

---

# Reproducing "Identifying through flows for recovering latent representations"

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 The authors claim to introduce a model for recovering the latent representation of observed data, that outperforms  
4 the state-of-the-art method for this task; namely the iVAE model. They claim that iFlow outperforms iVAE both in  
5 preservation of the original geometry of source manifold and correlation per dimension of the latent space.

### 6 **Methodology**

7 To reproduce the results of the paper, the main experiments are reproduced and the figures are recreated. To do so, we  
8 largely worked with code from the repository belonging to the original paper. We added plotting functionality as well as  
9 various bug fixes and optimisations. Additionally, attempts were made to improve the iVAE by making it more complex  
10 and fixing a mistake in its implementation. We also tried to investigate possible correlation between the structure  
11 of the dataset and the performance. All code used is publicly available at [https://github.com/HiddeLekanne/  
12 Reproducibility-Challenge-iFlow](https://github.com/HiddeLekanne/Reproducibility-Challenge-iFlow).

### 13 **Results**

14 The obtained mean and standard deviation of the MCC over 100 seeds are within 1 percent of the results reported in  
15 the paper. The iFlow model obtained a mean MCC score of 0.718 (0.067). Efforts to improve and correct the baseline  
16 increased the mean MCC score from 0.483 (0.059) to 0.556 (0.061). The performance, however, remains worse than the  
17 performance of iFlow, further supporting the authors' claim that the iFlow implementation is correct and more effective  
18 than iVAE.

### 19 **What was easy**

20 The GitHub repository associated with the paper provided most necessary code and ran with only minor changes. The  
21 code included all model implementations and data generation. The script that was used to obtain results was provided,  
22 which allowed us to determine which exact hyperparameters were used with experiments on the iFlow models. Overall,  
23 the code was well organised and the structure was easy to follow.

### 24 **What was difficult**

25 The specific versions of the Python libraries used were unknown, which made it infeasible to achieve the exact results  
26 from the paper when running on the same seeds. The code used to create figures 1-3 in the original paper was missing  
27 and had to be recreated. Furthermore, the long time the models needed to train made experimentation with e.g., different  
28 hyperparameters challenging. Finally, the code was largely undocumented.

### 29 **Communication with original authors**

30 Communication with the authors was attempted but could not be established.

# 31 1 Introduction

32 Nowadays, different types of deep generative models excel at generating new data by either explicitly or implicitly  
33 modelling the distribution of the training data. However, sometimes it is useful to recover the distribution that generated  
34 the observed data, i.e. the latent distribution, rather than the data distribution itself. It is easy to see that this is a more  
35 difficult task due to the unknown relation between the unobserved latent variables and the observed data. The concept  
36 of recovering the true latent distribution underlying the data is a form of *identifiability*.

37 Some research has been done in this area. Previously, models (notably  $\beta$ -VAE [1] and its variations) were created with  
38 the purpose of creating *disentangled* representations, where single latent units correspond to single generative factors.  
39 While related to identifiability, such models do not provide any proof or guarantee that they can recover the true latent  
40 representations.

41 More recently, an identifiable variation of the VAE called iVAE was proposed [4], which uses a factorised prior  
42 conditioned on an auxiliary variable to guarantee a basic form of identifiability. In practice however, the fact that this  
43 model optimises a lower bound on the posterior, rather than the actual posterior, could negatively affect the capability of  
44 the model to recover the true latent variables.

45 The paper "Identifying through flows for recovering latent representations" proposes *iFlow*, a model that aims to  
46 alleviate these problems by using Normalising Flow models rather than VAEs [8]. The fact that Normalising flows  
47 model exact distributions rather than approximating the posterior could make them more suitable for this task.

## 48 2 Scope of reproducibility

49 In this review, the work of the proposed iFlow model by Li et al. [8] is reproduced and examined. The aim is to  
50 reproduce the results obtained by the authors and to investigate the claims made in the paper. The claims made can be  
51 seen below. Each claim will be examined in a corresponding subsection in section 4.

- 52 1. Simulations on synthetic data validate the correctness and effectiveness of the proposed iFlow method and  
53 demonstrate its practical advantages over other existing methods.
- 54 2. iFlow outperforms iVAE in identifying the original sources while preserving the original geometry of source  
55 manifold.
- 56 3. iFlow exhibits much stronger correlation than iVAE does in each single dimension of the latent space.
- 57 4. Making iVAE more expressive does not help it approximate the real latent space further, justifying the  
58 discrepancy in parameters.

## 59 3 Methodology

60 Most of the original source code was available and used to test the reproducibility of the paper which can be found in  
61 the corresponding GitHub repository <sup>1</sup>. This repository itself contained code from the repository of the iVAE model<sup>2</sup>  
62 and *nflows*<sup>3</sup>.

63 The iFlow implementations were used largely as is, while the iVAE implementation was refactored to apply modifications  
64 more easily. The *nflows* code base has been removed from the repository and imported as a library instead. Furthermore,  
65 some small optimisations were made to make certain functions more efficient by vectorising them. For reproducing the  
66 results the models were trained on a GPU (see section 3.5). The code for creating the visualisations was not included in  
67 the repository and was therefore recreated. The implementation was made using the PyTorch and NumPy libraries with  
68 Python 3.7.9. TensorBoard was used for logging of variables during training.

### 69 3.1 Model descriptions

70 The paper compares two models: the proposed iFlow model and the iVAE model.

71 The iFlow model is a variation on the Normalising Flow model *rational-quadratic neural spline flows (featuring*  
72 *autoregressive layers)* (RQ-NSF (AR)) [2], where the prior has been replaced with a factorised exponential prior

---

<sup>1</sup><https://github.com/MathsXDC/iFlow>

<sup>2</sup><https://github.com/siamakz/iVAE>

<sup>3</sup><https://github.com/bayesiains/nflows>

73 distribution conditioning the latent variables  $z$  on auxiliary variables  $u$  to obtain identifiability up to an equivalence  
74 relation. The natural parameters of the prior are obtained through a trainable multi-layer perceptron (MLP) which takes  
75 the auxiliary variable  $u$  as input. Each iFlow model contains approximately 3 million trainable parameters.

76 The iVAE model is implemented as an extension of vanilla VAE models [6], using MLPs for both the encoder and  
77 decoder. The number of layers, hidden dimensions and activation functions are hyperparameters. The encoder uses two  
78 MLPs (one for mean and variance each), while the decoder uses just one. Additionally, the prior mapping the auxiliary  
79 variables  $u$  to the latent variables  $z$  is also implemented as an MLP. In total, the iVAE model with standard parameters  
80 has roughly 18,000 trainable parameters.

81 There is a significant difference in the complexities of iFlow and iVAE, seen in the number of trainable parameters the  
82 models have. The authors argue that this is not the cause of the inferior performance of the iVAE, showing that adding  
83 more layers/increasing the hidden dimensions of the model does not increase performance. However, only a limited  
84 range of parameters were used for this, resulting in only weak evidence to support the claim that the comparison is  
85 fair. We further investigate this claim by scaling up the complexity of iVAE through various methods, namely adding  
86 residual connections and layer normalisation in addition to changing hyperparameters.

87 When looking at the implementation of the iVAE model, there appears to be a difference with the theory: the mean of  
88 the prior distribution is not a function of the auxiliary variables  $u$ , as the theory states, but simply fixed to be 0 at all  
89 times. We aim to incorporate this change into the implementation of iVAE to see if it leads to better performance.

## 90 3.2 Dataset

91 A synthetic dataset is required in order to truly know the underlying latent distribution, which is necessary for quantitative  
92 analysis of the performance. The authors chose to use a dataset consisting of sources of non-stationary Gaussian  
93 time-series. Such data was previously used to introduce time-contrastive learning as a means of achieving identifiable  
94 non-linear independent component analysis [3] and was additionally used to assess the performance of the iVAE [4].

95 The latent representation (source) is created as non-stationary Gaussian time series. This data consists of  $M$  segments,  
96 which are modelled as Gaussian distributions with different, randomly selected mean and variance. The means are  
97 sampled from uniform distribution  $[-5, 5]$ , while the variances are sampled from uniform distribution  $[0.5, 3]$ . Each  
98 segment contains  $L$  samples drawn from the corresponding distribution of segment  $M$ . The segment labels serve as the  
99 auxiliary variables  $u$ .

100 A 3 layer invertible MLP is used to transform the samples in a non-linear manner to obtain the observable data. The  
101 invertible MLP consists of mixing matrices with the non-linear activation function  $h(x) = \tanh(x) + \alpha \cdot x$ . The last  
102 layer does not contain a non-linear activation function. Due to the constraints of Flow models, the dimensionality  $d$  of  
103 these observed data points has to be the same as the dimensionality of the latent representation  $n$ . The data generator  
104 allows for the addition of noise to the data points, but this is not utilised.

105 In the paper, results are reported on a dataset created using  $M = 40$ ,  $L = 1000$ ,  $n = d = 5$  and  $\alpha = 0.1$ . For  
106 visualisations of the sources and the estimations of the models,  $M = 5$ ,  $L = 1000$ ,  $n = d = 2$  and  $\alpha = 0.1$  are used.  
107 This differs from the reported  $M = 40$  from the original paper where the figure indicates that the true  $M = 5$ .

## 108 3.3 Hyperparameters

109 The authors of the original paper mention specific values for some of the hyperparameters. However, for other  
110 hyperparameters only a range is provided without a clear indication of what values were used for each evaluation.

111 As mentioned before, for generating the data, the parameters  $M = 40$ ,  $L = 1000$ ,  $n = d = 5$  and  $M = 5$ ,  $L = 1000$ ,  
112  $n = d = 2$  were used for experiments and visualisation respectively. A factorised Gaussian distribution was used  
113 as a prior for the source distributions. The means and variances for these distributions were sampled from uniform  
114 distributions  $[-5, 5]$  and  $[0.5, 3]$  respectively. The data was transformed with an invertible MLP of depth 3 with  $\tanh$   
115 activation function and a slope of 0.1.

116 Both the iVAE and iFlow used the same batch size ( $B = 64$ ) and learning rate of 0.001. A learning rate drop factor of  
117 0.25 was used and a learning patience of 10. An Adam optimiser without weight decay and with standard  $\beta$  values (0.9,  
118 0.999) and  $\epsilon$  ( $1e-8$ ) was used [5]. A learning rate scheduler to reduce the learning rate with a factor 0.1 on plateaus  
119 ensured that the learning rate decreased over time.

120 The iFlow models were initialised with a flow length of 10 with 8 bins. The Rational Quadratic Neural Spline Flows  
121 with Autoregressive transforms (RQNSF-AR) was used as flow type. The Softplus activation was exerted on the natural  
122 parameters.

123 To replicate the iVAE baseline, a model with a hidden dimensionality of 50, a latent dimension equal to that of the data  
124 ( $d = n = 5$ ) and 3 layer MLPs with leaky ReLU with  $\alpha = 0.1$  as activation function. These same hyperparameters  
125 were used for the additional experiments.

### 126 3.4 Experimental setup and code

127 The code for this reproducibility review is publicly available at [https://github.com/HiddeLekanne/](https://github.com/HiddeLekanne/Reproducibility-Challenge-iFlow)  
128 `Reproducibility-Challenge-iFlow`. As mentioned earlier, this code consist of a combination of the iFlow  
129 and iVAE codebases (see section 3).

130 The iFlow and iVAE models were trained with 100 different seeds to generate datasets and the aforementioned  
131 hyperparameters. As is standard for these types of models, the iFlow model was trained using negative log likelihood  
132 as a loss, and the iVAE was used using the ELBO as a loss [6]. Model performance was evaluated using the mean  
133 correlation coefficient (MCC) between the original source of the data and the estimated latent variables from the models.

### 134 3.5 Computational requirements

135 All experiments were run on the LISA system<sup>4</sup> provided by the University of Amsterdam. This system provides  
136 multi-core nodes for research projects. A GeForce GTX 1080 Ti was used to train the models.

137 Training of iVAE models for 100 different seeds took approximately 2 hours. Training of a single iFlow model with a  
138 flow length of 10 took approximately 45 minutes. To alleviate some of the computational cost, the 100 models were  
139 trained in two worker nodes instead of one. This totalled to approximately 1.5 days of training per 100 models.

140 Upgrading of computational resources would not garner better results, with respect to time, based on the fact that the  
141 bottleneck for the computations was the speed of a single thread CPU. Attempts were made to improve this performance  
142 but these did not decrease training time.

## 143 4 Results

144 In this section, results from the original paper are recreated. In addition, some further experiments were done of which  
145 the results can also be found below. These additional experiments consist of improving the existing base line proposed  
146 in the paper, as well as exploring the relation between the complexity of the synthetic data and the achieved MCC  
147 scores.

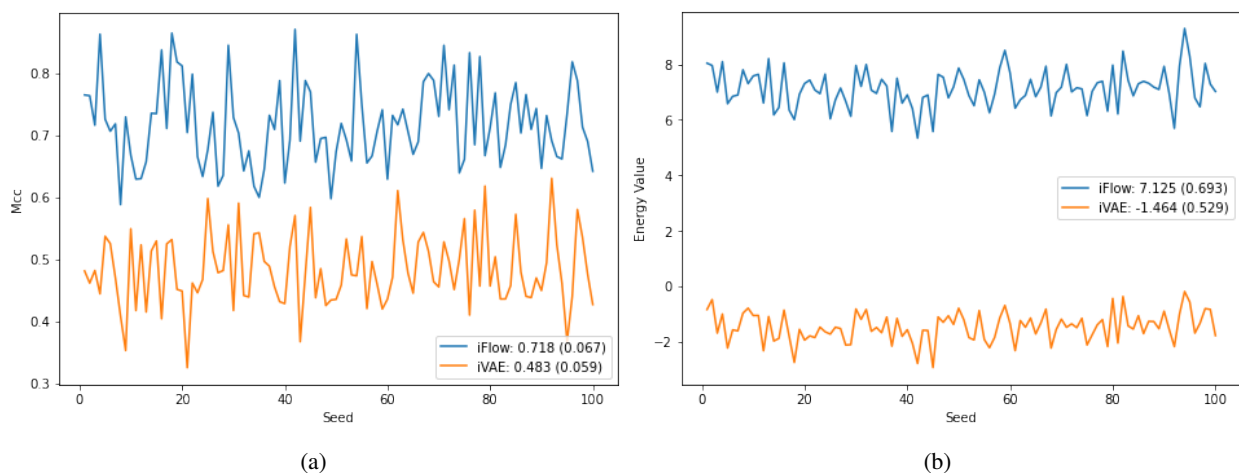


Figure 1: Comparison of identifying performance (MCC) and the energy value (log-likelihood) versus seed number respectively.

<sup>4</sup><https://userinfo.surfsara.nl/systems/lisa>

148 **4.1 Results reproducing original paper**

149 **4.1.1 Comparison of identifying performance**

150 The MCC scores and log-likelihood over 100 seeds are displayed in figure 1a and figure 1b respectively. The figures  
151 show that there is high variance in MCC scores for different datasets. The iFlow models obtained a mean accuracy of  
152 0.718 with a standard deviation of 0.067 whereas the iVAE models obtained a mean accuracy of 0.483 with a standard  
153 deviation of 0.059 which is in compliance with the results produced in the original paper. The results for the iVAE  
154 models are significantly worse than in the original iVAE paper. An improvement to the implementation was made to  
155 better emulate the performance of this paper which resulted in a fairer comparison (see section 4.2.1).

156 As can be seen in figure 1b, the energy values of the iVAE are significantly lower compared to those of iFlow, matching  
157 the results of the paper. The authors noted that the difference in energy values could indicate that the gap between the  
158 ELBO and the actual log likelihood is not negligible.

159 **4.1.2 Preservation of original source manifold geometry**

160 Figure 2 shows the 2D visualisation for different data seeds. The original paper stated that an  $M = 40$  was used but  
161 figures indicated that this should be  $M = 5$ .

162 The results largely support the claim of the author that the original geometry of the source manifold is preserved. The  
163 estimations from the iFlow model seem more similar to the original source than the estimations from the iVAE models,  
164 although it still contain artefacts from the observations. Figure 2a is an example of such where the latent dimensions are  
165 not successfully recovered. In other examples, the original Gaussian distributions are mostly recovered apart from some  
166 transformation as was the case in the original paper. The collapse of the latent space from iVAE models observed in the  
167 original paper was not prevalent during experiments. However, the preservation of the original geometry of the source  
168 manifold is better captured by the iFlow models.

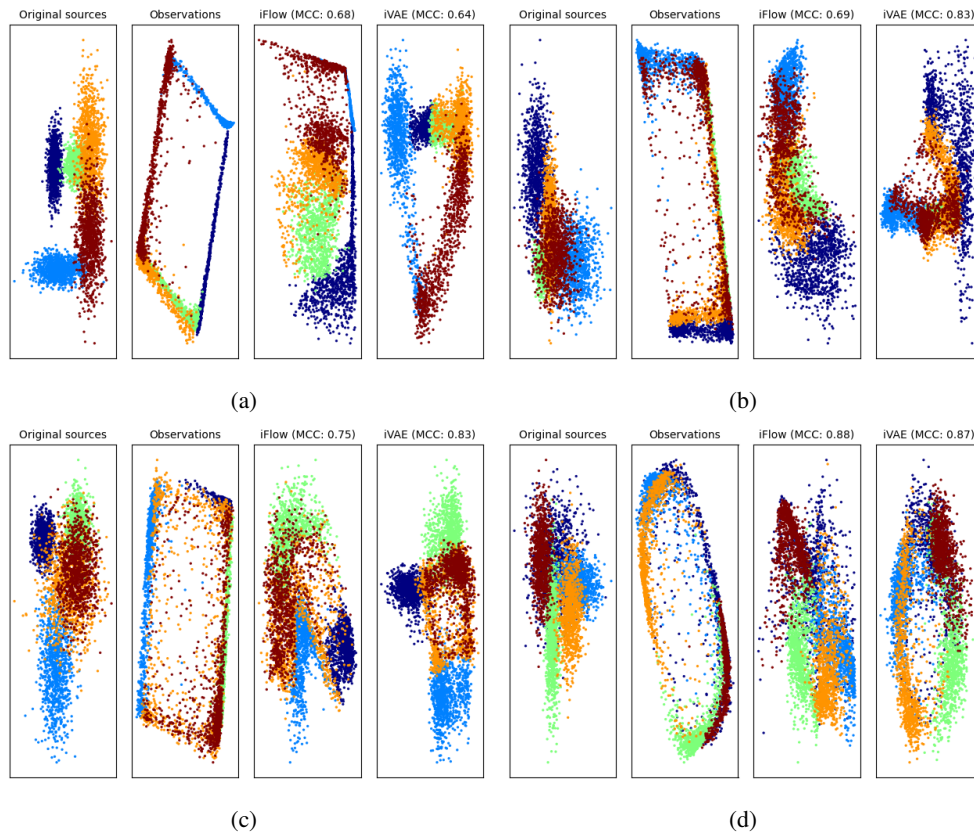


Figure 2: Visualisation of 2D-cases

169 **4.1.3 Separate latent dimension correlation**

170 Figure 3 and figure 4 show the correlation between the source signal used to generate the data and the latent variables  
 171 recovered by the iFlow and iVAE models. Figure 3 shows the results of the best performing iFlow, which we assume  
 172 is what figure 3 of the original paper also depicts. For fairness, we also show the results for the dataset that iVAE  
 173 performed best on in figure 4.

174 These results largely support the claim that iFlow exhibits stronger correlation than does iVAE in each single dimension  
 175 of the latent space: while this is generally the case, it does occur for some datasets that iVAE has a higher correlation  
 176 coefficient than iFlow on one or even two of the latent dimensions, as shown in figure 4.

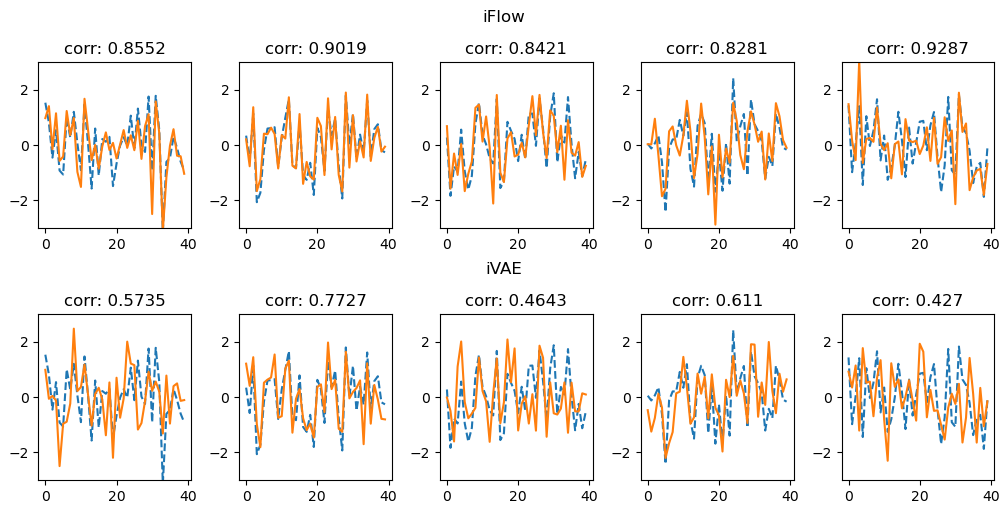


Figure 3: Comparison of the latent variables recovered by the models (orange lines) to the true latent variables (dashed blue lines) for individual dimensions. This figure shows results for the seed that resulted in the best iFlow performance.

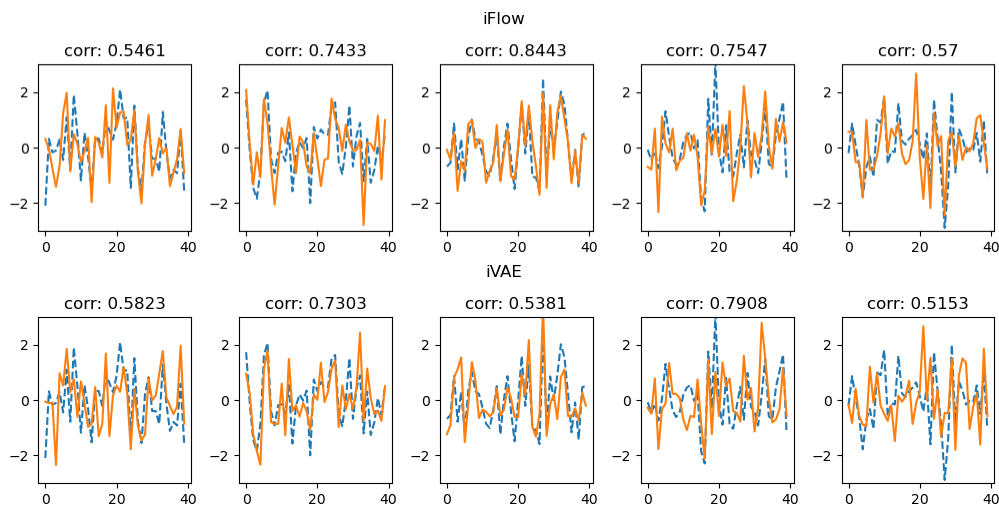


Figure 4: Comparison of the latent variables recovered by the models (orange lines) to the true latent variables (dashed blue lines) for individual dimensions. This figure shows results for the seed that resulted in the best iVAE performance.

## 177 4.2 Results beyond original paper

### 178 4.2.1 Improved baseline

179 In figure 5a and 5b the MCC scores and energy values over 100 seeds are displayed for the iFlow model, iVAE model  
180 and improved iVAE model. The addition of the trainable mean, based on auxiliary parameters, shows an increase in the  
181 mean MCC score from 0.483 (0.059) to 0.556 (0.061). The ELBO score improves almost with a constant value for  
182 every seed.

183 Other attempts had been made to improve this baseline by increasing the complexity of the iVAE model. However,  
184 were tried before the mistake in the iVAE implementation had been noticed, and are therefore not very useful. These  
185 results can be seen in appendix B.

186 In addition to figure 1, recreations of figures 2 3 using the improved baseline were made. These are not included in this  
187 report but can be found in appendix C

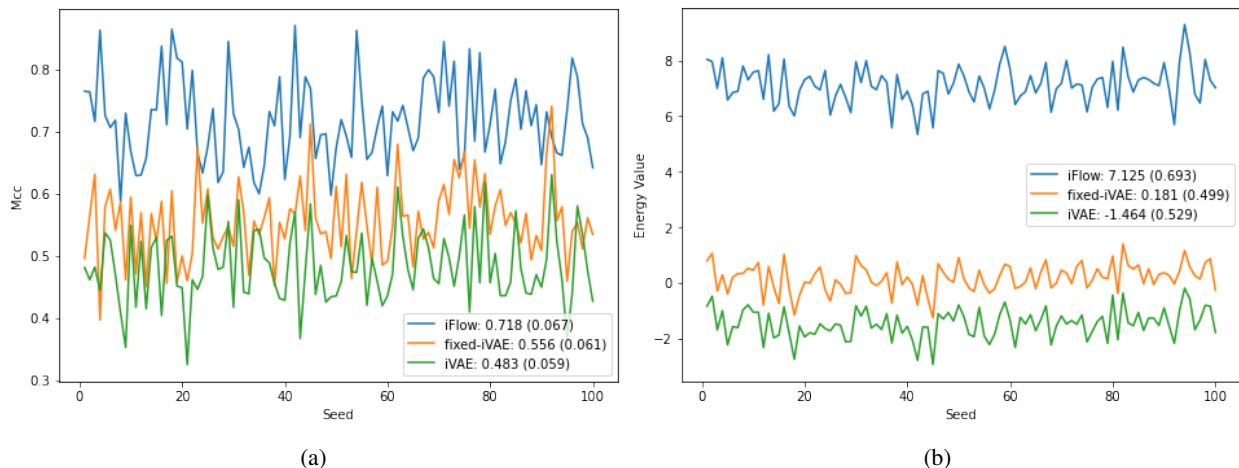


Figure 5: Comparison of identifying performance (MCC) and the energy value (log-likelihood) versus seed number respectively, with the fixed version iVAE included.

### 188 4.2.2 Synthetic data complexity

189 To measure the complexity of the dataset, the mean Kullback-Leibler divergence [7] of each source and its nearest  
190 neighbour was used.

191 This metric showed no correlation (0.06) with the MCC scores obtained by iVAE or iFlow, showing that the randomly  
192 sampled parameters of the source distribution were likely not to blame for the high variance in MCC scores.

## 193 5 Discussion

194 The results reproduced in the previous section largely support the claims by the original authors. Firstly, the MCC  
195 scores that we obtained after training the model on synthetic data are very similar to the ones reported. Secondly, the  
196 recreated visualisation of 2D latent sources seems to support the claim that the iFlow method outperforms iVAE in  
197 identifying the original sources. Finally, the claim that iFlow exhibits much stronger correlation than iVAE in each  
198 single dimension of the latent space is not fully supported by our results. In the original paper, the authors show this  
199 correlation only for the best iFlow results. When visualising the individual dimensions of the latent variables for the  
200 best iVAE results, iVAE outperforms iFlow in two of the latent dimensions. This shows that the claim does not strictly  
201 hold true for all seeds. Nevertheless, since iFlow outperforms even the best iVAE on most latent dimensions, it still  
202 seems to be a reasonable claim.

203 Our experiments to improve the performance of the iVAE model, by modelling the prior means as a trainable function  
204 of the auxiliary variable  $u$ , managed to increase its performance significantly. However, the performance remains worse  
205 than that of iFlow. This further cements the claim that the iFlow model is more suited for the task of identifiability than  
206 iVAE.

207 The strength of our approach was that we were generally faithful to the original implementation, using largely the same  
208 code which we examined thoroughly. Therefore, the chance of implementation differences with the original code is very  
209 small. Additionally, we rigorously compared the code with the underlying theory, allowing us to correct an important  
210 mistake in the baseline.

211 A weakness of our approach was that we did not do any work to examine the models on a more realistic dataset,  
212 meaning the generalisability of the model remains an open question. Furthermore, due to the high variance in the results  
213 of identifying models, all experiments had to be run with a large number of seeds (100), which took a long time given  
214 the fact that training of a single model took approximately 40 minutes. For this reason, experimentation done with  
215 hyperparameters was limited. The experiments in the appendix of the paper were not replicated for similar reasons.  
216 These experiments looked at the effect of different activation functions on the performance of iFlow and the effect of  
217 more and larger hidden layers on the performance of iVAE.

218 Overall, the authors provided a model which outperforms the previously best method for this problem in a quantifiable  
219 measure. Additionally, high variance in the results is addressed appropriately by running the experiments over a large  
220 number of seeds. Furthermore, the visualisation of the true sources and the estimations by the models makes it easier to  
221 interpret the MCC scores. Lastly, the model is theoretically well motivated.

222 Despite these strengths of the original paper, some improvements could be made to further substantiate the claims made  
223 in the paper. There is a clear advantage that iVAE has over iFlow, which is not mentioned by the authors: iVAE can be  
224 used when the dimensionality of the latent sources differs from the data dimensionality, while iFlow cannot. The fact  
225 that iFlow needs data with such corresponding dimensionalities also means that the iVAE had to be trained without a  
226 bottleneck. This is an important part of the VAE architecture, and the lack thereof could have contributed to the weaker  
227 performance of iVAE; compared to the paper introducing iVAE (MCC of above 0.95), the MCC scores of the iVAE  
228 reported by the authors are significantly worse (MCC of 0.496). This discrepancy is not addressed or explained by the  
229 authors.

## 230 5.1 What was easy

231 The code provided in the GitHub repository worked almost out of the box, with only small adjustments needed; the  
232 source code of the nflows library that was included in the repository was replaced with an import. This fixed an issue  
233 that prevented the code from running on a CPU. The code was well organised into separate files for e.g., the iFlow  
234 model, iVAE models or training, making it easy to quickly find specific parts of the code when needed. The code that  
235 generates the data the models are trained on also came with the implementation, and worked without any issues.

236 With the code, a shell script was provided that seems to be the one used for the experiments on iFlow in the paper  
237 (although this was not explicitly stated). This allowed for easy replication of these experiments, with all of the used  
238 hyperparameters provided.

## 239 5.2 What was difficult

240 There were difficulties in replicating some parts of the paper. The lack of a provided environment means that our code  
241 was likely run using different versions of some libraries such as PyTorch or NumPy. This could have contributed to the  
242 difference in outcomes of our experiments compared to the paper while using the same seeds.

243 While the script used to run iFlow experiments was provided, the same was not true for the iVAE experiments. This was  
244 not a large problem, however, since the authors do state that the hyperparameters used are the same as in the original  
245 iVAE paper [4]. The code for creating plots (Figures 1,2,3 in the iFlow paper) was also not provided and additional  
246 code had to be written to recreate these figures.

247 The training of the iFlow models for all 100 seeds took a significant amount of time. With the training for one seed  
248 taking approximately 40 minutes, the full training took roughly a day and a half (running two batches of 50 seeds  
249 simultaneously). This made it difficult to do full-scale experiments with different hyperparameters.

250 Lastly, there was a large portion of unused code present in the repository, which made it more difficult to understand the  
251 overall structure of the code. This includes the source code of the nflows library, code for planar flows, multiple different  
252 iVAE variations, an alternative dataloader, an unused dataset and an implementation of training using annealing.



253 **References**

254 [1] Christopher P Burgess et al. “Understanding disentangling in  $\beta$ -VAE”. In: *arXiv preprint arXiv:1804.03599*  
255 (2018).

256 [2] Conor Durkan et al. “Neural Spline Flows”. In: *Advances in Neural Information Processing Systems*. 2019,  
257 pp. 7511–7522.

258 [3] Aapo Hyvarinen and Hiroshi Morioka. “Unsupervised feature extraction by time-contrastive learning and nonlinear  
259 ica”. In: *arXiv preprint arXiv:1605.06336* (2016).

260 [4] Ilyes Khemakhem et al. “Variational autoencoders and nonlinear ica: A unifying framework”. In: *International  
261 Conference on Artificial Intelligence and Statistics*. PMLR. 2020, pp. 2207–2217.

262 [5] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint  
263 arXiv:1412.6980* (2014).

264 [6] Diederik P Kingma and Max Welling. “Auto-Encoding Variational Bayes”. In: *arXiv preprint arXiv:1312.6114*  
265 (2013).

266 [7] Solomon Kullback and Richard A Leibler. “On information and sufficiency”. In: *The annals of mathematical  
267 statistics* 22.1 (1951), pp. 79–86.

268 [8] Shen Li, Bryan Hooi, and Gim Hee Lee. “Identifying through Flows for Recovering Latent Representations”. In:  
269 *arXiv preprint arXiv:1909.12555* (2019).

270 **A Alternative iFlow model**

271 A comment in the files mentioned an error in the implementation because the Softplus function was applied to  $\xi$  as well  
272 as  $\eta$  from the natural parameters  $\lambda(\mathbf{u})$ . An alternative version of the implementation was also tested where the Softplus  
273 activation function was only exerted on  $\xi$ , as there are no constraints on the sign of  $\eta$ .

274 The results obtained using this method were approximately the same as the original performance of the iFlow. A mean  
275 MCC of 0.72 with a standard deviation of 0.057 was achieved. Because these results were not a significant improvement,  
276 it was decided to include this experiment as an appendix.

277 **B Baseline improvement experiments**

278 The table below shows the results of experiments with changing the iVAE architecture to increase complexity. As  
279 shown, adding skip connections or layer normalisation to the architecture did not increase performance with respect to  
280 the unchanged baseline. Due to the high training time, no additional experiments could be done.

281

addition	NUM_HIDDEN	NUM_LAYERS	AVG MCC
-	50	3	0.483 ( $\pm 0.059$ )
residual connections	50	3	0.474 ( $\pm 0.053$ )
layer normalisation	50	3	0.461 ( $\pm 0.051$ )

282 **C Visualizations for Fixed iVAE**

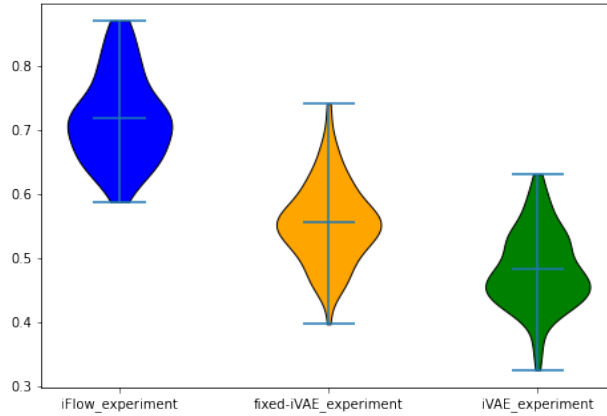


Figure 6: Alternative visualisation of the MCC scores obtained by the models, including the fixed iVAE.

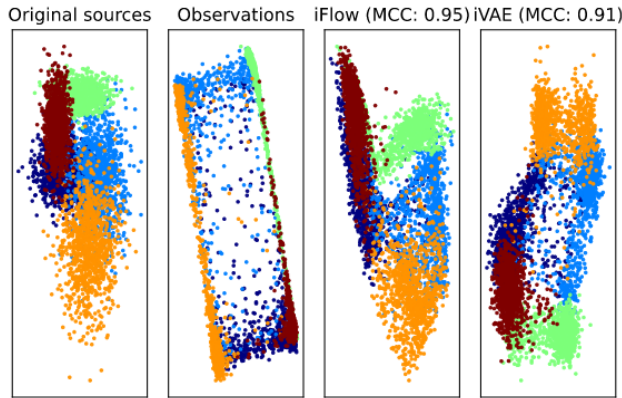


Figure 7: Visualisation of 2D-cases, comparing iFlow to the fixed version of iVAE.

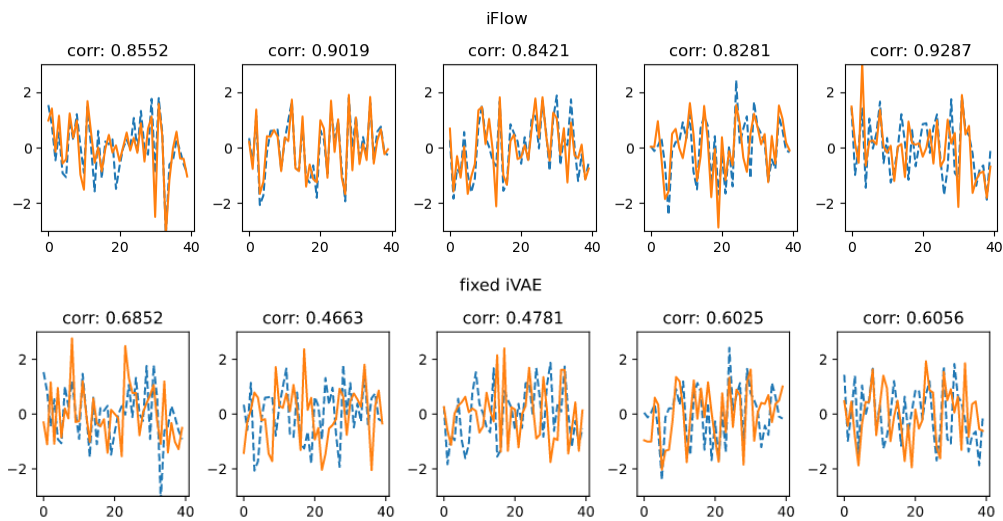


Figure 8: Comparison of the latent variables recovered by the models (orange lines) to the true latent variables (dashed blue lines) for individual dimensions. This figure shows results for the seed that resulted in the best iFlow performance.

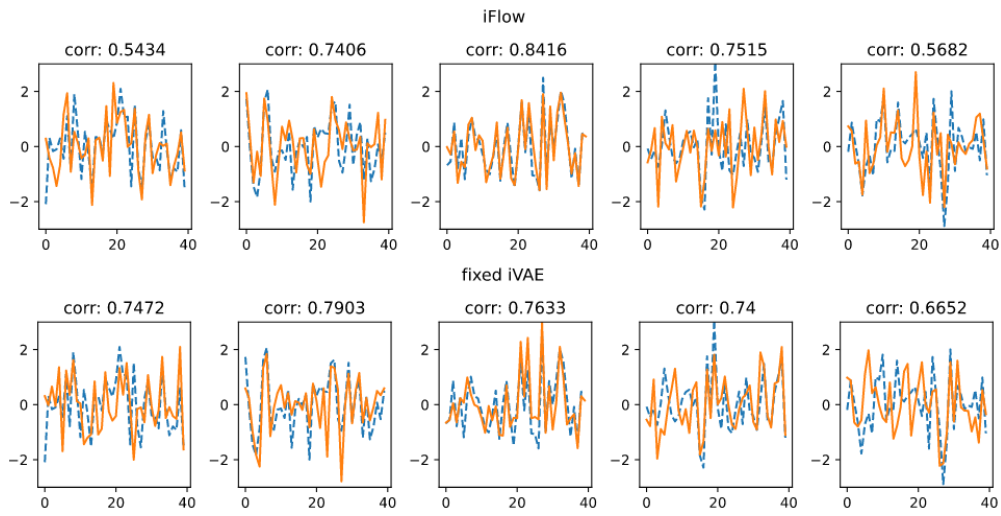


Figure 9: Comparison of the latent variables recovered by the models (orange lines) to the true latent variables (dashed blue lines) for individual dimensions. This figure shows results for the seed that resulted in the best fixed iVAE performance.