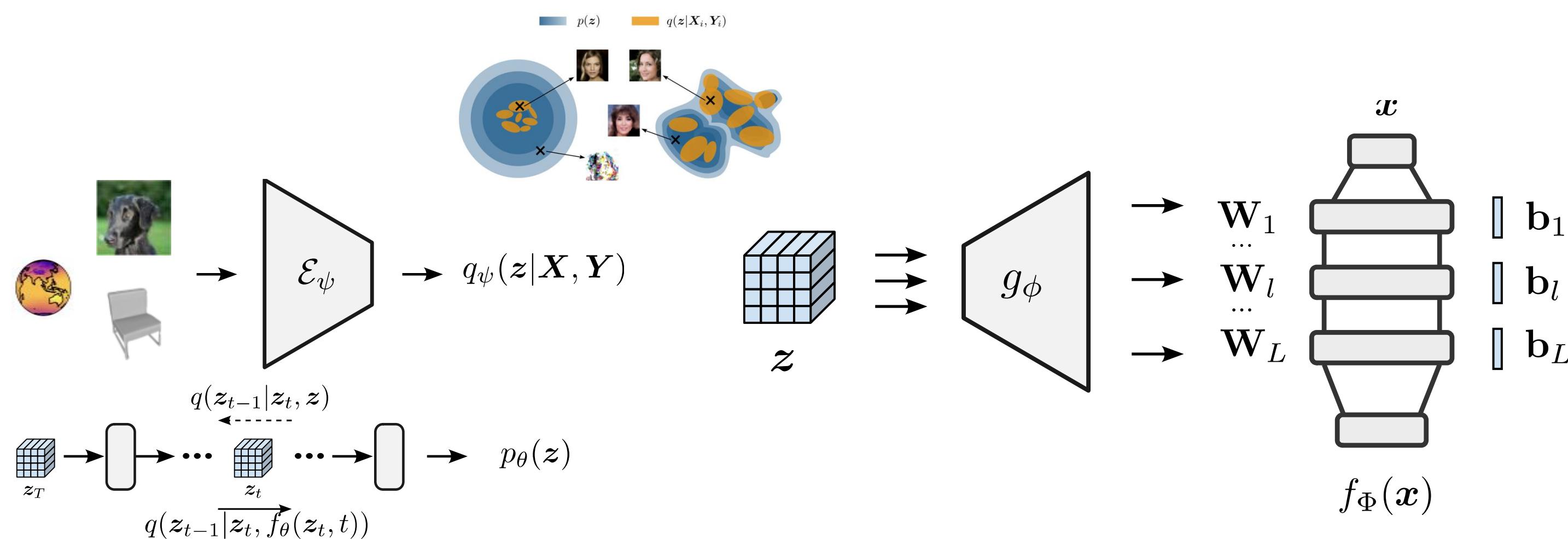
$$\begin{aligned}\mathcal{L}_{\text{VAE}}(\phi, \psi) = & \mathbb{E}_{q_{\psi}(\mathbf{z} | \mathbf{X}, \mathbf{Y})} [\log p_{\Phi}(\mathbf{Y} | \mathbf{X})] \\ & - \beta \cdot D_{\text{KL}} (q_{\psi}(\mathbf{z} | \mathbf{X}, \mathbf{Y}) \| p(\mathbf{z})) ,\end{aligned}$$



Latent Diffusion Models for Implicit Neural Representations

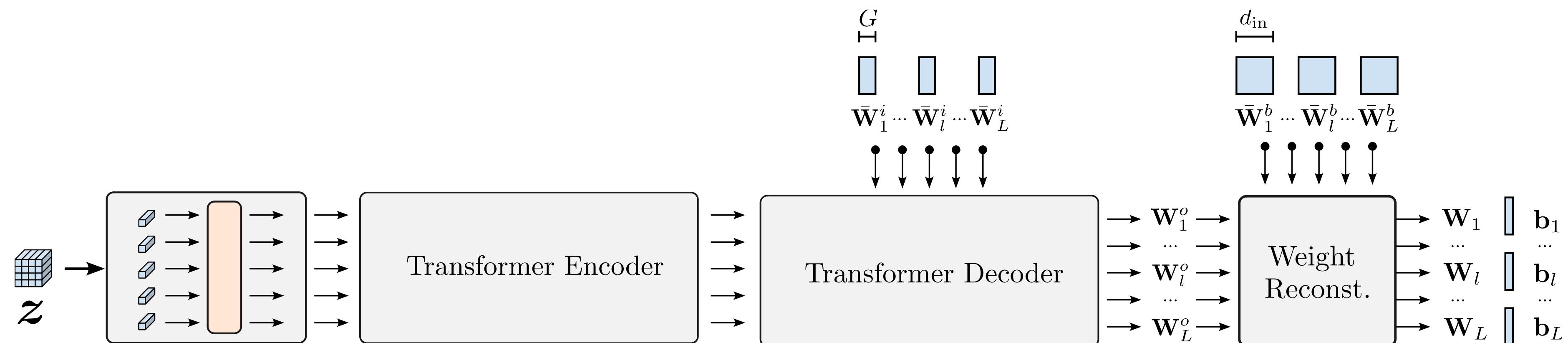
2. We will fit a Diffusion Model (DDPM) to the learned latent space.

$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \mathbf{z}, \epsilon, t} \left[\lambda(t) \|\epsilon - \epsilon_\theta(\mathbf{z}_t, t)\|^2 \right],$$



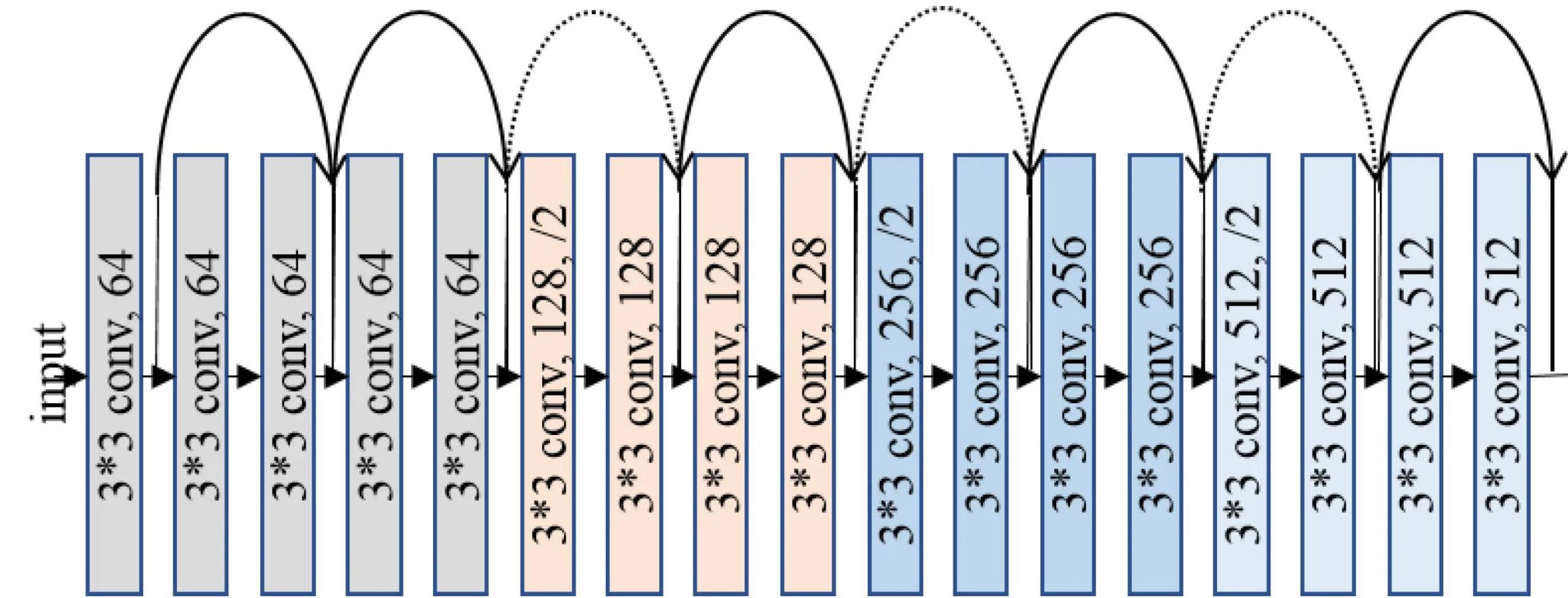
The Hyper-Transformer Decoder

- The latents are **tokenized** (following ViT [32]).
- Two sets of globally shared, learnable parameters:
 - Compressed weights that cross-attend the latent tokens.
 - Full weights to expand the compressed weights.



ResNet encoders

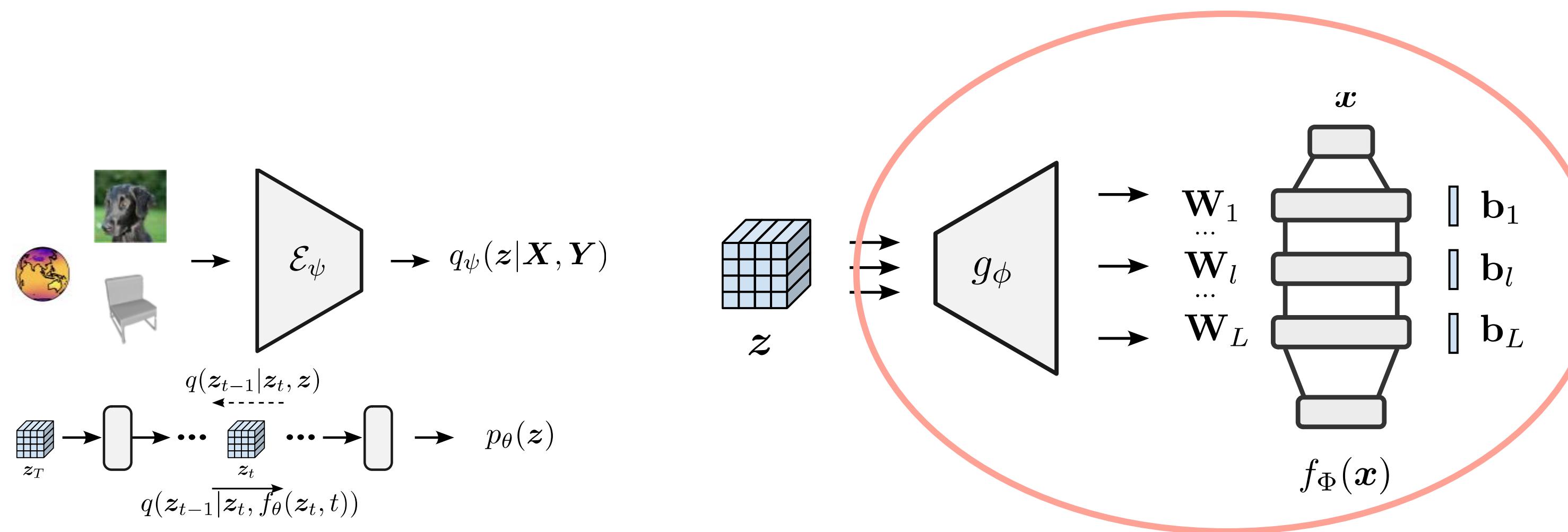
- The data is stored in a structured representation.
- We can make use of powerful encoders tailored to structured data.



Hyper-Transforming

- We can download pre-trained LDMs and just re-train only our decoder!

$$\mathcal{L}_{\text{HT}}(\phi) = \mathbb{E}_{q_{\psi}(\mathbf{z}|\mathbf{X}_m, \mathbf{Y}_m)} [\log p_{\Phi}(\mathbf{Y} \mid \mathbf{X})] + \mathcal{L}_{\text{perceptual}} + \mathcal{L}_{\text{GAN}}$$



Pretrained LDMs

Datset	Task	Model	FID	IS	Prec	Recall	
CelebA-HQ	Unconditional Image Synthesis	LDM-VQ-4 (200 DDIM steps, eta=0)	5.11 (5.11)	3.29	0.72	0.49	https://omr-diffusion/ci
FFHQ	Unconditional Image Synthesis	LDM-VQ-4 (200 DDIM steps, eta=1)	4.98 (4.98)	4.50 (4.50)	0.73	0.50	https://omr-diffusion/ff
LSUN-Churches	Unconditional Image Synthesis	LDM-KL-8 (400 DDIM steps, eta=0)	4.02 (4.02)	2.72	0.64	0.52	https://omr-diffusion/l8

Experiments

Datasets

CelebA (64x64)



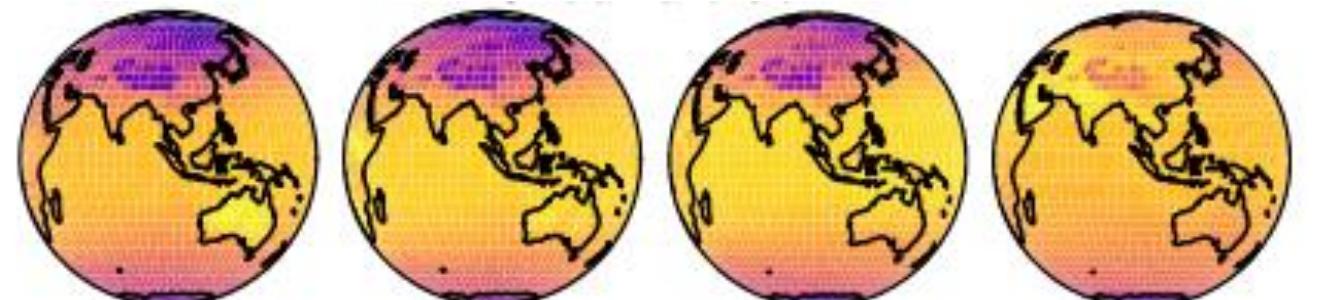
CelebA-HQ (256x256)



ImageNet (256x256)



ERA5 (Polar)



ShapeNET (Voxels)



Experiments

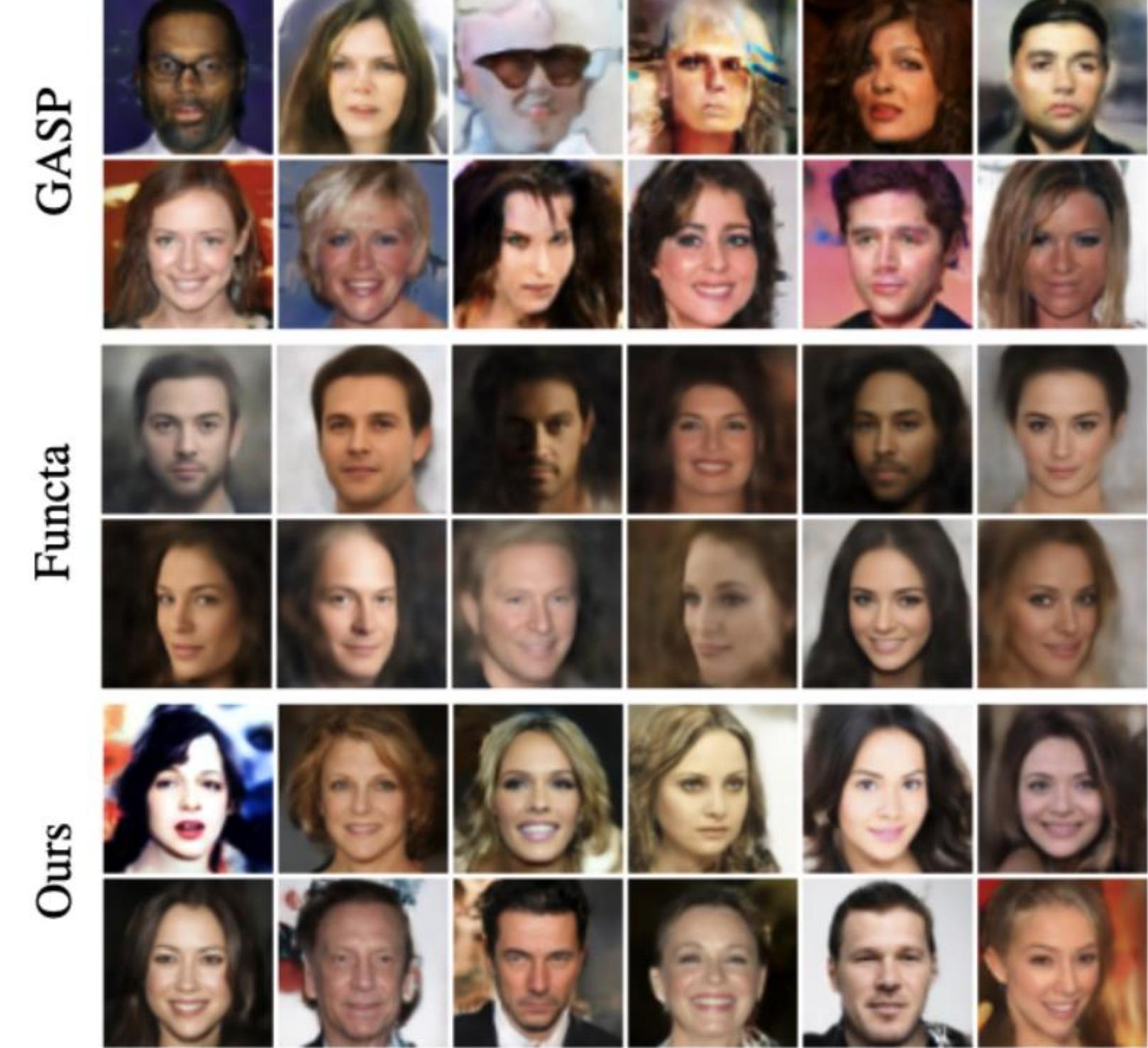
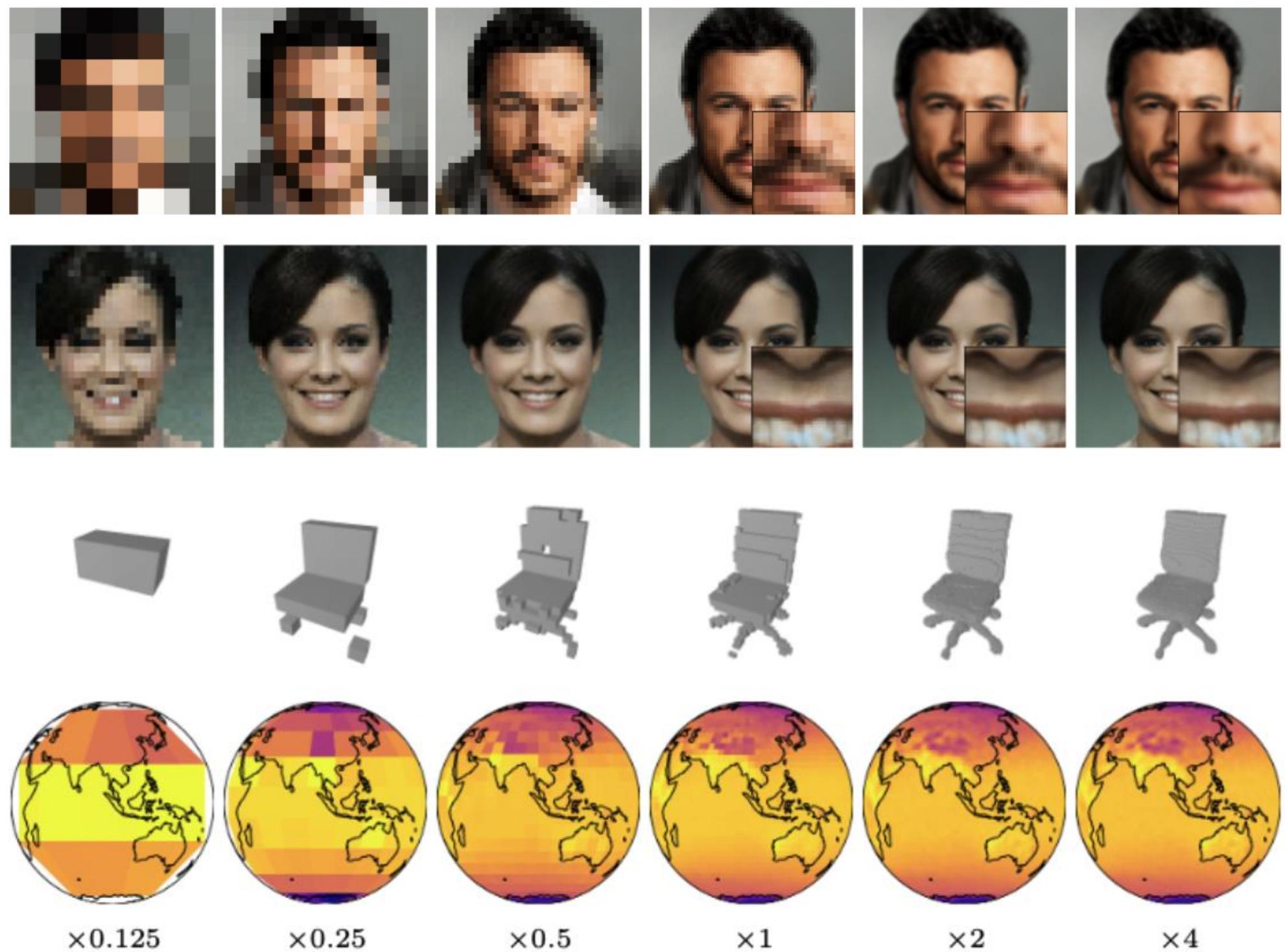
Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution	Scalable	Flexible
GASP (2021) [5]	GAN	Minimax	Forward Pass	✗	✗	✗
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample	✗	✗
VaMoH (ours)	VAE-based	Single optimization	Forward Pass	Forward pass	✗	✗
LDMI (ours)	LDM-based	Hyper-Transforming	Forward Pass	Forward pass	✓	✓

LDMI enhances efficiency, scalability quality of the learned representations.

Experiments

Generation: qualitative results



(a) CelebA-HQ

Experiments

Generation: quantitative results

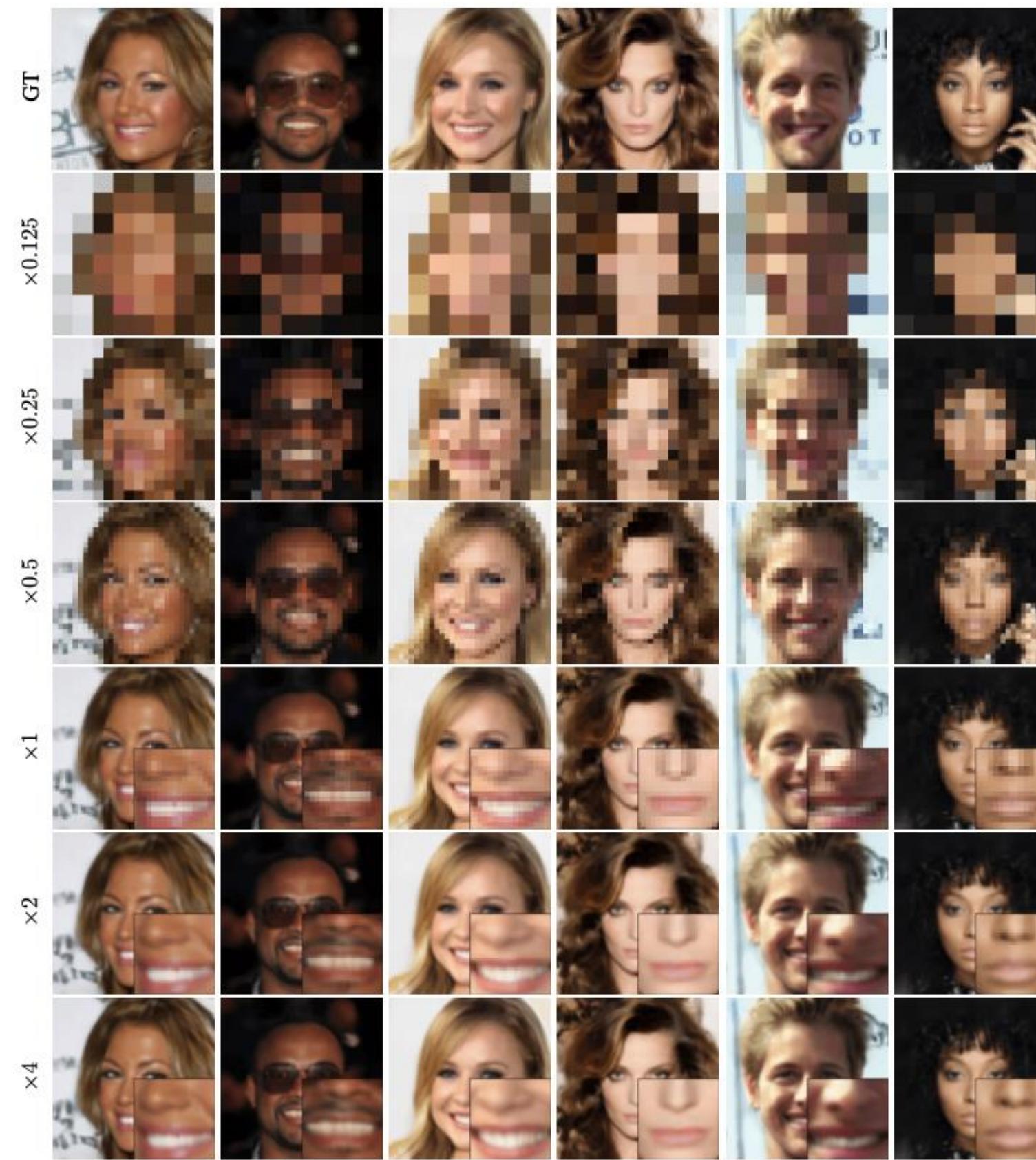
Model	PSNR (dB) \uparrow	FID \downarrow	HN Params \downarrow
CelebA-HQ (64×64)			
GASP [Dupont et al., 2022a]	-	7.42	25.7M
Functa [Dupont et al., 2022b]	≤ 30.7	40.40	-
VAMoH [Koyuncu et al., 2023]	23.17	66.27	25.7M
LDMI	24.80	18.06	8.06M
ImageNet (256×256)			
Spatial Functa [Bauer et al., 2023]	≤ 38.4	≤ 8.5	-
LDMI	20.69	6.94	102.78M

Table 1: Metrics on CelebA-HQ and ImageNet.

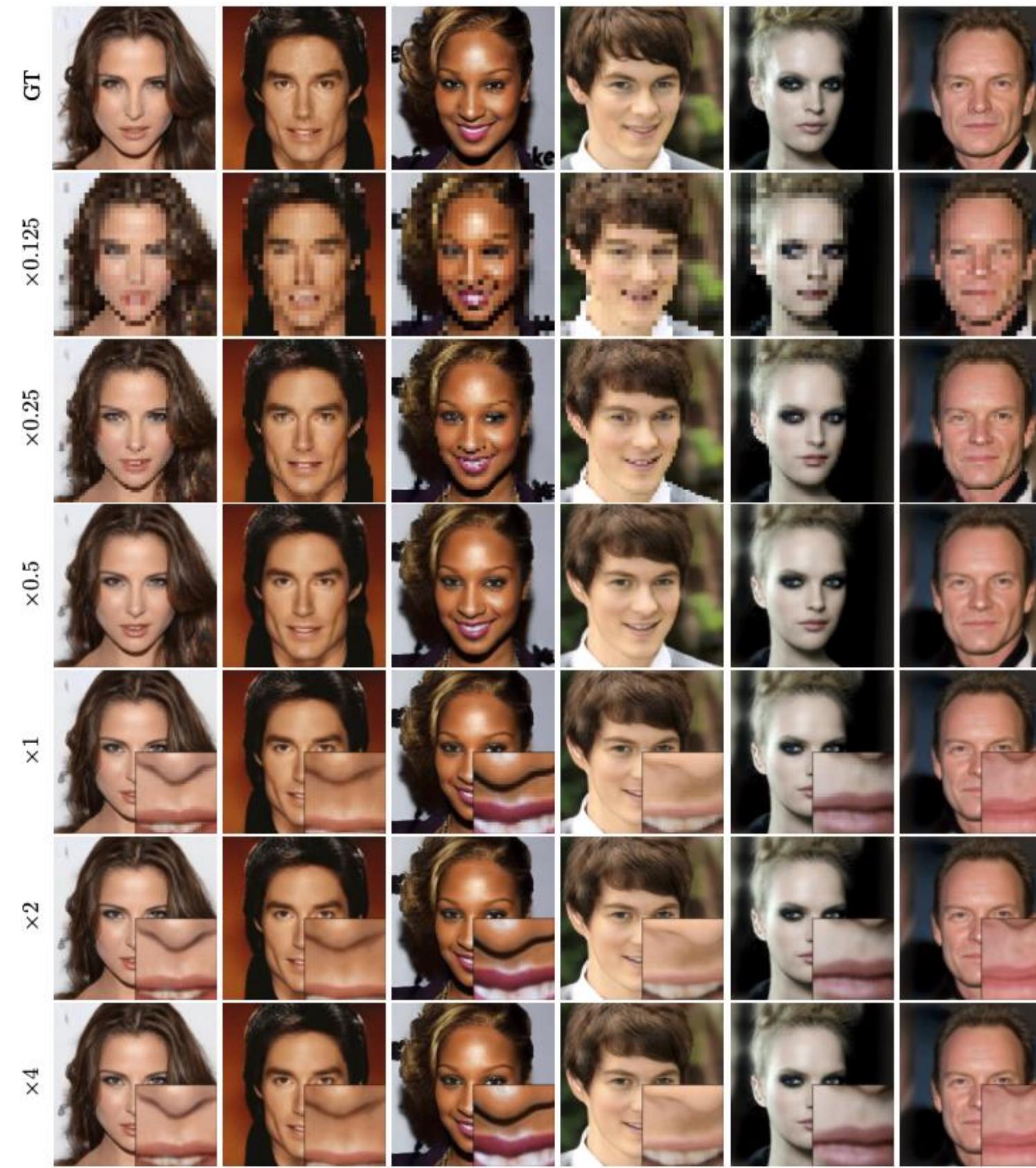
Experiments

Reconstruction

CelebA-HQ (64x64)



CelebA-HQ (256x256)



Model	Chairs (PSNR) \uparrow	ERA5 (PSNR) \uparrow
Functa [Dupont et al., 2022b]	29.2	34.9
VAMoH [Koyuncu et al., 2023]	38.4	39.0
LDMI	38.8	44.6

Table 2: Reconstruction quality (PSNR in dB) on ShapeNet Chairs and ERA5 climate data, demonstrating LDMI's strong generalization capabilities across modalities. Note that GASP is omitted as it is not applicable to INR reconstruction tasks.

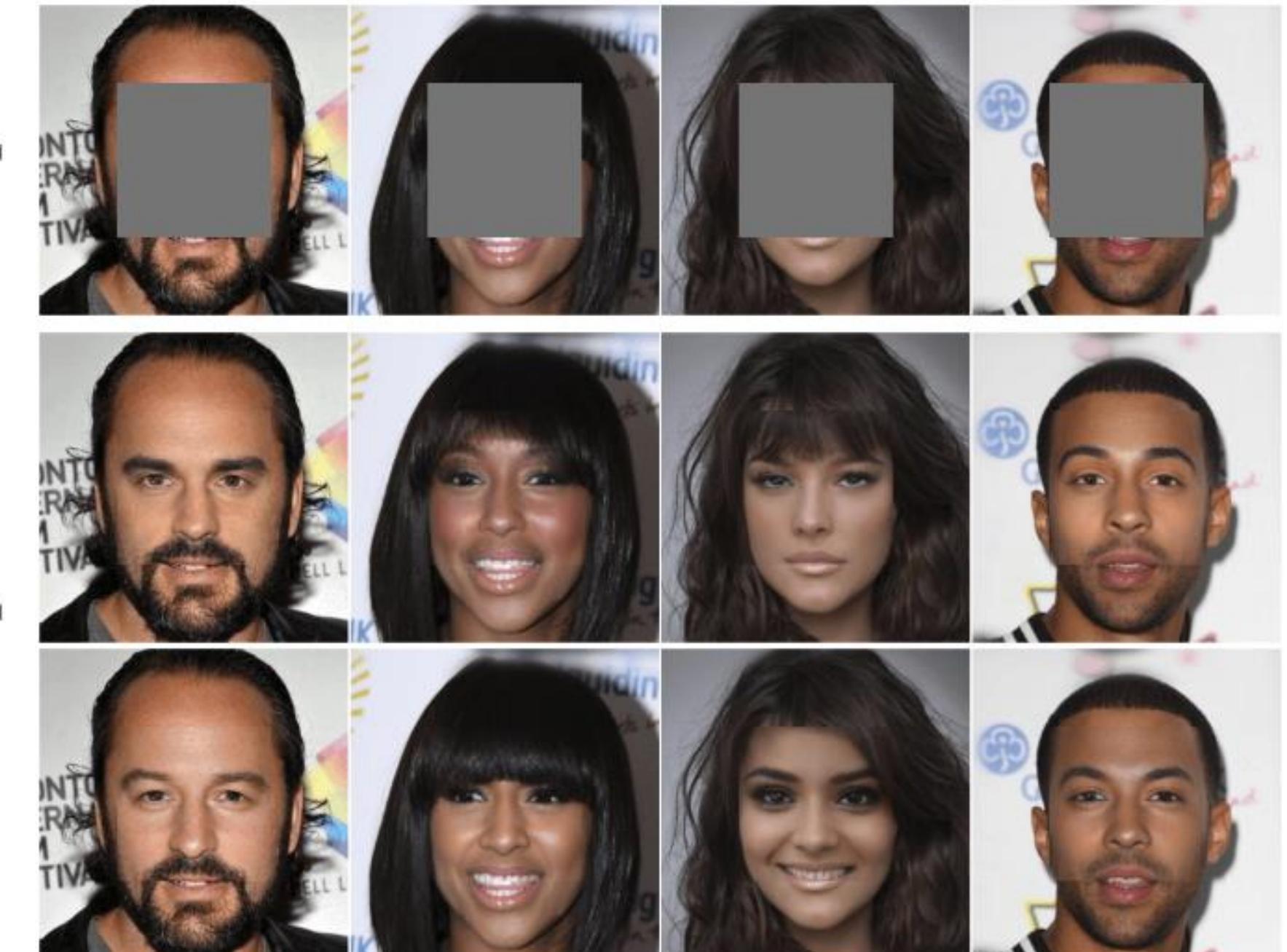
Experiments

Data completion

VAMoH



Input



Samples

LDMI

Experiments

Parameter efficiency

Method	HN Params	INR Weights	Ratio (INR/HN)
GASP/VAMoH	25.7M	50K	0.0019
LDMI	8.06M	330K	0.0409

Table 3: Parameter efficiency of hypernetworks (HN) in GASP/VaMoH and LDMI.

Method	HN Params	PSNR (dB)
LDMI-MLP	17.53M	24.93
LDMI-HD	8.06M	27.72

Table 4: Ablation study comparing MLP and hyper-transformer HD decoders on CelebA-HQ.

Conclusion

Thanks to learning **distributions of functions**, our proposed **VAMoH** can easily perform:

- Generation.
- Reconstruction.
- Conditional generation.
- Super resolution (interpolation).

While being:

- ✓ Robust to partially observed data.
- ✓ Expressive for generating high-quality data.
- ✓ Efficient in terms of inference.

Conclusion

Thanks to using **Latent Diffusion** and a **Transformer-based hypernetwork**, **LDMI** enhances

- Generation quality.
- Reconstruction accuracy.
- Conditional generation.
- Super resolution.

While:

- ✓ Being scalable.
- ✓ Being parameter efficient.
- ✓ Allowing for generation of **bigger INRs** and more complex data.

Further details

VARIATIONAL MIXTURE OF HYPERGENERATORS FOR LEARNING DISTRIBUTIONS OVER FUNCTIONS

Batuhan Koyuncu*
Saarland University
Saarbrücken, Germany

Pablo Sánchez-Martín
Max Planck Institute for Intelligent Systems
Tübingen, Germany

Ignacio Peis
Universidad Carlos III de Madrid
Madrid, Spain

Pablo M. Olmos
Universidad Carlos III de Madrid
Madrid, Spain

Isabel Valera
Saarland University
Saarbrücken, Germany



[\[Paper\]](#)

Further details

HYPER-TRANSFORMING LATENT DIFFUSION MODELS

Ignacio Peis*
Technical University of Denmark

Isabel Valera
Saarland University

Batuhan Koyuncu
Saarland University

Jes Frellsen
Technical University of Denmark



[\[Paper\]](#)

[27] Peis et al., 2025

References

- [1] Campbell, A., Chen, W., Stimper, V., Hernandez-Lobato, J. M., & Zhang, Y. (2021, July). A gradient based strategy for hamiltonian monte carlo hyperparameter optimization. In *International Conference on Machine Learning* (pp. 1238-1248). PMLR.
- [2] Caterini, A. L., Doucet, A., & Sejdinovic, D. (2018). Hamiltonian variational auto-encoder. *Advances in Neural Information Processing Systems*, 31.
- [3] Salimans, T., Kingma, D., & Welling, M. (2015, June). Markov chain monte carlo and variational inference: Bridging the gap. In *International conference on machine learning* (pp. 1218-1226). PMLR.
- [4] Ruiz, F. J., Titsias, M. K., Cemgil, T., & Doucet, A. (2021, December). Unbiased gradient estimation for variational auto-encoders using coupled Markov chains. In *Uncertainty in Artificial Intelligence* (pp. 707-717). PMLR.
- [5] Dupont, E., Whyte Teh, Y. & Doucet, A.. (2022). Generative Models as Distributions of Functions. *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics, in Proceedings of Machine Learning Research* 151:2989-3015.
- [6] Dupont, E., Kim, H., Eslami, S. A., Rezende, D. J., & Rosenbaum, D. (2022, June). From data to functa: Your data point is a function and you can treat it like one. In *International Conference on Machine Learning* (pp. 5694-5725). PMLR.

References

- [7] Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- [8] Cremer, C., Li, X., & Duvenaud, D. (2018, July). Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning* (pp. 1078-1086). PMLR.
- [9] Bernardo, J. M. (1979). Expected information as expected utility. *the Annals of Statistics*, 686-690.
- [10] Ma, C., Tschiatschek, S., Palla, K., Hernández-Lobato, J. M., Nowozin, S., & Zhang, C. (2018). Eddi: Efficient dynamic discovery of high-value information with partial vae. *arXiv preprint arXiv:1809.11142*.
- [11] Ma, C., Tschiatschek, S., Turner, R., Hernández-Lobato, J. M., & Zhang, C. (2020). VAEM: a deep generative model for heterogeneous mixed type data. *Advances in Neural Information Processing Systems*, 33, 11237-11247.
- [12] Child, R. (2020). Very deep vaes generalize autoregressive models and can outperform them on images. *arXiv preprint arXiv:2011.10650*.

References

- [13] Nazabal, A., Olmos, P. M., Ghahramani, Z., & Valera, I. (2020). Handling incomplete heterogeneous data using vaes. *Pattern Recognition*, 107, 107501.
- [14] Mattei, P. A., & Frellsen, J. (2019, May). MIWAE: Deep generative modelling and imputation of incomplete data sets. In *International conference on machine learning* (pp. 4413-4423). PMLR.
- [15] [Peis, I., Ma, C., & Hernández-Lobato, J. M. \(2022\). Missing Data Imputation and Acquisition with Deep Hierarchical Models and Hamiltonian Monte Carlo. arXiv preprint arXiv:2202.04599.](#)
- [16] Kraskov, A., Stögbauer, H., & Grassberger, P. (2004). Estimating mutual information. *Physical review E*, 69(6), 066138.
- [17] Gong, W., Li, Y., & Hernández-Lobato, J. M. (2020). Sliced kernelized Stein discrepancy. *arXiv preprint arXiv:2006.16531*.
- [18] Betancourt, M. (2017). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.

References

- [19] Betancourt, M., & Girolami, M. (2015). Hamiltonian Monte Carlo for hierarchical models. *Current trends in Bayesian methodology with applications*, 79(30), 2-4.
- [20] Sitzmann, V., Martel, J., Bergman, A., Lindell, D., & Wetzstein, G. (2020). Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33, 7462-7473.
- [21] Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., & Geiger, A. (2019). Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 4460-4470).
- [22] Sitzmann, V., Zollhöfer, M., & Wetzstein, G. (2019). Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32.
- [23] Ha, D., Dai, A. M., & Le, Q. V. HyperNetworks. In International Conference on Learning Representations.
- [24] Wu, W., Qi, Z., & Fuxin, L. (2019). Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition* (pp. 9621-9630).
- [25] [Koyuncu, B., Sanchez-Martin, P., Peis, I., Olmos, P. M., & Valera, I. \(2023\). Variational Mixture of HyperGenerators for Learning Distributions Over Functions. In Proceedings of the 40th International Conference on Machine Learning, 2023.](#)

References

- [26] Bauer, M., Dupont, E., Brock, A., Rosenbaum, D., Schwarz, J. R., & Kim, H. (2023). Spatial functa: Scaling functa to imagenet classification and generation. arXiv preprint arXiv:2302.03130.
- [27]: Peis, I., Koyuncu, B., Valera, I. & Frellsen, J. (2025). Hyper-Transforming Latent Diffusion Models. arXiv preprint arXiv:2504.16580.
- [28]: Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 10684-10695).
- [29]: Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. Score-Based Generative Modeling through Stochastic Differential Equations. In International Conference on Learning Representations.
- [30]: Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. Advances in neural information processing systems, 33, 6840-6851.
- [31]: Song, J., Meng, C., & Ermon, S. Denoising Diffusion Implicit Models. In International Conference on Learning Representations.

References

- [32]: Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In International Conference on Learning Representations.

Thank you!



ipeaz@dtu.dk

