

FissionVAE: Federated Non-IID Image Generation with Latent Space and Decoder Decomposition

Chen Hu¹, Hanchi Ren¹, Jingjing Deng², Xianghua Xie^{1,*} and Xiaoke Ma³

¹Swansea University, United Kingdom

²Durham University, United Kingdom

³Xi'Dian University, P. R. China

{chen.hu, hanchi.ren, x.xie}@swansea.ac.uk, jingjing.deng@durham.ac.uk

Abstract

Federated learning is a machine learning paradigm that enables decentralized clients to collaboratively learn a shared model while keeping all the training data local. While considerable research has focused on federated image generation, particularly Generative Adversarial Networks, Variational Autoencoders have received less attention. In this paper, we address the challenges of non-IID (independently and identically distributed) data environments featuring multiple groups of images of different types. Non-IID data distributions can lead to difficulties in maintaining a consistent latent space and can also result in local generators with disparate texture features being blended during aggregation. We thereby introduce FissionVAE that decouples the latent space and constructs decoder branches tailored to individual client groups. This method allows for customized learning that aligns with the unique data distributions of each group. Additionally, we incorporate hierarchical VAEs and demonstrate the use of heterogeneous decoder architectures within FissionVAE. We also explore strategies for setting the latent prior distributions to enhance the decoupling process. To evaluate our approach, we assemble two composite datasets: the first combines MNIST and FashionMNIST; the second comprises RGB datasets of cartoon and human faces, wild animals, marine vessels, and remote sensing images. Our experiments demonstrate that FissionVAE greatly improves generation quality on these datasets compared to baseline federated VAE models.

1 Introduction

Generative models have attracted increasing attention in recent years due to their impressive ability to generate new data across various modalities, including images [Ho *et al.*, 2020], texts [Touvron *et al.*, 2023], and audios [Borsos *et al.*, 2023]. As these models, like other deep learning systems, require substantial amounts of data, concerns regarding data

privacy have elevated among regulatory authorities and the public. Unlike the traditional centralized learning paradigm, which collects all data on a single computer system for training, federated learning allows private data to remain on the owner's device. In this paradigm, local devices train models independently, and a central server aggregates these models without accessing the individual data directly. Although this distributed approach enhances privacy protection, it also introduces unique challenges not encountered in centralized systems. Since data remains distributed across various client devices, the training samples are not guaranteed to be identically distributed. This can lead to inconsistencies in learning objectives among clients, resulting in degraded performance when these models are aggregated on the server.

In the context of FL with non-IID data, generative models such as Generative Adversarial Networks (GANs) [Goodfellow *et al.*, 2014] and Variational Autoencoders (VAEs) [Kingma and Welling, 2014] face additional challenges. These models involve sampling from a latent distribution, and the generator or decoder trained on client devices may develop differing interpretations of the same latent space. This discrepancy can lead to difficulties in maintaining a consistent and unified latent space, resulting in ambiguous latent representations. A further challenge arises from the role of the generator or decoder, which are tasked with mapping latent inputs to the sample space by synthesizing the shape, texture, and colors of images. Aggregating generative models trained on non-IID image data can produce artifacts that appear as a blend of disparate image types, because generators trained on non-IID local data capture the characteristics of varied visual features. Specifically for GANs, another problem arises from local discriminators, which may provide conflicting feedback that hinders model convergence. With the limited data available in FL settings, discriminators can quickly overfit to the training samples [Karras *et al.*, 2020]. If an updated generator from the server produces images of classes not present in a client's local dataset, the local discriminator might incorrectly label well-generated images as fake, simply because they do not match the local data distribution. This mislabeling can significantly impede the generator's ability to synthesize realistic images.

Existing research on generative models for non-IID data in federated learning (FL) has primarily focused on GANs. MDGAN [Hardy *et al.*, 2019] proposes exchanging local dis-

* Corresponding author.

84 criminators among clients during training. This strategy al- 143
85 lows discriminators to access a broader spectrum of local 144
86 data, thereby avoiding biased feedback to the generator. The 145
87 authors of [Yonetani *et al.*, 2019] uses the local discriminator 146
88 that gives the highest score to a generated sample to update 147
89 the global generator, promoting the idea that local discrimina- 148
90 tors should only judge samples from familiar distributions. In 149
91 [Xiong *et al.*, 2023], the authors aggregate generators at the 150
92 group level for client groups sharing similar data distributions 151
93 before performing a global aggregation, then the global gen- 152
94 erator is aggregated similar to [Yonetani *et al.*, 2019]. Both 153
95 [Yonetani *et al.*, 2019] and [Xiong *et al.*, 2023] involve send- 154
96 ing synthesized samples back to local clients, which could 155
97 potentially increase the risk of compromising client data pri- 156
98 vacy. 157

99 Studies employing VAEs solely for image generation pur- 158
100 poses are less common. The works in [Chen and Vikalo, 159
101 2023] and [Heinbaugh *et al.*, 2023] utilize VAEs to produce 160
102 synthetic images that assist in training global classifiers. In 161
103 [Chen and Vikalo, 2023], the global decoder generates mi- 162
104 nority samples for local classifiers by sampling from class 163
105 means with added noise. The approach in [Heinbaugh *et* 164
106 *al.*, 2023] treats converged local decoders as teacher mod- 165
107 els and uses knowledge distillation to train a global generator 166
108 on the server side without further local updates. While this 167
109 decoder can produce useful samples for classification tasks, 168
110 it risks overfitting to the potentially flawed output from local 169
111 decoders and lacks generative diversity, which is crucial for 170
112 high-quality image generation. Recent studies [Bohacek and 171
113 Farid, 2023] [Shumailov *et al.*, 2024] have shown that gen- 172
114 erative models trained on generated samples instead of real 173
115 data are prone to collapsing. VAEs are also widely used in 174
116 collaborative filtering tasks for recommendation systems [Po- 175
117 lato, 2021; Zhang *et al.*, 2024; Li *et al.*, 2025]. These mod- 176
118 els typically learn user embeddings from interaction vectors 177
119 using a standard Gaussian prior, and decode into item-score 178
120 distributions for ranking. In contrast, image generation tasks 179
121 require decoding into high-dimensional pixel space, where is- 180
122 sues such as latent space ambiguity and domain-specific tex- 181
123 ture blending and arise, which are not present in collaborative 182
124 filtering. As such, the architectural and modeling considera- 183
125 tions in our work are fundamentally different. 184

126 In response to the challenges posed by non-IID data in fed- 185
127 erated image generation, we introduce a model named Fis- 186
128 sionVAE. This model is specifically tailored to environments 187
129 featuring multiple groups of images of different types. To 188
130 mitigate the problem of mixed latent space interpretation, Fis- 189
131 sionVAE decomposes the latent space into distinctive priors, 190
132 hence adapting to the diverse data distributions across differ- 191
133 ent image types. We further refine this approach by investi- 192
134 gating strategies for encoding the prior Gaussians. Addition- 193
135 ally, to prevent the blending of unrelated visual features in 194
136 the generated outputs, FissionVAE employs specialized de- 195
137 coder branches for each client group. This method not only 196
138 accommodates the unique characteristics of each data sub- 197
139 set but also enhances the model’s generative capabilities in 198
140 highly heterogeneous environments. The primary contribu- 199
141 tions of our research are detailed as follows:

142 1. We introduce FissionVAE for federated non-IID image

generation. In FissionVAE, we decompose the latent space 143
according to the distinct data distributions of client groups. 144
This approach ensures that each client’s data are mapped to 145
its corresponding latent distribution without the adverse ef- 146
fects of averaging dissimilar distributions during aggregation. 147
Moreover, by implementing separate decoder branches for 148
different groups of data, FissionVAE allows for specialized 149
generation tailored to different image types, which is crucial 150
for preserving the distinct visual features of different image 151
types during the generative process. 152

2. We explore various strategies for encoding Gaussian pri- 153
ors to enhance the effectiveness of latent space decomposi- 154
tion. We further extends FissionVAE by introducing the hi- 155
erarchical inference architecture. We demonstrate that with 156
the decomposed decoder branches, it is feasible to employ 157
heterogeneous decoder architectures in FissionVAE, allowing 158
for more flexible model deployment on clients. 159

3. We validate FissionVAE with extensive experiments on 160
two composite datasets combining MNIST with FashionM- 161
NIST, and a more diverse set comprising cartoon and human 162
faces, animals, marine vessels, and remote sensing images. 163
Our results demonstrate improvements in generation quality 164
over the existing baseline federated VAE. 165

The remainder of the paper is organized as follows: In 166
Section 2, we describe the baseline FedVAE model and the 167
FissionVAE variants we propose. Section 3 presents the ex- 168
perimental setup, including the configuration details and an 169
analysis of the results. Finally, we conclude the paper in Sec- 170
tion 4 with a summary of our findings and a discussion on 171
potential future directions. 172

2 Investigating Strategies for Non-IID Image 173 Generation with VAEs 174

In this section, we describe our methodology for exploring 175
VAE configurations tailored for generating images under non- 176
IID conditions in a federated learning framework. For back- 177
ground on FL and VAEs, please refer to the supplementary 178
material. We specifically address scenarios where clients are 179
categorized based on distinct data distributions. For illustra- 180
tive purposes, we consider the case where some clients ex- 181
clusively possess hand-written digit images from the MNIST 182
dataset, while others maintain only clothing images from the 183
FashionMNIST dataset. We follow to the standard federated 184
learning framework, wherein a central server is tasked with 185
aggregating updates from the clients and subsequently dis- 186
tributing the updated model back to them. FedAvg [McMa- 187
han *et al.*, 2023] is employed for server-side aggregation. 188
Each client retains a subset of data representative of its re- 189
spective group and conducts local training independently. A 190
more practical scenario with RGB images and a larger num- 191
ber of client groups is explored and discussed in the experi- 192
ments section (Section 3). 193

2.1 FedVAE 194

A straightforward strategy for implementing VAEs in feder- 195
ated learning is using a unified encoder-decoder architecture. 196
In this configuration, all clients share a common latent space 197
(often predefined as the normal distribution $\mathcal{N}(0, 1)$) and the 198

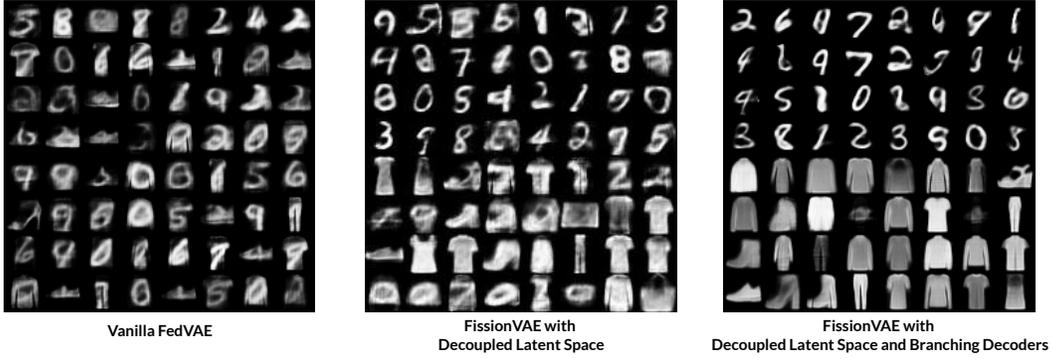


Figure 1: Qualitative results of the baseline FedVAE and proposed FissionVAEs. As we further decoupling the latent space and decoders in the federated environment, the quality of generated images is improved.

199 central server indiscriminately aggregates client models at the
 200 end of each training round. This approach is named FedVAE
 201 in [Jiang *et al.*, 2023] for trajectory data generation. Fig. 2
 202 illustrates this baseline training scheme.

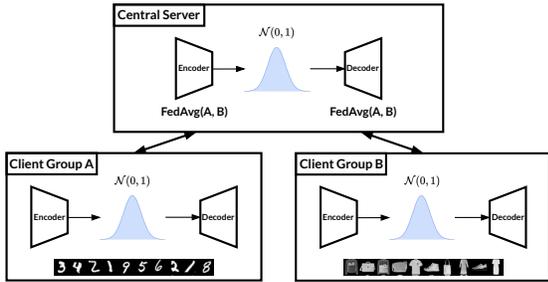


Figure 2: An illustration of baseline FedVAE. The encoder and the decoder of the VAE are aggregated through FedAvg regardless of their client groups.

203 Despite the simplicity of this strategy, it present significant
 204 challenges in the non-IID scenario. Specifically, employing a
 205 single prior distribution for the latent space does not account
 206 for the distinct data distributions across different clients. En-
 207 coders from different client groups may map their uniquely
 208 distributed data into the same region of the latent space. Con-
 209 sequently, client decoders might interpret this shared latent
 210 space differently, leading to inconsistencies or even conflicts
 211 among client models during aggregation at the server. Figure
 212 1 shows randomly generated samples produced after training
 213 the federated Vanilla VAE on the combined dataset of MNIST
 214 and FashionMNIST. These samples clearly exhibit artifacts
 215 that appear to blend features of handwritten digits with cloth-
 216 ing items, indicating the aggregation conflicts inherent in this
 217 method.

2.2 FissionVAE with Latent Space Decoupling

219 To address the conflicting latent space issue identified above,
 220 we propose decomposing the latent space according to differ-
 221 ent data groups, while maintaining a unified architecture for
 222 the encoder and decoder. This approach corresponds to the

architecture shown in Fig. 3.

223

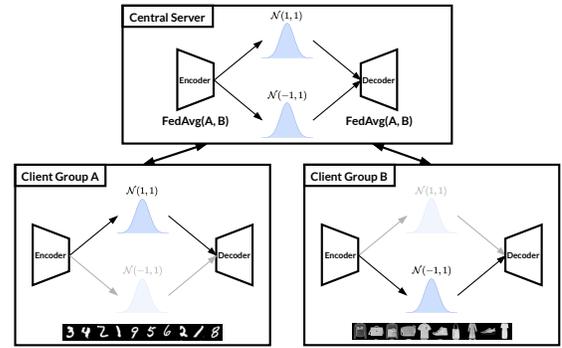


Figure 3: An illustration of FissionVAE with Latent Space Decoupling. The latent variables are forced to follow their respective group prior distributions. The model is aggregated the same way as the baseline FedVAE.

224 When decoupling the latent space, the encoder maps the input
 225 data to different distributions based on the client’s group.
 226 For instance, MNIST client may map to $\mathcal{N}(-1, 1)$ and Fash-
 227 ionMNIST clients to $\mathcal{N}(1, 1)$. The KL divergence in the
 228 ELBO for this model is given by:

$$D_{\text{KL}}(\mathcal{N}(\mu_q, \sigma_q) \parallel \mathcal{N}(\pm 1, 1)) = \frac{1}{2} \sum_{i=1}^k [\sigma_i + \mu_i^2 \mp 2\mu_i - \log \sigma_i] \quad (1)$$

229 Here, μ_q and σ_q represent the encoder’s estimates for the
 230 parameters of the latent code’s distribution, and k is the di-
 231 mension of the latent code.

232 Figure 1 shows randomly generated amples produced after
 233 training the FissionVAE with latent space decoupling on
 234 the Mixed MNIST dataset. While the quality of reconstructed
 235 images are improved compared to the baseline FedVAE, the
 236 generated images still exhibit a mixture of handwritten digits
 237 and clothing items, even when explicitly sampling from their
 238 respective latent distributions. This suggests that while de-
 239 composing latent encoding helps improving reconstructions,

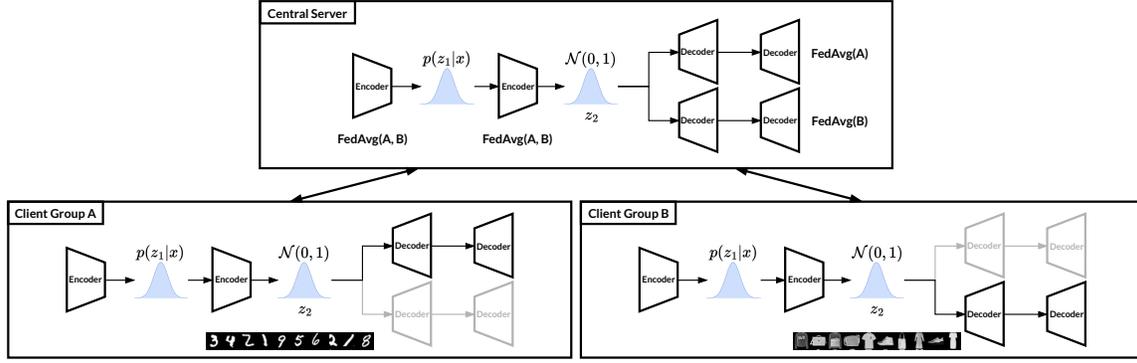


Figure 4: An illustration of Hierarchical FissionVAE. This FissionVAE architecture extends to allow two levels of latent variables. The latent variable z_1 can be either learned or predefined. As input from different groups has been separated by z_1 , the latent variable z_2 is set to follow the standard normal distribution.

240 the unified decoder still blends features due to the aggregation
 241 of model weights from diverse visual domains. This observa-
 242 tion motivates the architecture described in the next section,
 243 where the decoder is also split based on client groups.

2.3 FissionVAE with Group-specific Decoder Branches

246 **Non-Hierarchical FissionVAE** Building on the concept intro-
 247 duced by FissionVAE with latent space decoupling, we
 248 further refines non-IID data generation by incorporating de-
 249 coder branches specific to each data group while maintaining
 250 a unified encoder. This design allows the central server to ag-
 251 gregate the encoder updates agnostically of the client groups,
 252 whereas decoder branches are aggregated specifically accord-
 253 ing to their corresponding groups. In addition, this approach
 254 also offers flexibility in the choice of the prior latent distribu-
 255 tion $p(z)$ for each group to exert more explicit control over
 256 the data generation through the decoder. Figure 5 illustrates
 257 this branching architecture.

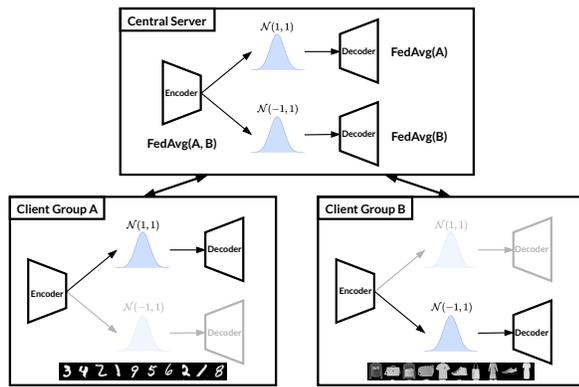


Figure 5: An illustration of FissionVAE with Decoder Branch Decoupling. This FissionVAE creates decoders specific to client groups and enforces constraints for latent variable priors. The encoder is aggregated across groups while the group-specific decoder is only aggregated from local models within the corresponding group.

258 Figure 1 also includes randomly generated samples pro-
 259 duced after training the FissionVAE with decoder branches.
 260 The results indicate a significant reduction in the blending
 261 feature issue in previously discussed VAE architectures.

262 **Hierarchical FissionVAE** Next, we show that the branching
 263 architecture can be enhanced by integrating hierarchical infer-
 264 ence [Kingma *et al.*, 2016] [Sønderby *et al.*, 2016] to
 265 the federated learning framework, which enables the use of
 266 deeper network structures to capture more complex data dis-
 267 tributions. Fig 4 depicts the FissionVAE with two levels of
 268 hierarchical inference. In this architecture, the first encoder
 269 module estimates $q(z_1|x)$ from the input data, then the sec-
 270 ond encoder module estimates $q(z_2|z_1)$ based on the first
 271 level latent code. The decoder reverses the encoding process,
 272 which estimates $p(z_1|z_2)$ based on z_2 to reconstruct z_1 , and
 273 subsequently reconstructs the original input x by estimating
 274 $p(x|z_1)$.

275 Following the convention in hierarchical VAEs, we assume
 276 conditional independence among the latent codes. Then the
 277 ELBO for this hierarchical VAE is expressed as (refer to sup-
 278 plimentary material for derivation),

$$\begin{aligned}
 \text{ELBO}_H &= \mathbb{E}_{q_\phi(z_1|x)}[\log p_\theta(x|z_1)] \\
 &\quad - \mathbb{E}_{q_\phi(z_1|x)}[D_{\text{KL}}(q_\phi(z_2|z_1)||p(z_2))] \\
 &\quad - \mathbb{E}_{q_\phi(z_2|z_1)}[D_{\text{KL}}(q_\phi(z_1|x)||p_\theta(z_1|z_2))] \quad (2)
 \end{aligned}$$

279 In the equation above, the first term is the reconstruction term
 280 as it is the expectation of the log-likelihood for the input sam-
 281 ples under the distribution estimated from the encoded z_1 , the
 282 second term is the prior matching term which is enforcing the
 283 encoded z_2 to conform the prior distribution $z_2 \sim \mathcal{N}(0, 1)$,
 284 and the last term is the consistency term which requires z_1
 285 from either the encoder or the decoder to be consistent. In
 286 practice, we find that adding the reconstruction loss from z_2
 287 to x is also crucial for generating meaningful samples. Op-
 288 tionally, perceptual losses such as the VGG loss [Ledig *et*
 289 *al.*, 2017] or the structural similarity index measure (SSIM)
 290 [Wang *et al.*, 2004] loss can be used to promote the fidelity of
 291 reconstructed images. However, no significant improvement
 292 is observed in our experiments. Therefore no perceptual loss

293 is included in our implementation. The final loss function for
294 the hierarchical and branching FissionVAE then becomes,

$$\mathcal{L} = \mathbb{E}_{q_\phi(z_1|x)}[D_{\text{KL}}(q_\phi(z_1|z_x)||p(z_1))] \\ - \mathbb{E}_{q_\phi(z_2|z_1)}[\log p_\theta(x|z_1, z_2)] - \text{ELBO}_H \quad (3)$$

295 Here we minimize the KL divergence for z_1 only when the
296 prior distribution for z_1 is explicitly defined, otherwise the
297 model learns the latent distribution by itself.

298 The proposed hierarchical FissionVAE also allows hetero-
299 geneous decoder architectures for each client groups, as each
300 decoder branch is trained and aggregated independently. This
301 flexibility is particularly advantageous in federated learning
302 environments, where clients often possess varying computa-
303 tional resources. Client groups with more resources can im-
304 plement deeper and more complex network structures, while
305 groups with limited computational capacity can utilize lighter
306 models.

307 **Complexity of FissionVAE** FissionVAE’s space complex-
308 ity grows linearly with the number of clients, due to group-
309 specific decoder branches. Time complexity per client fol-
310 lows standard feedforward model training. While we use
311 smaller batch sizes to encourage better latent space explo-
312 ration, this does not change asymptotic complexity.

313 3 Experiments

314 3.1 Datasets and Evaluation Metrics

315 We evaluated the proposed federated VAEs using two com-
316 posite datasets. Mixed MNIST combines MNIST [LeCun
317 and Cortes, 2010] and FashionMNIST [Xiao *et al.*, 2017], di-
318 viding samples into two client groups (one per dataset) with
319 10 clients each. Training samples were evenly distributed
320 within each group, and the default test sets served as evalua-
321 tion benchmarks. An equal number of images were generated
322 using the global model for comparison.

323 CHARM is a more diverse dataset combining five domains:
324 Cartoon faces [Churchill, 2019], Human faces [Karras *et al.*,
325 2018], Animals [Xian *et al.*, 2019], Remote sensing images
326 [Helber *et al.*, 2019], and Marine vessels [Gundogdu *et al.*,
327 2016], using preprocessed square images from Meta-Album
328 for Awa2 and MARVEL. Images were resized to 32×32 ,
329 and each domain was represented by 20 clients, with 20,000
330 images for training and 5,000 for evaluation. As with Mixed
331 MNIST, the global model generated evaluation samples.

332 For Mixed MNIST, encoders and decoders used Multi-
333 Layer Perceptrons (MLPs). On CHARM, encoders $q(z_1|x)$
334 and decoders $p(x|z_1)$ were convolutional, while $q(z_2|z_1)$
335 and $p(z_1|z_2)$ used MLPs. Client participation followed
336 a Bernoulli distribution: $B(0.5)$ for Mixed MNIST and
337 $B(0.25)$ for CHARM. Hyperparameters included learning
338 rates of 1×10^{-3} (Mixed MNIST) and 1×10^{-4} (CHARM),
339 with 70 and 500 training rounds, respectively. Clients per-
340 formed 5 local epochs per round with a batch size of 32. Cen-
341 tralized settings used 70 epochs for Mixed MNIST and 250
342 for CHARM.

343 Evaluation metrics included Fréchet Inception Distance
344 [Heusel *et al.*, 2017] and Inception Score [Salimans *et al.*,

2016] for generation quality, and the negative log-likelihood
(NLL) of the ELBO for reconstruction performance. IS
was computed using an ImageNet-pretrained Inception model
[Szegedy *et al.*, 2016].

349 3.2 Results and Analysis

350 Here we present the following experiments: we first evaluate
351 the overall generative performance of the proposed VAE ar-
352 chitectures in both federated and centralized settings, then we
353 explore strategies for encoding the prior distribution $p(z_1)$,
354 and lastly we showcase the use of heterogeneous decoder ar-
355 chitectures in our FissionVAEs. For experiments investigat-
356 ing different generation pathways of hierarchical VAEs and
357 the effect of reconstruction losses, please refer to our supple-
358 mentary material.

359 Overall Performance

360 The overall performance of the proposed FissionVAE mod-
361 els is summarized in Table 1, and generated examples are
362 shown in Fig. 6. In addition to the FedVAE baseline, a
363 Deep Convolutional GAN (DCGAN) [Radford *et al.*, 2016]
364 trained via FedGAN [Rasouli *et al.*, 2020] is used for com-
365 parison. Since GAN does not directly model the likelihood of
366 data, NLL is not evaluated for FedGAN. Also, FedGAN on
367 CHARM suffers from severe mode collapse, therefore per-
368 formance evaluation is not available on this dataset. Notably,
369 the performance of all models on the CHARM dataset is less
370 robust compared to the Mixed MNIST dataset. This discrep-
371 ancy arises because the CHARM dataset, encompassing RGB
372 images from diverse domains, presents a more complex and
373 realistic federated learning scenario. The dataset’s diversity,
374 coupled with a lower local data availability and participation
375 rate among clients, poses greater challenges to federated gen-
376 erative models.

377 **Latent Space Decoupling vs Decoder Branches** As shown
378 in Table 1, both latent space decoupling and group-specific
379 decoder branches improve image quality (lower FID, higher
380 IS). Decoder branches alone yield larger gains, highlighting
381 the negative impact of mixing decoders trained on non-IID
382 data.

383 FissionVAE+L moderately improves upon FedVAE by par-
384 titioning the latent space by client group, helping the decoder
385 better distinguish domain-specific features and reducing rep-
386 resentation overlap. Fig. 6 shows that while FissionVAE+L
387 enables group-specific sampling, shared decoder aggregation
388 still causes artifacts such as blended features.

389 FissionVAE+D, with a unified encoder and domain-
390 specific decoder branches, greatly reduces visual blending.
391 The encoder functions like a routing module akin to Mixture-
392 of-Experts, which directs inputs to group-specific latent dis-
393 tributions. As decoders remain distinct during aggregation,
394 texture mixing is avoided, producing cleaner outputs (Fig. 6).

395 FissionVAE+L+D combines both latent space decoupling
396 and decoder branches. As shown in Table 1, Fission-
397 VAE+L+D yields marginal gains on Mixed MNIST but out-
398 performs FissionVAE+D on CHARM. Enforcing latent space
399 decoupling yields different outcomes depending on the num-
400 ber of client groups. For Mixed MNIST (2 groups), the FID
401 is lowered due to the extra latent constraints. However, as the

Model	Mixed MNIST						CHARM					
	Federated			Centralized			Federated			Centralized		
	FID ↓	IS ↑	NLL ↓	FID ↓	IS ↑	NLL ↓	FID ↓	IS ↑	NLL ↓	FID ↓	IS ↑	NLL ↓
FedGAN	118.52	2.39	-	91.08	<u>3.18</u>	-	-	-	-	-	-	-
FedVAE	117.03	2.29	<u>0.23</u>	40.59	3.62	0.18	167.18	1.57	40.80	89.26	2.57	46.99
FissionVAE+L	64.99	2.83	0.22	39.27	3.03	0.18	155.81	1.73	43.49	86.19	2.53	51.45
FissionVAE+D	40.78	3.01	0.26	34.76	3.05	<u>0.25</u>	120.39	2.16	<u>33.07</u>	<u>63.25</u>	2.95	36.76
FissionVAE+L+D	<u>42.11</u>	3.04	0.25	<u>34.39</u>	3.08	0.20	<u>109.10</u>	<u>2.27</u>	33.29	50.30	<u>2.89</u>	40.14
FissionVAE+H+L+D	47.72	<u>2.98</u>	0.30	28.82	3.16	0.24	107.69	2.32	27.46	74.59	2.58	27.09

Table 1: Evaluation of proposed FissionVAEs on the Mixed MNIST and CHARM dataset. +L is for decoupled latent space. +D is for branching decoders. +H is for the hierarchical architecture. Best results in are in **bold**. Second best results are underlined. ↑ denotes the higher the better, while ↓ means the lower the better.

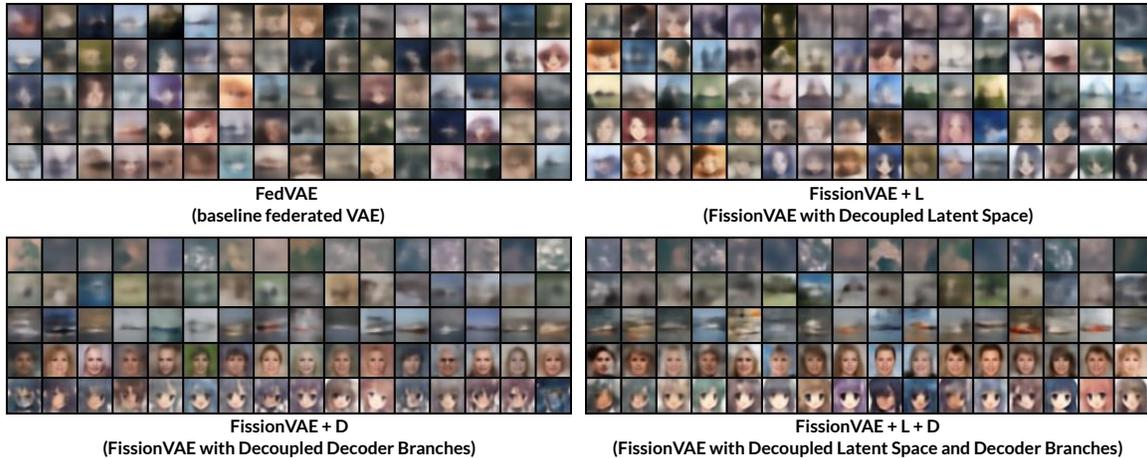


Figure 6: Qualitative results of image generation with FissionVAEs on the CHARM dataset. Best viewed in color.

402 number of client groups increases on CHARM (5 groups),
403 explicit latent space decoupling provides more direct signal
404 to the VAE to identify the intra-group difference, resulting an
405 improved FID. In Fig. 6 it can be observed that images gener-
406 ated by FissionVAE+L+D are sharper than the ones generat-
407 ed by FissionVAE+D.

408 **Hierarchical FissionVAE** As discussed in Section 2, here we
409 consider a hierarchical VAE with two levels of latent vari-
410 able. In Table 1, the architecture FissionVAE+H+L+D per-
411 forms the best on the CHARM dataset and falls behind its
412 non-hierarchical counterpart on the Mixed MNIST dataset.
413 The hierarchical VAE employs multiple levels of latent rep-
414 resentations, which refines the model’s ability to capture and re-
415 construct complex data distributions more faithfully. The per-
416 formance degradation on simpler datasets like Mixed MNIST
417 suggests that the hierarchical approach might introduce un-
418 necessary redundancy without proportional gains in perfor-
419 mance.

420 Decoupling the Prior of z_1

421 Explicitly decoupling the latent space for different client
422 groups improves the ability of VAEs to generate images that
423 align with the true data distribution (Table 1). We explore sev-
424 eral priors for the latent distribution, modeled as multivariate
425 Gaussians with customizable means and identity covariance
426 matrices and evaluate them in Table 2. Details regarding the

Model	Prior $p(z_1)$	Mixed MNIST		CHARM	
		FID ↓	IS ↑	FID ↓	IS ↑
FissionVAE+L+D	identical	40.78	3.01	120.39	2.16
	one-hot	42.01	<u>3.02</u>	113.82	2.25
	symmetrical	<u>41.79</u>	2.95	-	-
	random	43.26	3.00	<u>111.77</u>	2.47
	wave	42.11	3.04	109.10	<u>2.27</u>
FissionVAE+H+L+D	identical	55.91	2.96	122.16	2.30
	one-hot	53.22	<u>2.97</u>	<u>121.33</u>	<u>2.29</u>
	symmetrical	58.21	3.03	-	-
	random	53.99	2.94	124.91	2.23
	wave	<u>53.68</u>	2.94	118.56	2.24
	learnable	47.72	2.98	107.69	2.32

Table 2: Evaluation of Generation Performance with z_1 Priors

427 formal definition of priors can be found in the supplementary
428 material.

429 In non-hierarchical VAEs, z_1 represents the sole latent
430 variable, while in hierarchical VAEs, z_1 is controlled, with
431 z_2 following a standard normal distribution $N(0, 1)$. Base-
432 line priors are identical across client groups. Other prior
433 variations include one-hot encoding, symmetrical positive
434 and negative integers, random vectors, wave encodings (with
435 grouped 1’s in dimensions corresponding to client groups),
436 and a learnable approach unique to hierarchical VAEs. The
437 learnable approach dynamically aligns priors but sacrifices

Decoder Architecture on the FashionMNIST Branch	MNIST			FashionMNIST			Overall		
	FID ↓	IS ↑	NLL ↓	FID ↓	IS ↑	NLL ↓	FID ↓	IS ↑	NLL ↓
Homogeneous	46.73	2.41	0.38	61.81	2.92	0.61	47.72	2.98	0.30
Deeper MLP	49.54	2.38	0.33	60.95	2.90	0.78	48.79	2.95	0.39
Deeper MLP + Conv	<u>48.21</u>	2.38	0.38	65.82	2.99	0.60	50.16	3.00	<u>0.30</u>

Table 3: Evaluation of FissionVAE+H+L+D with Heterogeneous Decoder Architectures on the Mixed MNIST

438 direct sampling from $p(z_1)$.

439 Hierarchical FissionVAE often underperforms non-
440 hierarchical variants when predefined priors are used due to
441 increased uncertainty from additional latent layers. How-
442 ever, the learnable approach excels in capturing complex
443 distributions dynamically. In simpler datasets like Mixed
444 MNIST, identical priors suffice, but explicit latent encoding
445 becomes crucial as client group diversity increases, as seen
446 with CHARM.

447 Among prior definitions, symmetrical priors often lead to
448 divergence on CHARM, as their means may exceed neu-
449 ral network initialization ranges. One-hot and random ap-
450 proaches show comparable results but are less consistent
451 than wave encoding, which clearly distinguishes group pri-
452 ors without out-of-range values.

453 Group-level Privacy

454 In the presence of hierarchical VAEs, it is possible to incorpo-
455 rate the encoder $q_\phi(z_2|z_1)$ into the generation process, that
456 is, we can first sample the latent code z_1 from its prior dis-
457 tribution, then feed it to the subsequent encoder $q_\phi(z_2|z_1)$
458 and the decoders $p_\theta(z_1|z_2)$ and $p_\theta(x|z_1)$ to obtain the syn-
459 thesize a generated sample. On the Mixed MNIST dataset,
460 we observe that swapping the prior distributions of the two
461 client groups in the such a generation pathway leads to ev-
462 ident mode collapse, shown in Figure 7. This suggests that
463 the group-level privacy may be preserved by maintaining the
464 confidentiality of prior distributions. This strategy ensures
465 that high-quality samples are generated only when the cor-
466 rect prior distribution is used, while mismatched distributions
467 yield unrecognizable outputs. This phenomenon is more pro-
468 nounced in both hierarchical and non-hierarchical Fission-
469 VAEs on the Mixed MNIST dataset than on the CHARM
470 dataset, likely due to the simpler, more uniform nature of the
471 Mixed MNIST data compared to the diverse and colorful im-
472 age types in CHARM, which pose greater challenges in sat-
473 isfying complex latent distribution constraints. Evaluation on
474 other generation pathways are presented in the supplementary
475 material.

476 Heterogeneous Decoders in FissionVAE

477 As discussed in Section 2, the decoupling of decoders for
478 client groups allow for the use of heterogeneous architectures
479 in FissionVAE. The Mixed MNIST dataset, with its relatively
480 simple and grayscale colors, can be generated from both fully
481 connected (MLP) and convolutional layers. In contrast, the
482 more complex and colorful images in the CHARM dataset
483 predominantly require convolutional layers for effective gen-
484 eration.

485 Table 3 details the performance evaluation of various de-
486 coder architectures. The term ‘homogeneous’ refers to iden-

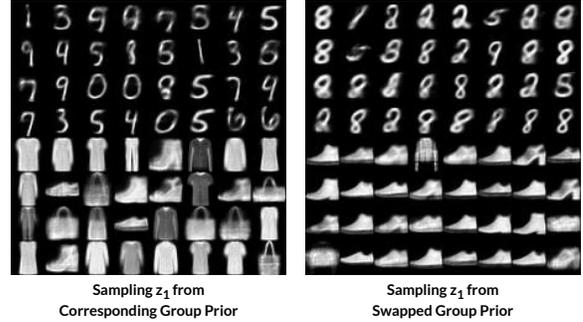


Figure 7: In hierarchical FissionVAE, when the prior distribu-
tion $p(z_1)$ of the MNIST and FashionMNIST groups are swapped,
the generation pathway $q(z_1) \rightarrow q_\phi(z_2|z_1) \rightarrow p_\theta(z_1|z_2) \rightarrow$
 $p_\theta(x|z_1)$ leads to severe mode collapse, suggesting potential group-
level privacy preserving through protected prior distribution.

487 tical architectural configurations across all decoder branches,
488 namely a three-layer MLP for each decoder modules. In the
489 ‘Deeper MLP’ configuration, we add two additional fully
490 connected layers to both $p_\theta(z_1|z_2)$ and $p_\theta(x|z_1)$. Mean-
491 while, we completely replace the decoder $p_\theta(x|z_1)$ from
492 MLP to a series of transpose convolution layers in the ‘Deeper
493 MLP + Conv’ configuration. The results indicate a gradual re-
494 duction in overall FID scores as the decoder architecture be-
495 comes more heterogeneous. However, the integration of con-
496 volutional layers does not improve generation performance
497 over the MLP models, underscoring that while heterogeneous
498 architectures are feasible, they can disrupt the convergence of
499 the VAE due to mismatches in architecture and the model’s
500 weight space.

501 4 Conclusion

502 We presented FissionVAE, a generative model for federated
503 image generation in non-IID data settings. By decoupling the
504 latent space and employing group-specific decoder branches,
505 FissionVAE enhances generation quality while preserving
506 the distinct features of diverse data subsets. Experiments
507 on Mixed MNIST and CHARM datasets demonstrated signifi-
508 cant improvements over baseline federated VAE models,
509 with heterogeneous decoder branches and wave-encoded pri-
510 ors proving particularly effective. Future work includes im-
511 proving the stability of heterogeneous decoder branches, en-
512 abling cross-modality data generation, and developing scal-
513 able strategies for handling an increasing number of client
514 groups in real-world federated learning scenarios.

515 Acknowledgments

516 This work is supported by the EPSRC National Edge AI Hub
517 (EP/Y007697/1).

518 References

- 519 [Bohacek and Farid, 2023] Matyas Bohacek and Hany Farid.
520 Nepotistically trained generative-ai models collapse, 2023.
- 521 [Borsos *et al.*, 2023] Zalán Borsos, Raphaël Marinier,
522 Damien Vincent, Eugene Kharitonov, Olivier Pietquin,
523 Matt Sharifi, Dominik Roblek, Olivier Teboul, David
524 Grangier, Marco Tagliasacchi, and Neil Zeghidour. Audi-
525 olm: a language modeling approach to audio generation.
526 In *arXiv*, 2023.
- 527 [Chen and Vikalo, 2023] Huancheng Chen and Haris Vikalo.
528 Federated learning in non-iid settings aided by differen-
529 tially private synthetic data. In *CVPRW*, 2023.
- 530 [Churchill, 2019] Spencer Churchill. Anime face dataset,
531 2019.
- 532 [Goodfellow *et al.*, 2014] Ian J. Goodfellow, Jean Pouget-
533 Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley,
534 Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Gen-
535 erative adversarial networks. In *NIPS*, 2014.
- 536 [Gundogdu *et al.*, 2016] Erhan Gundogdu, Berkan Solmaz,
537 Veysel Yucesoy, and Aykut Koc. Marvel: A large-scale
538 image dataset for maritime vessels. In *Asian Conference*
539 *on Computer Vision*, 2016.
- 540 [Hardy *et al.*, 2019] Corentin Hardy, Erwan Le Merrer, and
541 Bruno Sericola. Md-gan: Multi-discriminator generative
542 adversarial networks for distributed datasets. In *2019*
543 *IEEE International Parallel and Distributed Processing*
544 *Symposium (IPDPS)*, 2019.
- 545 [Heinbaugh *et al.*, 2023] Clare Elizabeth Heinbaugh, Emilio
546 Luz-Ricca, and Huajie Shao. Data-free one-shot federated
547 learning under very high statistical heterogeneity. In *ICLR*,
548 2023.
- 549 [Helber *et al.*, 2019] Patrick Helber, Benjamin Bischke, An-
550 dreas Dengel, and Damian Borth. Eurosat: A novel dataset
551 and deep learning benchmark for land use and land cover
552 classification. *IEEE Journal of Selected Topics in Applied*
553 *Earth Observations and Remote Sensing*, 2019.
- 554 [Heusel *et al.*, 2017] Martin Heusel, Hubert Ramsauer,
555 Thomas Unterthiner, Bernhard Nessler, and Sepp Hochre-
556 iter. Gans trained by a two time-scale update rule converge
557 to a local nash equilibrium. In *NIPS*, 2017.
- 558 [Ho *et al.*, 2020] Jonathan Ho, Ajay Jain, and Pieter Abbeel.
559 Denoising diffusion probabilistic models, 2020.
- 560 [Jiang *et al.*, 2023] Yuchen Jiang, Ying Wu, Shiyao Zhang,
561 and James J.Q. Yu. Fedvae: Trajectory privacy preserving
562 based on federated variational autoencoder. In *IEEE 98th*
563 *Vehicular Technology Conference (VTC2023-Fall)*, 2023.
- 564 [Karras *et al.*, 2018] Tero Karras, Timo Aila, Samuli Laine,
565 and Jaakko Lehtinen. Progressive growing of gans for im-
566 proved quality, stability, and variation. In *ICLR*, 2018.
- [Karras *et al.*, 2020] Tero Karras, Miika Aittala, Janne Hell- 567
sten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. 568
Training generative adversarial networks with limited 569
data. In *NIPS*, 2020. 570
- [Kingma and Welling, 2014] Diederik P Kingma and Max 571
Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 572
- [Kingma *et al.*, 2016] Diederik P. Kingma, Tim Salimans, 573
Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max 574
Welling. Improving variational inference with inverse au- 575
toregressive flow. In *NIPS*, 2016. 576
- [LeCun and Cortes, 2010] Yann LeCun and Corinna Cortes. 577
MNIST handwritten digit database, 2010. 578
- [Ledig *et al.*, 2017] Christian Ledig, Lucas Theis, Ferenc 579
Huszár, Jose Caballero, Andrew Cunningham, Alejandro 580
Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, 581
Zehan Wang, and Wenzhe Shi. Photo-realistic single im- 582
age super-resolution using a generative adversarial net- 583
work. In *CVPR*, 2017. 584
- [Li *et al.*, 2025] Zhiwei Li, Guodong Long, Tianyi Zhou, 585
Jing Jiang, and Chengqi Zhang. Personalized feder- 586
ated collaborative filtering: A variational autoencoder ap- 587
proach. In *AAAI*, 2025. 588
- [McMahan *et al.*, 2023] H. Brendan McMahan, Eider 589
Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera 590
y Arcas. Communication-efficient learning of deep net- 591
works from decentralized data, 2023. 592
- [Polato, 2021] Mirko Polato. Federated variational autoen- 593
coder for collaborative filtering. In *2021 International*
594 *Joint Conference on Neural Networks*, 2021. 595
- [Radford *et al.*, 2016] Alec Radford, Luke Metz, and 596
Soumith Chintala. Unsupervised representation learning 597
with deep convolutional generative adversarial networks. 598
In *ICLR*, 2016. 599
- [Rasouli *et al.*, 2020] Mohammad Rasouli, Tao Sun, and 600
Ram Rajagopal. Fedgan: Federated generative adversarial 601
networks for distributed data. In *arXiv:2006.07228*, 2020. 602
- [Salimans *et al.*, 2016] Tim Salimans, Ian Goodfellow, Wo- 603
jciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, 604
and Xi Chen. Improved techniques for training gans. In 605
NIPS, 2016. 606
- [Shumailov *et al.*, 2024] Ilia Shumailov, Zakhar Shumaylov, 607
Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin 608
Gal. Ai models collapse when trained on recursively gen- 609
erated data. *Nature*, 2024. 610
- [Szegedy *et al.*, 2016] Christian Szegedy, Vincent Van- 611
houcke, Sergey Ioffe, Jonathon Shlens, and Zbigniew 612
Wojna. Rethinking the inception architecture for computer 613
vision. In *CVPR*, 2016. 614
- [Sønderby *et al.*, 2016] Casper Kaae Sønderby, Tapani 615
Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole 616
Winther. Ladder variational autoencoders. In *NIPS*, 2016. 617
- [Touvron *et al.*, 2023] Hugo Touvron, Thibaut Lavril, Gau- 618
tier Izacard, Xavier Martinet, Marie-Anne Lachaux, Tim- 619
othée Lacroix, Baptiste Rozière, Naman Goyal, Eric Ham- 620
bro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, 621

- 622 Edouard Grave, and Guillaume Lample. Llama: Open and
623 efficient foundation language models. In *arXiv*, 2023.
- 624 [Wang *et al.*, 2004] Zhou Wang, A.C. Bovik, H.R. Sheikh,
625 and E.P. Simoncelli. Image quality assessment: from error
626 visibility to structural similarity. *IEEE Transactions on*
627 *Image Processing*, 2004.
- 628 [Xian *et al.*, 2019] Yongqin Xian, Christoph H. Lampert,
629 Bernt Schiele, and Zeynep Akata. Zero-shot learning -
630 a comprehensive evaluation of the good, the bad and the
631 ugly. *IEEE Transactions on Pattern Analysis and Machine*
632 *Intelligence*, 2019.
- 633 [Xiao *et al.*, 2017] Han Xiao, Kashif Rasul, and Roland
634 Vollgraf. Fashion-mnist: a novel image dataset
635 for benchmarking machine learning algorithms. In
636 *arXiv:1708.07747*, 2017.
- 637 [Xiong *et al.*, 2023] Zuobin Xiong, Wei Li, and Zhipeng
638 Cai. Federated generative model on multi-source hetero-
639 geneous data in iot. In *AAAI*, 2023.
- 640 [Yonetani *et al.*, 2019] Ryo Yonetani, Tomohiro Takahashi,
641 Atsushi Hashimoto, and Yoshitaka Ushiku. Decentralized
642 learning of generative adversarial networks from non-iid
643 data. In *CVPR Workshop on Challenges and Opportuni-*
644 *ties for Privacy and Security*, 2019.
- 645 [Zhang *et al.*, 2024] Lu Zhang, Qian Rong, Xuanang Ding,
646 Guohui Li, and Ling Yuan. Efvae: Efficient federated vari-
647 ational autoencoder for collaborative filtering. In *Proceed-*
648 *ings of the 33rd ACM International Conference on Infor-*
649 *mation and Knowledge Management*, 2024.