

Flow-based Generative Models for Learning Manifold to Manifold Mappings (Supplemental Material)

1 In the main paper, we showed some quantitative results
 2 of generating ODF conditioned on DTI. In this supplement,
 3 we present several additional quantitative results as well as
 4 an example of the heatmap of the p -value of the generated
 5 ODF compared with other measurements including FA, DTI,
 6 and ODF. Also, we provide some discussion regarding the
 7 motivation of the layers we choose and more details about
 8 the preliminary differential geometry concepts so that the
 9 presentation is more self-contained for the reader.

10 Quantitative evaluation for generating ODF 11 from DTI

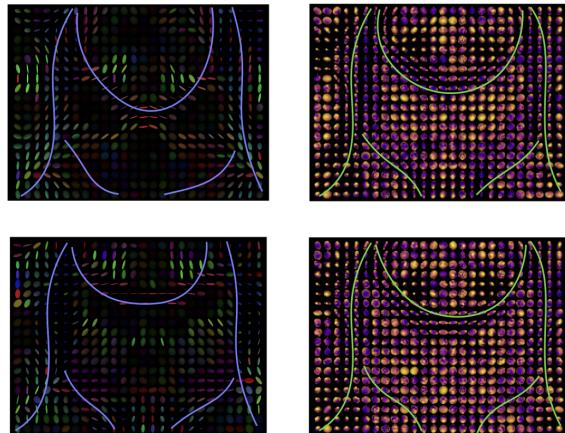


Figure 1: The samples of some fibers in the original DTI slice and in the generated ODF slice. Our model can generate the ODF structures.

12 In this section, we will show additional results of our model
 13 for ODF generation. As mentioned in the paper, the main
 14 focus for our paper is to generate ODF images based on DTI
 15 images. Note that in terms of acquisition time, a typical scan
 16 that may provide ODF data, will take 35 mins compared to
 17 7–12 mins for a typical DTI acquisition. Thus, successfully
 18 using our model to generate the ODF from DTI, e.g., for
 19 statistical analysis on groups, can facilitate a nice trade-off

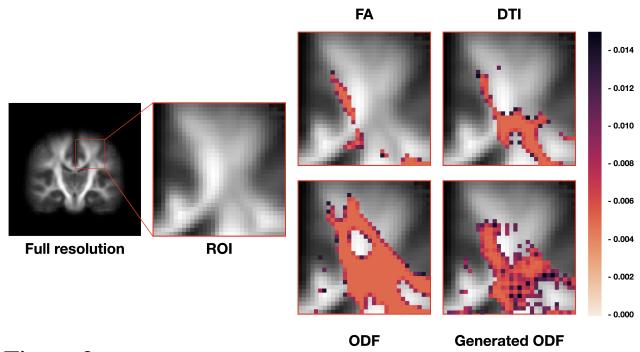


Figure 2: The p -value of one of the ROIs of the entire brain scan with full-resolution. We show that our proposed method can generate meaningful ODF with respect to the group level differences.

20 between the angular resolution, sensitive of statistical anal-
 21 ysis and the actual scanning time. A few examples of our
 22 results were included in the main paper Fig. 6(a). Here, we
 23 show some additional examples of DTI to ODF generation.
 24 In Fig. 3, each “pair” includes the input DTI (first and third
 25 lines) and the generated/synthetic ODF (second and fourth
 26 lines). We can see that the structure of the generated ODF is
 27 well preserved with our model. Similar to the main paper, we
 28 only show 2D slices of the 3D images.

29 Also, for the reader not familiar with diffusion weighted
 30 MR images, we highlight the white matter tracts in Fig. 1.
 31 The white matter bundles are the “path” or “roads/highways”
 32 connecting different gray matter regions (processing centers)
 33 of the brain. In our generation process, the white matter bun-
 34 des are the major Region of Interests (ROIs) that we seek
 35 to reconstruct, since they contain information pertinent to
 36 assessing brain connectivity, which may be different across
 37 clinically disparate groups. As an example, one might see
 38 group-level differences in brain connectivity between male
 39 and female or alternatively, Alzheimer Disease (AD) sub-
 40 jects and healthy controls. Since the HCP dataset that we use
 41 includes participants that are all healthy, perhaps a scientifi-
 42 cally more interesting experiment of evaluating group-level
 43 connectivity differences between healthy and diseased sub-
 44 jects, cannot be performed, and we used male versus female
 45 connectivity analysis for evaluations. For statistical testing
 46 between male and female, we presented the IoU of the entire
 47 brain with low-resolution in the main paper. Here, we show a

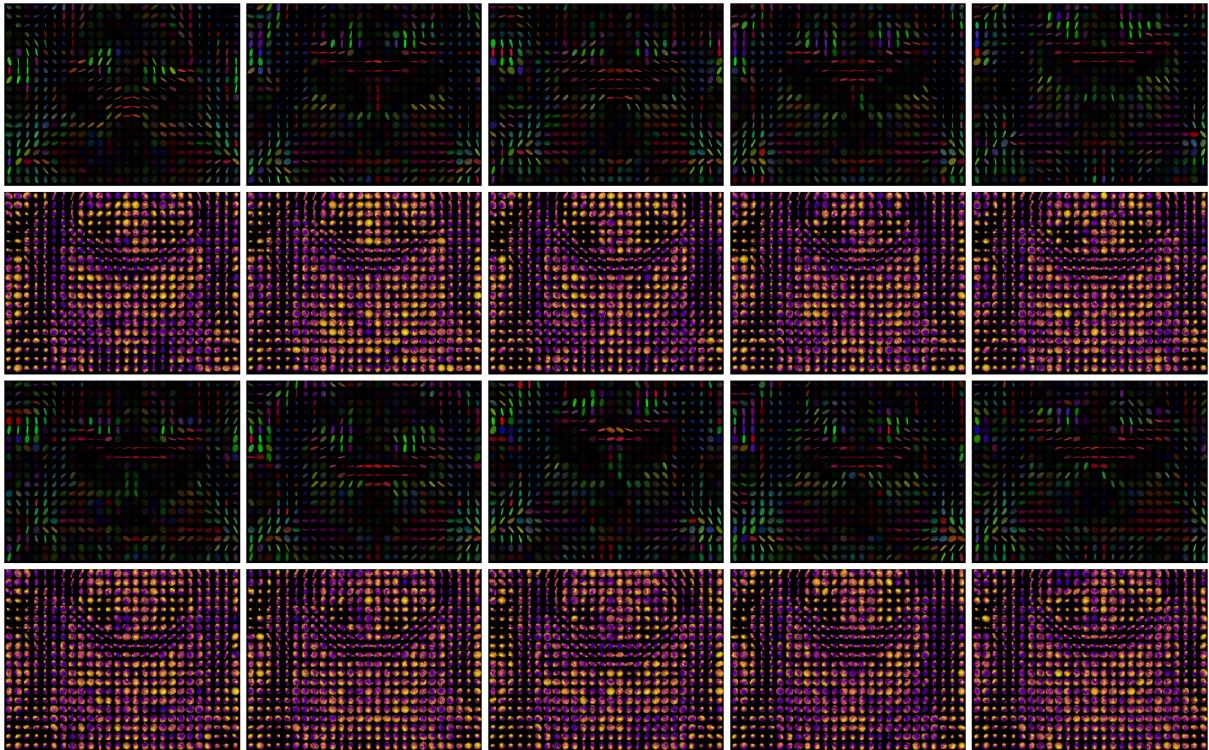


Figure 3: Generated ODF from corresponding DTI. Each pair here contains the input DTI (*top*) and the generated ODF (*bottom*).

zoomed in example of a ROI for the full-resolution images in Fig. 2. The p -values for different ROIs are all < 0.001 in both the original ODF and our generated ODF, indicating consistency of our results, at least in terms of regions identified in downstream statistical analysis. Note that the analysis on the real ODF images serves as the ground truth.

Other potential applications for our model

Since our model basically requires two parallel GLOW models (two-stream GLOW) on two manifolds, our model can be further extended or applied to other applications. We will introduce some additional aspects of our model in this supplement as illustrative examples, noting that these secondary experiments are only to show applicability rather than claiming state of the art performance.

Generate DTI from ODF

In the main paper, we only discussed generating ODF based on DTI data, because normally this is the direction which is meaningful in the real-world acquisitions, i.e., in other words, data from shorter acquisition time (DTI) providing images comparable to those acquired from a longer acquisition time (ODF). But our model imposes no constraint that the direction can only proceed in one direction. If we reverse the direction, we can generate the corresponding DTI based on ODF data. Even though ODF contains more angular information than DTI, the standard analysis is based on functional anisotropy (FA) (scalar-valued) summary images, which is calculated from DTI directly.

Suppose $\lambda_1, \lambda_2, \lambda_3$ are three eigen values of the SPD matrix at a specific voxel.

$$FA = \sqrt{\frac{1}{2} \frac{\sqrt{(\lambda_1 - \lambda_2)^2 + (\lambda_2 - \lambda_3)^2 + (\lambda_3 - \lambda_1)^2}}{\sqrt{\lambda_1^2 + \lambda_2^2 + \lambda_3^2}}}$$

An image consisting of FA values at each voxel can be constructed based on the ODF data.

In Fig. 4, we show some results for generating DTI based on ODF. Each pair contains the input ODF and the corresponding generated DTI. The colors of DTI indicates the eigenvalues of the SPD matrix at each voxel, and the orientation indicates the eigenvectors of SPD. Thus, qualitatively, the generated DTI preserves the DTI field structure fairly well.

Model setup: generating DTI based on ODF Similar to the model described in the main paper, the generation of DTI based on ODF contains two flow-based generation model and a transformation from latent space of ODF to the latent space of DTI. Two generation models remain identical while the transformation model is from ODF to DTI.

Motivation of the three basic layers

In the main paper, we introduced three basic layers: Actnorm, 1×1 convolution, and the Affine Coupling. Here, we will provide additional motivation behind these three layers. We will also include more details about the full structure of the GLOW model and the NanoFlow (Lee, Kim, and Yoon 2020) trick.

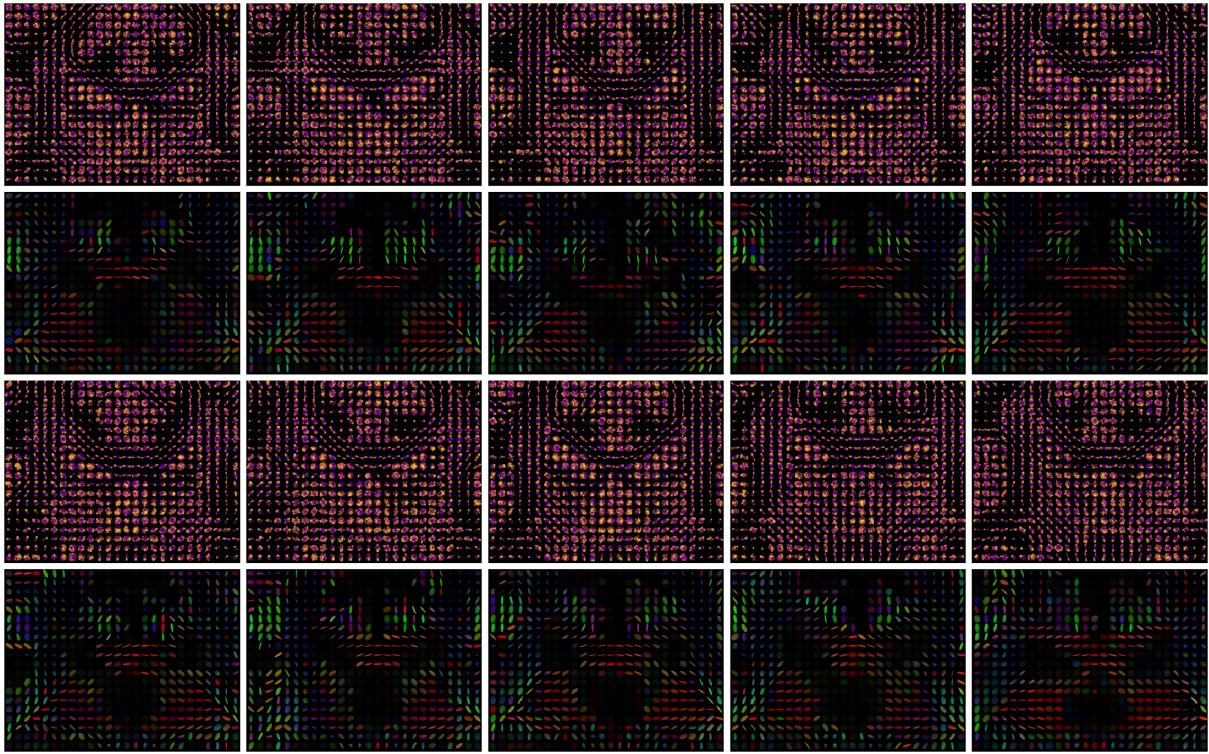


Figure 4: Generated DTI from corresponding ODF. Each pair here contains the input ODF (*top*) and the generated DTI (*bottom*).

97 Actnorm

98 The Actnorm layer plays a similar role as Batch Normalization
 99 ([Ioffe and Szegedy 2015](#)). In a traditional neural network,
 100 e.g., Resnet ([He et al. 2016](#)), batch normalization is used to
 101 stabilize the training procedure and make the distribution
 102 of the input to be simpler. However, in a normalizing-flow
 103 based architecture, one of the biggest drawbacks of the batch
 104 normalization is that it is not invertible. Batch normalization
 105 depends on the data within one mini-batch, which is not avail-
 106 able in the generating phase (inverse of the model). Thus, we
 107 need to modify this operation to be invertible and preserve
 108 stability in training. Thus, in ([Kingma and Dhariwal 2018](#)),
 109 the authors introduced the actnorm in Euclidean space, which
 110 is initialized with the data and then trained unconstrained.
 111 The initialization keeps the input close to the normal distri-
 112 bution. After the initialization, the parameters of the mean
 113 and variance are updated based on the gradient of the loss
 114 function. In our case, the target field is on the manifold. Thus,
 115 we use our modification of “adding” and “multiplication” op-
 116 erations on the manifold valued data to extend the Actnorm
 117 layer to the manifold setting. Mathematically, the order of
 118 the two operations “adding” and “multiplication” does not
 119 matter. But practically, applying “multiplication” first would
 120 reduce the cost to compute the gradient and thus is faster.

121 1×1 convolution

122 The 1×1 invertible convolution is one of the main contrib-
 123 utions in GLOW ([Kingma and Dhariwal 2018](#)). The 1×1
 124 convolution is used to permute the information, channel-wise,
 in the data with learnable weights. The rationale behind this

layer comes from ([Dinh, Sohl-Dickstein, and Bengio 2016](#))
 126 which used a fixed permutation to reverse the ordering of
 127 the channels to make each channel contribute equally. In the
 128 manifold-valued data, we can define the permutation after
 129 the chart map which projects the original manifold-valued
 130 data on to the Euclidean space. To make training more robust
 131 without exploding/vanishing as well as make it easier to com-
 132 pute the determine of Jacobian matrix, we impose an extra
 133 constraint that the kernel of 1×1 convolution is a rotation
 134 matrix. Thus, the determinant of the Jacobian is ± 1 which is
 135 a constant number.
 136

137 Affine Coupling

As the reader will notice, both operations we discussed above
 138 are linear without any activation functions. Indeed, we can
 139 build a linear neural network to be invertible easily (for exam-
 140 ple, a full-rank matrix multiplication). However, we will use
 141 the non-linearity to fit the distribution of the complicated data.
 142 In ([Dinh, Krueger, and Bengio 2014](#)), motivated by Advanced
 143 Encryption Standard (AES), the author introduced an invert-
 144 ible nonlinear layer called Affine Coupling. The goal of the
 145 Affine Coupling layer is to use half of the pixels of the data to
 146 estimate the mean and the variance of other half of the pixels
 147 (in the Euclidean space) and by adding and multiplication to
 148 make it more like a normal distribution. The first half of the
 149 pixels are preserved through this layer so that when inverting,
 150 the network can use the same nonlinear function to estimate
 151 the mean and variance again. In the next layer, two parts of
 152 the pixels are switched to iteratively make the final distri-
 153 bution of *all* pixels to be close to normal distribution. The

155 key functionality of this layer involves adding, multiplication,
 156 and the neural network to estimate the mean and variance.
 157 In the manifold setting, we use the chart map to project the
 158 first half of pixels into the Euclidean space to apply a simpler
 159 neural network (compared with the available neural network
 160 defined on manifold) to estimate the mean and variance of
 161 the other half. Also, we apply our modification of adding and
 162 multiplication on the manifold-valued data so that the final
 163 output will be similar as the normal distribution defined on
 164 the manifold.

165 The full structure of the network

166 In GLOW (Kingma and Dhariwal 2018), the authors designed
 167 a multi-scale architecture which includes split and squeeze
 168 functions. The basic idea of the split function is similar to
 169 Wavelets that divides a given function or continuous-time signal
 170 into different scale components. We call the combination
 171 of these three layers mentioned above as a “block”. After
 172 every few blocks, the model will split the channels of the
 173 data into two pieces equally, and directly pass one of these
 174 two pieces into the latent space. So, the model will learn to
 175 first pay more attention to the overall low-frequency distribution
 176 to make it a normal distribution. After splitting, the
 177 model will try to learn the distribution of the difference (the
 178 remaining part). The squeeze function is to reduce the spatial
 179 resolution and put it to the channel-wise to enable the next
 180 splitting. For example, for the natural images case, the input
 181 of the data will be $128 \times 128 \times 3$. In the first layer, we will
 182 squeeze the spatial resolution to the channel. Thus, it will
 183 become $(128/2) \times (128/2) \times (3 \times 2^2) = 64 \times 64 \times 12$. After
 184 this block, the output will also have the same size as the input,
 185 which is $64 \times 64 \times 12$. At this stage, we split the data into
 186 two pieces and directly pass $64 \times 64 \times (12/2) = 64 \times 64 \times 6$
 187 into the latent space. The remaining $64 \times 64 \times 6$ will again
 188 be squeezed into $(64/2) \times (64/2) \times (6 \times 2^2) = 32 \times 32 \times 24$
 189 and pass it into the next block. In our case, we use exactly
 190 the same idea as the split and squeeze operations. The major
 191 difference is that each pixel is a point on the manifold which
 192 may have different dimensions, 361 for ODF and 6 for DTI.
 193 The channel of each manifold-valued data is always 1 in the
 194 input layer.

195 NanoFlow

196 NanoFlow (Lee, Kim, and Yoon 2020) was introduced to
 197 reduce the number of parameters for sequential data processing.
 198 The assumption of NanoFlow is that the Affine Coupling
 199 layer, if fully trained, can estimate the distribution for any
 200 parts of the input data in a fixed order. There will be some per-
 201 formance drop compared with training different Affine Cou-
 202 pling layers for different parts of the data. But the gain from
 203 reducing the parameters is significant. Thus, in our setup, due
 204 to the large 3D input, we apply the NanoFlow trick for DTI
 205 and ODF separately. For example, for the DTI path, we first
 206 split the data into $2n$ parts called $\{x_1, x_2, \dots, x_{2n}\}$, $n \geq 1$.
 207 The input of two neural networks in Affine Coupling layer
 208 would be the x_{2k-1} , $k = 1, 2, \dots, n$, while the output will be
 209 the mean and variance of x_{2k} , $k = 1, 2, \dots, n$ respectively.
 210 The pixels to be fixed when passing through the Affine Cou-
 211 pling layer will be $1/2$ of the total pixels. But when $n > 1$

212 the input of the neural network will be smaller ($1/2n$), which
 213 can in turn reduce the complexity of the neural network to
 214 save the number of parameters needed. A special case will be
 215 $n = 1$ which is the original Affine Coupling layer introduced
 216 in GLOW.

217 Some more discussion about the differential 218 geometry background

219 Given a Riemannian manifold \mathcal{M} with dimension m , it can
 220 be equipped with an (countable) indexed tuple, denoted as
 221 *chart*. Each chart, denoted by the tuple (U, Φ) where $U \subset \mathcal{M}$
 222 and $\Phi : U \rightarrow \mathbf{R}^m$ is a diffeomorphism (an invertible differen-
 223 tiable mapping). Furthermore, the set of charts, $\{U_\alpha, \Phi_\alpha\}_{\alpha \in I}$
 224 covers \mathcal{M} , i.e., $\mathcal{M} \subset \cup_{\alpha \in I} U_\alpha$. The collection of charts that
 225 cover \mathcal{M} is called an *atlas*.

226 A chart (or atlas) is required to define differentiable op-
 227 erations on \mathcal{M} , as it alleviates the non-linearity of \mathcal{M} by
 228 defining locally-linear structure given by ΦU . The choice of
 229 chart map for the given Riemannian manifold is not unique.
 230 As long as the chosen tuple (U, Φ) has the covering property
 231 and diffeomirphic property as mentioned before, it forms an
 232 well-defined chart map and hence an atlas.

233 Observe that using chart map, we can get a local coor-
 234 dinate system as given by $\Phi : U \rightarrow \mathbf{R}^m$. This raises an
 235 important question *does a local parametrization of \mathcal{M} give
 236 a chart?* The answer is yes as long as the parameterization
 237 is a diffeomorphism as by definition of parametrization, it
 238 covers the space \mathcal{M} . For manifold of SPD matrices, Cholesky
 239 decomposition is a choice of parametrization and is used in
 240 our work as the choice of chart map.

241 Another choice of chart map is using Riemannian expo-
 242 nential and inverse exponential/log map. Given $p \in \mathcal{M}$, with
 243 the tangent space $T_p \mathcal{M}$, the Riemannian exponential map is
 244 defined as the diffeomorphic mapping: $\text{Exp}_p : V \subset T_p \mathcal{M} \rightarrow$
 245 $U \subset \mathcal{M}$, here $U = \text{Exp}(V)$ and hence, is also a well-defined
 246 choice for chart maps. Note that the inverse of exponential
 247 map is called Riemannian log map.

248 In general for a choice of atlas (an the corresponding chart
 249 maps), one may look for the *minimal atlas*, i.e., one with
 250 smaller cardinality index set, I . But in this work, we are not
 251 looking for using a minimal atlas as finding such is a compu-
 252 tationally hard problem. Below, we provide some examples
 253 of charts used in the main paper.

254 For example, for the Earth without height, there are several
 255 different methods can be used to describe the location. The
 256 most acceptable one is by the Latitude and the Longitude.
 257 However, we could also use (x, y, z) where $x^2 + y^2 + z^2 = 1$
 258 to describe the location. By splitting the Earth into two halves,
 259 we can use x, y to describe the location given the north part
 260 or the south part. All these choices are not Log-map, but can
 261 uniquely describe one point on the sphere. If we consider
 262 the tangent space at the north pole, we can always use a
 263 straight line connecting the south pole and the point on the
 264 sphere which passes through the tangent space at the north
 265 pole uniquely. Thus, we can also use the Log map and the
 266 Exp map of the north pole to transfer the points between the
 267 sphere and the tangent space uniquely.

268 We have shown several examples of the choice of the chart

map in the main paper in Table. 3. The rationale behind the choice is to choose a map that can uniquely map the points and be as simple as possible. Thus, for the SPD matrices, we use the Cholesky decomposition to project the original SPD matrix into the Euclidean space uniquely. As the DTI lying on $\text{SPD}(3)$, the Cholesky decomposition can be computed explicitly.

$$\begin{aligned} & \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{12} & x_{22} & x_{23} \\ x_{13} & x_{23} & x_{33} \end{bmatrix} \\ &= \begin{bmatrix} a_1 & 0 & 0 \\ b_1 & a_2 & 0 \\ b_2 & b_3 & a_3 \end{bmatrix} \times \begin{bmatrix} a_1 & b_1 & b_2 \\ 0 & a_2 & b_3 \\ 0 & 0 & a_3 \end{bmatrix} \\ &= \begin{bmatrix} a_1^2 & a_1 b_1 & a_1 b_2 \\ a_1 b_1 & b_1^2 + a_2^2 & b_1 b_2 + a_2 b_3 \\ a_1 b_2 & b_1 b_2 + a_2 b_3 & b_2^2 + b_3^2 + a_3^2 \end{bmatrix} \quad (1) \end{aligned}$$

Thus, one can solve all the components of the Cholesky decomposition.

$$\begin{aligned} a_1 &= \sqrt{x_{11}} & b_1 &= x_{12}/a_1 \\ b_2 &= x_{13}/a_1 & a_2 &= \sqrt{x_{22} - b_1^2} \\ b_3 &= (x_{23} - b_1 b_2)/a_2 & a_3 &= \sqrt{x_{33} - b_2^2 - b_3^2} \end{aligned}$$

268 Since the computation is simple and unique (when the matrix
269 is full rank), we finally choose the Cholesky decomposition
270 as the chart map for $\text{SPD}(3)$. Similar argument holds for \mathbf{S}^n
271 and \mathbf{R}_+ .

272 One more thing to notice is that the choice of the tangent
273 space T_p depends on the point p on the manifold. Thus, when
274 using the local chart map/ Log map, we will need to pass
275 the point p along the neural network to enable invertibility.
276 Thus, in this paper, we normally choose the global chart map
277 that can uniquely map any point on the manifold into the
278 Euclidean space. For those Riemannian manifolds without
279 global chart maps, the architecture will be far more compli-
280 cated and a problem of independent interest that we do not
281 address in this work.

282 Final remarks: Why making a simple 283 Euclidean model work is extremely 284 challenging

285 A natural question a reader may ask is why we cannot embed
286 the manifold valued data into the Euclidean space and simply
287 apply the traditional GLOW model. In other words, what is
288 the rationale of deriving such a model based on differential
289 geometric principles? This is a valid question and the expla-
290 nation below describes why this approach allows practical
291 implementation when it will not even be possible otherwise.

292 As we briefly mentioned in the main paper, the Euclidean
293 simplification will typically require a *much larger* memory
294 (almost infeasible) to store and compute the model. For ex-
295 ample, in the DTI to ODF example, if we directly embed
296 two data sources into the Euclidean space, we will have \mathbf{R}^6
297 and \mathbf{R}^{361} respectively. Interestingly, notice that in the tradi-
298 tional GLOW model, due to the use of “squeeze” after every

block, the number of channels at least doubles between layers. In the traditional GLOW paper, which only deals with the color images with channel 3, this setup already requires a huge amount of space but is at least feasible. However, things get much worse when we consider the ODF data, which has 361 channels: this is $120\times$ more than the color images. Thus, in *each* layer, the number of parameters will be $\{120\times\}^2 = 14400\times$ more than the traditional GLOW model! This requirement is not even feasible to run, even on some large computer clusters that we tried. On the other hand, if we simplify the model to contain fewer layers, the expressive power of the model drops to a level that it is not able to converge.

Further, with the Euclidean setup, there is no geometric structure on the pixel values. Thus, important hard-constraints will be violated. For example, the sum of the measurements at a ODF pixel is 1 and every value should be non-negative. Also, the pixel of DTI should be a semi-positive definite matrix. Without these constraints, the output will not be scientifically meaningful or interpretable.

We note that contemporary with our work, a paper recently proposed \mathcal{M}_e -flow (Brehmer and Cranmer 2020) which also argues considering the structure of the manifold. They treat the structure of the manifold as a “surface” in the Euclidean space. So the authors propose training an encoder-decoder network to project the sparse surface in the high-dimensional space into a dense low-dimensional space. Every operator in the low-dimensional space is similar to GLOW in the Euclidean space. This idea is effective for natural images and works well. But for **known** Riemannian manifolds, we already know the dimension of the manifold. Thus, there might be little space for the encoder to lower the dimension of the data unless one is willing to compromise and accept the loss of some details: this makes the final data more blurry, a phenomena we see in the Texture images shown in the main paper. Otherwise, since the \mathcal{M}_e -flow is even more general than GLOW, the dimension of the hidden space is even higher than the GLOW in the Euclidean space.

But in our setup, as the manifold constraint is inherent in the construction of the model, if we treat every pixel/ voxel to be manifold valued, the input channel will be 1. Thus, this will significantly reduce the number of the parameters. Since the manifold structure is “hard coded” in the neural network, we can train the model to converge better with sensible resource requirements, as we show in our main paper.

References

- Brehmer, J.; and Cranmer, K. 2020. Flows for simultaneous manifold learning and density estimation. *arXiv preprint arXiv:2003.13913*.
- Dinh, L.; Krueger, D.; and Bengio, Y. 2014. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*.
- Dinh, L.; Sohl-Dickstein, J.; and Bengio, S. 2016. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.

- 356 Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating
357 deep network training by reducing internal covariate shift. *arXiv*
358 preprint arXiv:1502.03167 .
- 359 Kingma, D. P.; and Dhariwal, P. 2018. Glow: Generative flow with
invertible 1x1 convolutions. In *Advances in Neural Information
Processing Systems*, 10215–10224. 360
361
- Lee, S.-g.; Kim, S.; and Yoon, S. 2020. NanoFlow: Scalable Normal-
izing Flows with Sublinear Parameter Complexity. *arXiv preprint
arXiv:2006.06280* . 362
363
364