

Conditional- t^3 VAE: Equitable Latent Space Allocation for Fair Generation

Aymene Mohammed Bouayed^{*1,3}, Samuel Deslauriers-Gauthier²
Adrian Iaccovelli³, David Naccache¹

¹DIÉNS, ÉNS, CNRS, PSL University, Paris, France

²Centre Inria d'Université Côte d'Azur, Nice, France

³Be-Ys Research, France

Abstract

Variational Autoencoders (VAEs) with global priors mirror the training set's class frequency in latent space, underrepresenting tail classes and reducing generative fairness on imbalanced datasets. While t^3 VAE improves robustness via heavy-tailed Student's t -distribution priors, it still allocates latent volume proportionally to the class frequency. In this work, we address this issue by explicitly enforcing equitable latent space allocation across classes. To this end, we propose Conditional- t^3 VAE, which defines a per-class Student's t joint prior over latent and output variables, preventing dominance by majority classes. Our model is optimized using a closed-form objective derived from the γ -power divergence. Moreover, for class-balanced generation, we derive an equal-weight latent mixture of Student's t -distributions. On SVHN-LT, CIFAR100-LT, and CelebA, Conditional- t^3 VAE consistently achieves lower FID scores than both t^3 VAE and Gaussian-based VAE baselines, particularly under severe class imbalance. In per-class F1 evaluations, Conditional- t^3 VAE also outperforms the conditional Gaussian VAE across all highly imbalanced settings. While Gaussian-based models remain competitive under mild imbalance ratio ($\rho \lesssim 3$), our approach substantially improves generative fairness and diversity in more extreme regimes.

1 Introduction

Class imbalance and long-tail distributions are prevalent in real-world datasets, yet generative models often fail to represent rare classes accurately. When trained on skewed data, these models tend to overfit dominant modes and underrepresent minority ones in latent and output spaces, resulting in biased or unfair generations. This issue is especially critical in sensitive applications such as facial synthesis [1] and medical imaging [2], where such biases can exacerbate social and diagnostic disparities [3].

Variational Autoencoders (VAEs) [4] are a widely used class of generative models, valued for their probabilistic formulation, stable training, and compatibility with latent-variable modeling frameworks leading in image quality such as Latent Diffusion Models (LDMs) [5]. Although GANs and diffusion models often achieve lower FID scores, VAEs offer unique advantages in class-conditional generation, interpretability, and efficient inference, making them strong candidates for improving fairness under class imbalance. Standard VAEs commonly use isotropic Gaussian priors, which inadequately model heavy-tailed structures and rare phenomena [6]. Prior efforts to address this have introduced non-Gaussian priors, particularly Student's t -distributions [7, 8, 9, 10], to enhance robustness. However, these approaches often rely on global priors, causing the latent space to be dominated by majority classes under skewed distributions.

We address this issue with Conditional- t^3 VAE (C- t^3 VAE), a conditional generative model that imposes a per-class Student's t -distribution prior over the joint latent-output space. This design ensures allocating an equal latent space volume per class, thereby mitigating majority class dominance,

^{*}Corresponding author : aymene.bouayed@ens.fr

while the heavy tails of the Student’s t-distribution more effectively capture intra-class variation. To enable class-balanced sampling, we introduce an equal-weight mixture of Student’s t-distributions with analytically derived component variances. Together, these components enable balanced class-conditional generation and mitigate bias present in unconditional models. We summarize our main theoretical and empirical contributions in the following points :

- We propose the C- t^3 VAE model with a training objective based on the γ -power divergence.
- We develop an equal-weight latent mixture sampling scheme with analytically derived optimal variance scaling for each component.
- We outperform relevant baselines in FID on SVHN-LT [11], CIFAR100-LT [12], and CelebA [13] under severe imbalance, and show via per-class evaluation that C- t^3 VAE better avoids mode collapse, exceeding a conditional VAE in per class Recall and F1 while remaining competitive on Precision.
- We identify the imbalance ratio threshold $\rho \approx 3$, beyond which Gaussian priors become sub-optimal, providing guidance for model selection on skewed datasets.

2 Related Work

Since the introduction of VAEs [4], many extensions have sought to improve latent representation by replacing the standard Gaussian prior with more expressive alternatives. These include Gaussian mixtures [14, 15], hyperspherical priors [16], normalizing flows [17], Riemannian priors [18], and implicit distributions [19]. While most retain the Evidence Lower Bound (ELBO) formulation, others adopt alternative objectives or divergence measures for added flexibility.

To address long-tailed or imbalanced data, Student’s t-distributions have been explored for their robustness and heavy tails [6]. Methods such as [7] and [8] model the latent space of the autoencoder through a t-distributed prior, and rely on KL-divergence-based ELBO objectives. However, since the KL divergence lacks a closed-form solution for the Student’s t-distributions, these methods resort to numerical approximation. The t^3 VAE [10] improves on this by modeling the joint latent-output distribution and optimizing a closed-form γ -divergence objective [9]. Nonetheless, it still employs a global latent prior, which results in a latent space volume allocation reflecting class frequency and leading to imbalance present in the generated samples.

Other works address fairness in generative modeling using normalizing flows [17] and diffusion models [20, 21]. However, since leading image generation models are based on Latent Diffusion Models [5] and VAEs are a corner-stone in the design of these models, we deem it crucial to improve VAEs’ ability to handle imbalanced data. Consequently, as fairness in latent space allocation remains under-explored in VAEs, we introduce a class-conditional, heavy-tailed prior to address this gap and allow for balanced latent space allocation across classes.

3 Background

This section introduces the theoretical background and baseline models relevant to our work. We assume access to a labeled, imbalanced dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, where $x_i \in \mathbb{R}^n$ is a data sample of dimension n , $y_i \in \{1, \dots, K\}$ its class label and m the latent space dimension.

3.1 VAEs and Conditional-VAEs

VAEs [4] are generative models trained via variational inference by maximizing the ELBO of the log-likelihood. The standard objective of this model is

$$\mathcal{L}_{\theta, \phi} := \mathbb{E}_{z \sim q_{\phi}(\cdot|x)}[\log p_{\theta}(x|z)] - \mathcal{D}_{KL}(q_{\phi}(z|x) \| p(z)), \quad (1)$$

where the first term is the reconstruction loss with $p_\theta(x|z)$ being the decoder. The second term is the Kullback–Leibler (KL) divergence between the approximate posterior $q_\phi(z|x)$ and the prior $p(z)$. The β -VAE is a weighted variant of the VAE model which introduces a β scaling term for the KL divergence [22]:

$$\mathcal{L}_{\theta,\phi} := \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\theta(x|z)] - \beta \mathcal{D}_{\text{KL}}(q_\phi(z|x) \| p(z)), \quad (2)$$

with $p(z) \sim \mathcal{N}_m(0, I)$, $q_\phi(\cdot|x) \sim \mathcal{N}_m(\mu_\phi(x), \Sigma_\phi(x))$, and $p_\theta(x|z) \sim \mathcal{N}_m(\mu_\theta(z), \sigma^2 I)$. $\mu_\phi(\cdot)$ and $\Sigma_\phi(\cdot)$ are the mean and the covariance matrices inferred through a neural network with parameter ϕ given the input x . This variant of the VAE model allows to place more weight on disentangling the latent space or on the reconstruction of the data points. To generate samples from the VAE or the β -VAE model, we sample a latent vector $z \sim \mathcal{N}_m(0, I)$. Then, the generated data point would be $\hat{x} \sim \mathcal{N}_m(\mu_\theta(z), \sigma^2 I)$.

Nevertheless, since Eq. (1) and Eq. (2) optimize the ELBO over the data distribution $p_{\text{data}}(x)$, which can be decomposed as $p_{\text{data}}(x) = \sum_{y_i} p(y_i) p_{\text{data}}(x | y_i)$, this optimization inherently biases the model toward head classes with larger $p(y_i)$. As a result, most generated samples come from overrepresented classes, while tail classes' samples are underrepresented and of lower quality, a phenomenon commonly referred to as *mode collapse*. Therefore, when labels are available, it is preferable to define class-conditional approximate posterior and prior distributions: $q_\phi(z|x, y)$ and $p(z|y)$. This yields the Conditional-VAE (CVAE) model trained using the objective [23]:

$$\sum_y \mathbb{E}_{z \sim q_\phi(\cdot|x, y)} [\log p_\theta(x|z, y)] - \beta \mathcal{D}_{\text{KL}}(q_\phi(z|x, y) \| p(z|y)). \quad (3)$$

Here, in Eq. (3) we constrain all $p(y_i)$ to be equal and omit them from the loss function in order not to exacerbate the issue of class imbalance in the latent space. Also, we define $p(z|y) \sim \mathcal{N}_m(\mu_y, I)$ with learnable class-wise means μ_y . To generate a data point \hat{x}_y from class y , we sample $z_y \sim \mathcal{N}_m(\mu_y, I)$, then we get $\hat{x}_y \sim p_\theta(x|z_y, y)$. Nevertheless, despite conditioning, this formulation remains Gaussian. Unlike Student's t-distributions, Gaussian priors poorly approximate heavy-tailed data distributions [6].

3.2 Multivariate Student's t-Distribution and γ -Power Divergence

A d -dimensional Student's t-distribution with mean $\mu \in \mathbb{R}^d$, covariance $\Sigma \in \mathbb{R}^{d \times d}$, and degrees of freedom $\nu > 2$ is a heavy-tail, super-Gaussian distribution defined as

$$t_d(x) = C_{\nu,d} |\Sigma|^{-\frac{1}{2}} \left(1 + \frac{(x - \mu)^\top \Sigma^{-1} (x - \mu)}{\nu} \right)^{-\frac{\nu+d}{2}}, \quad C_{\nu,d} = \frac{\Gamma\left(\frac{\nu+d}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) (\nu\pi)^{\frac{d}{2}}}. \quad (4)$$

The power form of this distribution prevents a closed-form KL divergence between two Student's t-distributions. Instead, the γ -power divergence $\mathcal{D}_\gamma(q||p)$ is used [9, 10]. This divergence is defined for $q \sim t_d(\mu_0; \Sigma_0; \nu)$, $p \sim t_d(\mu_1; \Sigma_1; \nu)$ starting from the γ -entropy $\mathcal{H}_\gamma(p)$ and cross-entropy $\mathcal{C}_\gamma(q, p)$

$$\mathcal{H}_\gamma(p) := -\|p\|_{1+\gamma} = - \left(\int p(x)^{1+\gamma} dx \right)^{\frac{1}{1+\gamma}}, \quad \mathcal{C}_\gamma(q, p) := - \int q(x) \left(\frac{p(x)}{\|p\|_{1+\gamma}} \right)^\gamma dx.$$

$$\mathcal{D}_\gamma(q||p) := \gamma^{-1} (\mathcal{C}_\gamma(q, p) - \mathcal{H}_\gamma(p)) \quad (5)$$

with $\gamma = -\frac{2}{\nu+d}$. Then, substituting the definition of a Student's t-distribution from Eq. (4) into Eq. (5), the following closed-form formula for the γ -power divergence can be derived (See Appendix A):

$$\begin{aligned} \mathcal{D}_\gamma(q||p) = & -\frac{C_{\nu,d}^{\frac{\gamma}{1+\gamma}}}{\gamma} \left(1 + \frac{d}{\nu-2} \right)^{-\frac{\gamma}{1+\gamma}} \left[-|\Sigma_0|^{-\frac{\gamma}{2(1+\gamma)}} \left(1 + \frac{d}{\nu-2} \right) \right. \\ & \left. + |\Sigma_1|^{-\frac{\gamma}{2}} |\Sigma_0|^{\frac{\gamma^2}{2(1+\gamma)}} \left(1 + \frac{\text{Tr}(\Sigma_1^{-1} \Sigma_0)}{\nu-2} + \frac{(\mu_0 - \mu_1)^\top \Sigma_1^{-1} (\mu_0 - \mu_1)}{\nu} \right) \right]. \end{aligned} \quad (6)$$

3.3 t^3 -Variational Autoencoder

3.3.1 Definition

The t^3 VAE model [10] is a non-ELBO-based autoencoder which models the joint prior distribution $p_\theta(x, z)$ using multivariate Student's t-distributions

$$p_\theta(x, z) \propto \sigma^{-n} \left[1 + \frac{1}{\nu} \left(\|z\|^2 + \frac{\|x - \mu_\theta(z)\|^2}{\sigma^2} \right) \right]^{-\frac{\nu+m+n}{2}},$$

where $\mu_\theta(\cdot)$ is the decoder neural network with parameter θ whereas σ is a parameter controlling the decoder's output covariance. From this joint distribution, the marginal latent prior $p(z)$ and decoder distribution $p_\theta(x|z)$ can be defined. Furthermore, the approximate posterior distribution is defined as :

$$q_\phi(z|x) = t_m \left(x \middle| \mu_\phi(x), \frac{\Sigma_\phi(x)}{1 + \nu^{-1}n}, \nu + n \right).$$

Hence, the data distribution would be $q_\phi(x, z) = p_{\text{data}}(x)q_\phi(z|x)$. As a result, relying on the γ -divergence in Eq. (6) applied to the $p_\theta(x, z)$ and $q_\phi(x, z)$ distributions, the following loss function is derived to optimize the t^3 VAE's parameters :

$$\mathcal{L}_\gamma = \mathbb{E}_x \left[\frac{\mathbb{E}_z [\|x - \mu_\theta(z)\|^2]}{\sigma^2} + \|\mu_\phi(x)\|^2 + \frac{\nu \text{Tr}(\Sigma_\phi(x))}{\nu + n - 2} - \frac{\nu C_1}{C_2} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right], \quad (7)$$

with $\gamma = -\frac{2}{\nu+n+m}$ and C_1 and C_2 being theoretically derived constants. We note that the first term in this loss function represents the standard reconstruction term in VAE models and the rest of the terms are regularization terms over the latent space. To sample from the latent space of the t^3 VAE, [10] propose the $p_\nu^*(z) = t_m(0, \tau^2 I, \nu + n)$ distribution with

$$\tau^2 = \frac{1}{1 + \nu^{-1}n} \left(\frac{C_{\nu,n}}{\sigma^n} \cdot \frac{\nu - 2}{\nu + n - 2} \right)^{\frac{2}{\nu+n-2}}. \quad (8)$$

Moreover, sampling from a multi-variate Student's t-distribution $T \sim t_d(\mu, \Sigma, \nu)$ both in the learning (Eq. (7)) and sampling (Eq. (8)) phases is performed through the standard reparameteration trick for Student's t-distributions $T := \mu + Z\sqrt{\nu V^{-1}}$ where $Z \sim \mathcal{N}(0, \Sigma)$ and $V \sim \mathcal{X}^2(\nu)$.

3.3.2 β - t^3 VAE

From Eq. (7) we can also define a β - t^3 VAE model by multiplying all the regularization terms by a β factor. Similarly to β -VAE models, this improves the versatility of the model and allows either a focus on generation or disentangling.

3.3.3 τ^2 Improvement

Closely analyzing the proposed derivation of τ^2 of the t^3 VAE model, we discovered a subtle issue in its mathematical formulation. The employed γ -power divergence presents a discrepancy to the correct formula in Eq. (6). We revised the formulation and the corrected τ^2 is (See Appendix D):

$$\tau^2 = \frac{1}{1 + \nu^{-1}n} \left(|\Sigma_\phi(x)|^{\frac{\gamma}{2}} \frac{C_{\nu,n}}{\sigma^n} \cdot \frac{\nu - 2}{\nu + n - 2} \right)^{\frac{2\gamma}{(1+\gamma)(2+\gamma m)}} \approx \frac{1}{1 + \nu^{-1}n} \left(\frac{C_{\nu,n}}{\sigma^n} \cdot \frac{\nu - 2}{\nu + n - 2} \right)^{-\frac{2\gamma}{2+\gamma m}}. \quad (9)$$

The corrected exact form of τ^2 is applicable when $|\Sigma_\phi(x)|$ is known and when the dimension of the data is low. However, for high dimensional data, as handled in this work, we get $\gamma \approx 0$. Hence, one can use the approximation without any loss in accuracy. We note that the new form of τ^2 leads to a similar empirical value of standard deviation compared to the previous form. Nevertheless, for

correctness, in our sampling from the latent space of the t^3 VAE model we use the approximation in Eq. (9).

In summary, although the t^3 VAE effectively models heavy-tailed distributions through Student's t-distributions and γ -power divergence, it does not explicitly address class imbalance in the latent space by not allocating equal volume for each class. In the next section, we introduce a class-conditional variant of the t^3 VAE, designed to ensure fair and balanced generation across all classes.

4 Conditional t^3 -Variational Autoencoder

We propose the Conditional t^3 -Variational Autoencoder (C- t^3 VAE), present its formulation, training objective, and sampling strategy. C- t^3 VAE models the latent space as a mixture of Student's t-distributions, one per class, ensuring equal latent volume allocation and promoting fairness in generation. Intra-class imbalance is further handled by the heavy-tailed nature of the Student's t-distribution prior.

4.1 Model definition

The C- t^3 VAE we propose is based on the following class conditional joint prior distribution

$$p_\theta(x, z|y) = \frac{C_{\nu, m+n}}{|\Sigma_x|^{\frac{1}{2}} |\Sigma_y|^{\frac{1}{2}}} \left[1 + \frac{(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) + (x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z))}{\nu} \right]^{-\frac{\nu+m+n}{2}},$$

with ν , n and m being the degrees of freedom of the Student's t-distribution, the dimension of the input data and the dimension of the latent space respectively. $\mu_y \in \mathbb{R}^m$ is a learnable mean vector representing class centers in latent space of dimension m . Moreover, Σ_x and Σ_y are the covariance matrices of the prior distributions over the latent and output variables.

From this joint distribution, we can derive the conditional latent prior $p(z|y) = t_m(z|\mu_y, \Sigma_y, \nu)$ and decoder distribution (See Appendix B)

$$p_\theta(x|z, y) = t_n\left(x \middle| \mu_\theta(z), \frac{(1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y))}{(1 + \nu^{-1}m)} \Sigma_x, \nu + m\right).$$

Furthermore, as in t^3 VAE, we define the approximate posterior $q_\phi(z|x)$ as a multivariate Student's t-distribution capturing heavy-tailed structure in the latent space :

$$q_\phi(z|x) = t_m\left(z \middle| \mu_\phi(x), \frac{\Sigma_\phi(x)}{1 + \nu^{-1}n}, \nu + n\right).$$

4.2 Objective function

Harnessing Eq. (5) and the defined prior and posterior distributions of the proposed C- t^3 VAE, we derive in Appendix C the following class-wise objective

$$\begin{aligned} \mathcal{L}(\gamma, y) = \mathbb{E}_x \left[\mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] + (\mu_\phi(x) - \mu_y)^\top \Sigma_y^{-1} (\mu_\phi(x) - \mu_y) \right. \\ \left. + \frac{\nu \text{Tr}(\Sigma_y^{-1} \Sigma_\phi(x))}{\nu + n - 2} - \frac{\nu C_1}{C_2} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right], \end{aligned}$$

with $C_1 = \left(C_{\nu+n, m}^\gamma \left(1 + \frac{n}{\nu}\right)^{\frac{\gamma m}{2}} \frac{\nu+n+m-2}{\nu+n-2} \right)^{\frac{1}{1+\gamma}}$ and $C_2 = \left(\frac{C_{\nu, m+n}^\gamma}{|\Sigma_x|^{\frac{\gamma}{2}} |\Sigma_y|^{\frac{2\gamma+1}{2}}} \left(1 + \frac{m+n}{\nu-2}\right)^{-\gamma} \right)^{\frac{1}{1+\gamma}}$. By taking $\Sigma_x = \sigma^2 I$ and $\Sigma_y = I$, $\mathcal{L}(\gamma, y)$ objective function simplifies to :

$$\mathcal{L}(\gamma, y) = \mathbb{E}_x \left[\frac{\mathbb{E}_z [\|x - \mu_\theta(z)\|^2]}{\sigma^2} + \|\mu_\phi(x) - \mu_y\|^2 + \frac{\nu \text{Tr}(\Sigma_\phi(x))}{\nu + n - 2} - \frac{\nu C_1}{C_2} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right]. \quad (10)$$

Therefore, we express the final loss function $\mathcal{L}(\gamma)$ over the whole dataset as : $\mathcal{L}(\gamma) = \sum_y \mathcal{L}(\gamma, y)$. As in Eq. (3), here too we consider all $p(y_i)$ to be equal and hence rule out their contribution to the loss function. This is done to avoid emphasizing the imbalance present in the data.

4.3 Sampling distribution

Similarly to the objective function of t^3 VAE, $\mathcal{L}(\gamma, y)$ in Eq. (10) can be decomposed into a reconstruction term and regularization terms. To sample from the latent space of the C- t^3 VAE, we focus on the regularization terms and define the following sampling distribution :

$$p_v^*(z) = \sum_{y=1}^K \alpha_y \cdot t_m(\mu_y, \tau^2 I, v + n), \quad \forall y, \quad \alpha_y = \frac{1}{K}. \quad (11)$$

The theoretical derivation of the variance τ^2 leads to the form expressed in Eq. (9) (Derivation in Appendix D).

The mixture-based sampling distribution we define in Eq. (11) with equal α_y ensures a uniformly sampled synthetic data across all classes, regardless of their frequency in the original training data. As a result, C- t^3 VAE equipped with this sampling distribution mitigates the common problem in generative models where head-class samples dominate due to their density in latent space. Furthermore, by modifying the mixture weights α_y , one can prioritize specific classes. This makes our method flexible for targeted data augmentation or fairness-aware sampling strategies.

4.3.1 β -C- t^3 VAE

As with t^3 VAE, the class-wise objective defined in Eq. (10) can be split into a reconstruction and regularization terms. By preceding the regularization term with a β scalar, we can define a β -C- t^3 VAE model thereby improving the domain of applicability of the model.

Overall, C- t^3 VAE provides a principled, flexible, and tractable framework for fair generative modeling, particularly under class-imbalanced conditions. In the following, we study the C- t^3 VAE’s performance across multiple datasets with varying imbalance degrees.

5 Experiments

In this section we outline the generative performance of the proposed C- t^3 VAE model on labeled datasets compared to relevant VAE baselines¹. We conduct experiments on three datasets notably SVHN-LT [11], CIFAR100-LT [24, 12] and CelebA [13] each chosen to highlight different challenges.

5.1 Evaluation procedure

All models are evaluated using Fréchet Inception Distance (FID) [25], computed against a *balanced* test set for each dataset. This setup measures how effectively a model overcomes training set imbalance by assessing its ability to generate high-quality samples across all classes. To evaluate robustness, we impose varying degrees of imbalance during training.

For SVHN-LT and CIFAR100-LT, we introduce class imbalance by applying an exponential decay to the number of samples per class after equalizing class sizes in the original train and test sets. The imbalance ratio ρ defines the ratio between the most and least frequent classes, with class-wise sample counts M_{y_i} given by: $M_{y_i} = M \cdot \rho^{-\frac{y_i-1}{K-1}}$, where M is the original per-class sample count.

For CelebA, we compute FID per attribute, treating each attribute (eg. Mustache) and its negation (eg. no Mustache) as separate binary classes. The training set uses CelebA’s inherent imbalance, while the test set is balanced by downsampling to the smaller class size. Multi-attribute generation is not considered in this work and is left for future exploration.

¹We restrict our comparisons to VAE-based generative models, as Latent Diffusion Models [5], the state-of-the-art in image generation, depend on VAEs as a core component. Enhancing fairness in VAEs under class imbalance is therefore a necessary precursor to advancing more complex models.

5.2 Results

We present both quantitative and qualitative results of the $C\text{-}t^3\text{VAE}$ model and the models it directly improves upon notably the VAE, $C\text{-VAE}$ and $t^3\text{VAE}$ models with their β variants. This controlled comparison helps isolate the contributions of key design choices:

- **VAE** : ELBO trained standard Gaussian-based VAEs.
- **C-VAE** : VAE supplemented by conditional Gaussian priors to assess the class conditioning effect without changing the prior family.
- **$t^3\text{VAE}$** : Student’s t-distribution latent prior and γ -power divergence objective, does not use class-conditional priors, nor allows for class conditional generation. Through a comparison with this model we assess the role of conditional modeling on Student’s t-distribution priors.

In the following, we analyze the latent sampling standard deviation τ employed in $t^3\text{VAE}$ and $C\text{-}t^3\text{VAE}$ models, varying it to assess alignment between empirical and theoretical values. Then, we present FID comparisons across baselines with optimized hyper-parameters (The tuning of β , ν and τ hyper-parameters is reported in Appendix F), followed by per-class generative evaluation.

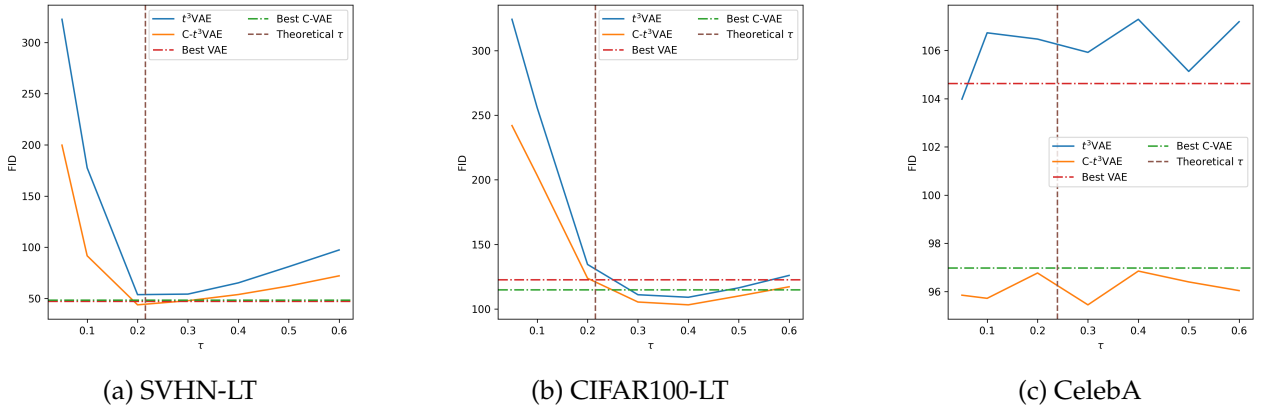


Figure 1: FID score as a function of τ for the $t^3\text{VAE}$ and $C\text{-}t^3\text{VAE}$ models. Results are for the imbalance ratio $\rho = 100$ for the SVHN-LT and CIFAR100-LT, and for the Mustache attribute ($\rho = 25$) in the case of the CelebA dataset. Other imbalance ratios’ results paint a similar picture and are provided in Appendix F.3. The horizontal dashed lines is the FID value of the best performing VAE and C-VAE on each dataset and the vertical dashed line is the value of τ as derived in Eq. (9). We note that the used models in these figures have optimized β and ν hyper-parameters.

5.2.1 τ parameter study

Figure 1 illustrates the sampling standard deviation’s τ impact on the FID metric. We observe that models based on the Student’s t-distribution benefit from higher standard deviation on CIFAR100-LT compared to SVHN-LT. Specifically, $C\text{-}t^3\text{VAE}$ outperforms the $C\text{-VAE}$ for $\tau \in [0.25; 0.55]$ on CIFAR100-LT and $\tau \in [0.19; 0.28]$ on SVHN-LT and for all τ values on the CelebA dataset. Moreover, it surpasses the $t^3\text{VAE}$ models’ FID for all τ values and across all datasets.

For both Student’s t-distribution based models, the optimal FID score for the SVHN-LT occurs near the theoretically derived τ value. However, for the more complex CIFAR100-LT dataset, the optimal τ is higher than the theoretical value $\tau = 0.4$. We hypothesize that the higher variance reflects the model’s need during training to accommodate the greater complexity and entropy of the dataset, which makes it harder to align the learned approximate posterior with the prior and to produce confident reconstructions. For the CelebA dataset, the standard deviation parameter τ has minimal to no impact on model performance likely due to the lower variability in the dataset’s images as they all represent faces.

5.2.2 Optimized Model Results Discussion

After optimizing the hyperparameters of the various models tested in this work, we present their generation FID scores in Table 1. We provide results for β models after optimization and non β models ($\beta = 1$) to underscore the importance of this parameter on the resulting performance as it was not explored in the t^3 VAE work [10].

Table 1: Generation FID results on the SVHN-LT, CIFAR100-LT and CelebA datasets. For the SVHN-LT and CIFAR100-LT datasets we use different imbalance ratios $\rho \in \{100, 50, 10, 1\}$. However, for the CelebA dataset we use the Mustache, Young, Male and Smiling attributes which have imbalance ratios of 25, 3.5, 1.4 and 1 respectively. The β models undertook an optimization of the β hyper-parameter while non- β models have $\beta = 1$. All models have optimized ν and τ hyper-parameters. The attributes for the CelebA dataset column indicate which attribute is used to condition the conditional models and balance the test set.

Models	SVHN-LT				CIFAR100-LT				CelebA			
	$\rho = 100$	50	10	1	$\rho = 100$	50	10	1	Mustache	Young	Male	Smiling
VAE	93.89	91.91	91.66	92.16	163.66	162.91	165.47	166.46	110.58	92.01	110.58	82.05
β -VAE	47.11	49.81	45.70	43.48	122.62	123.07	123.72	124.43	104.63	92.87	87.96	83.15
C-VAE	74.75	70.40	72.30	74.16	157.90	163.67	162.09	163.24	96.98	89.17	86.17	78.35
β -C-VAE	48.39	46.39	43.97	43.87	114.88	118.89	114.89	118.21	98.35	85.53	79.76	78.46
t^3 VAE	57.07	54.30	52.10	51.52	136.63	137.24	138.92	135.23	105.80	88.07	83.62	78.90
β - t^3 VAE	51.62	49.55	48.93	45.37	109.11	107.93	108.97	111.00	105.86	88.21	83.83	78.89
C- t^3 VAE	47.09	46.29	47.43	51.32	125.48	127.96	130.28	129.40	101.18	87.07	81.92	80.97
β -C- t^3 VAE	44.02	42.60	42.01	44.49	103.25	102.99	105.92	112.37	95.82	82.61	81.65	80.08

From Table 1, t^3 VAE consistently improves over the VAE model. This underlines the generative advantage of the Student’s t-distribution prior over a Gaussian one, in addition to the reconstruction advantage highlighted by [10]. Moreover, optimizing the β hyper-parameter improves t^3 VAE’s FID over the VAE on the CIFAR100-LT and CelebA, while remaining competitive on SVHN-LT. Also, qualitative results of Figure 2 show that the t^3 VAE model compared to the VAE model produces significantly sharper synthetic images on the CelebA dataset.



Figure 2: Sample synthetic images from the optimized VAE and t^3 VAE models trained on the CelebA dataset. No class conditioning is possible for these models.

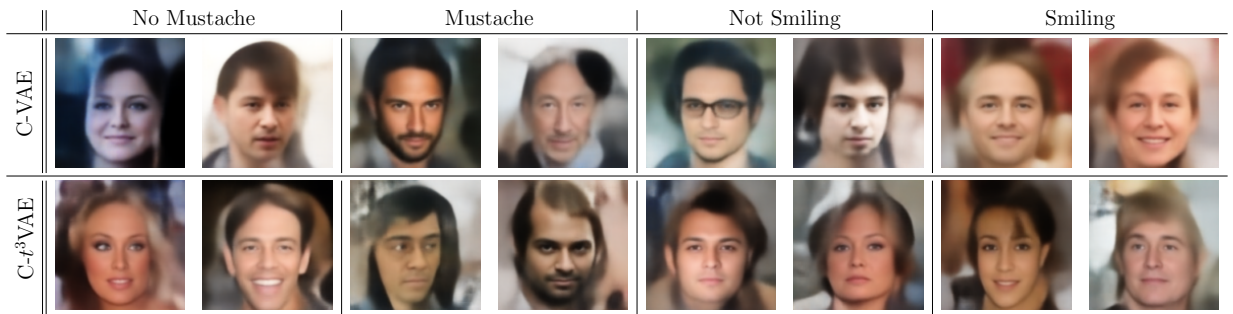


Figure 3: Sample synthetic images for the optimized C-VAE and C- t^3 VAE trained on specific attributes of the CelebA dataset.

For class-conditional models, β optimized C- t^3 VAE yields strong FID improvements across all imbalance settings, surpassing all baselines. As per Table 1, it achieves a gain of up to 4, 5 and 10 FID points over the β - t^3 VAE model on the imbalanced settings of the SHVN-LT, CIFAR100-LT and CelebA, respectively. This underlines the importance of an equal per-class volume in an imbalanced dataset setting. Moreover, β -C- t^3 VAE reduces the FID by up to 4 and 15 points over C-VAE on

SVHN-LT and CIFAR100-LT respectively as per Table 1. For the CelebA dataset, Table 1 shows that β -C- t^3 VAE achieves the best results on heavily imbalanced attributes like Mustache, demonstrating better generation for underrepresented groups. This gain in performance compared to the C-VAE goes back to the use of Student’s t-distribution latent prior and its ability to better capture intra-class long-tail distributions. Additionally, qualitative samples of conditional optimized models (Figure 3) reveal sharper facial features in C- t^3 VAE compared to the C-VAE. These results establish it as the most reliable model across imbalance ratios and resolutions.

In summary, C- t^3 VAE after optimization of all hyper-parameters notably β , ν and τ consistently achieves the lowest FID, validating the use of class-conditional Student’s t-distribution priors and balanced latent sampling. Improvements on CIFAR100-LT are especially pronounced, highlighting the model’s robustness to severe imbalance. However, for balanced settings, C- t^3 VAE remains competitive.

5.2.3 Per-class evaluation

In this section, we evaluate the conditional models on a per-class basis. However, since FID is biased on small datasets and offers limited insight as a single scalar metric, we rely on Precision, Recall, and F1 generative metrics² [26]. Our results on CelebA are shown in Figure 4, with additional results for SVHN-LT and CIFAR100-LT included in Appendix G.

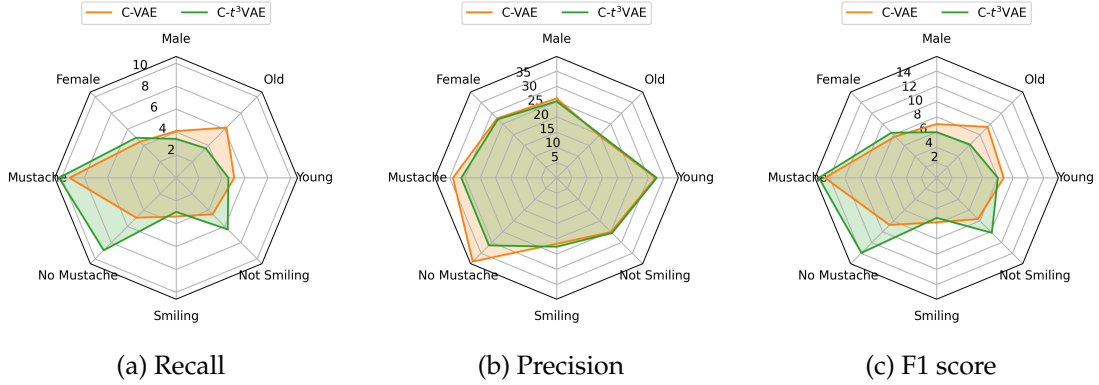


Figure 4: Per-class generative metrics on CelebA after optimization of all hyper-parameters notably β , ν and τ . We note that the imbalance ratio of the Mustache, Young, Male and Smiling factors ρ are 25, 3.5, 1.4 and 1 respectively.

On the CelebA dataset (Figure 4), C- t^3 VAE improves Recall and F1 on the most imbalanced attribute (Mustache), but not on more balanced ones (Male, Smiling) which is an expected behavior. However, surprisingly, for the attribute Young (imbalance ratio 3.5), C-VAE performs slightly better this suggests that there is a regime where Gaussian priors may still suffice. To explore this, we vary the imbalance ratio on SVHN-LT from 100 to 1 and plot the results in Figure 5. This figure shows there is a threshold $\rho \approx 3$ before which C-VAE performs better; beyond that, C- t^3 VAE has the advantage. Also, we notice that the more the imbalance ratio increases, the more FID gap becomes more significant and in favor of the C- t^3 VAE.

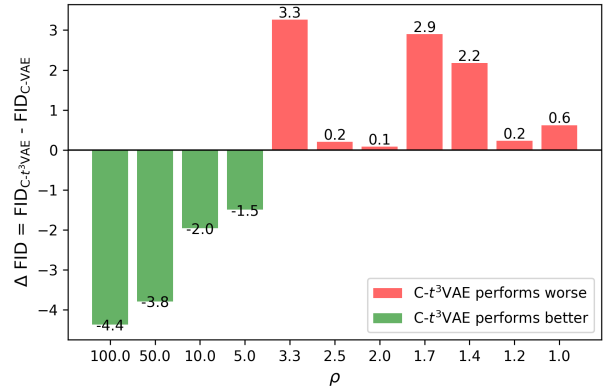


Figure 5: Fine-grained comparison of C- t^3 VAE and C-VAE models under varying imbalance ratios on SVHN-LT.

²Precision quantifies sample quality (sharpness), Recall reflects mode coverage and F1 is their harmonic mean. Together they allow us to gain deeper insight into the models’ behavior.

For the SVHN-LT dataset (Figures in Appendix G), when comparing the $C-t^3$ VAE to the C-VAE, the $C-t^3$ VAE achieves significantly higher Recall and competitive Precision. This indicates a better coverage of data modes and maintaining images quality which yields higher F1 scores especially on tail classes. For the CIFAR100-LT dataset (Figures in Appendix G), C-VAE often achieves high Precision but zero Recall, indicating mode collapse. In contrast, the $C-t^3$ VAE preserves Recall at the cost of slightly lower Precision, also resulting in improved F1 scores. Consequently, on the SVHN-LT CIFAR100-LT and CelebA datasets, $C-t^3$ VAE sets itself as the most reliable method for fair, high quality image generation.

6 Conclusion

We introduced $C-t^3$ VAE, a class-conditional generative model that leverages Student’s t-distributions in the latent space with a theoretically derived balanced sampling scheme. This approach improves fairness and sample quality under class imbalance. Extensive experiments on SVHN-LT, CIFAR100-LT, and CelebA demonstrate that, after the optimization of the β , ν and τ hyper-parameters, $C-t^3$ VAE consistently outperforms the t^3 VAE and conditional VAE baselines, achieving up to a 15-point FID improvement on highly imbalanced datasets. Per-class Precision, Recall, and F1 evaluations confirm improved mode coverage, particularly in tail classes. Moreover, we identify $\rho > 3$ as the regime where $C-t^3$ VAE outperforms Gaussian-prior models, with the performance gap widening as ρ increases. For $\rho \lesssim 3$, $C-t^3$ VAE remains competitive with C-VAE.

Beyond generative performance, $C-t^3$ VAE contributes to fairness-aware modeling by promoting balanced latent sampling and better mode coverage in underrepresented classes. Future work will explore extending $C-t^3$ VAE to multi-label problems and integration with Latent Diffusion models.

Acknowledgement

The authors were granted access to the HPC resources of MesoPSL financed by the Région Île-de-France and the Equip@Meso project (reference ANR-10-EQPX-29-01) of the *programme investissements d’avenir* supervised by France’s Agence nationale pour la recherche.

References

- [1] D. Mehta, A. Mehta, and P. Narang, “Ldfacenet: Latent diffusion-based network for high-fidelity deepfake generation,” in *International Conference on Pattern Recognition*, pp. 386–400, Springer, 2024.
- [2] W. H. Pinaya, P.-D. Tudosiu, J. Dafflon, P. F. Da Costa, V. Fernandez, P. Nachev, S. Ourselin, and M. J. Cardoso, “Brain imaging generation with latent diffusion models,” in *MICCAI Workshop on Deep Generative Models*, pp. 117–126, Springer, 2022.
- [3] R. Naik and B. Nushi, “Social biases through the text-to-image generation lens,” in *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, p. 786–808, 2023.
- [4] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” 2013.
- [5] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, 2022.
- [6] E. Tam and D. B. Dunson, “On the statistical capacity of deep generative models,” *arXiv preprint arXiv:2501.07763*, 2025.

- [7] H. Takahashi, T. Iwata, Y. Yamanaka, M. Yamada, and S. Yagi, “Student-t variational autoencoder for robust density estimation,” in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pp. 2696–2702, International Joint Conferences on Artificial Intelligence Organization, 7 2018.
- [8] N. Abiri and M. Ohlsson, “Variational auto-encoders with student’s t-prior,” *arXiv preprint arXiv:2004.02581*, 2020.
- [9] S. Eguchi, “Pythagoras theorem in information geometry and applications to generalized linear models,” in *Information Geometry* (A. Plastino, A. S. Srinivasa Rao, and C. Rao, eds.), vol. 45 of *Handbook of Statistics*, pp. 15–42, Elsevier, 2021.
- [10] J. Kim, J. Kwon, M. Cho, H. Lee, and J.-H. Won, “ t^3 -variational autoencoder: Learning heavy-tailed data with student’s t and power divergence,” in *The Twelfth International Conference on Learning Representations*, 2024.
- [11] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with unsupervised feature learning,” in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
- [12] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, “Learning imbalanced datasets with label-distribution-aware margin loss,” in *Advances in Neural Information Processing Systems*, 2019.
- [13] Z. Liu, P. Luo, X. Wang, and X. Tang, “Deep learning face attributes in the wild,” in *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [14] A. Saseendran, K. Skubch, S. Falkner, and M. Keuper, “Shape your space: A gaussian mixture regularization approach to deterministic autoencoders,” in *Advances in Neural Information Processing Systems* (M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, eds.), vol. 34, pp. 7319–7332, Curran Associates, Inc., 2021.
- [15] N. Dilokthanakul, P. A. Mediano, M. Garnelo, M. C. Lee, H. Salimbeni, K. Arulkumaran, and M. Shanahan, “Deep unsupervised clustering with gaussian mixture variational autoencoders,” *arXiv preprint arXiv:1611.02648*, 2016.
- [16] T. R. Davidson, L. Falorsi, N. De Cao, T. Kipf, and J. M. Tomczak, “Hyperspherical variational auto-encoders,” *34th Conference on Uncertainty in Artificial Intelligence (UAI-18)*, 2018.
- [17] P. Jaini, I. Kobyzev, Y. Yu, and M. Brubaker, “Tails of lipschitz triangular flows,” in *International Conference on Machine Learning*, pp. 4673–4681, PMLR, 2020.
- [18] C. Chadebec, E. Thibeau-Sutre, N. Burgos, and S. Allasonnière, “Data augmentation in high dimensional low sample size setting using a geometry-based variational autoencoder,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 2879–2896, 2023.
- [19] H. Takahashi, T. Iwata, Y. Yamanaka, M. Yamada, and S. Yagi, “Variational autoencoder with implicit optimal priors,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 5066–5073, 2019.
- [20] Y. Qin, H. Zheng, J. Yao, M. Zhou, and Y. Zhang, “Class-balancing diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18434–18443, 2023.
- [21] K. Pandey, J. Pathak, Y. Xu, S. Mandt, M. Pritchard, A. Vahdat, and M. Mardani, “Heavy-tailed diffusion models,” in *The Thirteenth International Conference on Learning Representations*, 2025.
- [22] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, “beta-VAE: Learning basic visual concepts with a constrained variational framework,” in *International Conference on Learning Representations*, 2017.

- [23] D. P. Kingma, S. Mohamed, D. Jimenez Rezende, and M. Welling, "Semi-supervised learning with deep generative models," *Advances in neural information processing systems*, vol. 27, 2014.
- [24] A. Krizhevsky, "Learning multiple layers of features from tiny images," tech. rep., 2009.
- [25] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in neural information processing systems*, vol. 30, 2017.
- [26] T. Kynkäänniemi, T. Karras, S. Laine, J. Lehtinen, and T. Aila, "Improved precision and recall metric for assessing generative models," *Advances in neural information processing systems*, vol. 32, 2019.

Appendix for "Conditional- t^3 VAE: Equitable Latent Space Allocation for Fair Generation"

A γ -power divergence corrected derivation

In this section, we present our correction of the previously derived form of the γ -power divergence in [10]. We start from the integrals derived in [10] and keep a similar notation :

$$\begin{aligned}\int q(x)p(x)^\gamma dx &= C_{\nu,d}^\gamma |\Sigma_1|^{-\frac{\gamma}{2}} \left(1 + \frac{1}{\nu-2} \text{Tr}(\Sigma_1^{-1}\Sigma_0) + \frac{1}{\nu}(\mu_0 - \mu_1)^\top \Sigma_1^{-1}(\mu_0 - \mu_1) \right) \\ \int q(x)^{1+\gamma} dx &= C_{\nu,d}^\gamma |\Sigma_0|^{-\frac{\gamma}{2}} \left(1 + \frac{d}{\nu-2} \right) \\ \int p(x)^{1+\gamma} dx &= C_{\nu,d}^\gamma |\Sigma_1|^{-\frac{\gamma}{2}} \left(1 + \frac{d}{\nu-2} \right)\end{aligned}$$

Combining these formulas gives a similar entropy \mathcal{H}_γ to the one derived in [10]:

$$\mathcal{H}_\gamma(q) = -C_{\nu,d}^{\frac{\gamma}{1+\gamma}} |\Sigma_0|^{-\frac{\gamma}{2(1+\gamma)}} \left(1 + \frac{d}{\nu-2} \right)^{\frac{1}{1+\gamma}}.$$

However, the cross-entropy \mathcal{C}_γ takes the following form :

$$\mathcal{C}_\gamma(q, p) = -C_{\nu,d}^{\frac{\gamma}{1+\gamma}} |\Sigma_0|^{\frac{\gamma^2}{2(1+\gamma)}} |\Sigma_1|^{-\frac{\gamma}{2}} \left(1 + \frac{1}{\nu-2} \text{Tr}(\Sigma_1^{-1}\Sigma_0) + \frac{1}{\nu}(\mu_0 - \mu_1)^\top \Sigma_1^{-1}(\mu_0 - \mu_1) \right),$$

In red, we highlight the main difference to the formula derived in [10]. $\mathcal{H}_\gamma(q)$ and $\mathcal{C}_\gamma(q, p)$ combine to give :

$$\begin{aligned}\mathcal{D}_\gamma(q||p) &= -\frac{C_{\nu,d}^{\frac{\gamma}{1+\gamma}}}{\gamma} \left(1 + \frac{d}{\nu-2} \right)^{-\frac{\gamma}{1+\gamma}} \left[-|\Sigma_0|^{-\frac{\gamma}{2(1+\gamma)}} \left(1 + \frac{d}{\nu-2} \right) \right. \\ &\quad \left. + |\Sigma_1|^{-\frac{\gamma}{2}} |\Sigma_0|^{\frac{\gamma^2}{2(1+\gamma)}} \left(1 + \frac{\text{Tr}(\Sigma_1^{-1}\Sigma_0)}{\nu-2} + \frac{(\mu_0 - \mu_1)^\top \Sigma_1^{-1}(\mu_0 - \mu_1)}{\nu} \right) \right].\end{aligned}$$

B Priors derivations

In this section, we present our derivations of the different prior distributions defining our proposed C- t^3 -VAE model. Starting from the proposed joint distribution :

$$p_\theta(x, z|y) = \frac{C_{\nu, m+n}}{|\Sigma_x|^{\frac{1}{2}} |\Sigma_y|^{\frac{1}{2}}} \left[1 + \frac{(z - \mu_y)^\top \Sigma_y^{-1}(z - \mu_y) + (x - \mu_\theta(z))^\top \Sigma_x^{-1}(x - \mu_\theta(z))}{\nu} \right]^{-\frac{\nu+m+n}{2}}.$$

To calculate the prior distribution on the latent space we marginalize out x as follows :

$$\begin{aligned}
p(z|y) &= \int p_\theta(x, z|y) dx \\
&= \int C_{\nu, m+n} |\Sigma_x|^{-\frac{1}{2}} |\Sigma_y|^{-\frac{1}{2}} \left[1 + \frac{(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)}{\nu} \right. \\
&\quad \left. + \frac{(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z))}{\nu} \right]^{-\frac{\nu+m+n}{2}} dx \\
&= C_{\nu, m+n} |\Sigma_x|^{-\frac{1}{2}} |\Sigma_y|^{-\frac{1}{2}} \left[1 + \frac{1}{\nu} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right]^{-\frac{\nu+m+n}{2}} \\
&\quad \times \int \left(1 + \frac{(1 + \nu^{-1}m)(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z))}{(1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)) (\nu + m)} \right)^{-\frac{\nu+m+n}{2}} dx.
\end{aligned}$$

Given that :

$$\begin{aligned}
&\int C_{\nu+m, n} |\Sigma|^{-\frac{1}{2}} \left(1 + \frac{(x - \mu)^\top \Sigma^{-1} (x - \mu)}{\nu + m} \right)^{-\frac{\nu+m+n}{2}} dx = 1 \\
&\Rightarrow \int \left(1 + \frac{(x - \mu)^\top \Sigma^{-1} (x - \mu)}{\nu + m} \right)^{-\frac{\nu+m+n}{2}} dx = C_{\nu+m, n}^{-1} |\Sigma|^{\frac{1}{2}},
\end{aligned}$$

and when setting :

$$\Sigma^{-1} = \frac{(1 + \nu^{-1}m) \Sigma_x^{-1}}{1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)},$$

We get :

$$\begin{aligned}
&\int \left(1 + \frac{(1 + \nu^{-1}m)(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z))}{(1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)) (\nu + m)} \right)^{-\frac{\nu+m+n}{2}} dx \\
&= C_{\nu+m, n}^{-1} \left| \left(\frac{(1 + \nu^{-1}m) \Sigma_x^{-1}}{1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)} \right)^{-1} \right|^{\frac{1}{2}} \\
&= C_{\nu+m, n}^{-1} \left| \frac{1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)}{(1 + \nu^{-1}m)} \Sigma_x \right|^{\frac{1}{2}} \\
&= C_{\nu+m, n}^{-1} \left(\frac{1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)}{1 + \nu^{-1}m} \right)^{\frac{n}{2}} |\Sigma_x|^{\frac{1}{2}}.
\end{aligned}$$

Therefore, $p(z|y)$ simplifies to :

$$\begin{aligned}
p(z|y) &= C_{\nu, m+n} |\Sigma_x|^{-\frac{1}{2}} |\Sigma_y|^{-\frac{1}{2}} \left[1 + \frac{1}{\nu} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right]^{-\frac{\nu+m+n}{2}} C_{\nu+m, n}^{-1} |\Sigma_x|^{\frac{1}{2}} \\
&\quad \times \left(\frac{1 + \nu^{-1}(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)}{1 + \nu^{-1}m} \right)^{\frac{n}{2}} \\
&= C_{\nu, m+n} C_{\nu+m, n}^{-1} \left(1 + \frac{m}{\nu} \right)^{-\frac{n}{2}} |\Sigma_y|^{-\frac{1}{2}} \left[1 + \frac{1}{\nu} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right]^{-\frac{\nu+m}{2}} \\
&= C_{\nu, m} |\Sigma_y|^{-\frac{1}{2}} \left[1 + \frac{1}{\nu} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right]^{-\frac{\nu+m}{2}} \\
&= t_m(z | \mu_y, \Sigma_y, \nu).
\end{aligned}$$

Here and in the following, we use the fact

$$C_{\nu, m+n} = C_{\nu+m, n} C_{\nu, m} \left(1 + \frac{m}{\nu}\right)^{\frac{n}{2}}.$$

Besides, the prior distribution over the output of the decoder model $p(x|z, y)$ can be derived as follows :

$$\begin{aligned} p_{\theta}(x|z, y) &= \frac{p_{\theta}(x, z|y)}{p(z|y)} \\ &= \frac{C_{\nu, m+n}}{|\Sigma_x|^{\frac{1}{2}} |\Sigma_y|^{\frac{1}{2}}} \left[1 + \frac{(z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y) + (x - \mu_{\theta}(z))^{\top} \Sigma_x^{-1} (x - \mu_{\theta}(z))}{\nu} \right]^{-\frac{\nu+m+n}{2}} \\ &\quad \times C_{\nu, m}^{-1} |\Sigma_y|^{\frac{1}{2}} \left[1 + \frac{1}{\nu} (z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y) \right]^{\frac{\nu+m}{2}} \\ &= C_{\nu+m, n} |\Sigma_x|^{-\frac{1}{2}} \left(1 + \frac{m}{\nu}\right)^{\frac{n}{2}} \left[1 + \frac{1}{\nu} (z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y) \right]^{-\frac{n}{2}} \\ &\quad \times \left(1 + \frac{(1 + \nu^{-1}m)(x - \mu_{\theta}(z))^{\top} \Sigma_x^{-1} (x - \mu_{\theta}(z))}{\left(1 + \nu^{-1}(z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y)\right) (\nu + m)} \right)^{-\frac{\nu+m+n}{2}} \\ &= t_n \left(x \middle| \mu_{\theta}(z), \frac{\left(1 + \nu^{-1}(z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y)\right)}{(1 + \nu^{-1}m)} \Sigma_x, \nu + m \right). \end{aligned}$$

C Loss function derivation

In this section, we derive the loss function of C- t^3 -VAE. We start by calculating the different double integrals $\iint p_{\theta}(x, z|y)^{1+\gamma} dx dz$, $\iint q_{\phi}(x, z|y) p_{\theta}(x, z|y)^{\gamma} dx dz$, and $\iint q_{\phi}(x, z|y)^{1+\gamma} dx dz$.

Firstly,

$$\begin{aligned} \iint p_{\theta}(x, z|y)^{1+\gamma} dx dz &= \mathbb{E}_{z \sim p(z|y)} \mathbb{E}_{x \sim p_{\theta}(x|z, y)} [p_{\theta}(x, z|y)^{\gamma}] \\ &= C_{\nu, m+n}^{\gamma} |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_z \mathbb{E}_x \left[1 + \frac{(z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y)}{\nu} \right. \\ &\quad \left. + \frac{(x - \mu_{\theta}(z))^{\top} \Sigma_x^{-1} (x - \mu_{\theta}(z))}{\nu} \right] \\ &= C_{\nu, m+n}^{\gamma} |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_z \left[1 + \nu^{-1} (z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y) \right. \\ &\quad \left. + \nu^{-1} \mathbb{E}_x \left[\text{Tr}(\Sigma_x^{-1} (x - \mu_{\theta}(z)) (x - \mu_{\theta}(z))^{\top}) \right] \right] \\ &= C_{\nu, m+n}^{\gamma} |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_z \left[1 + \nu^{-1} (z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y) \right. \\ &\quad \left. + \nu^{-1} \text{Tr} \left(\Sigma_x^{-1} \Sigma_x \frac{\nu + m}{\nu + m - 2} \frac{\left(1 + \nu^{-1}(z - \mu_y)^{\top} \Sigma_y^{-1} (z - \mu_y)\right)}{(1 + \nu^{-1}m)} \right) \right] \end{aligned}$$

Here, we use the following identities

$$(k - p)^{\top} H^{-1} (k - p) = \text{Tr} \left(H^{-1} (k - p) (k - p)^{\top} \right); \quad \mathbb{E}[\text{Tr}(\cdot)] = \text{Tr}(\mathbb{E}[\cdot])$$

and the covariance of a multivariate Student's t distribution $p \sim t(\mu; \Sigma; \nu)$ is $\frac{\nu}{\nu-2} \Sigma$. Consequently, and after a few simplifications we get

$$\begin{aligned}
\iint p_\theta(x, z|y)^{1+\gamma} dx dz &= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_z \left[1 + \nu^{-1} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right. \\
&\quad \left. + \frac{n}{\nu + m - 2} \left(1 + \nu^{-1} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right) \right] \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_z \left[\left(1 + \frac{n}{\nu + m - 2} \right) \right. \\
&\quad \left. \times \left(1 + \nu^{-1} (z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right) \right] \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \left(1 + \frac{n}{\nu + m - 2} \right) \\
&\quad \times \left(1 + \nu^{-1} \mathbb{E}_z \left[(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right] \right) \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \left(1 + \frac{n}{\nu + m - 2} \right) \\
&\quad \times \left(1 + \nu^{-1} \mathbb{E}_z \left[\text{Tr}(\Sigma_y^{-1} (z - \mu_y)(z - \mu_y)^\top) \right] \right) \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \left(1 + \frac{n}{\nu + m - 2} \right) \left(1 + \frac{m}{\nu - 2} \right) \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \left(1 + \frac{m+n}{\nu - 2} \right).
\end{aligned}$$

Secondly,

$$\begin{aligned}
\iint q_\phi(x, z|y) p_\theta(x, z|y)^\gamma dx dz &= \mathbb{E}_{x \sim p_{\text{data}}} \mathbb{E}_{z \sim q(z|x)} [p_\theta(x, z|y)^\gamma] \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_x \mathbb{E}_z \left[1 + \frac{(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y)}{\nu} \right. \\
&\quad \left. + \frac{(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z))}{\nu} \right] \\
&= C_{\nu, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_x \left[1 + \frac{1}{\nu} \mathbb{E}_z \left[(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right] \right. \\
&\quad \left. + \frac{1}{\nu} \mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] \right].
\end{aligned}$$

Simplifying $\mathbb{E}_z \left[(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right]$:

$$\begin{aligned}
\mathbb{E}_z \left[(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right] &= \mathbb{E}_z \left[\text{Tr} \left(\Sigma_y^{-1} (z - \mu_y)(z - \mu_y)^\top \right) \right] \\
&= \mathbb{E}_z \left[\text{Tr} \left(\Sigma_y^{-1} (z - \mu(x) + \mu(x) - \mu_y)(z - \mu(x) + \mu(x) - \mu_y)^\top \right) \right] \\
&= \mathbb{E}_z [\text{Tr}(\Sigma_y^{-1} ((z - \mu(x))(z - \mu(x))^\top + (z - \mu(x))(\mu(x) - \mu_y)^\top \\
&\quad + (\mu(x) - \mu_y)(z - \mu(x))^\top + (\mu(x) - \mu_y)(\mu(x) - \mu_y)^\top))] \\
&= \frac{\nu}{\nu + n - 2} \text{Tr} \left(\Sigma_y^{-1} \Sigma_\phi(x) \right) + (\mu(x) - \mu_y)^\top \Sigma_y^{-1} (\mu(x) - \mu_y).
\end{aligned}$$

Then, $\iint q(x, z|y) p_\theta(x, z|y)^\gamma dx dz$ simplifies to :

$$\begin{aligned} \iint q_\phi(x, z|y) p_\theta(x, z|y)^\gamma dx dz &= C_{v, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \mathbb{E}_x \left[1 + \frac{1}{v} \mathbb{E}_z \left[(z - \mu_y)^\top \Sigma_y^{-1} (z - \mu_y) \right] \right. \\ &\quad \left. + \frac{1}{v} \mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] \right] \\ \iint q_\phi(x, z|y) p_\theta(x, z|y)^\gamma dx dz &= C_{v, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{1}{2}} \mathbb{E}_x \left[1 + \frac{1}{v} \frac{v \text{Tr}(\Sigma_y^{-1} \Sigma_\phi(x))}{v + n - 2} \right. \\ &\quad + \frac{(\mu(x) - \mu_y)^\top \Sigma_y^{-1} (\mu(x) - \mu_y)}{v} \\ &\quad \left. + \frac{\mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right]}{v} \right]. \end{aligned}$$

Finally, the third term $\iint q(x, z|y)^{1+\gamma} dx dz$ is

$$\iint q_\phi(x, z|y)^{1+\gamma} dx dz = C_{v+n, m}^\gamma \left(1 + \frac{n}{v}\right)^{\frac{\gamma m}{2}} \left(1 + \frac{m}{v+n-2}\right) \int |\Sigma_\phi(x)|^{-\frac{\gamma}{2}} p_{data}(x)^{1+\gamma} dx,$$

where this last double integral is equal to the one computed for the t^3 -VAE.

Equipped with these formulas we can calculate the entropy \mathcal{H}_γ , cross-entropy \mathcal{C}_γ and the γ -divergence $\mathcal{D}(q||p)$ of our model. Firstly,

$$\begin{aligned} \mathcal{H}_\gamma &= - \left(\iint q(x, z)^{1+\gamma} dx dz \right)^{\frac{1}{1+\gamma}} \\ &= - C_{v+n, m}^{\frac{\gamma}{1+\gamma}} \left(1 + \frac{n}{v}\right)^{\frac{\gamma m}{2(1+\gamma)}} \left(1 + \frac{m}{v+n-2}\right)^{\frac{1}{1+\gamma}} \left(\int |\Sigma_\phi(x)|^{-\frac{\gamma}{2}} p_{data}(x)^{1+\gamma} dx \right)^{\frac{1}{1+\gamma}}, \end{aligned}$$

Which is similar to the one calculated in the t^3 VAE model.

Secondly,

$$\begin{aligned} \mathcal{C}_\gamma &= - \left(\iint q(x, z|y) p_\theta(x, z|y)^\gamma dx dz \right) \left(\iint p_\theta(x, z|y)^{1+\gamma} \right)^{-\frac{\gamma}{1+\gamma}} \\ &= - C_{v, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2} - \frac{1}{2}} \mathbb{E}_x \left[1 + \frac{1}{v} \frac{v \text{Tr}(\Sigma_y^{-1} \Sigma_\phi(x))}{v + n - 2} + \frac{(\mu(x) - \mu_y)^\top \Sigma_y^{-1} (\mu(x) - \mu_y)}{v} \right. \\ &\quad \left. + \frac{1}{v} \mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] \right] \left(C_{v, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \left(1 + \frac{m+n}{v-2}\right) \right)^{-\frac{\gamma}{1+\gamma}} \\ &= - \left(C_{v, m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{\gamma}{2}} \right)^{\frac{1}{1+\gamma}} |\Sigma_y|^{-\frac{1}{2}} \left(1 + \frac{m+n}{v-2}\right)^{-\frac{\gamma}{1+\gamma}} \mathbb{E}_x \left[1 + \frac{1}{v} \left(\frac{v \text{Tr}(\Sigma_y^{-1} \Sigma_\phi(x))}{v + n - 2} \right. \right. \\ &\quad \left. \left. + (\mu(x) - \mu_y)^\top \Sigma_y^{-1} (\mu(x) - \mu_y) + \mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] \right) \right]. \end{aligned}$$

Hence, we can define our divergence as :

$$\begin{aligned}
\mathcal{D}_\gamma(q||p) &= \frac{C_1}{\gamma} \left(\int |\Sigma_\phi(x)|^{-\frac{\gamma}{2}} p_{data}(x|y)^{1+\gamma} dx \right)^{\frac{1}{1+\gamma}} - \frac{C_2}{\gamma} \mathbb{E}_x \left[1 + \frac{1}{\nu} \left(\frac{\nu \text{Tr} \left(\Sigma_y^{-1} \Sigma_\phi(x) \right)}{\nu + n - 2} \right. \right. \\
&\quad \left. \left. + (\mu(x) - \mu_y)(\mu(x) - \mu_y)^\top + \mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] \right) \right] \\
&= \mathbb{E}_x \left[\frac{C_1}{\gamma} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} - \frac{C_2}{\gamma} \left(1 + \frac{1}{\nu} \left(\frac{\nu}{\nu + n - 2} \text{Tr} \left(\Sigma_y^{-1} \Sigma_\phi(x) \right) \right. \right. \right. \\
&\quad \left. \left. + (\mu(x) - \mu_y)^\top \Sigma_y^{-1} (\mu(x) - \mu_y) + \mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] \right) \right) \right],
\end{aligned}$$

with C_1 and C_2 being :

$$\begin{aligned}
C_1 &= C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}} \left(1 + \frac{n}{\nu} \right)^{\frac{\gamma m}{2(1+\gamma)}} \left(1 + \frac{m}{\nu + n - 2} \right)^{\frac{1}{1+\gamma}} \\
C_2 &= \left(C_{\nu,m+n}^\gamma |\Sigma_x|^{-\frac{\gamma}{2}} |\Sigma_y|^{-\frac{2\gamma+1}{2}} \right)^{\frac{1}{1+\gamma}} \left(1 + \frac{m+n}{\nu-2} \right)^{-\frac{\gamma}{1+\gamma}}.
\end{aligned}$$

On that account, the loss function for a class y is :

$$\begin{aligned}
\mathcal{L}(\gamma, y) &= -\frac{\nu\gamma}{C_2} \cdot \mathcal{D}_\gamma(q||p) \\
&= \mathbb{E}_x \left[\mathbb{E}_z \left[(x - \mu_\theta(z))^\top \Sigma_x^{-1} (x - \mu_\theta(z)) \right] + (\mu(x) - \mu_y)^\top \Sigma_y^{-1} (\mu(x) - \mu_y) \right. \\
&\quad \left. + \frac{\nu}{\nu + n - 2} \text{Tr} \left(\Sigma_y^{-1} \Sigma_\phi(x) \right) - \frac{\nu C_1}{C_2} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right],
\end{aligned}$$

and by taking $\Sigma_x = \sigma^2 I$ and $\Sigma_y = I$, we obtain :

$$\mathcal{L}(\gamma, y) = \mathbb{E}_x \left[\frac{\mathbb{E}_z [||x - \mu_\theta(z)||^2]}{\sigma^2} + ||\mu(x) - \mu_y||^2 + \frac{\nu \text{Tr} (\Sigma_\phi(x))}{\nu + n - 2} - \frac{\nu C_1}{C_2} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right].$$

D Sampling distribution variance derivation

In this section, we present the derivation of τ^2 used in the sampling of t^3 VAE and C- t^3 VAE model. We present only the derivation for the C- t^3 -VAE and it is identical to the one for the t^3 -VAE since the former model is a generalization of the later.

First, we simplify the divergence $\mathcal{D}(q||p^*)$:

$$\begin{aligned}
\mathcal{D}(q||p^*) &= -\frac{C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}}}{\gamma} \left(1 + \frac{m}{\nu + n - 2} \right)^{-\frac{\gamma}{1+\gamma}} \left[- \left(1 + \nu^{-1}n \right)^{\frac{\gamma m}{2(1+\gamma)}} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right. \\
&\quad \times \left(1 + \frac{m}{\nu + n - 2} \right) + |\tau^2 I|^{-\frac{\gamma}{2}} \left(1 + \nu^{-1}n \right)^{-1} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \\
&\quad \left. \times \left(1 + \frac{\text{Tr} \left(\tau^{-2} (1 + \nu^{-1}n)^{-1} \Sigma_\phi(x) \right)}{\nu + n - 2} + \frac{\tau^{-2}}{\nu + n} ||\mu(x) - \mu_y||^2 \right) \right]
\end{aligned}$$

Here, we use the fact that $|\alpha A|^\delta = \alpha^{\delta n} |A|^\delta$ where n is the dimension of the square A matrix. Also, we use $\text{Tr}(\alpha A) = \alpha \text{Tr}(A)$. After simplification and rearranging we get :

$$\begin{aligned}
\mathcal{D}(q||p^*) &= -\frac{1}{\gamma} C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}} \left(1 + \frac{m}{\nu+n-2}\right)^{-\frac{\gamma}{1+\gamma}} \left[- \left(1 + \nu^{-1}n\right)^{\frac{\gamma m}{2(1+\gamma)}} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right. \\
&\quad \times \left(1 + \frac{m}{\nu+n-2}\right) + \tau^{-\gamma m} \left(1 + \nu^{-1}n\right)^{-\frac{\gamma^2 m}{2(1+\gamma)}} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \\
&\quad \left. \left(1 + \frac{\tau^{-2} (1 + \nu^{-1}n)^{-1}}{\nu+n-2} \text{Tr}(\Sigma_\phi(x)) + \frac{\tau^{-2}}{\nu+n} \|\mu(x) - \mu_y\|^2\right) \right] \\
&= -\frac{1}{\gamma} C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}} \left(1 + \frac{m}{\nu+n-2}\right)^{-\frac{\gamma}{1+\gamma}} \left[- \left(1 + \nu^{-1}n\right)^{\frac{\gamma m}{2(1+\gamma)}} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \right. \\
&\quad \times \left(1 + \frac{m}{\nu+n-2}\right) + \nu^{-1} \tau^{-2-\gamma m} \left(1 + \nu^{-1}n\right)^{-\frac{\gamma^2 m}{2(1+\gamma)}-1} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \\
&\quad \left. \left(\kappa + \frac{\nu}{\nu+n-2} \text{Tr}(\Sigma_\phi(x)) + \|\mu(x) - \mu_y\|^2\right) \right] \\
&= -\frac{1}{\gamma} C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}} \left(1 + \frac{m}{\nu+n-2}\right)^{-\frac{\gamma}{1+\gamma}} \nu^{-1} \tau^{-2-\gamma m} \left(1 + \nu^{-1}n\right)^{-\frac{\gamma^2 m}{2(1+\gamma)}-1} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \\
&\quad \times \left[- \nu \tau^{2+\gamma m} \left(1 + \nu^{-1}n\right)^{\frac{\gamma^2 m + \gamma m}{2(1+\gamma)}+1} |\Sigma_\phi(x)|^{-\frac{\gamma^2}{2(1+\gamma)}} |\Sigma_\phi(x)|^{-\frac{\gamma}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right) \right. \\
&\quad \left. + \kappa + \frac{\nu}{\nu+n-2} \text{Tr}(\Sigma_\phi(x)) + \|\mu(x) - \mu_y\|^2 \right],
\end{aligned}$$

with:

$$\kappa = \nu \tau^2 \left(1 + \nu^{-1}n\right).$$

Then, we match the result to the loss function in Eq. (10) to get :

$$\tau^{2+\gamma m} \left(1 + \nu^{-1}n\right)^{\frac{\gamma m}{2}+1} |\Sigma_\phi(x)|^{-\frac{\gamma^2}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right) = \frac{C_1}{C_2}.$$

Moreover, we have :

$$\begin{aligned}
\frac{C_1}{C_2} &= C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}} \left(1 + \frac{n}{\nu}\right)^{\frac{\gamma m}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right)^{\frac{1}{1+\gamma}} C_{\nu,m+n}^{-\frac{\gamma}{1+\gamma}} \sigma^{\frac{n\gamma}{1+\gamma}} \left(1 + \frac{m+n}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}} \\
&= \sigma^{\frac{n\gamma}{1+\gamma}} C_{\nu+n,m}^{\frac{\gamma}{1+\gamma}} C_{\nu,m+n}^{-\frac{\gamma}{1+\gamma}} \left(1 + \frac{n}{\nu}\right)^{\frac{\gamma m}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right)^{\frac{1}{1+\gamma}} \left(1 + \frac{m+n}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}} \\
&= \sigma^{\frac{n\gamma}{1+\gamma}} C_{\nu,n}^{\frac{-\gamma}{1+\gamma}} \left(1 + \frac{m}{\nu+n-2}\right)^{\frac{1}{1+\gamma}} \left(1 + \frac{m+n}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}}.
\end{aligned}$$

Consequently we obtain :

$$\begin{aligned}
\tau^{2+\gamma m} \left(1 + \nu^{-1}n\right)^{\frac{\gamma m}{2}+1} |\Sigma_\phi(x)|^{-\frac{\gamma^2}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right) \\
= \sigma^{\frac{n\gamma}{1+\gamma}} C_{\nu,n}^{\frac{-\gamma}{1+\gamma}} \left(1 + \frac{m}{\nu+n-2}\right)^{\frac{1}{1+\gamma}} \left(1 + \frac{m+n}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}}
\end{aligned}$$

$$\begin{aligned}
\tau^{2+\gamma m} &= \sigma^{\frac{n\gamma}{1+\gamma}} C_{\nu,n}^{-\frac{\gamma}{1+\gamma}} \left(1 + \nu^{-1}n\right)^{-\frac{\gamma m}{2}-1} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right)^{-1} \left(1 + \frac{m}{\nu+n-2}\right)^{\frac{1}{1+\gamma}} \\
&\quad \times \left(1 + \frac{m+n}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}} \\
&= \sigma^{\frac{n\gamma}{1+\gamma}} C_{\nu,n}^{-\frac{\gamma}{1+\gamma}} \left(1 + \nu^{-1}n\right)^{-\frac{\gamma m}{2}-1} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \left(1 + \frac{m}{\nu+n-2}\right)^{-\frac{\gamma}{1+\gamma}} \left(1 + \frac{m+n}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}} \\
&= \sigma^{\frac{n\gamma}{1+\gamma}} C_{\nu,n}^{-\frac{\gamma}{1+\gamma}} \left(1 + \nu^{-1}n\right)^{-\frac{\gamma m}{2}-1} |\Sigma_\phi(x)|^{\frac{\gamma^2}{2(1+\gamma)}} \left(\frac{\nu+n-2}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}} \\
&= \left(1 + \nu^{-1}n\right)^{-\frac{\gamma m}{2}-1} \left(\sigma^n C_{\nu,n}^{-1} |\Sigma_\phi(x)|^{\frac{\gamma}{2}} \frac{\nu+n-2}{\nu-2}\right)^{\frac{\gamma}{1+\gamma}}.
\end{aligned}$$

Hence, we get :

$$\begin{aligned}
\tau^2 &= \left(1 + \nu^{-1}n\right)^{-\frac{\gamma m+2}{\gamma m+2}} \left(\sigma^n C_{\nu,n}^{-1} |\Sigma_\phi(x)|^{\frac{\gamma}{2}} \frac{\nu+n-2}{\nu-2}\right)^{\frac{2\gamma}{(1+\gamma)(2+\gamma m)}} \\
&= \left(1 + \nu^{-1}n\right)^{-1} \left(\sigma^n C_{\nu,n}^{-1} |\Sigma_\phi(x)|^{\frac{\gamma}{2}} \frac{\nu+n-2}{\nu-2}\right)^{\frac{2\gamma}{(1+\gamma)(2+\gamma m)}},
\end{aligned}$$

which is the form of τ^2 we report in Eq. (9).

E Experimental Setup Details

E.1 Datasets

We conduct experiments on three datasets notably SVHN-LT [11], CIFAR100-LT [24, 12] and CelebA [13] each chosen to highlight different challenges related to generative modeling under class imbalance and varying visual complexity.

- **SVHN-LT** : This dataset is comprised of colored images of digits from 0 to 9 of size $32 \times 32 \times 3$. It serves as our controlled experimental setting. While simple enough for all models to converge, it is rich enough to reflect performance differences. However, as this dataset is naturally imbalanced, we balance the number of images across classes to have more control over the imposed imbalance ratio. In Table 2 we provide the number of images present in the dataset before balancing.

Table 2: The Number of images in the SVHN dataset for the train and test sets before balancing. The value in bold is the one used to balance the dataset.

Class	0	1	2	3	4	5	6	7	8	9
Train set	4948	13861	10585	8497	7458	6882	5727	5595	5045	4656
Test set	1744	5099	4149	2882	2523	2384	1977	2019	1660	1595

- **CIFAR100-LT** : This balanced dataset comprised of colored images of 100 classes of different natural objects. The images are of size $32 \times 32 \times 3$ presents high variability with 100 fine-grained classes. In the most extreme imposed imbalance setting, tail classes may contain as few as five examples. Though the images remain low-resolution, this setting stresses the models' robustness in the face of sparse data and various categories.
- **CelebA** : This dataset enables evaluation in real-world class-imbalance scenarios using attribute labels (e.g., Mustache, Young). It also introduces the complexity of higher resolution images ($178 \times 218 \times 3$), testing the scalability of our models. We select four binary attributes (Mustache,

Young, Male, Smiling) for evaluation, with test sets balanced by down-sampling the larger class. In Table 3, we provide the number of images for each chosen attribute in this dataset.

Table 3: The Number of images in the CelebA dataset for the train and test sets for the Mustache, Young, Male and Smiling attributes.

	Mustache	No Mustache	Young	Old	Male	Female	Smiling	Not Smiling
Train set	6642	156128	126788	35982	68261	94509	78080	84690
Test set	722	19190	15114	4848	7715	12247	9987	9975

E.2 Data Preprocessing

- **SVHN-LT Dataset** For both training and testing, we crop each class to the minimum number of samples available across all classes. The only data augmentation applied is a random horizontal flip with 50% probability.
- **CIFAR100-LT Dataset** We use the dataset in its entirety without class filtering. As with SVHN-LT, we apply a random horizontal flip with 50% probability for data augmentation.
- **CelebA Dataset** We use the full dataset and apply a center crop of 160×160 to each image and then resize to 128×128 . A random horizontal flip with 50% probability is also applied.

E.3 Model Architecture

Our encoder-decoder models follow a modular block design. Each encoder block consists of a convolutional layer, followed by 2D batch normalization and ReLU activation. Decoder blocks mirror this structure but replace convolutional layers with transposed convolutions.

- **SVHN-LT and CIFAR100-LT** : Encoders consist of four convolutional blocks with channels $\{64, 128, 256, 512\}$, followed by two linear layers for estimating mean and covariance. The decoder uses three transposed convolutional blocks with channel sizes $\{128, 64, 32\}$, ending with a three-channel convolution and Sigmoid activation.
- **CelebA** : The CelebA encoder includes six convolutional blocks with channels $\{64, 128, 256, 512, 512, 512\}$, ending with two linear layers. The decoder has six transposed convolutional layers with channels $\{512, 512, 256, 128, 64, 32\}$, followed by a final convolutional layer and Sigmoid activation.

E.4 Training Details

All models are trained using the AdamW optimizer with a learning rate of 10^{-3} for 150 epochs. We use a batch size of 64 for SVHN-LT and CIFAR100-LT, and 128 for CelebA.

F Hyperparameter Tuning

We present the hyperparameter tuning process used across all evaluated models. We first optimize β , then ν , and finally τ , yielding the models' results reported in Table 1.

F.1 β Optimization

We perform a hyperparameter study over β for all tested models. Unless otherwise noted, we use the theoretically derived τ^2 and set $\nu = 10$.

F.1.1 On the SVHN-LT dataset

As shown in Figure 6, the optimal β values for Student's t models lies in the range $\beta \in [0.4, 0.6]$ whereas it lies in the $\beta \in [0.05, 0.07]$ range for Gaussian-based models. This is because the regularization term in the γ -power divergence loss is ten times larger than the KL divergence. Figure 6 also shows that FID performance is highly sensitive to β in the Gaussian setting, requiring careful tuning which is not the case for Student's t based models. Finally, C- t^3 VAE achieves the best FID surpassing the t^3 VAE and the C-VAE for all imbalance settings.

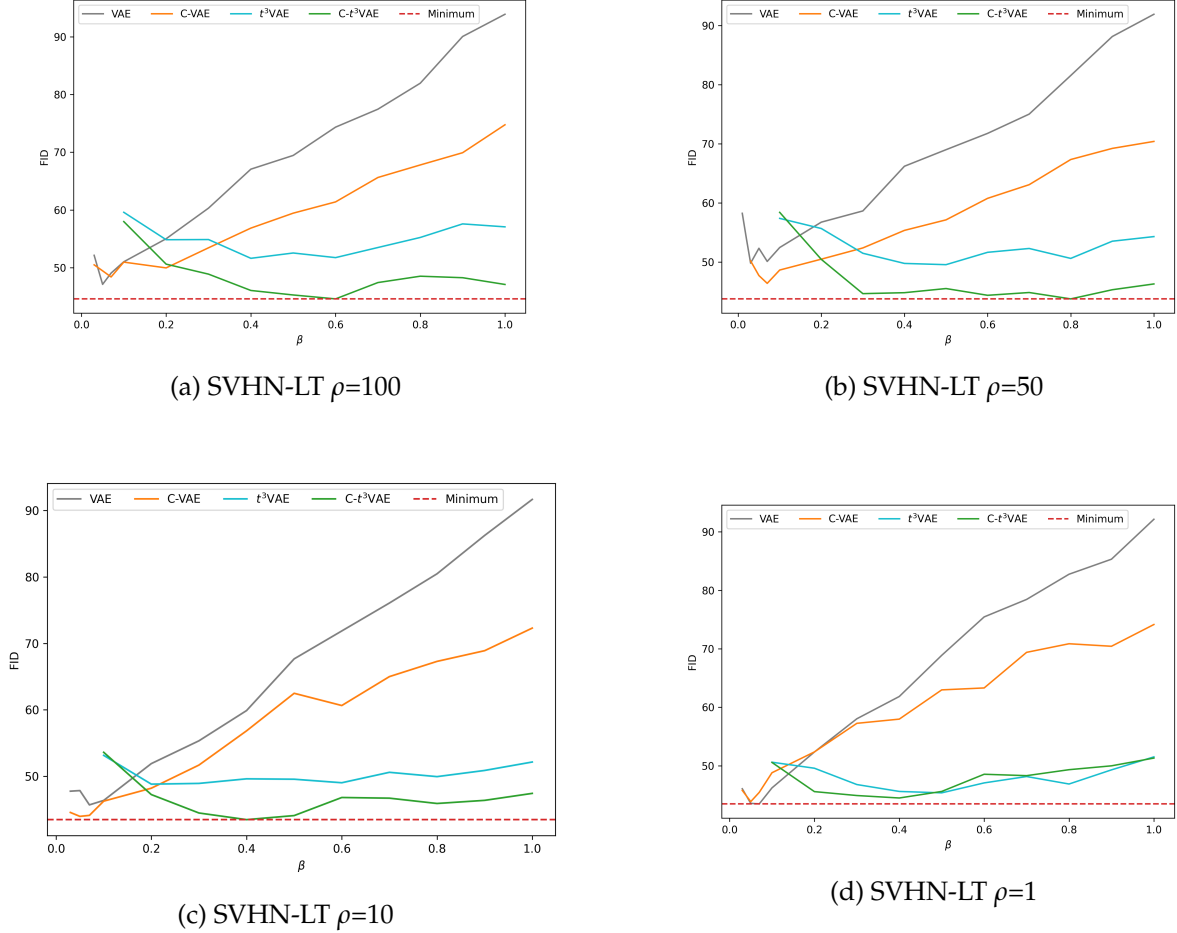
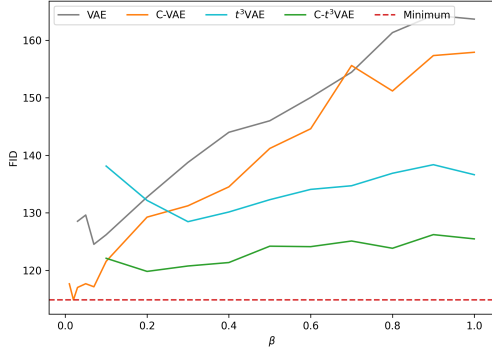


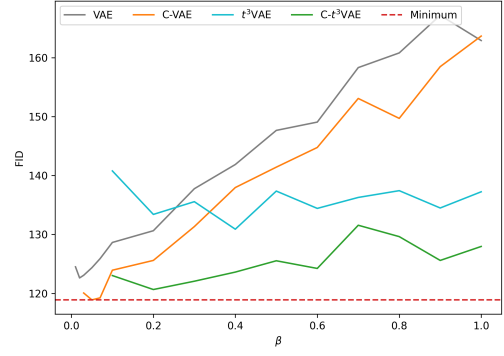
Figure 6: Variability of the FID as a function of the β hyperparameter for the VAE, C-VAE, t^3 VAE and C- t^3 VAE on the SVHN-LT dataset.

F.1.2 On the CIFAR100-LT dataset

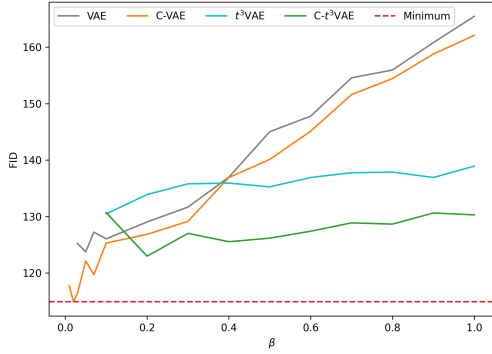
From Figure 7, we observe that Student's t models obtain the best performance in terms of FID at $\beta = 0.2$ for CIFAR100-LT dataset. However, for the Gaussian-based models, the optimal value is much lower with $\beta \in [0.02, 0.05]$. The reason for this is that on this dataset too the KL regularization term is ten times smaller than the regularization terms present in the γ -power divergence loss. Additionally, we notice that C-VAE performs slightly better, likely due to the complexity of the dataset preventing full convergence to the imposed latent distribution. We further investigate this hypothesis in the τ analysis below.



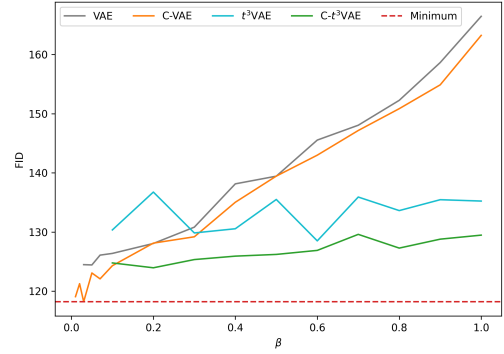
(a) CIFAR100-LT $\rho=100$



(b) CIFAR100-LT $\rho=50$



(c) CIFAR100-LT $\rho=10$

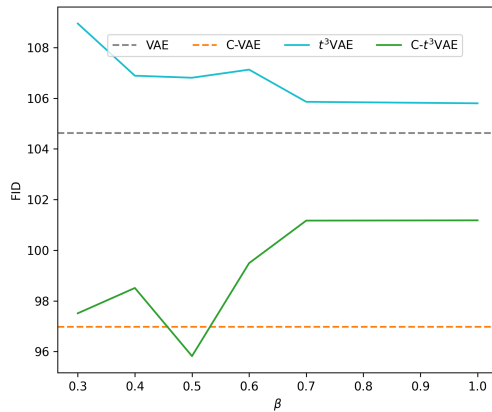


(d) CIFAR100-LT $\rho=1$

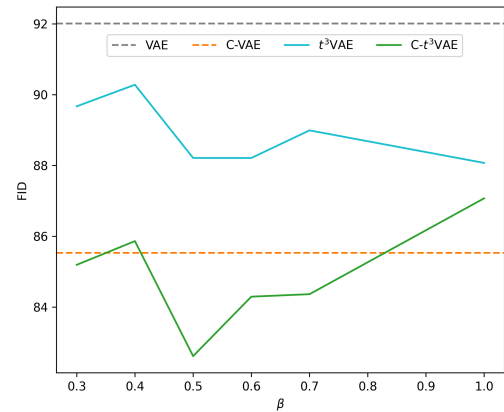
Figure 7: Variability of the FID as a function of the β hyperparameter for the VAE, C-VAE, t^3 VAE and C- t^3 VAE on the CIFAR100-LT dataset.

F.1.3 On the CelebA dataset

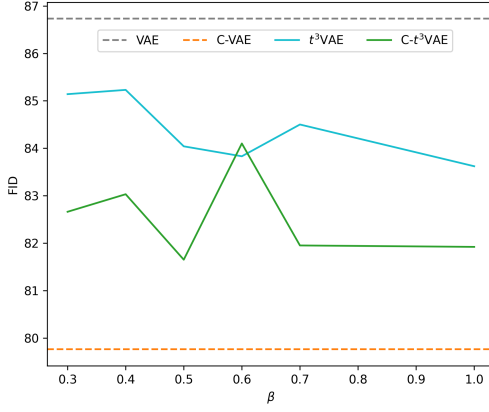
For the CelebA dataset, we focus only on β optimization for Student's t models. As for the Gaussian-based models, we set $\beta = 0.1$, because as shown in Table 1 the optimization of the β hyperparameter has marginal impact on the FID for the CelebA dataset, which is not the case for the SVHN-LT and CIFAR100-LT dataset. Hence, we did not deem necessary to perform a hyper-parameter optimization over β for the Gaussian based models. Nevertheless, from Figure 8 we see that β has marginal effect on the Student's t based models too, likely due to the low intra-class variability.



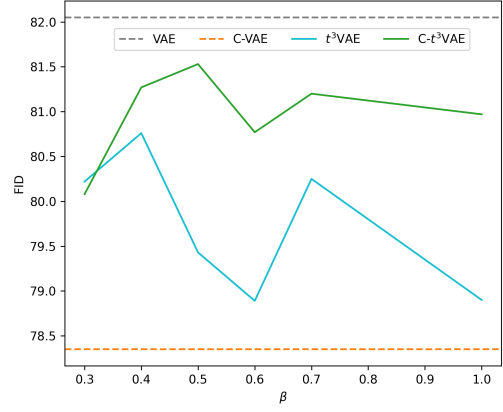
(a) CelebA - Mustache



(b) CelebA - Young



(c) CelebA - Male



(d) CelebA - Smiling

Figure 8: Variability of the FID as a function of the β hyperparameter for the t^3 VAE and $C-t^3$ VAE on the CelebA dataset. The horizontal lines for the VAE and C-VAE models are for the best performing model between $\beta = 0.1$ and $\beta = 1$.

F.2 ν Optimization

Table 4 shows results from tuning the degrees of freedom parameter ν in $C-t^3$ VAE on both SVHN-LT and CIFAR100-LT across all imbalance ratios, using the optimal β from the previous study. On average, $\nu = 10$ performs well, consistent with prior work [10]. However, performance can be further improved by selecting ν within the range $[2.5, 20]$. Still, no major influence of hyper-parameter ν on the generative FID can be observed similarly to what was observed by [10] for the reconstruction FID.

ν	SVHN-LT				CIFAR100-LT			
	100	50	10	1	100	50	10	1
2.1	45.50	44.51	42.96	46.23	121.28	122.03	121.93	123.41
2.5	45.76	43.96	45.81	45.40	119.15	120.19	120.10	124.83
5	44.89	42.60	45.03	46.33	120.52	123.21	124.29	123.71
10	44.59	44.37	43.48	44.49	119.83	120.65	122.96	123.95
20	44.02	43.89	42.01	44.75	121.48	118.41	124.58	126.13
50	48.03	46.39	43.59	45.57	119.58	126.36	124.38	127.48
100	45.97	44.63	43.74	47.52	123.26	122.90	127.42	125.67

Table 4: Variability of the FID as a function of the standard deviation ν for the $C-t^3$ VAE model on the SVHN-LT and CIFAR100-LT datasets.

F.3 τ Optimization

In this section, we evaluate the effect of the τ parameter on the SVHN-LT, CIFAR100-LT and CelebA datasets for all imbalance ratios while setting β and ν to their previously optimized values. As shown in Figure 9, the optimal τ for SVHN-LT aligns closely with our theoretical prediction. In contrast, CIFAR100-LT consistently benefits from a larger $\tau = 0.4$, yielding improved FID across all imbalance settings and outperforming C-VAE. On CelebA, τ has minimal impact and the most likely value is $\tau \approx 0.3$. This value is between the value observed in for the SVHN-LT and CIFAR100-LT, as is mostly linked to the complexity and the entropy of the dataset where the more complex the data is the higher the value of τ is required.

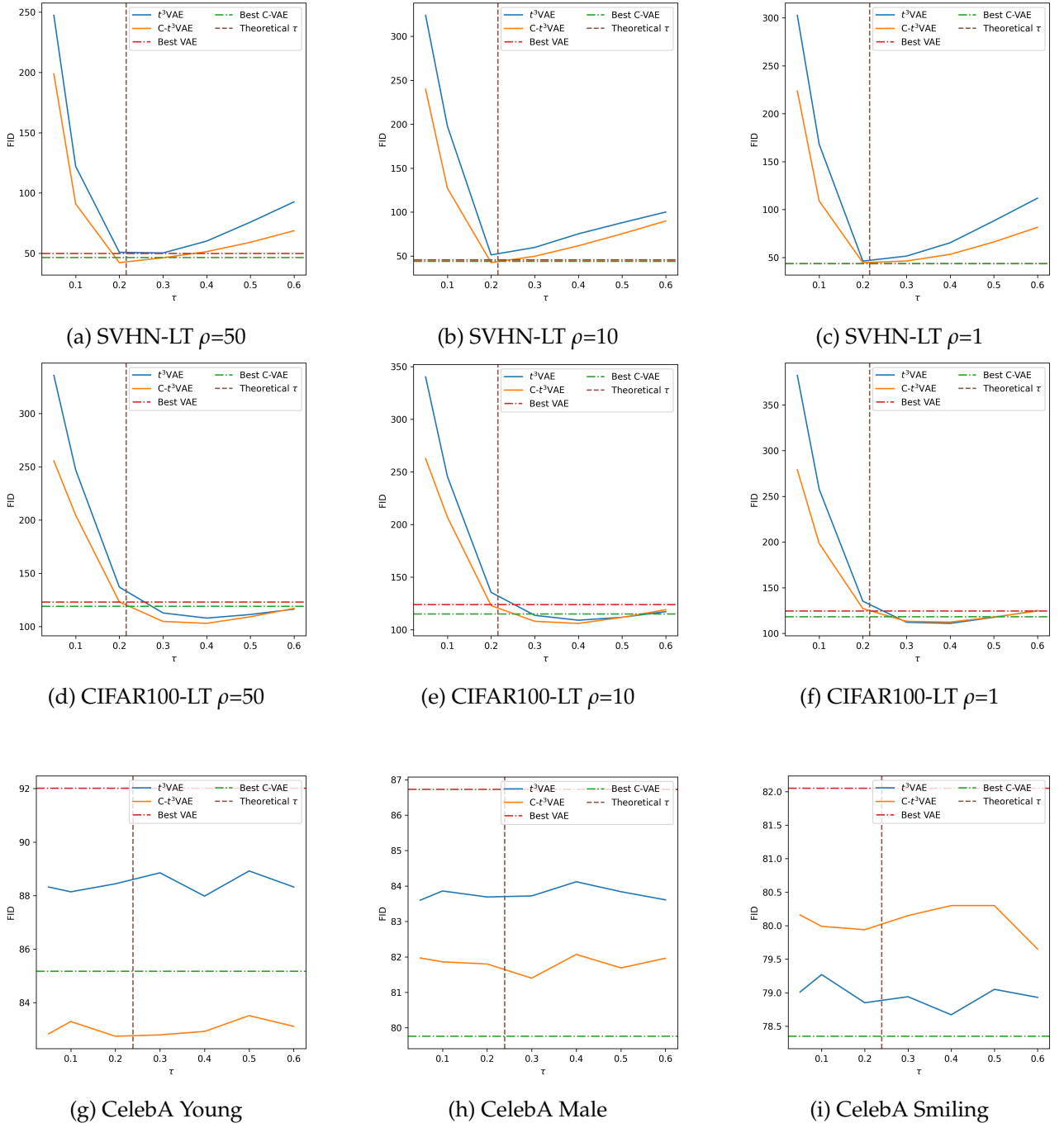


Figure 9: Variability of the FID as a function of the standard deviation τ^2 for the t^3 VAE and C- t^3 VAE. In horizontal dashed lines is the FID value of the best performing VAE and C-VAE on each dataset. In vertical dashed lines is the theoretically identified value of τ .

G Per-Class Evaluation

In this section, we assess the conditional models' per-class Recall, Precision, and F1 metrics under all imbalance settings and for all tested datasets after optimization of all hyper-parameters.

From the following figures in Table 5 and 6, we see that the $C\text{-}t^3\text{VAE}$ consistently improves Recall and mode coverage in highly imbalanced settings with $\rho = 100$ and $\rho = 50$. This comes at a minor Precision cost but results in significantly better F1 scores across most classes. However, on balanced or mildly imbalanced datasets, its performance remains competitive with Gaussian-based models. This observation is valid for both the SVHN-LT and CIFAR100-LT but is more pronounced on the later.

Table 5: Per-class generative metrics on SVHN-LT after optimization of β , ν and τ hyper-parameters.

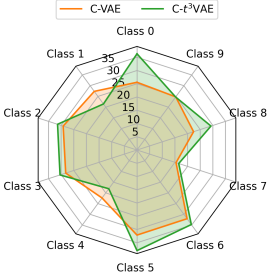
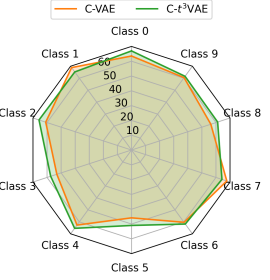
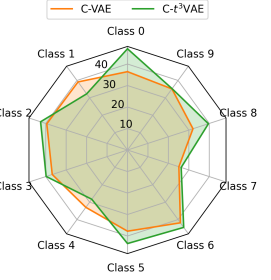
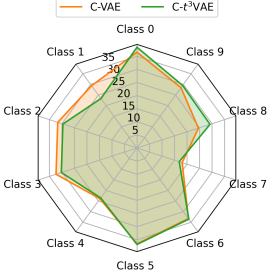
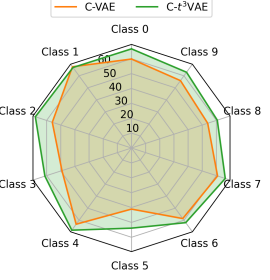
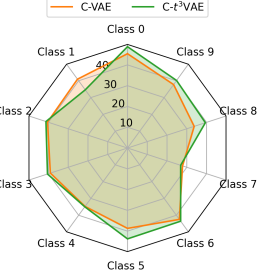
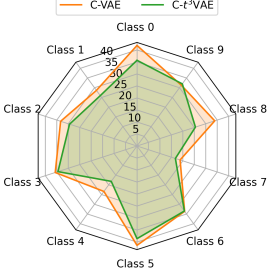
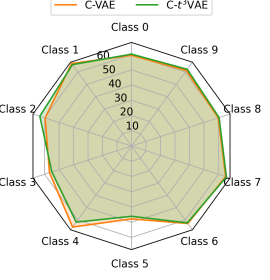
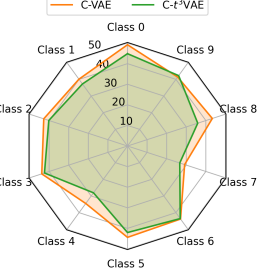
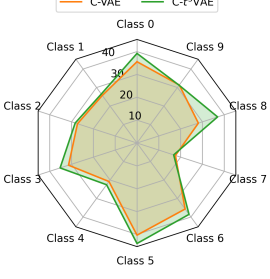
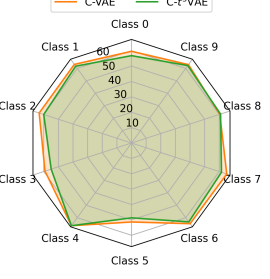
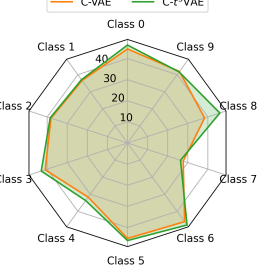
ρ	Recall	Precision	F1 score
100			
50			
10			
1			

Table 6: Per-class generative metrics on CIFAR100-LT after optimization of β , ν and τ hyper-parameters, we focus on the top 5 head and tail classes.

