

Diverse super-resolution with pretrained hierarchical VAEs

*Jean Prost*¹ Antoine Houdard¹ Andrés Almansa² Nicolas Papadakis¹

Univ. Bordeaux, Bordeaux INP, CNRS, IMB, UMR 5251, F-33400 Talence, France

Université de Paris, MAP5, CNRS, F-75006 Paris, France

December 7, 2022

Outline

- 1 Variational autoencoders
- 2 Plugging the VAE prior in a restoration model
- 3 Application: Super-resolution with VDVAE
- 4 Results

Variational autoencoders

Variational autoencoders (VAEs)¹ allows to learn latent variable models of the form:

$$p_{\theta}(\mathbf{x}) = \int \underbrace{p_{\theta}(\mathbf{x}|\mathbf{z})}_{\text{decoder}} \underbrace{p_{\theta}(\mathbf{z})}_{\text{(latent) prior}} d\mathbf{z},$$

jointly with an approximate posterior distribution (encoder):

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p_{\theta}(\mathbf{z}|\mathbf{x})$$

VAE are trained to maximize the evidence lower-bound (ELBO):

$$\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] - KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{x}|\mathbf{z})) \quad (1)$$

$$\geq \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]}_{\mathcal{L}_{\theta, \phi}(\mathbf{x})} \quad (2)$$

¹Diederik P Kingma and Max Welling. “Auto-encoding variational bayes”. In: *arXiv preprint arXiv:1312.6114* (2013).

Hierarchical VAE

Hierarchical prior

$$p_{\theta}(\mathbf{z}) = p_{\theta}(\mathbf{z}_0) \prod_{l=1}^{L-1} p_{\theta}(\mathbf{z}_l | \mathbf{z}_{<l}) \quad (3)$$

with $\mathbf{z}_l \in \mathbb{R}^{C_l \times H_l \times W_l}$, and

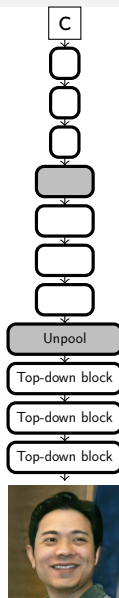
$$p_{\theta}(\mathbf{z}_l | \mathbf{z}_{<l}) \sim \mathcal{N}(\mathbf{z}_l | \mu_{\theta}(\mathbf{z}_l), \Sigma_{\theta}(\mathbf{z}_l)) \quad (4)$$

Hierarchical encoder²:

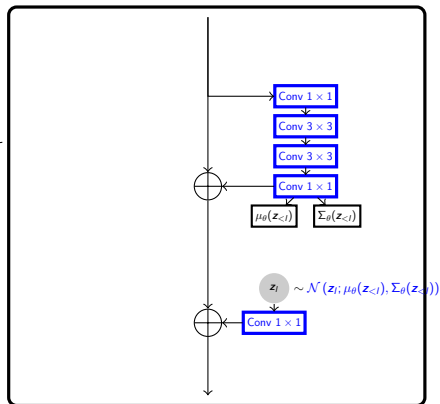
$$q_{\phi}(\mathbf{z} | \mathbf{x}) = q_{\psi}(\mathbf{z}_0 | \mathbf{x}) \prod_{l=1}^{L-1} q_{\phi}(\mathbf{z}_l | \mathbf{z}_{<l}, \mathbf{x}) \quad (5)$$

²Casper Kaae Sønderby et al. “How to train deep variational autoencoders and probabilistic ladder networks”. In: *Advances in Neural Information Processing Systems*. Vol. 29. 2016.

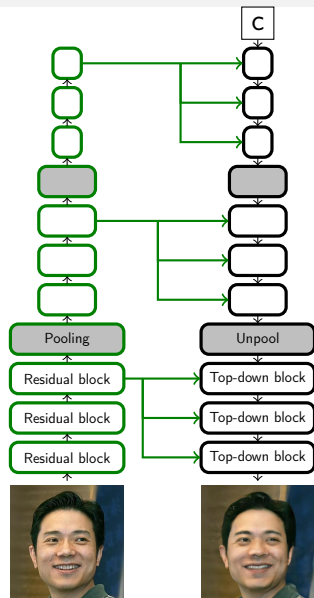
Hierarchical VAE: Architecture



$$p_{\theta}(\mathbf{z}, \mathbf{x}) = p_{\theta}(\mathbf{z}_0) \prod_{l=1}^{L-1} p_{\theta}(\mathbf{z}_l | \mathbf{z}_{<l}) p_{\theta}(\mathbf{x} | \mathbf{z})$$

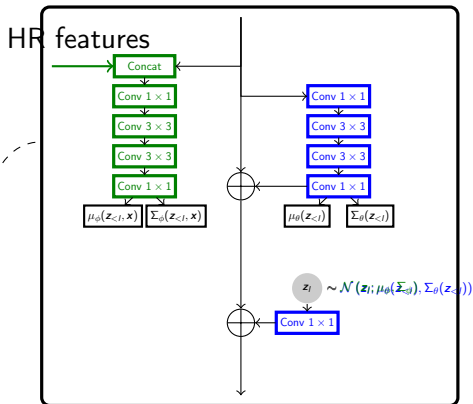


Hierarchical VAE: Architecture



$$p_{\theta}(\mathbf{z}, \mathbf{x}) = p_{\theta}(\mathbf{z}_0) \prod_{l=1}^{L-1} p_{\theta}(\mathbf{z}_l | \mathbf{z}_{<l}) p_{\theta}(\mathbf{x} | \mathbf{z})$$

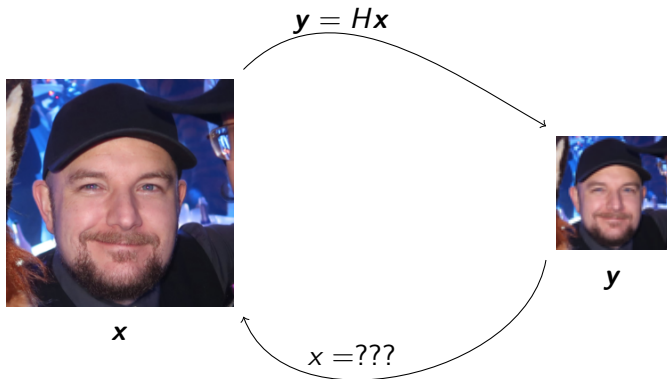
$$q_{\phi}(\mathbf{z} | \mathbf{x}) = q_{\phi}(\mathbf{z}_0 | \mathbf{x}) \prod_{l=1}^{L-1} q_{\phi}(\mathbf{z}_l | \mathbf{z}_{<l}, \mathbf{x})$$



Outline

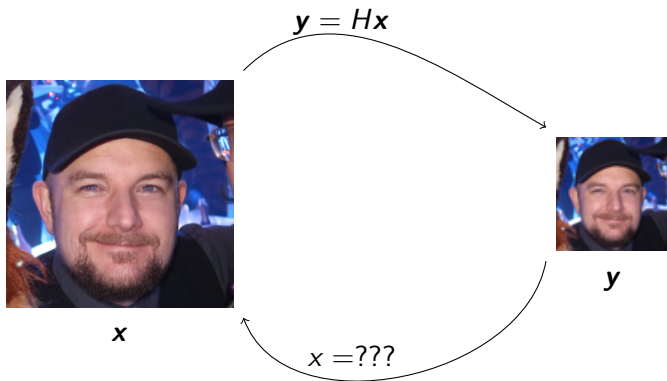
- 1 Variational autoencoders
- 2 Plugging the VAE prior in a restoration model
- 3 Application: Super-resolution with VDVAE
- 4 Results

Super-resolution



- High-resolution (HR) image (x)
- Low-resolution (LR) image (y)
- Forward operator H (low-pass + subsampling)

Super-resolution



- One-to-many problem : there can be many good solutions (consistent and realistic)
- Goal : model the distributions of solutions : $p_{SR}(x|y)$

Plugging the VAE prior in a restoration model

Target distribution:

$$p_{\theta}(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}) \underbrace{p_{\theta}(\mathbf{x})}_{\text{VAE prior}} \quad (6)$$

$$p_{\theta}(\mathbf{x}|\mathbf{y}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}, \mathbf{y}) p_{\theta}(\mathbf{z}|\mathbf{y}) d\mathbf{z} \quad (7)$$

$$= \int \underbrace{p_{\theta}(\mathbf{x}|\mathbf{z})}_{\text{VAE decoder}} p_{\theta}(\mathbf{z}|\mathbf{y}) d\mathbf{z} \quad (8)$$

→ When $p_{\theta}(\mathbf{x}|\mathbf{z}, \mathbf{y}) = p_{\theta}(\mathbf{x}|\mathbf{z})$, we only need to model $p_{\theta}(\mathbf{x}|\mathbf{y})$

Learning the low-resolution encoder

Model distribution :

$$p_{SR}(\mathbf{x}|\mathbf{y}) = \int \underbrace{p_{\theta}(\mathbf{x}|\mathbf{z})}_{\text{decoder}} \underbrace{q_{\psi}(\mathbf{z}|\mathbf{y})}_{\text{LR encoder}} d\mathbf{z} \quad (9)$$

Conditional log likelihood of the model³:

$$\log p_{SR}(\mathbf{x}|\mathbf{y}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || q_{\psi}(\mathbf{z}|\mathbf{y})) \quad (10)$$

→ Loss function:

$$\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x},\mathbf{y})} \left[KL \left(\underbrace{q_{\phi}(\mathbf{z}|\mathbf{x})}_{\text{HR encoder}} || \underbrace{q_{\psi}(\mathbf{z}|\mathbf{y})}_{\text{LR encoder}} \right) \right] \quad (11)$$

³William Harvey, Saeid Naderiparizi, and Frank Wood. “Conditional Image Generation by Conditioning Variational Auto-Encoders”. In: *International Conference on Learning Representations*. 2022.

Outline

- 1 Variational autoencoders
- 2 Plugging the VAE prior in a restoration model
- 3 Application: Super-resolution with VDVAE**
- 4 Results

For deep HVAE, the low-frequency information is only encoded at the beginning of the hierarchy, in the low-resolution latents:



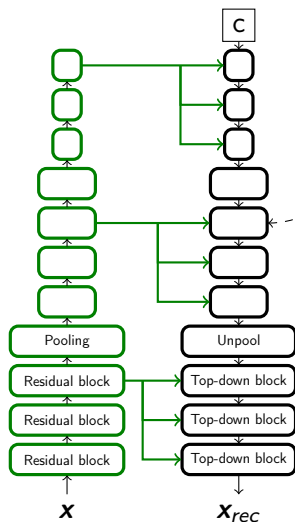
Figure: Samples from VDVAE⁴ generative model $p_{\theta}(\mathbf{x}|\mathbf{z}_{<43})$. Right: pixel-wise standard deviation.

$$\rightarrow p_{\theta}(\mathbf{x}|\mathbf{z}_{<k}, y) = p_{\theta}(\mathbf{x}|\mathbf{z}_{<k}).$$

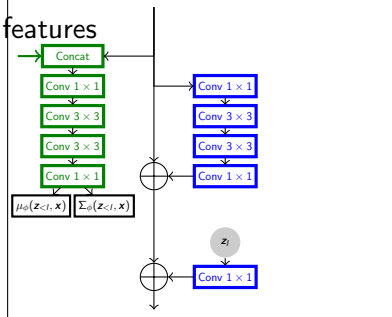
$$\rightarrow \text{We just have to learn } q_{\psi}(\mathbf{z}_{<k}|\mathbf{y})$$

⁴Rewon Child. “Very Deep VAEs Generalize Autoregressive Models and Can Outperform Them on Images”. In: *arXiv preprint arXiv:1201.10650* (2021).

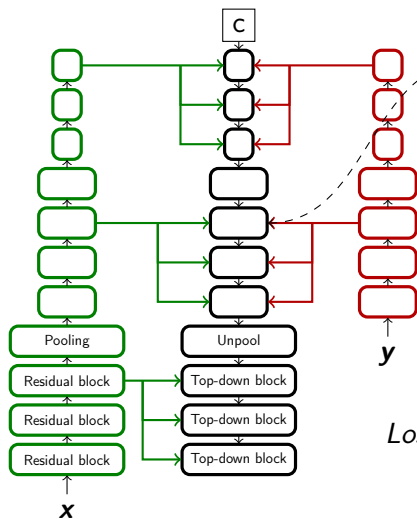
Low-resolution encoder (super-resolution)



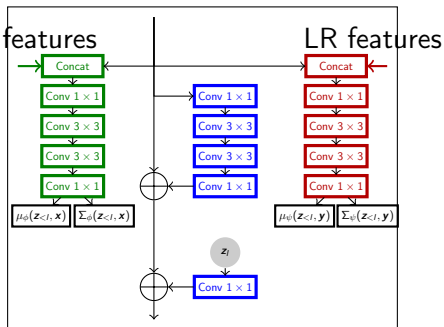
HR features



Low-resolution encoder (super-resolution)

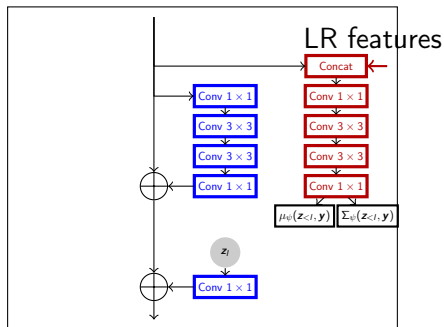
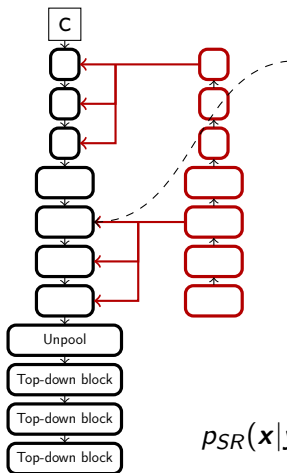


HR features



$$\text{Loss}(x, y, \psi) = KL(q_\phi(z|x) || q_\psi(z|y))$$

Low-resolution encoder (super-resolution)



$$p_{SR}(\mathbf{x}|\mathbf{y}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}_{<k}) q_{\psi}(\mathbf{z}_{<k}|\mathbf{x}) d\mathbf{z}_{<k}$$

Outline

- 1 Variational autoencoders
- 2 Plugging the VAE prior in a restoration model
- 3 Application: Super-resolution with VDVAE
- 4 Results

Experiments

- Super-resolution on FFHQ256 dataset
- $\times 4$, $\times 8$, $\times 16$ upscaling
- Comparaison with HCFlow⁵ and bicubic interpolation

⁵ [Jingyun Liang et al.](#) “Hierarchical Conditional Flow: A Unified Framework for Image Super-Resolution and Image Rescaling”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 4076–4085.

Results



LR

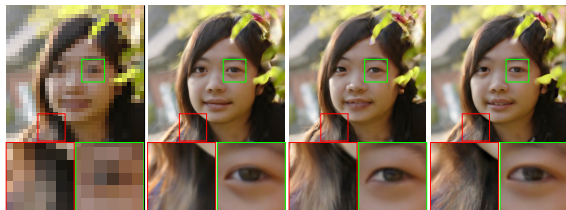
CVD-VAE ($\tau = 0.8$)



Bicubic

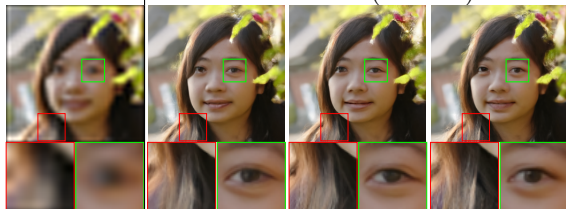
HCFflow ($\tau = 0.8$)

Figure: $\times 4$



LR

CVD-VAE ($\tau = 0.8$)



Bicubic

HCFlow ($\tau = 0.8$)

Figure: $\times 8$

Results

		Visual Quality	Consistency	Diversity	
		BRISQUE↓	LR-PSNR ↑	APD (MSE) ($\times 10^4$) ↑	APD (LPIPS) ($\times 10^3$) ↑
×4	Bicubic	61.79	36.99	0	0
	HCFlow ($\tau = 0.1$)	48.21	52.66	1.3	1.1
	HCFlow ($\tau = 0.8$)	37.21	52.81	161.8	62.6
	CVD-VAE ($\tau = 0.1$)	<u>36.47</u>	75.70	64.6	<u>104.5</u>
	CVD-VAE ($\tau = 0.8$)	32.30	<u>75.20</u>	<u>88.8</u>	123.0
×8	Bicubic	78.42	33.61		
	HCFlow ($\tau = 0.1$)	69.05	51.35	4.9	4.2
	HCFlow ($\tau = 0.8$)	<u>36.25</u>	51.13	575.5	155.3
	CVD-VAE ($\tau = 0.1$)	50.34	71.63	140.4	<u>179.0</u>
	CVD-VAE ($\tau = 0.8$)	32.26	<u>70.15</u>	<u>248.2</u>	236.4
×16	Bicubic	97.53	30.67		
	HCFlow ($\tau = 0.1$)	74.28	51.93	3.3	14.6
	HCFlow ($\tau = 0.8$)	30.83	53.81	<u>323.7</u>	<u>268.6</u>
	CVD-VAE ($\tau = 0.1$)	52.75	60.68	299.4	242.4
	CVD-VAE ($\tau = 0.8$)	<u>32.01</u>	<u>57.02</u>	613.8	341.2

APD = Average Pairwise Distance

Effect of the temperature

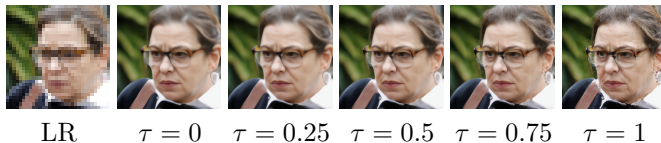


Figure: Effect of the sampling temperature τ



Figure: Reducing the temperature can help to reduce artifacts for difficult images.

More samples



Figure: $\times 4$

More samples



Figure: $\times 8$

Conclusion

- Hierarchical VAE can model complex image priors that can be used to solve inverse problems.
- For future works:
 - Apply on different datasets
 - Flexible method (for general inverse problems)

Expected consistency of the super-resolution

Expected consistency of $p_{SR}(\mathbf{x}|\mathbf{y})$

$$CE(k) = \mathbb{E}_{p_{\mathcal{D}}(\mathbf{y})} \mathbb{E}_{p_{SR}(\mathbf{x}|\mathbf{y})} \left[\frac{1}{\sqrt{m}} \|H\mathbf{x} - \mathbf{y}\|_2 \right]. \quad (12)$$

Proposition

If :

- 1 The low-resolution encoder has enough capacity and is trained to optimality
- 2 The VAE encoder $q_{\phi}(\mathbf{x}|\mathbf{z})$ and generative model $p_{\theta}(\mathbf{x}, \mathbf{z})$ have enough capacity and are trained to optimality

Then:

$$CE(k) = \mathbb{E}_{p_{\theta}(\mathbf{z}_{<k})} \mathbb{E}_{p_{\theta}(\mathbf{x}|\mathbf{z}_{<k})} \mathbb{E}_{p_{\theta}(\tilde{\mathbf{x}}|\mathbf{z}_{<k})} \left[\frac{1}{\sqrt{m}} \|H\tilde{\mathbf{x}} - H\mathbf{x}\|_2 \right] := U_k^s. \quad (13)$$

Expected consistency of the super-resolution model

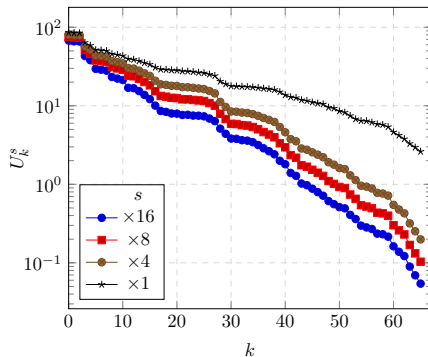


Figure: Average low-resolution pairwise distance, U_k^s between samples from the conditional generative model $p_\theta(\mathbf{x}|\mathbf{z}_{<k})$ of VD-VAE, for downscaling factors $s = 1, 4, 8, 16$. Image with pixel values in $[0, 255]$.

Enforcing the consistency with the input

Projection to space of consistent solution $\{\mathbf{x} | H\mathbf{x} = \mathbf{y}\}$:

$$\hat{\mathbf{x}} \leftarrow (I - H^T(HH^T)^{-1}H)\mathbf{x} + H^T(HH^T)^{-1}\mathbf{y} \quad (14)$$

We use the CEM implementation of⁶

⁶Yuval Bahat and Tomer Michaeli. “Explorable super resolution”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020, pp. 2716–2725.

Effect of the projection

	model	Distortion			Visual Quality	Consistency	Diversity	
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	BRISQUE \downarrow	LR-PSNR \uparrow	APD (MSE) ($\times 10^4$) \uparrow	APD (LPIPS) ($\times 10^3$) \uparrow
$\times 4$	CVD-VAE ($\tau = 0.1$)	29.09	0.839	0.177	34.56	35.31	77.5	<u>123.6</u>
	CVD-VAE ($\tau = 0.8$)	28.66	0.825	0.177	30.37	35.09	105.7	142.9
	CVD-VAE ($\tau = 0.1$) (cons)	30.75	0.863	0.152	36.47	75.70	64.6	104.5
	CVD-VAE ($\tau = 0.8$) (cons)	<u>30.24</u>	<u>0.850</u>	<u>0.157</u>	<u>32.30</u>	<u>75.20</u>	<u>88.8</u>	123.0
$\times 8$	CVD-VAE ($\tau = 0.1$)	25.46	0.735	0.312	47.48	33.93	181.6	200.7
	CVD-VAE ($\tau = 0.8$)	24.71	0.697	<u>0.284</u>	30.40	33.42	308.0	256.2
	CVD-VAE ($\tau = 0.1$) (cons)	26.27	0.747	0.299	50.34	71.63	140.4	179.0
	CVD-VAE ($\tau = 0.8$) (cons)	<u>25.47</u>	<u>0.708</u>	0.275	<u>32.26</u>	<u>70.15</u>	<u>248.2</u>	<u>236.4</u>
$\times 16$	CVD-VAE ($\tau = 0.1$)	20.16	<u>0.601</u>	0.440	47.34	25.20	711.1	283.1
	CVD-VAE ($\tau = 0.8$)	19.21	0.541	<u>0.401</u>	29.17	24.22	1278.8	374.0
	CVD-VAE ($\tau = 0.1$) (cons)	22.06	0.622	0.427	52.75	60.68	299.4	242.4
	CVD-VAE ($\tau = 0.8$) (cons)	<u>21.14</u>	0.564	0.388	<u>32.01</u>	<u>57.02</u>	<u>613.8</u>	<u>341.2</u>

Figure: Effect of the projection

Effect of the projection

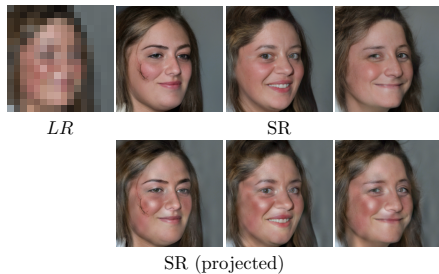
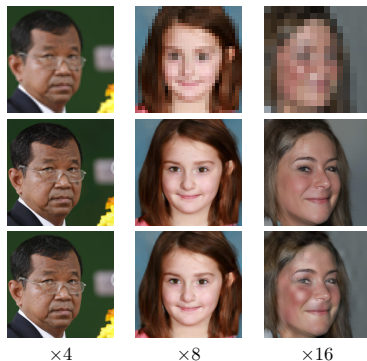


Figure: $\times 16$

From top row to bottom row : LR, SR,
SR projected

More samples ($\times 16$)



Figure: $\times 16$