

Hyper-Transforming Latent Diffusion Models

Ignacio Peis, Batuhan Koyuncu, Isabel Valera, Jes Frellsen

Contributions

1. Novel Hyper-Transformer Decoder (HD) - first full Transformer-based probabilistic decoder for INR parameter generation.
2. Integration into Latent Diffusion Models with support for both full training and **hyper-transforming** paradigms.
3. Latent Diffusion Models for INRs (LDMI) offers scalable framework overcoming MLP-based hypernetwork bottlenecks while increasing expressiveness INRs.

Motivation

Challenges: Existing generative frameworks rely on structured representations that constrain resolution and generalization. MLP-based hypernetworks suffer from scalability bottlenecks when generating high-dimensional INRs, limiting flexibility and expressiveness for complex data.

Solution: LDMI combines Transformer-based hypernetworks with latent diffusion models for scalable, probabilistic INR generation.

Preliminaries

Implicit Neural Representations (INRs): Neural networks representing continuous functions with parameters Φ :

$$f_{\Phi}(\mathbf{x}) = \mathbf{y}, \quad \mathbf{x} \in \mathbb{R}^d, \mathbf{y} \in \mathbb{R}^c \quad (1)$$

Hypernetworks: Networks generating parameters for other networks:

$$\Phi = g_{\phi}(\mathbf{z}) \mapsto f_{\Phi}(\mathbf{x}) = \hat{\mathbf{y}} \quad (2)$$

Latent Diffusion Models: Generative models applying diffusion in compressed latent space:

$$\mathcal{L}_{DDPM} = \mathbb{E}_{\mathbf{z}, \epsilon, t} [\|\epsilon - \epsilon_{\theta}(\mathbf{z}_t, t)\|^2] \quad (3)$$

Fast sampling via DDIM:

$$p_{\theta}(\mathbf{z}_{t-1}|\mathbf{z}_t) = \begin{cases} \mathcal{N}\left(f_{\theta}^{(1)}(\mathbf{z}_1), \sigma_1^2 \mathbf{I}\right) & \text{if } t = 1 \\ q_{\sigma}\left(\mathbf{z}_{t-1}|\mathbf{z}_t, f_{\theta}^{(t)}(\mathbf{z}_t)\right) & \text{otherwise,} \end{cases} \quad (4)$$

LDMI

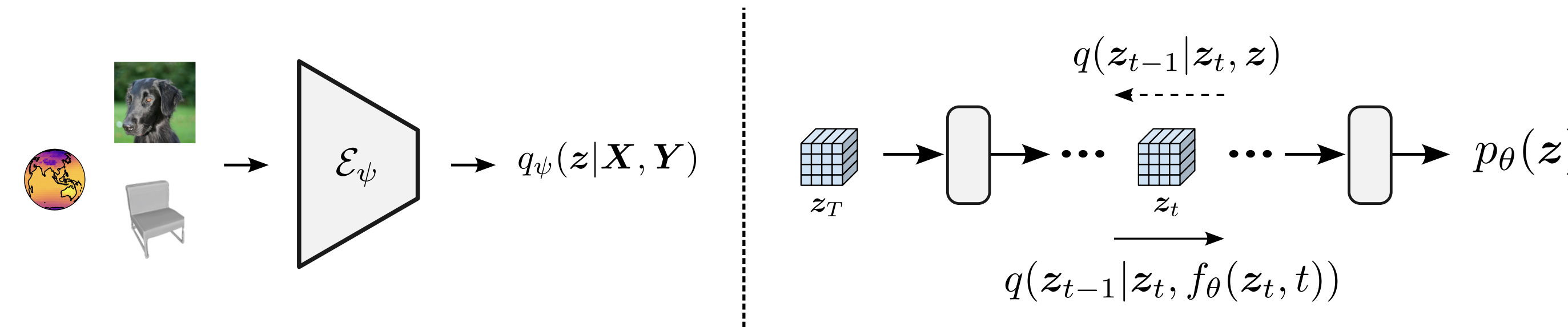


Figure 2: LDMI Encoder (left) and Latent Diffusion trained on latent space (right).

Approach: LDMI combines Transformer-based hypernetworks with latent diffusion models for scalable, probabilistic INR generation.

- **Encoder:** maps data to variational parameters \mathbf{z} .
- **Decoder:** full Transformer generates INR parameters via cross-attention and weight reconstruction.

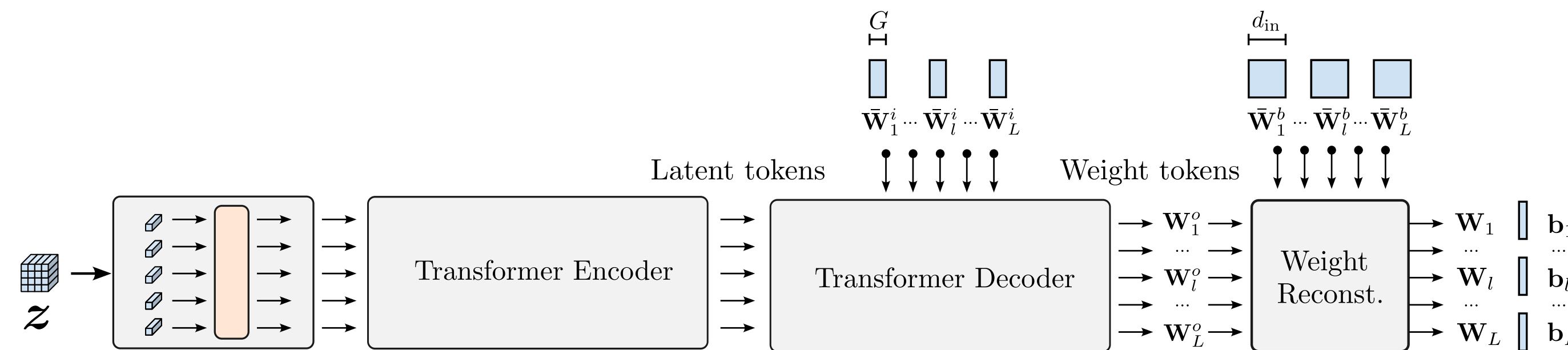


Figure 3: The HD Decoder

$$\mathcal{R}^{(\text{scale})}(\mathbf{w}_{[c/k]}^o, \bar{\mathbf{w}}_c^b) = (1 + \mathbf{w}_{[c/k]}^o) \odot \bar{\mathbf{w}}_c^b. \quad (5)$$

Training Paradigms:

1. Full Training:

First stage:

$$\mathcal{L}_{\text{VAE}}(\phi, \psi) = \mathbb{E}_{q_{\psi}(\mathbf{z}|\mathbf{X}, \mathbf{Y})} [\log p_{\Phi}(\mathbf{Y}|\mathbf{X})] - \beta \cdot D_{KL}(q_{\psi}(\mathbf{z}|\mathbf{X}, \mathbf{Y}) \| p(\mathbf{z})) \quad (6)$$

Second stage:

$$\mathcal{L}_{\text{DDPM}} = \mathbb{E}_{q(\mathbf{z}_t|\mathbf{z}_0)} [\|\epsilon - \epsilon_{\theta}(\mathbf{z}_t, t)\|^2] \quad (7)$$

2. Hyper-Transforming: Adapt pre-trained LDM by freezing $\{\psi, \theta\}$ and training decoder:

$$\mathcal{L}_{\text{HT}}(\phi) = \mathbb{E}_{q_{\psi}(\mathbf{z}|\mathbf{X}, \mathbf{Y})} [\log p_{\Phi}(\mathbf{Y}|\mathbf{X})], \quad (8)$$

Hyper-transforming enables efficient adaptation of existing diffusion models without full retraining, leveraging pre-trained latent spaces.

Results

Generation

LDMI achieves high-quality unconditional and conditional generation across multiple modalities and arbitrary resolutions.

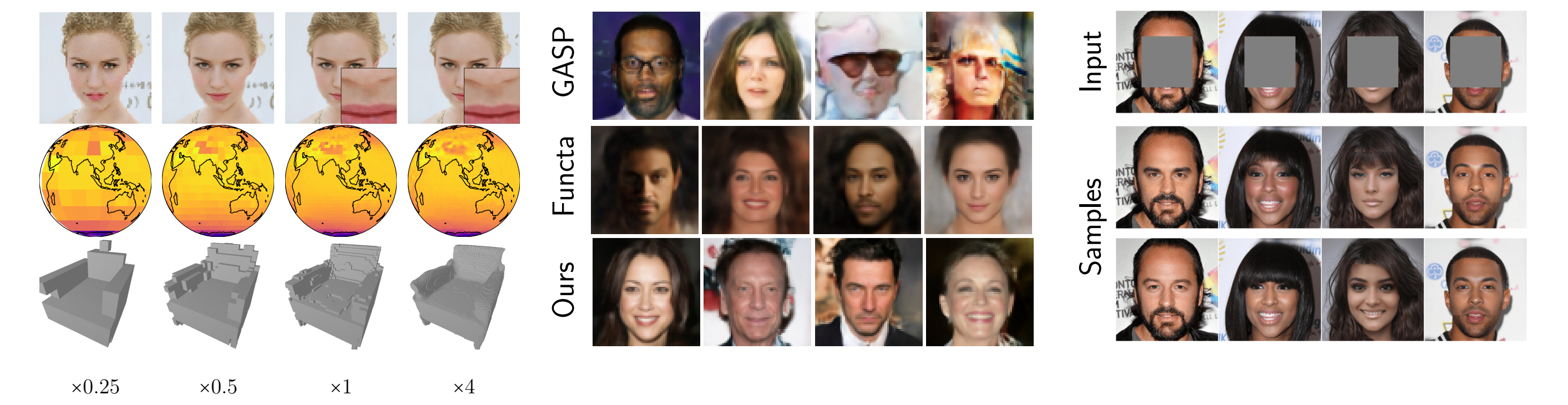


Figure 4: LDMI samples at multiple resolutions and modalities.

Figure 5: CelebA-HQ (64 x 64) samples from baselines and LDMI. Figure 6: Inpainting with LDMI on CelebA-HQ (256 x 256).

Reconstruction:

Our framework models the space of INRs to represent data at significantly higher resolutions, beyond the capabilities of existing methods.

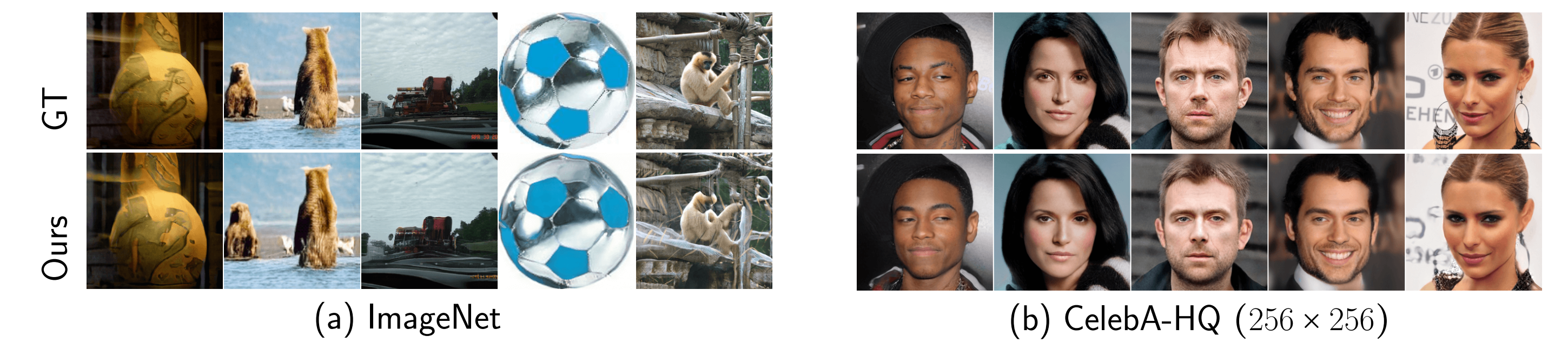


Figure 7: Reconstructions by our LDMI trained by hyper-transforming pre-trained LDMs.

Quantitative Results:

We generate INRs that are approximately $7\times$ larger in size using hyper-networks with less than $1/3$ the number of parameters.

Model	PSNR (dB) \uparrow	FID \downarrow	HN Params \downarrow	Model	Chairs (acc %) \uparrow	ERA5 (PSNR dB) \uparrow
CelebA-HQ (64 x 64)				Functa	99.51	34.9
GASP	-	7.42	25.7M	VAMoH	96.75	39.0
Functa	≤ 30.7	40.40	-	LDMI	97.25	44.6
VAMoH	23.17	66.27	25.7M			
LDMI	24.80	18.06	8.06M			

Table 3: Reconstruction quality on ShapeNet Chairs and ERA5.

Model	PSNR (dB) \uparrow	FID \downarrow	HN Params \downarrow
ImageNet (256 x 256)			
Spatial Functa	≤ 38.4	≤ 8.5	-
LDMI	20.69	6.94	102.78M

Table 1: Metrics on CelebA-HQ and ImageNet.

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

Table 2: Scalability analysis.

