

Uncertainty quantification in non-rigid image registration via stochastic gradient Markov chain Monte Carlo

Daniel Grzech

Department of Computing, Imperial College London, London, UK

d.grzech17@imperial.ac.uk

Mohammad Farid Azampour

Department of Computing, Imperial College London, London, UK

Computer Aided Medical Procedures, Technische Universität München, Munich, Germany

Department of Electrical Engineering, Sharif University of Technology, Tehran, Iran

f.azampour20@imperial.ac.uk

Huaqi Qiu

Department of Computing, Imperial College London, London, UK

huaqi.qiu15@imperial.ac.uk

Ben Glocker

Department of Computing, Imperial College London, London, UK

b.glocker@imperial.ac.uk

Bernhard Kainz

Department of Computing, Imperial College London, London, UK

Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany

b.kainz@imperial.ac.uk

Loïc Le Folgoc

Department of Computing, Imperial College London, London, UK

l.le-folgoc@imperial.ac.uk

Abstract

We develop a new Bayesian model for non-rigid registration of three-dimensional medical images, with a focus on uncertainty quantification. Probabilistic registration of large images with calibrated uncertainty estimates is difficult for both computational and modelling reasons. To address the computational issues, we explore connections between the *Markov chain Monte Carlo by backpropagation* and the *variational inference by backpropagation* frameworks, in order to efficiently draw samples from the posterior distribution of transformation parameters. To address the modelling issues, we formulate a Bayesian model for image registration that overcomes the existing barriers when using a dense, high-dimensional, and diffeomorphic transformation parametrisation. This results in improved calibration of uncertainty estimates. We compare the model in terms of both image registration accuracy and uncertainty quantification to VoxelMorph, a state-of-the-art image registration model based on deep learning.

Keywords: deformable image registration, uncertainty quantification, SG-MCMC, SGLD

1. Introduction

Image registration is the problem of aligning images into a common coordinate system such that the discrete pixel locations have the same semantic information. It is a common pre-processing step for many applications, e.g. the statistical analysis of imaging data and computer-aided diagnosis through comparison with an atlas. Image registration methods based on deep learning tend to incorporate task-specific knowledge from large datasets, whereas traditional methods are more general purpose. Many established models are based on the iterative optimisation of an energy function consisting of task-specific similarity and regularisation terms, which has to be done independently for every pair of images in order

to calculate the deformation field (Schnabel et al., 2001; Klein et al., 2009; Avants et al., 2014).

DLIR (de Vos et al., 2019) and VoxelMorph (Balakrishnan et al., 2018, 2019; Dalca et al., 2018, 2019) changed this paradigm by learning a function that maps a pair of input images to a deformation field. This gives a speed-up of several orders of magnitude at inference time and maintains an accuracy comparable to traditional methods. An overview of state-of-the-art models for image registration based on deep learning can be found in Lee et al. (2019).

Due to the perceived conceptual difficulty and computational overhead, Bayesian methods tend to be shunned when designing medical image analysis algorithms. However, in order to fully explore the parameter space and lessen the impact of ad-hoc hyperparameter choices, it is desirable to use Bayesian models. In addition, with help of open-source libraries with automatic differentiation like PyTorch, the implementation of even complex Bayesian models for image registration is very similar to that of non-probabilistic models.

In this paper, we make use of the stochastic gradient Markov chain Monte Carlo (SG-MCMC) algorithm to design an efficient posterior sampling algorithm for 3D non-rigid image registration. SG-MCMC is based on the idea of stochastic gradient descent interpreted as a stochastic process with a stationary distribution centred on the optimum and whose covariance can be used to approximate the posterior distribution (Chen et al., 2016; Mandt et al., 2017). SG-MCMC methods have been useful for training generative models on very large datasets ubiquitous in computer vision, e.g. Du and Mordatch (2019); Nijkamp et al. (2019); Zhang et al. (2020). We show that they are also applicable to image registration.

This work is an extended version of Grzech et al. (2020), where we first proposed use of the SG-MCMC algorithm for non-rigid image registration. The code to reproduce the results is available in a public repository: <https://github.com/dgrzech/ir-sgmcmc>. The following is a summary of the main contributions of the previous work:

1. We proposed a computationally efficient SG-MCMC algorithm for three-dimensional diffeomorphic non-rigid image registration;
2. We introduced a new regularisation loss, which allows to carry out inference of the regularisation strength when using a transformation parametrisation with a large number of degrees of freedom;
3. We evaluated the model both qualitatively and quantitatively by analysing the output transformations, image registration accuracy, and uncertainty estimates on inter-subject brain magnetic resonance imaging (MRI) data from the UK Biobank dataset.

In this version, we extend the previous work:

- We provide more details on the Bayesian formulation, including a comprehensive analysis of the learnable regularisation loss, as well as a more in-depth analysis of the model hyperparameters and hyperpriors;
- We conduct additional experiments in order to compare the uncertainty estimates output by variational inference (VI), SG-MCMC, and VoxelMorph qualitatively, as well as quantitatively by analysing the Pearson correlation coefficient between the displacement and label uncertainties;

- We analyse the differences between uncertainty estimates when the SG-MCMC algorithm is initialised to different transformations and when using different parametrisations of the transformation, including non-parametric stationary velocity fields (SVFs) and SVFs based on B-splines; and
- We include a detailed evaluation of the computational speed and of the output transformation smoothness.

2. Related work

The problem of uncertainty quantification in non-rigid image registration is controversial because of ambiguity regarding the definition of uncertainty as well as the accuracy of uncertainty estimates (Luo et al., 2019). Uncertainty quantification in probabilistic image registration relies either on variational Bayesian methods (Simpson et al., 2012, 2013; Wassermann et al., 2014), which are fast and approximate, and popular within models based on deep learning (Dalca et al., 2018; Liu et al., 2019; Schultz et al., 2019), or Markov chain Monte Carlo (MCMC) methods, which are slower but enable asymptotically exact sampling from the posterior distribution of the transformation parameters. The latter include e.g. Metropolis-Hastings used for intra-subject registration of brain MRI scans (Risholm et al., 2010, 2013) and estimating delivery dose in radiotherapy (Risholm et al., 2011), reversible-jump MCMC used for cardiac MRI (Le Folgoc et al., 2017), and Hamiltonian Monte Carlo used for atlas building (Zhang et al., 2013).

Uncertainty quantification for image registration has also been done via kernel regression (Zöllei et al., 2007; Janoos et al., 2012) and deep learning (Dalca et al., 2018; Krebs et al., 2019; Heinrich, 2019; Sedghi et al., 2019). More generally, Bayesian frameworks have been used e.g. to characterize image intensities (Hachama et al., 2012) and anatomic variability (Zhang and Fletcher, 2014).

One of the main obstacles to a more widespread use of MCMC methods for uncertainty quantification is the computational cost. This was recently tackled by embedding MCMC in a multilevel framework (Schultz et al., 2018). SG-MCMC was previously used for *rigid* image registration (Karabulut et al., 2017). It has also been employed in the context of unsupervised non-rigid image registration based on deep learning, where it allowed to sample from the posterior distribution of the network weights, rather than directly the transformation parameters (Khawaled and Freiman, 2020).

Previous work on data-driven regularisation focuses on transformation parametrisations with a relatively low number of degrees of freedom, e.g. B-splines (Simpson et al., 2012) and a sparse parametrisation based on Gaussian radial basis functions (RBFs) (Le Folgoc et al., 2017). Limited work exists also on spatially-varying regularisation, again with B-splines (Simpson et al., 2015). Deep learning has been used for spatially-varying regularisation learnt using more than one image pair (Niethammer et al., 2019). Shen et al. (2019) introduced a related model which could be used for learning regularisation strength based on a single image pair but suffered from non-diffeomorphic output transformations and slow speed.

3. Registration model

We denote an image pair by $\mathcal{D} = (F, M)$, where $F: \Omega_F \rightarrow \mathbb{R}$ and $M: \Omega_M \rightarrow \mathbb{R}$ are a fixed and a moving image respectively. The goal of image registration is to align the underlying domains Ω_F and Ω_M with a transformation $\varphi(w): \Omega_F \rightarrow \Omega_M$, i.e. to calculate parameters w such that $F \simeq M(w) := M \circ \varphi^{-1}(w)$. The transformation is often expected to possess desirable properties, e.g. diffeomorphic transformations are smooth and invertible, with a smooth inverse.

We parametrise the transformation using the SVF formulation (Arsigny et al., 2006; Ashburner, 2007), which we briefly review below. The ordinary differential equation (ODE) that defines the evolution of the transformation is given by:

$$\frac{\partial \varphi^{(t)}}{\partial t} = w(\varphi^{(t)}) \quad (1)$$

where $\varphi^{(0)}$ is the identity transformation and $t \in [0, 1]$. If the velocity field w is spatially smooth, then the solution to Equation (1) is a diffeomorphic transformation. Numerical integration is done by scaling and squaring, which uses the following recurrence relation with 2^T steps (Arsigny et al., 2006):

$$\varphi^{(1/2^{t-1})} = \varphi^{(1/2^t)} \circ \varphi^{(1/2^t)} \quad (2)$$

The Bayesian image registration framework that we present is not limited to SVFs. Moreover, there is a very limited amount of research on the impact of the transformation parametrisation on uncertainty quantification. Previous work on uncertainty quantification in image registration characterised uncertainty using a single transformation parametrisation, e.g. a small deformation model using B-splines in Simpson et al. (2012), the finite element (FE) method in Risholm et al. (2013), and multi-scale Gaussian RBFs in Le Folgoc et al. (2017), or a large deformation diffeomorphic metric mapping (LDDMM) in Wassermann et al. (2014).

To help understand the potential impact of the transformation parametrisation on uncertainty quantification, we also implement SVFs based on cubic B-splines (Modat et al., 2012). In this case, the SVF consists of a grid of B-spline control points, with regular spacing $\delta \geq 1$ voxel. The dense SVF at each point is a weighted combination of cubic B-spline basis functions (Rueckert et al., 1999). To calculate the transformation based on the dense velocity field, we again use the scaling and squaring algorithm in Equation (2).

3.1 Likelihood model

The likelihood $p(\mathcal{D} | w)$ specifies the relationship between the data and the transformation parameters by means of a similarity metric. In probabilistic image registration, it usually takes the form of a Boltzmann distribution (Ashburner, 2007):

$$\log p(\mathcal{D} | w; \mathcal{H}) \propto -\mathcal{E}_{\text{data}}(\mathcal{D}, w; \mathcal{H}) \quad (3)$$

where $\mathcal{E}_{\text{data}}$ is the similarity metric and \mathcal{H} an optional set of hyperparameters.

Local cross-correlation (LCC), which is invariant to linear intensity scaling, is a popular similarity metric but not meaningful in a probabilistic context. For this reason, instead of

the sum of the voxel-wise product of intensities, like in standard LCC, we opt for the sum of voxel-wise squared differences of images standardised to zero mean and unit variance inside a local neighbourhood of five voxels. This way, we can benefit from robustness under linear intensity transformations, as well as desirable properties of a Gaussian mixture model (GMM) of intensity residuals, i.e. robustness to outlier values caused by acquisition artefacts and misalignment over the course of registration (Le Folgoc et al., 2017).

Let \bar{F} and $\bar{M}(w)$ be respectively the fixed and the warped moving image with intensities standardised to zero mean and unit variance inside a neighbourhood of five voxels. For each voxel, the intensity residual $r_i = \bar{F} - \bar{M}(w)$, $i \in \{1, \dots, N^3\}$, is assigned to the l -th component of the mixture, $1 \leq l \leq L$, if the categorical variable $c_i \in \{1, \dots, L\}$ is equal to l , in which case it follows a normal distribution $\mathcal{N}\left(0, \beta_l^{-1}\right)$ ¹. The component assignment c_i follows a categorical distribution and takes value l with probability ϱ_l . We use the same GMM of intensity residuals on a global basis rather than per neighbourhood. In all experiments it has $L = 4$ components, which we determine to be sufficient for a good model fit.

We also use the scalar virtual decimation factor α to account for the fact that voxel-wise residuals are not independent. This prevents over-emphasis on the data term and allows to better calibrate uncertainty estimates (Groves et al., 2011; Simpson et al., 2012). The full expression of the image similarity term is given by:

$$\mathcal{E}_{\text{data}}(\mathcal{D}, w; \beta, \varrho) = -\alpha \times \sum_{i=1}^{N^3} \log \sum_{l=1}^L \varrho_l \sqrt{\frac{\beta_l}{2\pi}} \exp\left(-\frac{\beta_l}{2} r_i^2\right) \quad (4)$$

3.2 Transformation priors

In Bayesian models, the transformation parameters are typically regularised with use of a multivariate normal prior that ensures smoothness:

$$\log p(w; \lambda_{\text{reg}}) \propto -\frac{1}{2} \lambda_{\text{reg}} (\mathbf{L}w)^\top \mathbf{L}w \quad (5)$$

where λ_{reg} is a scalar parameter that controls the regularisation strength, and \mathbf{L} is the matrix of a differential operator. Here we assume that \mathbf{L} represents the gradient operator, which penalises the magnitude of the 1st derivative of a velocity field. Note that $(\mathbf{L}w)^\top \mathbf{L}w = \|\mathbf{L}w\|^2 := \chi^2$.

The regularisation weight λ_{reg} can either be fixed or estimated from data. The latter has been done successfully only for transformation parametrisations with a relatively low number of degrees of freedom, e.g. B-splines (Simpson et al., 2012) and a sparse parametrisation (Le Folgoc et al., 2017). In case of an SVF, where the number of degrees of freedom is orders of magnitude higher, the problem is more difficult. However, a reliable method to adjust regularisation strength based on data is crucial, as both the output transformation and registration uncertainty are highly sensitive to regularisation. In order to infer the regularisation strength, we specify a log-normal prior on the scalar regularisation energy

1. In order to reduce the notation clutter we omitted the voxel index for the fixed and moving images.

$\chi^2 \sim \text{Lognormal}(\mu_{\chi^2}, \sigma_{\chi^2}^2)$, and derive a prior on the underlying SVF:

$$\log p(\chi^2) \propto -\log \chi^2 - \log \sigma_{\chi^2} - \frac{(\log \chi^2 - \mu_{\chi^2})^2}{2\sigma_{\chi^2}^2} \quad (6)$$

$$\log p(w) \propto -\left(\frac{\nu}{2} - 1\right) \log \chi^2 + \log p(\chi^2) \quad (7)$$

where $\nu = 3N^3$ is the number of degrees of freedom, i.e. the count of transformation parameters in all three directions. Given (semi-)informative hyperpriors on μ_{χ^2} and $\sigma_{\chi^2}^2$, which we discuss in the next section, we can adjust the regularisation strength to the input images. The full expression of the regularisation term is given by:

$$\mathcal{E}_{\text{reg}}(w) = \frac{\nu}{2} \log \chi^2 + \log \sigma_{\chi^2} + \frac{(\log \chi^2 - \mu_{\chi^2})^2}{2\sigma_{\chi^2}^2} \quad (8)$$

It is worth noting that the traditional L_2 regularisation with a fixed regularisation weight in Equation (5) actually belongs to this family of regularisation losses. If we specify a gamma prior instead of a log-normal prior on the scalar regularisation energy $\chi^2 \sim \Gamma(\nu/2, \lambda_{\text{reg}}/2)$, we get:

$$\log p(\chi^2) \propto \left(\frac{\nu}{2} - 1\right) \log \chi^2 - \frac{1}{2} \lambda_{\text{reg}} \cdot \chi^2 \quad (9)$$

$$\log p(w) \propto \left(\frac{\nu}{2} - 1\right) \log \chi^2 + \left(\frac{\nu}{2} - 1\right) \log \chi^2 - \frac{1}{2} \lambda_{\text{reg}} \cdot \chi^2 \quad (10)$$

$$\propto -\frac{1}{2} \lambda_{\text{reg}} (\mathbf{L}w)^\top \mathbf{L}w \quad (11)$$

3.3 Hyperpriors

We set the likelihood GMM hyperpriors similarly to Le Folgoc et al. (2017), with the mixture precision parameters $\beta = (\beta_1, \dots, \beta_L)$ assigned independent log-normal priors $\beta_l \sim \text{Lognormal}(\mu_{\beta_l}, \sigma_{\beta_l}^2)$ and the mixture proportions $\varrho = (\varrho_1, \dots, \varrho_L)$ with an uninformative Dirichlet prior $\varrho \sim \text{Dir}(\kappa)$, where $\kappa = (\kappa_1, \dots, \kappa_L)$.

Regularisation parameters require informative priors due to the difficulty of learning the regularisation strength based on a single image pair. Because of a gamma prior on the regularisation energy $\exp(\mu_{\chi^2}) \sim \Gamma(\nu/2, \lambda_{\text{init}}/2)$, we can rely on the familiar regularisation weight λ_{init} to initialise the logarithm of the regularisation energy μ_{χ^2} to the expected value of the logarithm of the gamma distribution, i.e. $\mathbb{E}[\mu_{\chi^2}] = \psi(\nu/2) - \log(\lambda_{\text{init}}/2)$, where ψ is the digamma function. The value of this expression is sharply peaked if the number of degrees of freedom ν is large, which yields a very informative prior on μ_{χ^2} . More details on how to calculate the expected value of the logarithm of a gamma distribution can be found in Appendix A.

The choice of a hyperprior on the scale parameter σ_{χ^2} , which controls the amount of deviation of $\log \chi^2$ from the location parameter μ_{χ^2} , is more intuitive. Here we use a log-normal prior $\sigma_{\chi^2}^2 \sim \text{Lognormal}(\eta, \varsigma^2)$.

4. Variational inference

Image registration methods often rely on VI for uncertainty quantification. We only use VI to initialise the SG-MCMC algorithm, which also lets us compare the uncertainty estimates output by approximate and asymptotically exact methods for sampling from the posterior distribution of transformation parameters.

We assume that the posterior distribution $p(w | \mathcal{D})$ is a multivariate normal distribution $q(w) \sim \mathcal{N}(\mu_w, \Sigma_w)$. To find parameters μ_w and Σ_w , we maximise the evidence lower bound (ELBO) (Jordan et al., 1999):

$$\mathcal{L}(q) = \mathbb{E}_q [\log p(\mathcal{D} | w)] - \text{D}_{\text{KL}}(q || p) = -\langle \mathcal{E}_{\text{data}} + \mathcal{E}_{\text{reg}} \rangle_q + H(q) \quad (12)$$

where $\text{D}_{\text{KL}}(q || p)$ is the Kullback-Leibler divergence (KL-divergence) between the approximate posterior q and the prior p . Like in traditional image registration, the energy function consists of the sum of similarity and regularisation terms, with an additional term for the entropy of the posterior distribution $H(q)$. We show how to calculate this term in Appendix B.

It is not possible to calculate every element of the covariance matrix Σ_w due to high dimensionality of the problem. Instead, we approximate the covariance matrix as a sum of diagonal and low-rank parts, i.e. $\Sigma_w \approx \text{diag}(\sigma_w^2) + u_w u_w^T$, with $\sigma_w \in \mathbb{R}^{3N^3 \times 1}$ and $u_w \in \mathbb{R}^{3N^3 \times R}$, where R is a hyperparameter which determines the parametrisation rank. Using a multivariate normal distribution as the approximate posterior distribution of transformation parameters is common in probabilistic image registration. The key difference in our work is the diagonal + low-rank parametrisation of the covariance matrix. Most recent image registration models with an SVF transformation parametrisation are based on deep learning and make the assumption of a diagonal covariance matrix (Dalca et al., 2018; Krebs et al., 2019).

We use the reparametrisation trick with two samples per update to backpropagate with respect to the variational posterior parameters:

$$\begin{aligned} w &= \mu_w \pm (\text{diag}(\sigma_w) \cdot \epsilon + u_w \cdot x) \\ \epsilon &\sim \mathcal{N}(0, I_{3N^3}), \quad x \sim \mathcal{N}(0, I_R) \end{aligned} \quad (13)$$

In order to make the optimisation less susceptible to undesired local maxima of the ELBO, we take advantage of Sobolev gradients (Neuberger et al., 1997). Samples from the posterior are convolved with a Sobolev kernel. We approximate the 3D kernel by three separable 1D kernels to lower the computational overhead. Using the notation in Slavcheva et al. (2018), we set the kernel width to $s_{H^1} = 7$ and the smoothing parameter to $\lambda_{H^1} = 0.5$.

The GMM and regularisation hyperparameters are fit using the stochastic approximation expectation maximisation (SAEM) algorithm (Richard et al., 2009; Zhang et al., 2013). The mixture precision hyperparameters β and proportion hyperparameters ϱ are updated by solving the optimisation problem:

$$\beta^{(k)}, \varrho^{(k)} = \arg \max_{\beta, \varrho} \mathbb{E}_q \left[\log p(\mathcal{D}, w; \beta, \varrho) \mid \beta^{(k-1)}, \varrho^{(k-1)} \right] + \log p(\beta) + \log p(\varrho) \quad (14)$$

This is done at each step of the iterative optimisation algorithm. We update the regularisation hyperparameters in an analogous way. Even though the hybrid VI and SAEM

approach is computationally efficient and requires minimal implementation effort, it disregards the uncertainty caused by hyperparameter variability.

5. Stochastic gradient Markov chain Monte Carlo

Image registration algorithms based on VI restrict parametrisation of the posterior distribution to a specific family of probability distributions, which may not include the true posterior. To avoid this problem and sample the transformation parameters in an efficient way, we use stochastic gradient Langevin dynamics (SGLD) (Besag, 1993; Welling and Teh, 2011). The update equation is given by:

$$w_{k+1} \leftarrow w_k + \tau A \nabla \log \pi(w_k) + \sqrt{2\tau A} \xi_k \quad (15)$$

where τ is the step size, A is an optional preconditioning matrix, $\nabla \log \pi(w_k)$ is the gradient of the logarithm of the posterior probability density function (PDF), and $\xi_k \sim \mathcal{N}(0, I_{3N^3})$. SGLD does not require a particular initialisation, so we study several different possibilities, including a sample $w_0 \sim \mathcal{N}(\mu_w, \Sigma_w)$ from the approximate variational posterior, in which case we set $A = \text{diag}(\sigma_w^2)$. The preconditioning helps with the MCMC mixing time in case the target distribution is strongly anisotropic.

It is worth noting that, except for the preconditioning matrix A and the noise term ξ_k , Equation (15) is equivalent to a gradient descent update when minimising the maximum a posteriori (MAP) objective function $-\log p(w | \mathcal{D}) = -\log p(\mathcal{D} | w) - \log p(w)$. When drawing samples from SGLD, we continue to update the GMM as well as regularisation hyperparameters like in Equation (14), except that the expected value is calculated with respect to the new posterior $\pi(w)$.

In the limit as $k \rightarrow \infty$ and $\tau \rightarrow 0$, SGLD can be used to draw exact samples from the posterior of the transformation parameters without Metropolis-Hastings accept-reject tests, which are computationally expensive. Indeed, these costs prevent the use of other MCMC algorithms for the registration of large 3D images. In practice, the step size needs to be adjusted to avoid high autocorrelation between samples yet remain smaller than the width of the most constrained direction in the local energy landscape (Neal, 2011). The step size can also be used to control the trade-off between accuracy and computation time. We can quantify uncertainty either quickly in a coarse manner or slowly, with more detail.

Despite the fact that the term $\nabla \log \pi(w)$ allows to traverse the energy landscape in an efficient way, SGLD suffers from high autocorrelation and slow mixing between modes (Hill, 2020). However, simplicity of the formulation makes it better suited than other MCMC methods for high-dimensional problems like three-dimensional image registration.

6. Experiments

6.1 Setup

The model is implemented in PyTorch. For all experiments we use three-dimensional T2-FLAIR MRI brain scans and subcortical structure segmentations from the UK Biobank dataset (Sudlow et al., 2015). Input images are pre-registered with the affine component of

drop2 (Glocker et al., 2008)² and resampled to $N = 128$ isotropic voxels of length 1.82 mm along every dimension.

We use 2^{12} steps to integrate SVFs. In order to start optimisation with small displacements, μ_w is initialised to zero, which corresponds to an identity transformation, σ_w to half a voxel length in each direction and u_w to a tenth of a voxel length in each direction. We are mainly interested in the approximate variational posterior in order to initialise the SG-MCMC algorithm, so the rank parameter is set to $R = 1$. We use the Adam optimiser with a step size of 1×10^{-2} for the approximate variational posterior parameters μ_w , $\log \sigma_w^2$, and u_w . Training is run until the ELBO value stops to increase, which requires approximately 1,024 iterations.

In the likelihood model, we use $\kappa = 0.5$ for an uninformative Jeffreys prior on the mixture proportions, while the mixture precision hyperparameters are set to $\mu_{\beta_l} = 0.0$ and $\sigma_{\beta_l} = 2.3$. The model is much less sensitive to the value of the likelihood hyperparameters than the regularisation hyperparameters, which are calibrated to guarantee diffeomorphic transformations sampled from the approximate variational posterior $q(w) \sim \mathcal{N}(\mu_w, \Sigma_w)$. The local transformation is diffeomorphic only in locations where the Jacobian determinant is positive (Ashburner, 2007), so we aim to keep the number of voxels where the Jacobian determinant is non-positive $|\det J_{\varphi^{-1}}| \leq 0$ close to zero. We calibrate the location hyperparameter λ_{init} in every experiment, while the scale hyperparameters are set to $\eta = 2.8$ and $\varsigma = 5.0$.

SAEM convergence is known to be conditional on decreasing step sizes (Delyon et al., 1999). For this reason, we use a small step size decay of 1×10^{-3} and the Adam optimiser with a step size of 2×10^{-1} for the GMM hyperparameters $\log \beta^{-0.5}$ and $\log \varrho$, and 1×10^{-2} for the regularisation hyperparameters μ_{χ^2} and $\log \sigma_{\chi^2}$. In case of parameters whose value is constrained to be positive, we state the step size used on the logarithms. In practice, we did not observe the result to be dependent on these step sizes.

6.2 Regularisation strength

First we evaluate the proposed regularisation. We compare it to a gamma prior on λ , i.e. $\lambda \sim \Gamma(s, r)$, where s and r are the shape and the rate parameters respectively, set to uninformative values $s = r = \nu/2$ (Simpson et al., 2012).

We compare the output of VI when using fixed regularisation weights $\lambda_{\text{reg}} \in \{0.1, 1.2\}$, the baseline method for learnable regularisation strength, and our regularisation loss. The result on a sample pair of input images is shown in Figure 1. For the baseline method, the learnt regularisation strength is too high, which effectively prevents the alignment of images. This indicates that previous schemes for inference of regularisation strength from data are inadequate when the transformation parametrisation involves a very large number of degrees of freedom. In case of $\lambda_{\text{reg}} = 0.1$, the resulting transformation is not diffeomorphic. The output when using our regularisation loss with $\lambda_{\text{init}} = 1.2$ strikes a balance between the baseline and $\lambda_{\text{reg}} = 0.1$, where there is an overemphasis on the data term.

In Figure 2, we show the output of VI for two pairs of images which require different regularisation strengths for accurate alignment. We choose a fixed image and two moving images M_1 and M_2 , with one visibly different and the other similar to the fixed image. We

2. <https://github.com/biomed-mira/drop2>

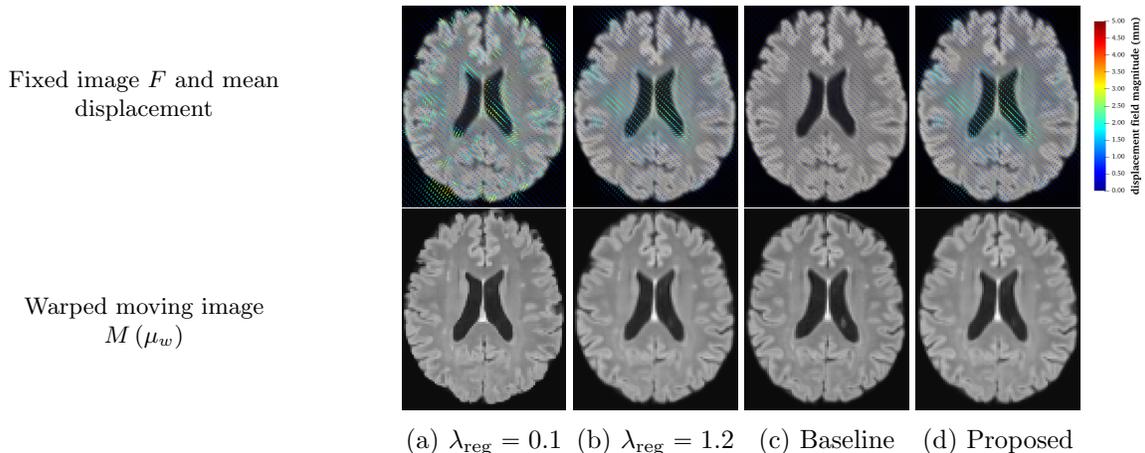


Figure 1: Output when using two fixed regularisation weights $\lambda_{\text{reg}} \in \{0.1, 1.2\}$, the baseline method for learnable regularisation strength, and the proposed learnable regularisation loss with $\lambda_{\text{init}} = 1.2$. For fixed regularisation weight $\lambda_{\text{reg}} = 0.1$, the sampled transformations are not diffeomorphic. In case of the baseline method, the learnt regularisation strength is too high, which effectively prevents the alignment of images. When using the proposed learnable regularisation loss, we strike a balance between the baseline method and fixed regularisation weight $\lambda_{\text{reg}} = 0.1$, where the regularisation strength is too low. The figure shows the middle axial slice of 3D images.

also analyse the result when using fixed regularisation weights $\lambda_{\text{reg}} = 0.2$, which leads to non-diffeomorphic transformations, and $\lambda_{\text{reg}} = 2.0$, which produces smooth transformations but, in case of M_1 , at the expense of accuracy. The proposed regularisation, initialised with $\lambda_{\text{init}} = 2.0$, helps to prevent oversmoothing.

6.3 Uncertainty quantification

We run a number of experiments to evaluate the uncertainty estimates and better understand the differences between uncertainty output by various non-rigid registration methods in practice:

1. We compare the uncertainty estimates output by VI, SG-MCMC, and VoxelMorph on inter-subject brain MRI data from UK Biobank qualitatively and quantitatively, by calculating the Pearson correlation coefficient between the displacement and label uncertainties;
2. We compare the uncertainty estimates when the SG-MCMC algorithm is initialised to different transformations;
3. We compare the result when using non-parametric SVFs and SVFs based on B-splines to parametrise the transformation.

In order to make sampling from SG-MCMC efficient, we determine the largest step size that guarantees diffeomorphic transformations as defined in Section 6.1 and set it to

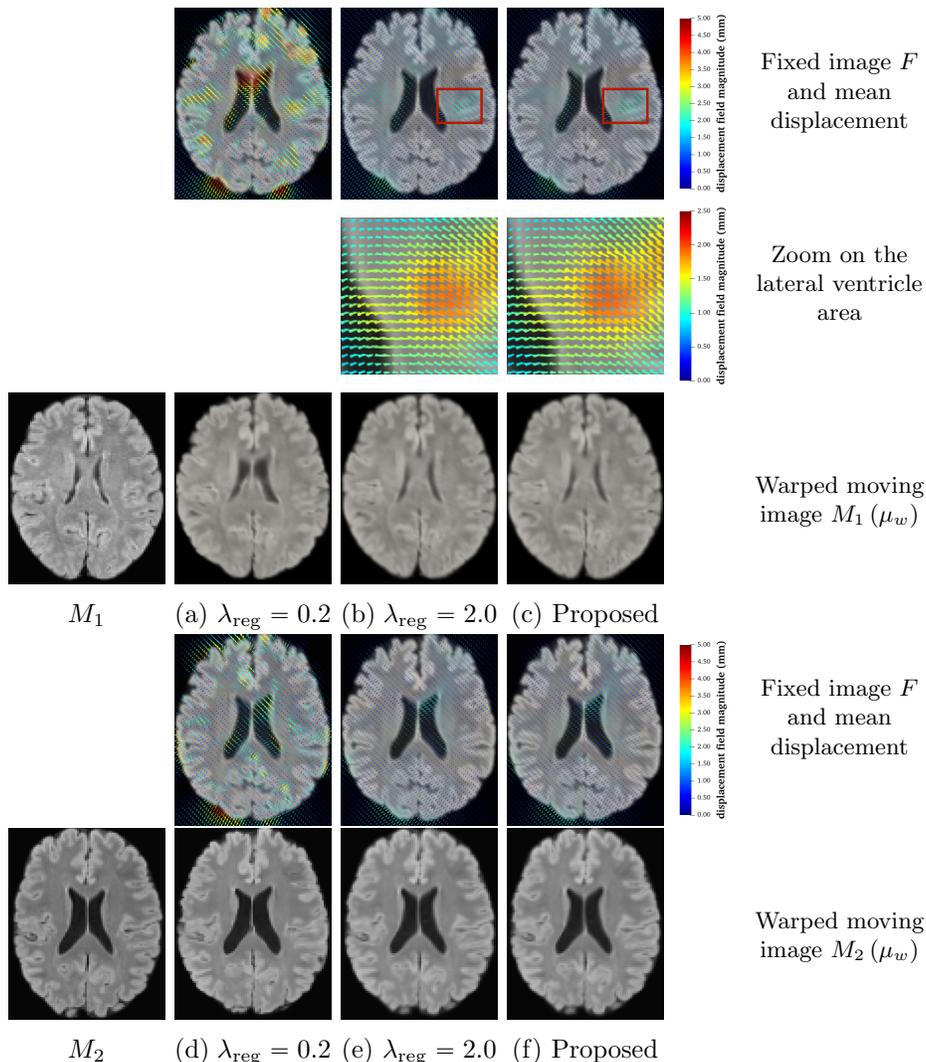


Figure 2: Output of VI for two image pairs which require different regularisation strengths. M_1 is visibly different to the fixed image and M_2 is similar to it. The proposed regularisation is initialised with $\lambda_{\text{init}} = 2.0$. For both image pairs, the alignment is best in case of the fixed regularisation strength $\lambda_{\text{reg}} = 0.2$ but the output transformation is not diffeomorphic. For M_2 , the proposed learnable regularisation loss is almost identical to $\lambda_{\text{reg}} = 2.0$, which ensures good accuracy and smoothness of the transformation. For M_1 , the proposed loss outputs somewhat higher accuracy than $\lambda_{\text{reg}} = 2.0$. The figure shows the middle axial slice of 3D images.

$\tau = 4 \times 10^{-1}$. A single Markov chain requires only 4 GB of memory, so two Markov chains are run in parallel, each initialised to a different sample from the approximate variational posterior. We discard the first 100,000 samples from each chain to allow the Markov chains to reach the stationary distribution.

6.3.1 COMPARISON OF UNCERTAINTY ESTIMATES OUTPUT BY DIFFERENT MODELS

VoxelMorph. To fill the knowledge gaps on differences between uncertainty estimates produced by non-rigid image registration algorithms, we train the probabilistic VoxelMorph (Dalca et al., 2018) in atlas mode on a random 80/20 split of 13,401 brain MRI scans in the UK Biobank³. We use the same fixed image as in the experiments and exclude the moving images from the training data. To enable a fair comparison, the chosen similarity metric is the sum of squared differences (SSD). We also study the differences between uncertainty estimates output by VI and SG-MCMC, based on 500 samples output by each model. In order to reduce auto-correlation, samples output by SG-MCMC are selected at regular intervals from the one million samples drawn from each chain.

Like in VI, which we use to initialise the SG-MCMC algorithm, the approximate variational posterior of transformation parameters output by VoxelMorph is assumed to be a multivariate normal distribution $q_{\text{VXM}}(w) \sim \mathcal{N}(\mu_{\text{VXM}}, \Sigma_{\text{VXM}})$. The only difference is the covariance matrix Σ_{VXM} , which is assumed to be diagonal rather than diagonal + low-rank. In order to set the model hyperparameters, we analyse the average surface distances (ASDs) and Dice scores (DSCs) on subcortical structure segmentations and the Jacobian determinants of sample transformations, with the aim of striking a balance between accuracy and smoothness. The most important hyperparameters are the fixed regularisation strength parameter, set to $\lambda_{\text{VXM}} = 10.0$, and the initial value of the diagonal covariance matrix, set to $\Sigma_{\text{VXM}} = \text{diag}(0.01)$.

Qualitative comparison of uncertainty estimates. The output of VI, SG-MCMC, and VoxelMorph on a sample image pair is shown in Figure 3. It should be noted that, due to the fact that direct correspondence between regions is hard to determine, the problem of registering inter-subject brain MRI scans is more challenging than problems where non-rigid registration uncertainty had been studied previously, e.g. intra-subject brain MRI (Risholm et al., 2010; Simpson et al., 2012) and intra-subject cardiac MRI (Le Folgoc et al., 2017). Nonetheless, the uncertainty estimates output by VI and SG-MCMC are consistent with previous findings, e.g. higher uncertainty in homogeneous regions (Simpson et al., 2012; Dalca et al., 2018).

Unlike in Le Folgoc et al. (2017), where Gaussian RBFs were used to parametrise the transformation on intra-subject cardiac MRI scans, VI outputs higher uncertainty than MCMC. SGLD is known to overestimate the posterior covariance due to non-vanishing learning rates at long times (Mandt et al., 2017), which further suggests that the uncertainty output by SGLD might be underestimated. The uncertainty estimates produced by VI and SG-MCMC are consistent but different in magnitude, while those output by VoxelMorph are noticeably different, with the values much smaller. This indicates the need for further research into calibration of uncertainty estimates for image registration methods based on deep learning.

Image registration accuracy. To evaluate image registration accuracy, we calculate ASDs and DSCs on the subcortical structure segmentations using the fixed and the moving segmentation warped with transformations sampled from the models. The metric comparison between our model and VoxelMorph on the image pair in Figure 3 is shown in Table 1

3. The official VoxelMorph implementation used in the experiments is available on GitHub: <https://github.com/voxelmorph/voxelmorph>.

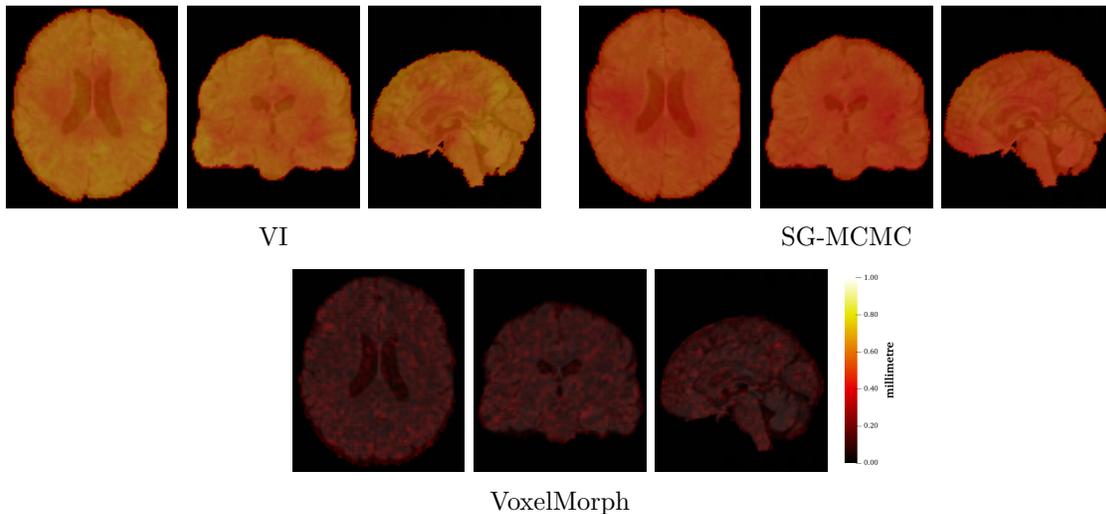


Figure 3: Uncertainty output by the models on the input image pair shown in Figure 1. The standard deviation of the displacement field magnitude is calculated using 500 samples. In case of SG-MCMC, samples are selected at regular intervals from the one million samples output by each of the two Markov chains, which is needed to prevent autocorrelation between samples. SG-MCMC outputs lower uncertainty than VI. The uncertainty estimates output by VI and VoxelMorph are very dissimilar, even though the two models assume a similar approximate variational posterior. In case of SGLD, visualising standard deviation of the displacement field magnitude is valid under the assumption that the posterior is approximately Gaussian and mono-modal. The standard deviations of displacement field magnitudes sampled from VI and SG-MCMC are very different, so the visualisation in case of SG-MCMC should not be treated as an accurate description of the true uncertainty, but rather as evidence that the true posterior distribution of transformation parameters likely is not Gaussian with a diagonal + low-rank covariance matrix.

for segmentations grouped by volume and in Figure 4 for individual segmentations. The metrics show significant improvement over affine registration on most subcortical structures. The differences between label uncertainties are less pronounced than between transformation uncertainties. ASDs and DSCs are generally marginally better when using samples from SG-MCMC than from the approximate variational posterior. Better accuracy in case of SGLD is expected, given restrictive assumptions of the approximate variational model. VI, SG-MCMC, and VoxelMorph produce similar accuracy, despite different uncertainty estimates.

Quantitative comparison of uncertainty estimates. It is difficult to define ground truth uncertainty with regards to image registration. Luo et al. (2019) suggested that well-calibrated image registration uncertainty estimates need to be informative of anatomical features, which are important in neurosurgery. To compare the uncertainty output by different models quantitatively and evaluate whether the uncertainty estimates are clinically useful, we adopt a method similar to that used by Luo et al. (2019) and calculate the Pearson correlation coefficient $r_{u_d u_l}$ between the displacement uncertainties u_d and label

Table 1: Mean and standard deviation of ASD and DSC on small, medium, and large subcortical structures, calculated using 500 output samples. Large structures include the brain stem and thalamus, medium structures—the caudate, hippocampus, and putamen, and small structures—the accumbens, amygdala, and pallidum. The values for the best performing model are underlined.

model	average surface distance (mm)			Dice score		
	small structures	medium structures	large structures	small structures	medium structures	large structures
VI	1.12 (0.16)	0.97 (0.08)	<u>1.04</u> (0.26)	0.68 (0.07)	0.79 (0.03)	<u>0.88</u> (0.02)
SG-MCMC	1.10 (0.15)	<u>0.95</u> (0.07)	1.05 (0.24)	0.68 (0.07)	0.79 (0.03)	<u>0.88</u> (0.02)
VoxelMorph	<u>1.06</u> (0.14)	0.96 (0.09)	1.05 (0.11)	<u>0.70</u> (0.05)	<u>0.80</u> (0.02)	0.87 (0.01)

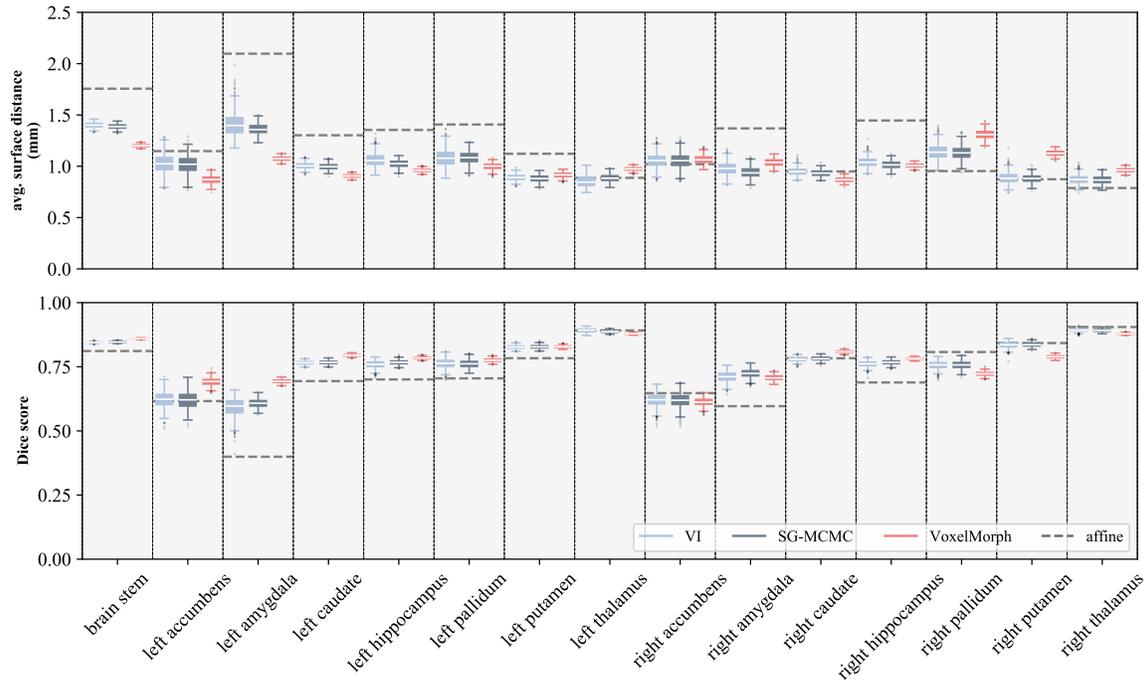


Figure 4: ASD and DSC on each subcortical structure. Dashed lines show the metric values prior to non-rigid registration. The label uncertainty of VI and SG-MCMC is comparable, despite different transformation uncertainty.

uncertainties u_l . Here we define label uncertainty u_l to be the standard deviation of the DSC for each subcortical structure, and displacement uncertainty u_d to be the voxel-wise mean of the displacement field standard deviation magnitude within the region that corresponds to a given subcortical structure.

Table 2: Pearson correlation coefficient $r_{u_d u_l}$ between the displacement uncertainties u_d and label uncertainties u_l . We define label uncertainty u_l to be the standard deviation of the DSC for each subcortical structure, and displacement uncertainty u_d to be the voxel-wise mean of the displacement field standard deviation magnitude within the region which corresponds to a given subcortical structure. The value of the correlation coefficient is highest in case of VoxelMorph but qualitative evidence suggests that the uncertainty estimates output by the model are not well calibrated. The correlation coefficient is also higher for SG-MCMC than for VI.

model	$r_{u_d u_l}$
VI	0.07
SG-MCMC	0.10
VoxelMorph	0.12

The correlation coefficient is highest in case of VoxelMorph but qualitative evidence suggests that image registration uncertainty output by the model is not well calibrated. The displacement field uncertainty is more informative of label uncertainty in case of SG-MCMC than VI, so the uncertainty estimates output by SG-MCMC are likely to be more useful for clinical purposes than those output by VI or VoxelMorph.

Transformation smoothness. Finally, in order to evaluate the quality of the model output, in Table 3 we report the number of voxels where the sampled transformations are not diffeomorphic. Each model produces transformations where the number of non-positive Jacobian determinants is nearly zero. SG-MCMC slightly reduced the number of non-positive Jacobian determinants compared to VI. The transformations output by VoxelMorph appear more smooth, which is directly related to the lower transformation uncertainty shown in Figure 3.

6.3.2 COMPARISON OF THE OUTPUT OF SG-MCMC FOR DIFFERENT INITIALISATIONS

To study the potential impact of initialisation on the output uncertainty estimates, we analyse the transformations sampled from SGLD run with different initial velocity fields w_0 . We experiment with a sample $w_0 \sim \mathcal{N}(\mu_w, \Sigma_w)$ from VI, a zero velocity field which corresponds to an identity transformation, and a random velocity field $w_0 \sim \mathcal{N}(0, I_{3N^3})$ sampled from a standard multivariate normal distribution. The first 100,000 samples from each chain are discarded to allow MCMC to reach the stationary distribution, and the two Markov chains are run for one million transitions each, from which we extract 500 total samples at regular intervals.

We observed no strong dependence of the result on the initialisation. The uncertainty values are visually identical to those in Figure 3. The mean voxel-wise discrepancy between the magnitudes of uncertainty estimates is approximately 0.1 mm, while the maximum is approximately 0.2 mm. This suggests that the Markov chains mix well.

Table 3: Mean and standard deviation of the number and percentage of voxels where the sampled transformation is not diffeomorphic as defined in Section 6.1. The values are based on 500 samples.

model	$ \det J_{\varphi^{-1}} \leq 0$	% ($\times 10^{-6}$)
VI	0.00 (0.04)	0.0 (0.2)
SG-MCMC	0.00 (0.00)	0.0 (0.0)
VoxelMorph	0.00 (0.00)	0.0 (0.0)

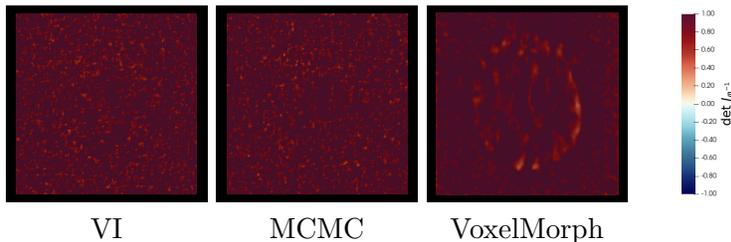


Figure 5: Jacobian determinant of a sample transformation from each model. Even though the output transformations are diffeomorphic, they are not convincingly smooth due to parametrisation of the approximate variational posterior of transformation parameters, whose covariance matrix is diagonal + low-rank in case of VI and diagonal in case of VoxelMorph. Middle slice of a 3D image in the axial plane.

6.3.3 COMPARISON OF THE OUTPUT OF SG-MCMC FOR NON-PARAMETRIC SVFs AND SVFs BASED ON B-SPLINES

One of the common features of recent state-of-the-art models for non-rigid registration that enable uncertainty quantification, e.g. Dalca et al. (2018) and Krebs et al. (2019), is an SVF transformation parametrisation as well as a diagonal covariance matrix of the approximate variational posterior of transformation parameters, which ignores spatial correlations. We showed in Section 6.3.1 that this assumption can lead to diffeomorphic transformations that are not smooth even in case of image registration that is not based on deep learning. Furthermore, previous work on uncertainty quantification in image registration made the assumption of independence between control points but not between directions and used transformation parametrisations that guaranteed smoothness in an implicit manner, e.g. B-splines (Simpson et al., 2012) or Gaussian RBFs (Le Folgoc et al., 2017).

To better understand the impact of transformation parametrisations on uncertainty quantification, we analyse the output transformations and uncertainty estimates when using sparse SVFs based on cubic B-splines. The control point spacing is set to two or four voxels along each dimension, which gives 3.64 mm or 7.28 mm. The SGLD step size is set to 5×10^{-2} , using the heuristic in Section 6.3.

In Figure 6, we show the output uncertainties and sample Jacobian determinants. Using SVFs based on B-splines results in smoother transformations. In fact, the implicit smoothness of B-splines helps both with the crude assumption of a diagonal covariance matrix in VI, which translates to a diagonal pre-conditioning matrix in SGLD, and with computational and memory efficiency. We also observed a positive impact on the image registration accuracy. For both control point spacings, SVFs based on B-splines are more accurate than non-parametric SVFs on the majority of structures. Despite comparable accuracy, the uncertainty estimates differ.

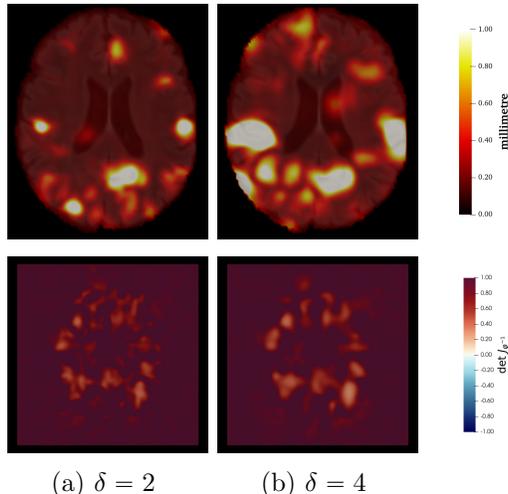


Figure 6: Uncertainty and the Jacobian determinant of transformations sampled from the models using SVFs based on cubic B-splines. The transformations are visibly smoother than in case of non-parametric SVFs. The uncertainty estimates are also shown to be highly dependent on the transformation parametrisation. The figure shows the middle axial slice of 3D images.

7. Discussion

7.1 Modelling assumptions

To draw samples from the true posterior of transformation parameters remains a difficult problem even with a large number of simplifying model assumptions. If the true posterior were approximately Gaussian, VI would provide a good approximation thereof and the lower uncertainty output by SG-MCMC would indicate that the output samples are auto-correlated even with subsampling. However, if the true posterior is not Gaussian, then the posterior output by VI is ill-fitting, and the samples output by SG-MCMC cover multiple modes near the mode that corresponds to VI, which makes a good case for the use of SGLD.

In practice, the quality of uncertainty estimates is also sensitive to the validity of model assumptions. These include coinciding image intensities up to the expected spatial noise offsets and ignoring spatial correlations between residuals. The first assumption is valid in case of mono-modal registration but the model can be easily adapted to other settings

through use of a different data loss. We manage the second assumption by use of virtual decimation (Groves et al., 2011; Simpson et al., 2012), which calculates the effective degrees of freedom of the residual map with stationary covariance by proxy of the correlation between neighbouring voxels in each direction.

We showed how good modelling choices, such as a careful choice of priors and transformation parametrisation can mitigate some of the issues caused by approximations necessary in practical applications. The main problem in inter-subject brain MRI registration remains accuracy but the trade-off between the quality of the transformation and registration accuracy can be managed effectively.

7.2 Runtime

The experiments were run on a system with an Intel i9-10920X CPU and a GeForce RTX 3090 GPU. Table 4 shows the runtime of the models. VI takes approximately 3 min to register a pair of images and produces 140 samples/sec, while SG-MCMC produces 25 autocorrelated samples/sec. The registration time for SG-MCMC is absent because it is difficult to pin down the mixing time, so a fair comparison of the relative efficiency of VI and SG-MCMC can be done only on the basis of the number of samples per second.

Due to lack of publicly available information we cannot directly compare the efficiency of our model to other Bayesian image registration methods. The speed of our model is an order of magnitude better than reported by Le Folgoc et al. (2017), while also being three- rather than two-dimensional. Thus, the proposed method based on SG-MCMC is very efficient given the Bayesian constraint. It is not as efficient as feed-forward neural networks, such as VoxelMorph. However, the proposed model is fully Bayesian and enables asymptotically exact sampling from the posterior of transformation parameters.

Table 4: Comparison of VI, SG-MCMC, and VoxelMorph vis-à-vis computational efficiency. In contrast to VI and VoxelMorph, SG-MCMC requires burn-in and produces samples which need to be further subsampled in order to avoid autocorrelation. For this reason, in case of SG-MCMC we report the sampling speed based on the time needed to draw 4,000 samples. Note that the ratio used to subsample the output of SG-MCMC may vary depending on the application (cf. Section 5).

model	training time	registration time	samples/sec
VI	—	3 min	1.4×10^2
SG-MCMC	—	—	< 1.0
VoxelMorph	38 h	55 ms	2.6×10^1

8. Conclusion

In this paper we presented a new Bayesian model for three-dimensional medical image registration. The proposed regularisation loss allows to adjust regularisation strength to the data when using an SVF transformation parametrisation, which involves a very large number of degrees of freedom. Sampling from the posterior distribution via SG-MCMC makes it possi-

ble to quantify registration uncertainty even for large images. The computational efficiency and theoretical guarantees regarding samples output by SG-MCMC make our model an attractive alternative for uncertainty quantification in non-rigid image registration compared to methods based on VI.

Acknowledgments

This research used UK Biobank resources under the application number 12579. Daniel Grzech is funded by the EPSRC CDT for Medical Imaging EP/L015226/1 and Loïc Le Folgoc by EP/P023509/1. We are very grateful to Prof. Wells and the Reviewers for their feedback during the review process.

Ethical Standards

The work follows appropriate ethical standards in conducting research and writing the manuscript, following all applicable laws and regulations regarding treatment of animals or human subjects.

Conflicts of Interest

We declare we do not have conflicts of interest.

References

- Vincent Arsigny, Olivier Commowick, Xavier Pennec, and Nicholas Ayache. A Log-Euclidean Framework for Statistics on Diffeomorphisms. In *MICCAI*, 2006.
- John Ashburner. A fast diffeomorphic image registration algorithm. *NeuroImage*, 2007.
- Brian B. Avants, Nicholas J. Tustison, Michael Stauffer, Gang Song, Baohua Wu, and James C. Gee. The Insight ToolKit image registration framework. *Frontiers in Neuroinformatics*, 8, 2014.
- Guha Balakrishnan, Amy Zhao, Mert R. Sabuncu, John Guttag, and Adrian V. Dalca. An Unsupervised Learning Model for Deformable Medical Image Registration. In *CVPR*, 2018.
- Guha Balakrishnan, Amy Zhao, Mert R. Sabuncu, John Guttag, and Adrian V. Dalca. VoxelMorph: A Learning Framework for Deformable Medical Image Registration. *IEEE Transactions on Medical Imaging*, 2019.
- Julian Besag. Comments on "Representations of knowledge in complex systems" by U. Grenander and MI Miller. *Journal of the Royal Statistical Society*, 56:591–592, 1993.

- Changyou Chen, David Carlson, Zhe Gan, Chunyuan Li, and Lawrence Carin. Bridging the Gap between Stochastic Gradient MCMC and Stochastic Optimization. In *AISTATS*, 2016.
- Adrian V. Dalca, Guha Balakrishnan, John Guttag, and Mert R. Sabuncu. Unsupervised Learning for Fast Probabilistic Diffeomorphic Registration. In *MICCAI*, 2018.
- Adrian V. Dalca, Guha Balakrishnan, John Guttag, and Mert R. Sabuncu. Unsupervised Learning of Probabilistic Diffeomorphic Registration for Images and Surfaces. *MedIA*, 57:226–236, 2019.
- Bob D. de Vos, Floris F. Berendsen, Max A. Viergever, Hessam Sokooti, Marius Staring, and Ivana Išgum. A deep learning framework for unsupervised affine and deformable image registration. *MedIA*, 52:128–143, 2019.
- Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the EM algorithm. *Annals of Statistics*, 27(1):94–128, 1999.
- Yilun Du and Igor Mordatch. Implicit Generation and Generalization in Energy-Based Models. In *NeurIPS*, 2019.
- Ben Glocker, Nikos Komodakis, Georgios Tziritas, Nassir Navab, and Nikos Paragios. Dense image registration through MRFs and efficient linear programming. *MedIA*, 12(6):731–741, 2008.
- Adrian R. Groves, Christian F. Beckmann, Steve M. Smith, and Mark W. Woolrich. Linked independent component analysis for multimodal data fusion. *NeuroImage*, 54(3):2198–2217, 2011.
- Daniel Grzech, Bernhard Kainz, Ben Glocker, and Loïc Le Folgoc. Image registration via stochastic gradient Markov chain Monte Carlo. In *UNSURE MICCAI*, 2020.
- Mohamed Hachama, Agnès Desolneux, and Frédéric J.P. Richard. Bayesian technique for image classifying registration. *IEEE Transactions on Image Processing*, 21(9):4080–4091, 2012.
- Mattias P. Heinrich. Closing the Gap Between Deep and Conventional Image Registration Using Probabilistic Dense Displacement Networks. In *MICCAI*, 2019.
- Mitchell Krupiarz Hill. *Learning and Mapping Energy Functions of High-Dimensional Image Data*. PhD thesis, UCLA, 2020.
- Firdaus Janoos, Petter Risholm, and William Wells. Bayesian characterization of uncertainty in multi-modal image registration. In *WBIR*, 2012.
- Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. Introduction to variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.

- Navdar Karabulut, Ertuğ Erdil, and Çetin Müjdat. A Markov Chain Monte Carlo based Rigid Image Registration Method. In *Signal Processing and Communications Applications Conference*, 2017.
- Samah Khawaled and Moti Freiman. Unsupervised deep-learning based deformable image registration: A Bayesian framework. In *Medical Imaging Meets NeurIPS*, 2020.
- Stefan Klein, Marius Staring, Keelin Murphy, Max A. Viergever, and Josien P.W. Pluim. Elastix: a toolbox for intensity-based medical image registration. *IEEE Transactions on Medical Imaging*, 29(1):196–205, 2009.
- Julian Krebs, Hervé Delingette, Boris Mailhé, Nicholas Ayache, and Tommaso Mansi. Learning a Probabilistic Model for Diffeomorphic Registration. *IEEE Transactions on Medical Imaging*, 38(9), 2019.
- Loïc Le Folgoc, Hervé Delingette, Antonio Criminisi, and Nicholas Ayache. Quantifying Registration Uncertainty With Sparse Bayesian Modelling. *IEEE Transactions on Medical Imaging*, 36(2):607–617, 2017.
- Matthew C. H. Lee, Ozan Oktay, Andreas Schuh, Michiel Schaap, and Ben Glocker. Image-and-Spatial Transformer Networks for Structure-Guided Image Registration. In *MICCAI*, 2019.
- Lihao Liu, Xiaowei Hu, Lei Zhu, and Pheng Ann Heng. Probabilistic Multilayer Regularization Network for Unsupervised 3D Brain Image Registration. In *MICCAI*, 2019.
- Jie Luo, Alireza Sedghi, Karteek Popuri, Dana Cobzas, Miaomiao Zhang, Frank Preiswerk, Matthew Toews, Alexandra Golby, Masashi Sugiyama, William M. Wells III, and Sarah Frisken. On the Applicability of Registration Uncertainty. In *MICCAI*, 2019.
- Stephan Mandt, Matthew D. Hoffman, and David M. Blei. Stochastic Gradient Descent as Approximate Bayesian Inference. *Journal of Machine Learning Research*, 18:1–35, 2017.
- Marc Modat, Pankaj Daga, M. Jorge Cardoso, Sebastien Ourselin, Gerard R. Ridgway, and John Ashburner. Parametric non-rigid registration using a stationary velocity field. In *Proceedings of the Workshop on Mathematical Methods in Biomedical Image Analysis*, 2012.
- Radford M. Neal. *MCMC Using Hamiltonian Dynamics*, chapter 5. CRC Press, 2011.
- J. W. Neuberger, A. Dold, and F. Takens. *Sobolev Gradients and Differential Equations*. Lecture Notes in Computer Science. Springer, 1997.
- Marc Niethammer, Roland Kwitt, and François-Xavier Vialard. Metric Learning for Image Registration. In *CVPR*, 2019.
- Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Learning Non-Convergent Non-Persistent Short-Run MCMC Toward Energy-Based Model. In *NeurIPS*, 2019.

- Frédéric J.P. Richard, Adeline M.M. Samson, and Charles A. Cuénod. A SAEM algorithm for the estimation of template and deformation parameters in medical image sequences. *Statistics and Computing*, 19(4):465–478, 2009.
- Petter Risholm, Steve Pieper, Eigil Samset, and William M. Wells. Summarizing and visualizing uncertainty in non-rigid registration. In *MICCAI*, 2010.
- Petter Risholm, James Balter, and William M. Wells III. Estimation of Delivered Dose in Radiotherapy: The Influence of Registration Uncertainty. In *MICCAI*, 2011.
- Petter Risholm, Firdaus Janoos, Isaiah Norton, Alex J. Golby, and William M. Wells. Bayesian characterization of uncertainty in intra-subject non-rigid registration. *MedIA*, 2013.
- Daniel Rueckert, L. I. Sonoda, C. Hayes, D. L. G. Hill, M. O. Leach, and David J. Hawkes. Nonrigid registration using free-form deformations: Application to breast MR images. *IEEE Transactions on Medical Imaging*, 18(8):712–721, 1999.
- Julia A. Schnabel, Daniel Rueckert, Marcel Quist, Jane M. Blackall, Andy D. Castellano-Smith, Thomas Hartkens, Graeme P. Penney, Walter A. Hall, Haiying Liu, Charles L. Truwit, Frans A. Gerritsen, Derek L. G. Hill, and David J. Hawkes. A Generic Framework for Non-rigid Registration Based on Non-uniform Multi-level Free-Form Deformations. In *MICCAI*, 2001.
- Sandra Schultz, Heinz Handels, and Jan Ehrhardt. A multilevel Markov Chain Monte Carlo approach for uncertainty quantification in deformable registration. In *SPIE Medical Imaging*, 2018.
- Sandra Schultz, Julia Krüger, Heinz Handels, and Jan Ehrhardt. Bayesian inference for uncertainty quantification in point-based deformable image registration. In *SPIE Medical Imaging*, 2019.
- Alireza Sedghi, Tina Kapur, Jie Luo, Parvin Mousavi, and William M. Wells. Probabilistic Image Registration via Deep Multi-class Classification: Characterizing Uncertainty. In *UNSURE MICCAI*, 2019.
- Zhengyang Shen, François Xavier Vialard, and Marc Niethammer. Region-specific diffeomorphic metric mapping. In *NeurIPS*, 2019.
- Jack Sherman and Winifred J. Morrison. Adjustment of an Inverse Matrix Corresponding to Changes in the Elements of a Given Column or a Given Row of the Original Matrix. *The Annals of Mathematical Statistics*, 21(1):124–127, 1950.
- I. J.A. Simpson, M. J. Cardoso, M. Modat, D. M. Cash, M. W. Woolrich, J. L.R. Andersson, J. A. Schnabel, and S. Ourselin. Probabilistic non-linear registration with spatially adaptive regularisation. *MedIA*, 26(1):203–216, 2015.
- Ivor J.A. Simpson, Julia A. Schnabel, Adrian R. Groves, Jesper L.R. Andersson, and Mark W. Woolrich. Probabilistic inference of regularisation in non-rigid registration. *NeuroImage*, 59(3):2438–2451, 2012.

- Ivor J.A. Simpson, Mark W. Woolrich, Jesper R. Andersson, Adrian R. Groves, and Julia Schnabel. Ensemble learning incorporating uncertain registration. *IEEE Transactions on Medical Imaging*, 32(4):748–756, 2013.
- Miroslava Slavcheva, Maximilian Baust, and Slobodan Ilic. SobolevFusion: 3D Reconstruction of Scenes Undergoing Free Non-rigid Motion. In *CVPR*, 2018.
- Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, et al. Uk biobank: An open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLOS Medicine*, 12(3), 2015.
- Demian Wassermann, Matt Toews, Marc Niethammer, and William Wells III. Probabilistic Diffeomorphic Registration: Representing Uncertainty. In *WBIR*, 2014.
- Max Welling and Yee Whye Teh. Bayesian Learning via Stochastic Gradient Langevin Dynamics. In *ICML*, 2011.
- whuber. What is the expected value of the logarithm of gamma distribution?, Oct 2018. URL <https://stats.stackexchange.com/questions/370880/what-is-the-expected-value-of-the-logarithm-of-gamma-distribution>.
- Miaomiao Zhang and P. Thomas Fletcher. Bayesian principal geodesic analysis in diffeomorphic image registration. In *MICCAI*, 2014.
- Miaomiao Zhang, Nikhil Singh, and P. Thomas Fletcher. Bayesian Estimation of Regularization and Atlas Building in Diffeomorphic Image Registration. In *IPMI*, 2013.
- Ruqi Zhang, Chunyuan Li, Jianyi Zhang, Changyou Chen, and Andrew Gordon Wilson. Cyclical stochastic gradient MCMC for bayesian deep learning. In *ICLR*, 2020.
- Lilla Zöllei, Mark Jenkinson, Samson Timoner, and William Wells. A marginalized MAP approach and EM optimization for pair-wise registration. In *IPMI*, 2007.

Appendix A.

In this appendix we show how to calculate the mean of the logarithm of the gamma distribution (whuber, 2018).

Let X be a random variable which follows the gamma distribution with the shape and rate parameters α and β , i.e. $X \sim \Gamma(\alpha, \beta)$. We are interested in the expected value of $Y = \log X$. If we assume that $\beta = 1$, then the PDF of X is given by:

$$f_X(x) = \frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x} \quad (16)$$

Note that $\Gamma(\alpha)$ is a constant and the integral of f_X must equal 1, so we have:

$$\Gamma(\alpha) = \int_{\mathbb{R}_+} x^{\alpha-1} e^{-x} dx \quad (17)$$

Let $x = \exp y$. This means that $\frac{dx}{x} = dy$, so the PDF of Y is given by:

$$f_Y(y) = \frac{1}{\Gamma(\alpha)} e^{\alpha y - e^y} dy \quad (18)$$

Again, because $\Gamma(\alpha)$ is a constant and the integral of the PDF of Y must equal 1, we have:

$$\Gamma(\alpha) = \int_{\mathbb{R}} e^{\alpha y - e^y} dy \quad (19)$$

Now, using Feynman's trick of differentiating under the integral sign, we see that:

$$\mathbb{E}(Y) = \int_{\mathbb{R}} y f_Y(y) dy = \frac{1}{\Gamma(\alpha)} \int_{\mathbb{R}} \Gamma(\alpha) y f_Y(y) dy \quad (20)$$

$$= \frac{1}{\Gamma(\alpha)} \int_{\mathbb{R}} \frac{d}{d\alpha} e^{\alpha y - e^y} dy = \frac{1}{\Gamma(\alpha)} \frac{d}{d\alpha} \int_{\mathbb{R}} e^{\alpha y - e^y} dy \quad (21)$$

$$= \frac{1}{\Gamma(\alpha)} \frac{d}{d\alpha} \Gamma(\alpha) = \frac{d}{d\alpha} \log \Gamma(\alpha) \quad (22)$$

$$= \psi(\alpha) \quad (23)$$

where ψ is the digamma function. Finally, the rate parameter β shifts the logarithm by $-\log \beta$. Therefore, the expected value of $\log X$ is given by:

$$\mathbb{E}[\log X] = \psi(\alpha) - \log \beta \quad (24)$$

Appendix B.

In this appendix we show how to calculate the KL-divergence term between the approximate variational posterior $q(w) \sim \mathcal{N}(\mu_w, \Sigma_w)$ and the prior $p(w)$, as well as the entropy $H(q)$ of the approximate variational posterior, both of which are needed to maximise the ELBO in Equation (12).

The KL-divergence can be evaluated in terms of the entropy $H(q)$ of the approximate variational posterior and the regularisation energy \mathcal{E}_{reg} :

$$\text{D}_{\text{KL}}(q \parallel p) = \int_w q_w(w) \log q(w) dw - \int_w q(w) \log p(w) dw = -H(q) + \langle \mathcal{E}_{\text{reg}} \rangle_q \quad (25)$$

where $\langle \cdot \rangle$ denotes the expected value.

The entropy $H(q)$ of the approximate variational posterior is calculated as follows:

$$H(q) = - \int_w q(w) \log q(w) dw = \frac{1}{2} \log \det(\Sigma_w) + \frac{1}{2} \langle (w - \mu_w)^\top \Sigma_w^{-1} (w - \mu_w) \rangle_q + \text{const.} \quad (26)$$

The first term on the left on the RHS is calculated using the matrix determinant lemma:

$$\det(\Sigma_w) = \det\left(\text{diag}(\sigma_w^2) + u_w u_w^\top\right) = \left(1 + u_w^\top \text{diag}(\sigma_w^{-2}) u_w\right) \times \det\left(\text{diag}(\sigma_w^2)\right) \quad (27)$$

To evaluate the other term, and in particular the precision matrix Σ_w^{-1} , we can use the Sherman-Morrison formula (Sherman and Morrison, 1950), which states that:

$$\Sigma_w^{-1} = \left(\text{diag}(\sigma_w^2) + u_w u_w^\top\right)^{-1} = \text{diag}(\sigma_w^{-2}) - \frac{\text{diag}(\sigma_w^{-2}) u_w u_w^\top \text{diag}(\sigma_w^{-2})}{1 + u_w^\top \text{diag}(\sigma_w^{-2}) u_w} \quad (28)$$