# HypoSVI: Hypocenter inversion with Stein variational inference and Physics Informed Neural Networks

Jonathan D. Smith [*], Zachary E. Ross [*], Kamyar Azizzadenesheli [†], Jack B. Muir [*]

## Abstract

We introduce a scheme for probabilistic hypocenter inversion with Stein variational inference. Our approach uses a differentiable forward model in the form of a physics-informed neural network, which we train to solve the Eikonal equation. This allows for rapid approximation of the posterior by iteratively optimizing a collection of particles against a kernelized Stein discrepancy. We show that the method is well-equipped to handle highly non-convex posterior distributions, which are common in hypocentral inverse problems. A suite of experiments is performed to examine the influence of the various hyperparameters. Once trained, the method is valid for any network geometry within the study area without the need to build travel time tables. We show that the computational demands scale efficiently with the number of differential times, making it ideal for large-N sensing technologies like Distributed Acoustic Sensing.

## 1  Introduction

Earthquake hypocenters represent the points in space and time at which earthquakes occur. They are a fundamental component of many downstream analyses in seismology, from seismic tomography to earthquake source properties. They are also used for real-time earthquake forecasting, such as during active sequences. Thus the ability to reliably estimate hypocenters and characterize their uncertainty is of major importance in seismology.

Determining a hypocenter from observations of seismic waves is a classic inverse problem in geophysics (Geiger, 1912; Thurber, 1985). More recently, Bayesian inference has been used for hypocenter inversion (Tarantola, 2004; Lomax *et al.*, 2000), in which prior information is combined with some observations to infer posterior distributions over the hypocentral parameters. In their most general form, the travel time solutions from ray theory are nonlinear in the hypocentral coordinates, which for this particular problem adds non-convexity to the posterior. Furthermore, and perhaps a more serious issue, is that the observations of seismic wave arrival times often contain significant errors resulting from the widespread adoption of automated picking algorithms (Ross *et al.*, 2018; Mousavi *et al.*, 2020). One strategy for adding robustness to hypocenter inversions has been to incorporate non-standard likelihood functions into the inverse problem, which has significantly improved the results but at the cost of creating highly non-convex posteriors. This makes many common techniques for performing Bayesian inference, such as Markov chain Monte Carlo or variational inference, ill-suited for this particular problem with the non-convex nature of the posterior expected to be computational slow.

---

[*]Seismological Laboratory, California Institute of Technology, Pasadena, CA, USA
[†]Lawson Computer Science Building, Purdue University, West Lafayette, IN, USA

Recent advances in deep learning have led to the development of physics-informed neural networks (PINNs), which are designed to learn solutions to partial differential equations (PDEs). Such approaches have a number of appealing properties that are not present with conventional approaches like finite difference methods, for example the solutions can be made differentiable, are often mesh-free, and can be rapidly calculated upon demand. These properties make PINNs well-suited as the forward model in an inverse problem, in particular since it often is desirable to take gradients of some objective function. Sampling and ensemble methods such as MCMC or variational inference typically require many evaluations of the forward model, so solving this with a PINN makes it more computationally tractable.

The rise of deep learning over the last decade has accelerated the development of novel approaches for performing Bayesian inference. One notable example is Stein variational inference (SVI), in which a collection of particles is iteratively optimized to approximate a target posterior (Qiang & Dilin, 2016). It is better suited than standard variational techniques at handling multi-modal distributions because the number of modes does not need to be known a priori; this results from a kernelized objective function that creates a natural repulsive force between the particles. SVI requires evaluating many gradients and as such, benefits from the differentiability of PINNs.

Our contributions to this paper are as follows: (1) we develop a framework for earthquake hypocenter inversion using Stein variational inference; (2) we incorporate a PINN trained to solve the Eikonal equation as a forward model; (3) we perform experiments on the hyperparameters of the inverse problem to characterize their effect on the solution; and (4) we benchmark the method against a catalog of earthquakes from Southern California.

## 2 Background

In this section, we provide background information on Stein variational inference as well as physics-informed neural networks.

### 2.1 Stein Variational Gradient Descent

For two random variables $x$ and $y$, let $p(x)$ denote the prior on $x$, $p(y|x)$ the likelihood function, and $p(x|y)$ the posterior over $x$ after observing (conditioning on) $y$. Using celebrated Bayes rule, these quantities are related as, $p(x|y) = \bar{p}(x)/Z$, where $\bar{p}(x) = p(y|x)p(x)$, and the normalization constant $Z = \int_x \bar{p}(x)dx$.

Let $\mathcal{H}$ denote a reproducing kernel Hilbert space on the domain $x$, with a positive definite reproducing kernel $\kappa$, endowed with the inter product $\langle \cdot, \cdot \rangle$ and the norm $\| \cdot \|_{\mathcal{H}}$. We further define $\mathcal{H}^d$, as a set of multivalued functions, with $d$ values, with the corresponding norm $\| \cdot \|_{\mathcal{H}^d}$, where for any $\boldsymbol{f} = [f_1, f_2, \ldots, f_d] \in \mathcal{H}^d$ we have $f_i \in \mathcal{H} \ \forall i \in [1, 2, \ldots, d]$.

For a function $\boldsymbol{f} \in \mathcal{H}^d$, we define Stein's operator endowed with $\mathcal{H}^d$ and $p$ as,

$$(\mathcal{A}\boldsymbol{f})(x) = \boldsymbol{f}(x)\nabla_x \log p(x)^\top + \nabla_x \boldsymbol{f}(x). \tag{1}$$

We further define the kernelized Stein's discrepancy between two distributions $p$ and $q$ using $\mathcal{H}^d$ is as follows,

$$\mathcal{D}(q, p) := \max_{\boldsymbol{f} \in \mathcal{H}^d \ s.t., \ \|\boldsymbol{f}\|_{\mathcal{H}^d} \leq 1} E_{x \sim q} \left[ trace \left( \mathcal{A}\boldsymbol{f}\left(x\right) \right) \right]^2. \tag{2}$$

This discrepancy equates to zero when $p = q$. Fortunately, the maximization in Eq. 2 has a closed-form solution $\mathcal{D}(q, p) = \|\boldsymbol{f}_q^*\|_{\mathcal{H}^d}$ where $\boldsymbol{f}_q^* := E_{x \sim q}[\mathcal{A}\kappa(x, \cdot)]$ is the maximizer.

Now consider the Kullback–Leibler between $q$ and $p$, i.e., $KL(q, p)$ divergence. We aim to find a gradient direction $g$, a function, such that $g$ maximally reduces the KL divergence. Using $g$, we can use gradient descent with learning rate $\alpha$ and update $q \leftarrow q + \alpha g$ to reduce the $KL(q, p)$, and make the $q$ closer to $p$. It is known that for the kernelized Stein's discrepancy, the direction $g \in \mathcal{H}$ that provides the direction of maximal change is $g := \boldsymbol{f}_q^*$ (Qiang & Dilin, 2016). In the following, we provide an update rule to update $q$ and approximate the posterior $p$ given observed data.

We represent $q$ with a set of particles, i.e., an average of many delta Dirac measures. $\{x_i\}_{i=1}^n$ where q is the distribution of this particle. In the following, we update $q$, and make it closer to $p$ by moving the particles. Therefore, for the update direction $\boldsymbol{f}_q^* = E_{x \sim q}[\mathcal{A}\kappa(x, \cdot)]$, at each point $x$, we have,

$$\boldsymbol{f}_q^*(x) = \sum_{i=1}^n \mathcal{A}\kappa(x_i, x)$$
$$= \sum_{i=1}^n [\kappa(x_i, x)\nabla_{x'} \log p(x')|_{x'=x_i} + \nabla_{x'}\kappa(x', x)|_{x'=x_i},$$

with the updating rule given by,

$$x_i^{l+1} \leftarrow x_i^l + \alpha_l \boldsymbol{f}_q^*(x_i^l). \tag{3}$$

Here $\alpha_l$ is the step size at the $l$th epoch. For the choice of kernel, we deploy the celebrated Radial Basis Function (RBF), $\kappa(x', x) = \exp(-\frac{1}{h}\|x - x'\|^2)$, with $h$ representing the width of kernel, for its empirical and universal approximation properties. As discussed above, the update in Eq. 3, updates $q$ (through updating the particles distribution) at each time step to make it closer to $p$ in the Stein's discrepancy sense.

## 2.2   Physics-informed Neural Networks for Ray Tracing

In solving inverse problems for earthquake hypocenters, the most common approach is to use a ray theoretical forward model to calculate the expected travel times, $T$, for seismic waves propagating from a given source location to a receiver location. In heterogeneous 3D Earth models, the Eikonal equation is often solved to determine $T$ (Rawlinson & Sambridge, 2005),

$$\|\nabla_r T_{s \to r}\|^2 = \frac{1}{V(\vec{x}_r)^2} = S(\vec{x}_r)^2 \tag{4}$$

where $\|\cdot\|^2$ is the Euclidean norm, $T_{s \to r}$ is the travel-time through the medium from a source location $s$ to a receiver location $r$, $V_r$ is the velocity of the medium at the receiver location, $S_r$ is the slowness of the medium at the receiver location, and $\nabla_r$ the gradient at the receiver location.

Smith *et al.* (2020) developed a PINN approach to solving the factored Eikonal equation (EikoNet), which trains a deep neural network to calculate the travel-time between any two points in a 3D medium for a given velocity model, satisfying the additional boundary condition that the travel-time at the source location equals zero, $T_{s \to s} = 0$. We leverage a factored eikonal formulations to mitigate the strong singularity affects at the source location, representing the travel-time as a deviation from

a homogeneous medium with $V = 1$ (Treister *et al.*, 2016). The factored travel-time form can then be represented by:

$$T_{s\rightarrow r} = T_0 \cdot \tau_{s\rightarrow r} \tag{5}$$

where $T_0 = \|\vec{x_r} - \vec{x_s}\|$, representing the distance function from the source location, and $\tau$ the deviation of the travel-time field from a model travel-time with homogeneous unity velocity. Substituting the formulation of equation 5 into equation 4 and expanding using the chain rule, then the velocity can be represented by;

$$V(\vec{x_r}) = \left[ T_0^2 \| \underset{r}{\nabla} \tau_{s\rightarrow r} \|^2 + 2\tau_{s\rightarrow r} (\vec{x_r} - \vec{x_s}) \cdot \underset{r}{\nabla} \tau_{s\rightarrow r} + \tau_{s\rightarrow r}^2 \right]^{-\frac{1}{2}}. \tag{6}$$

They leveraged the analytical differentiability of neural networks to solve the factored Eikonal equation from scratch, without the use of finite-difference solutions during training. Once trained, a network describing the travel-time between any source-receiver pair can be represented by:

$$T_{s\rightarrow r} = f_\theta (\vec{x_s}, \vec{x_r}) \tag{7}$$

where $T_{s\rightarrow r}$ is the travel-time between the source location $\vec{x_s}$ and receiver location $\vec{x_r}$, and $f$ is the neural network with weights and biases given by $\theta$.

EikoNet has several properties that are mathematically advantageous in solving inverse problems. First, the solutions to the Eikonal equation are mesh-independent, i.e. they are not discretized on a fixed grid and can be evaluated at truly any point within the 3D medium. Second, the network is a forward model that is analytically differentiable, which allows for gradient-based methods to be efficiently employed to calculate a downstream objective function. Third, by approximating the Eikonal equation with a deep neural network, the optimization part of the inverse problem is easily solved with graphical processing units (GPUs).

## 3 Methods

### 3.1 Overview

We now present an approach for probabilistic hypocenter inversion that uses a PINN as a forward model and SVI to approximate the posterior distribution. The method consists of several primary steps:

1. An EikoNet model is trained for a given Earth velocity model to solve the Eikonal equation. This is performed for both P-waves and S-waves.

2. A collection of particles is randomly initialized throughout the geographic study area. These represent preliminary hypocenter locations.

3. Travel times are calculated with EikoNet from each particle to every receiver with an observation.

4. The synthetic travel times are used together with the data to calculate a kernelized Stein discrepancy (loss function).

5. The gradients of the loss are calculated with automatic differentiation and used to collectively update the particles' locations.
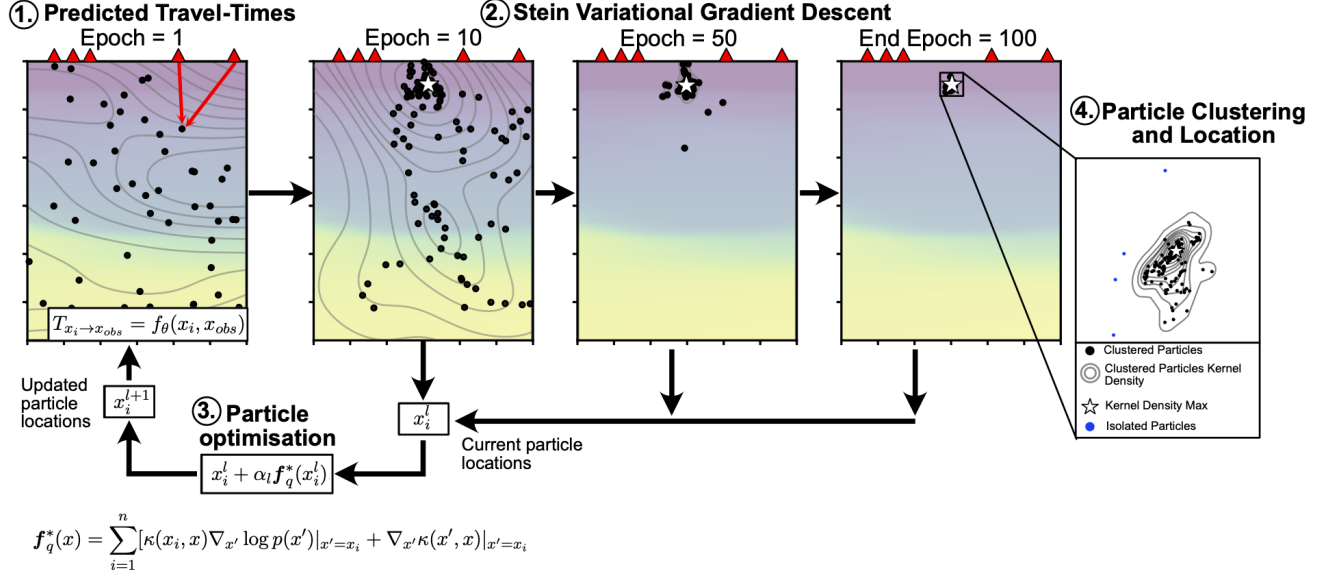
4

**①. Predicted Travel-Times**
Epoch = 1

**②. Stein Variational Gradient Descent**
Epoch = 10          Epoch = 50          End Epoch = 100

**④. Particle Clustering and Location**

$T_{x_i \to x_{obs}} = f_\theta(x_i, x_{obs})$

Updated particle locations    $x_i^{l+1}$

**③. Particle optimisation**

$x_i^l + \alpha_l \boldsymbol{f}_q^*(x_i^l)$

$x_i^l$    Current particle locations

$$\boldsymbol{f}_q^*(x) = \sum_{i=1}^{n} [\kappa(x_i, x)\nabla_{x'}\log p(x')|_{x'=x_i} + \nabla_{x'}\kappa(x', x)|_{x'=x_i}$$

- ● Clustered Particles
- ◎ Clustered Particles Kernel Density
- ☆ Kernel Density Max
- ● Isolated Particles

Figure 1: Overview of the inversion procedure. **Panel 1** represents the travel-time between the observational locations, red triangles and given by $x_{obs}$, and the particle locations, black dots and given by $x_i$. The particle kernel density are given by a series of gray contours with the particle with the maximum kernel density given by the white star. **Panel 2** shows the distribution of the particle locations changing with the Stein Variational Gradient descent. **Panel 3** represents the particle location optimisation with the step-size given by $\alpha_i$ and optimisation direction $\boldsymbol{f}_q^*(x)$. **Panel 4** represents the clustering procedure to determine optimised location, represented by the kernel density max, and the location uncertainty given by the clustered particle kernel density.

6. Steps 3-5 are repeated until convergence. The final collection of particle positions will approximate the posterior distribution of the hypocenter.

7. Uncertainty estimates are extracted from the particles using kernel density methods.

Next, we provide a detailed discussion of each stage of the procedure, with the outline of the inversion given in Figure 1.

## 3.2   Constructing the forward model

Throughout this study we train EikoNet travel-time models using a set of constant training parameters and network architecture as described in Smith *et al.* (2020) and supplied in Table 1. A model region is defined spanning our Longitude, Latitude, depth regional of interest, with xmin and xmax locations as $[117^o30'W, 32^o30'N, -2km]$ and $[115^o30'W, 34^o30'N, 50km]$ respectively. The grid is projected to a UTM coordinate system, with random source-reciever locations selected within the UTM model space. These points represent the training locations, with different velocity models discussed below.

In many earthquake location procedures the complex geometry of the subsurface is poorly understood, with the assumption that lateral variations in velocity are negligible compared to

Table 1: EikoNet training paradigm used to learn velocity models

| Parameter | Value |
|---|---|
| Dataset Size | $1 \times 10^6$ |
| Validation Fraction | 0.1 |
| Batch Size | 752 |
| Optimizer | ADAM (+ scheduler) |
| Learning Rate | $1 \times 10^{-5}$ |
| Sampling Type | Weighted Random Distance |
| Sampling Type Bounds | $[0.1, 0.9]$ |
| Domain Normalization | Offset Min-Max Normalization |
| Network Architecture | Dense $6 \rightarrow 32$ + Dense $32 \rightarrow 512$ + $10\times$ Residual Blocks $512 \rightarrow 512$ + Dense $512 \rightarrow 32$ + Dense $31 \rightarrow 1$ ELU Activation Function |

velocity variations in depth. As such one-dimensional velocity structure describing how the velocity changes with depth are specified. These models typically have independent velocity structure defined for both the P-wave and S-wave arrivals, or a scaling relationship of Vp/Vs. It is important to understand how reliabiable these methods are for location procedures such as HypoSVI, as this would be a typical starting model for many use cases. In addition, understanding of the computational demand for training more simplistic travel-time models, informs the feasibility of the method on typical computational systems. We investigate these problems for our region of interest by training EikoNet travel-time models from the Vp and Vs velocity structure shown by the blue dots in Figure 2a. We interpolate the velocity at the point locations as the linear interpolation of the observed velocity values. Two independent EikoNet neural networks are trained independently for the Vp and Vs velocity structure using the network parameters specified in Table 1. The training of each model took 10 epochs, with roughly a 10 minutes training time on a Nvidia V100 GPU and $\sim 20$ minutes on a free Colab GPU (either a Nvidia K80,T4 or P100). Once trained the travel-time models can be validated by comparing the imposed observed velocity to predicted velocity, determined as the analytical gradient of the travel-time over the neural network, for a series of $1 \times 10^5$ source-receiver pairs within the three-dimensional domain. Figure 2 outlines the comparison of the observed velocity structure and the predicted velocity, with the variance of the predicted velocity within $0.05 km/s$ of the observed values. The consistent velocity structure and low computational overhead shows that this method is viable regardless of the available computational infrastructure.

In more well studied regions prior geophysical datasets and analysis could have been leverage to gain a better insight into the complex subsurface and therefore the velocity structure. The velocity structure for our region of interest, that of Southern California, has been densely studied with a compilation of direct velocity model estimates (Shaw *et al.*, 2015; Lee *et al.*, 2014), used to construct detailed subsurface velocity models defined under the group project 'Southern California Earthquake Centre Community Velocity Model' (SCEC-CVM, https://www.scec.org/research/cvm). The current versions of these velocity models encompass the SCEC-CVM-H (Shaw *et al.*, 2015, version 15.1.0)) and SCEC-CVM-S (Lee *et al.*, 2014, version 4.26). Incremental improvements are made to the CVM-H, but the CVM-S has seen no advances since the current version. The SCEC-CVMs are computed using numerous datasets (Süss & Shaw, 2003), encompassing a detailed subsurface three-
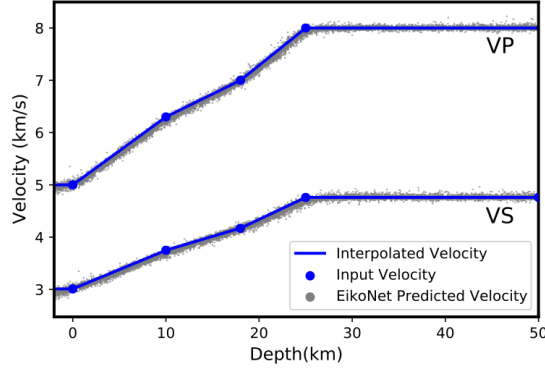
Figure 2: EikoNet trained travel-time formulation for a one-dimensional velocity model only changing in depth (Z). Each curve represent a different model computed for both the P-wave velocity structure (VP) and S-wave velocity structure (VS). Blue points represent the user defined velocity values at depths, blue lines the linear interpolation of velocity between points. Gray points represent the predicted velocity from EikoNet for $1 \times 10^5$ randomly selected points for each of the velocity models.

dimensional velocity structure from moho surface, basement surface and topological/bathymetric surfaces. Using the SCEC-CVM-H model we train EikoNet models to determine the travel-time within the complex 3D velocity structures. The models are trained on $1 \times 10^6$ randomly selected source-receiver points within the domain, with example slice at Longitude= $115^o30'W \pm 1.8'$ given for the P-wave and S-wave velocity structure in Figure 3a and 3b respectively. The EikoNet models once trained represent the travel-time and predicted velocity between any points, as such we show the recovered velocity model colourmap and travel-time contours (at $2s$ spacing) for a earthquake source at $[115^o30'W, 31^o12', 25km]$ on a receiver grid as separation $[Latitude, Depth] = [0.05^o, 0.5km]$ with Longitude= $115^o30'W$. This example shows consistent agreement between the observed and predicted velocity models, able to reconcile the sharp velocity contrasts which create deflection in the travel-time fields. This example demonstrates the viability of this method in complex 3D velocity structures.

## 3.3   Inverse problem formulation

An earthquake hypocenter, $m$, is composed of three spatial coordinates, $[x, y, z]$, and the origin time, $t_o$. Most commonly, the data used to locate earthquakes are measured times of seismic P- and S-wave arrivals ("phase picks") over a network seismic instruments. These phase picks define a set of absolute arrival time observations $d = T_{obs}$, where $d \in \mathbb{R}^N$. In a Bayesian framework (Tarantola, 2004; Lomax $et\ al.$, 2000), inference on $m$ is performed by combining prior knowledge together with the observations,

$$p(m|d) = Z^{-1}p(d|m)p(m) \tag{8}$$

where $p(m|d)$ is the posterior distribution, $p(m)$ is the prior distribution, $p(d|m)$ is the likelihood, and $Z$ is a constant.
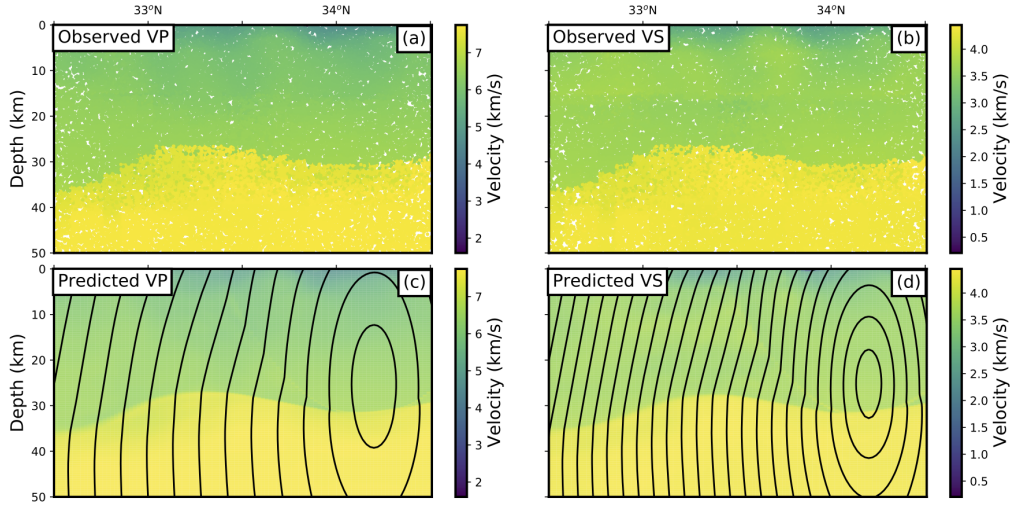
Figure 3: EikoNet trained travel-time formulation for the complex three-dimensional velocity of SCEC-CVM-H. Plots represent a slice in the three-dimensional structure taken at Longitude= $115^o30'W$. (a) and (b) represent the P-wave (VP) and S-wave (VS) velocity structures for the training points within $pm1.8'$ of the longitude slice and within the Latitude and Depth domain of the model space. (c) and (d) represent the predicted velocity structure colourmap and predicted travel-time contours, at $2s$ intervals, for the P-wave and S-wave EikoNet models.

A simple example of $p(d|m)$ for hypocenter inversion is,

$$p(d|m) = \exp\left(-\frac{1}{2}\sum_{obs_i}\frac{[T_{obs} - T_{pred}]^2}{\sigma_i^2}\right) \tag{9}$$

where $\sigma_i$ is an estimate of uncertainty and,

$$T_{pred} = t_o + f_\theta\left(\vec{x_s}, \vec{x_r}\right), \tag{10}$$

is a nonlinear forward model, i.e. a solution to the Eikonal equation plus the origin time. Thus, the forward model in this problem is a physics-informed neural network. Since $\vec{x_s}$ is included as an input to the neural network, this allows for downstream gradients to be taken with respect to it.

More recently, a likelihood function based on the Equal Differential Time method (Lomax $et$ $al.$, 2000, EDT) has seen increasing usage. The EDT likelihood builds differential times from all pairs of phases, and in the process, decouples origin time, $t_o$ from the spatial coordinates of the hypocenter. The formulation is given by;

$$p(d|m) = \left[\sum_a\sum_b\frac{1}{\sqrt{\sigma_a^2 + \sigma_b^2}}\exp(A)\right]^N, \tag{11}$$

$$A = -\frac{\left[\left(T_{obs(a)} - T_{obs(b)}\right) - \left(T_{pred(a)} - T_{pred(b)}\right)\right]^2}{\sigma_a^2 + \sigma_b^2}, \tag{12}$$

where $a$ and $b$ are different phase arrival time observations, $\sigma$ is a phase-dependent estimate of uncertainty, and $N$ is the total number of differential times. In addition to reducing the number of latent variables by one, this formulation acts to minimise the effects of outliers, which are particularly common with automated picking algorithms. This robustness results from the fact that in the EDT likelihood, the errors are combined in an additive manner, rather than multiplicative as typical in Bayesian inference problems. Each term in Eq. 11 produces a hyperbolic error surface that decays like a Gaussian in the direction normal to each point on the hyperbola. Thus, Eq. 11 can be viewed as producing a stack of hyperbolas with relatively limited intersection, which creates robustness in the presence of strong outliers. However, the downside is that it results in posterior distributions that are highly non-convex, making MCMC methods and standard variational inference schemes difficult to use for this problem (Lomax $et$ $al.$, 2000). The origin time is reintroduced by using the optimised earthquake location to determine the predicted origin times to each of the observational locations, determining the origin time as the median of the predicted origin times. The uncertainty is then defined by the median absolute deviation (MAD) from the predicted origin time. We use a uniform prior, $p(m)$, with samples selected within the model domain specified in the Eikonal physics informed neural network.

The uncertainty in the posterior distribution is assigned as a combination of the observational, $\sigma_{obs}$, and forward model uncertainty, $\sigma_{pred}$, given as

$$\sigma^2 = \sigma_{obs}^2 + \sigma_{pred}^2. \tag{13}$$

The observational uncertainty represents uncertainty in each of the observational times, with an expected standard deviation for each observation time supplied by the user. This value is then converted to a variance to define $\sigma_{obs}$ for each observation. The forward model uncertainty is constructed as a function of the predicted travel-time for each of the observational locations (similar to that given in Lomax $et$ $al.$ 2000 for LOCGAU2), given by

$$\sigma_{pred} = \begin{cases} \sigma_{min}, & \text{for } \sigma_f T_P < \sigma_{min} \\ \sigma_{frac}T_{pred}, & \text{for } \sigma_{min} \leq \sigma_f T_P \leq \sigma_{max} \\ \sigma_{max}, & \text{for } \sigma_f T_P > \sigma_{max} \end{cases} \tag{14}$$

where $\sigma_f$ is the fraction of the travel time to use as uncertainty, bounded within the max and min uncertainties specified by $\sigma_{min}$ and $\sigma_{max}$ respectively. Throughout this work we use the $[\sigma_f, \sigma_{min}, \sigma_{max}] = [0.1, 0.1s, 2.0s]$, discussing the effects of these parameters on synthetic testing within Section 4.5.

A Stein variational gradient descent procedure is used to optimise for the Equal-Differential Time posterior. We use a Gaussian kernel for the SVI procedure, also known as radial basis function (RBF) kernel, for its practical and universal approximation properties. First, we initialize $N$ particles randomly using a uniform prior over the 3D study area . For each of these particle locations, we calculate corresponding travel times using EikoNet forward model (Section 1. of Figure 1), evaluating the posterior (to within the normalization constant $Z$), and determine the kernelized Stein discrepancy (Section 2 of Figure 1). Then, we calculate the gradients of this loss function particle-wise with respect to the hypocentral coordinates using automatic differentiation, which is possible due to the differentiablity of the PINN (Section 3 of Figure 1). We use these gradients together with the ADAM optimizer (Kingma & Ba, 2014) to update the particle locations until convergence, where the optimal hypocentral location is consistent across multiple epochs. Supplementary Video 1 demonstrates the convergence for the example outlined in Figure Figure 1.

The next step is to extract summary statistics from the posterior distribution. As mentioned previously, the posterior is typically strongly non-convex due to the EDT likelihood function, although many of the local extrema are effectively negligible in amplitude. Therefore, we aim to determine the dominant cluster of particles representing the main peak of the posterior. This is achieved by using the DBSCAN clustering technique (Hahsler *et al.*, 2019, Section 4 of Figure 1) to identify high-density clusters of particles, with the cluster with the largest number of particles is used as the dominant cluster. Once the dominant cluster is identified, we apply kernel density methods using Gaussian kernels to estimate the MAP and quantify the location uncertainty from its covariance matrix. We discuss hyperparameter tuning for DBSCAN (Hahsler *et al.*, 2019) in a later section.

# 4  Experiments

## 4.1  Method validation

In this section we first demonstrate the earthquake inversion scheme on a series of synthetic tests. We construct a catalogue of synthetic earthquake locations across the region, determining the travel-time to a grid of observation points at fixed elevation of $0km$, before applying a $0.05s$ uncertainty in the synthetic phase arrival and inverting to determine the earthquake location and uncertainty. The earthquake locations are at a fixed latitude and depth of $33^o36'N$ and $5km$ respectively, with longitude varying from $116^o36'W$ to $116^o24'W$ at $6'$ separations. The recovered optimal hypocentre and location uncertainty are then compared with the imposed earthquake locations and an expected 2-std contour from a grid-search approach. We vary the possible user defined parameters with the optimised parameters given in Table 2 and earthquake locations in Figure 4. However, we expect that these parameters will need to be varied somewhat depending on the exact application, for example if the error models or network geometry are changed significantly. As such we recommend that initial synthetic testing is undertaken before real data is inverted. Outlined below are discussions on how each hyperparameter affects the recovered locations for this study, with corresponding Supplementary Figures S9-S11.

## 4.2  Number of particles

The number of particles used in the Stein Variational Gradient Descent is of great importance for the resolution of the resolved earthquake location. If the number of particles is too small then the particles density is unable to adequately represent the posterior distribution. However, a large number of particles would have increasing computational demand on the inversion procedure and is intractable for large earthquake catalogues. We specify a optimal number of samples equal to 150 and find that an increase in the number of particles does not provide additional information on the the earthquake location, but reducing the number
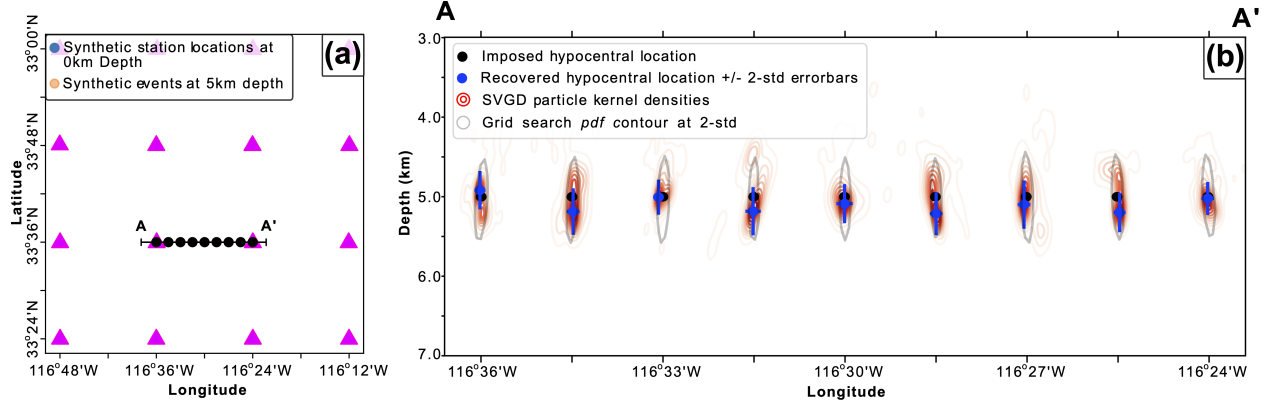
Figure 4: Synthetics earthquake location recovery for a synthetic seismic array. (a) represents a map view of the synthetic earthquake locations and synthetic stations locations. Black points represent the synthetic earthquake location Latitude/Longitudes, at a fixed depth of $5km$. Pink triangles represent the synthetic station locations, at a fixed depth of $0km$. Black line represent a cross section at a fixed latitude, with the cross section given in (b). (b) represents the imposed earthquake locations, black points, recovered optimal location with errors bars, blue dots, posterior determined by the particle density, red contours, and grid search derived posterior at 2-std, gray line.

Table 2: HypoSVI parameters used in earthquake location techniques

| Parameter | Value |
|---|---|
| Number of Particles | 150 |
| Number of Epochs | 175 |
| Observation Weights | $[0.1, 0.1s, 2.0s]$ |
| Radial Basis Function | 15 |
| Location Uncertainty | $[0.8km, 3]$ |

of particles greatly effects posterior. Additional plots for variations of number of particles, with remaining parameters set equal to Table 2, is given in Supplementary Figure S9.

## 4.3 Number of Epochs

The number of epochs effects how well the particles density represents the posterior for the earthquake location. Figure 1b demonstrates that the posterior drastically changes in the early epochs, but once it converges there is little to no change in the recovered posterior. However, an increase in number of epochs effects the computational time, which over a large number of earthquake locations could have a large effect. We optimise the number of epochs to determine when the earthquake locations and location uncertainty is consistent between epochs.

## 4.4 Influence of the kernel

The RBF kernel can be represented by $\kappa(x, x') = \exp(-\frac{1}{h}\|x - x'\|^2)$, where $h$ is the shape parameter and $x$ the pairwise particle difference. As this term acts as a repulsive force in the SVI procedure increasing the $\frac{1}{h}$ term has the effect of increasing the minimum distance between particles locations. Understanding the trade off for the shape parameter is important as larger values could effect on the recovered posterior. Qiang & Dilin (2016) defined a dynamic shape parameter with the value changing depending on $h = med^2/\log n$, where $med$ is the median distance between pairwise particles, with the definition $\sum_j k(x_i, x_j) \approx n \exp(-\frac{1}{h}med^2) = 1$ demonstrating that for each $x_i$ the contribution from its own gradient and the influence from other points balances out. We investigate the variation of the RBF shape parameters on the recovered synthetic earthquake locations finding that parameter has little effect the recovered optimal hypocentral location, with minor variations of the recovered posterior for static values between $2 - 20$ and that of the dynamic shape parameter (Supplementary Figure S10). We decided to use a static shape parameter of 15, to mitigate any difference that could occur to the posterior from mulitple run of the same observations for a dynamic shape parameter.

## 4.5 Error models

The total uncertainty assigned to the inverse problem is a combination of the picking uncertainty and the forward model uncertainty due to the velocity structure. As described previously, we follow Lomax *et al.* (2000) and characterize the uncertainty in the forward model as a fraction of the travel time. This is a reasonable choice as the uncertainty in the predicted travel times is expected grow in proportion to the travel time. In our hyperparameter investigation we found that a fraction of 0.1 should be used, as lower values lead to significant mis-location of the recovered events (Supplementary Figure S11). The upper and lower bounds to the allowed error has less of an effect on our synthetic testing, which we attribute to the synthetic station locations being regularly spaced. For observational data that is clustered spatially the upper and lower bounds could be of great importance and should be investigated with synthetic examples for the specific network geometry.

## 4.6 Clustering hyperparameter

The posterior for earthquake location is non-convex due to the EDT likelihood function and we aim to determine the dominant cluster of particles representing the main peak of the posterior. This is achieved by using the DBSCAN clustering technique (Hahsler *et al.*, 2019) to identify high-density clusters of particles. We investigate the variation of the two dominant parameters defining the DSCAN clustering technique, the maximum distance between two samples for them to be considered as in the neighbourhood of the other and the minimum number of samples per cluster (Supplementary figure S12). The definition of the minimum distance is crucial for effectively sampling the dominant peak of the distribution, with a minimum distance too small possibly subsampling the peak as multiple clusters and a value too large including additional local peaks in the the posterior. We found that the minimum distance must be defined large enough to remove

Table 3: HypoSVI computational cost on a Nvidia V100 GPU with different number of observations and corresponding differential time pairs. The remaining parameters used in this synthetic test are given Table 2

| # of Observations | # of Differential Times | Time per Event(s) |
|---|---|---|
| 32 | 496 | 6 |
| 128 | 8128 | 17 |
| 512 | 130816 | 64 |
| 1024 | 523776 | 155 |
| 1408 | 990528 | 247 |
| 1728 | 1492128 | 336 |
| 2028 | 2055378 | 439 |

the effects of subsampling the dominant peak in the posterior, in this use case we found $0.5km$ sufficient, although has little variation in the recovered kernal density function function until the mimum distance is expanded to at least $10km$. In addition, we find little effect of the kernel density function with changes in the minimum number of particles per cluster, something that is to be expected when the dominent cluster typically comprises $> 90\%$ of the particles for these synthetic inversions. We conclude that a minimum cluster separation of $0.5km$ should be used for these large regional scale problems and is insenstive to the low values in the minimum number of samples per cluster, which we use a minimum cluster comprising 3 particles.

## 4.7  Computational demands

The number of observations going into a inversion affects the compute time, as each observation requires predicted travel-time formulations from EikoNet and gradients to be computed . Here, we investigate the computational cost of the inversion procedure while increasing the number of observations. We replicate an increasing number of observations by copying the synthetic station deployment locations multiple times, labelling them as different station names but comprising the same arrival times. This synthetic testing was chosen to minimising the changing effect on the location estimate, which would occur if additional synthetic station locations are provided. All other location hyperparameters are fixed at values given in Table 2. The earthquake locations are then determined for the varying number of observations and the total number of pairwise differential times, with the average computational time for a Nvidia V100 shown in Table 3. The computational time even for the 2048 observations, 2055378 differential times, only takes $439s$ per event. These synthetic tests demonstrate that this approach is computationally scale-able with computational time increasing as a linearly in a log-log space of computational time vs number of observational differential times.

# 5   Case Study: Application to earthquake swarms in Southern California

## 5.1  Background

To further validate the developed method, we apply it to real earthquakes occurring within the Southern California region, with region defined in Section 3.2. This study area was chosen as it encompasses a large seismic network and complex 3D regional velocity structures (Allam & Ben-Zion, 2012). We used the detections and phase picks from the open source Southern California Earthquake Data Centre (SCEDC) phase arrival observational catalogue, for the fist $10k$ events starting 2019-01-01. The events and phase picks used have all been manually reviewed by analysts at the Southern California Seismic Network (Hutton *et al.*, 2010).

## 5.2    Earthquake Location comparisons with NonLinLoc

We infer hypocenters for the $10k$ earthquakes using two different velocity models (1D and 3D cases, described in Section 3.2). The hyperparameters used for the inversions are outlined in Table 2 with detailed explanation of the reasoning behind the parameter definition outlined in Section 4. The catalogues are generated on a Nvidia V100 GPU with an average of $5s$ per event, varying depending on the number of observations in the inversion procedure, with on average $\sim 30$ observations per event. Since the calculation of travel-times from EikoNet is independent on the complexity of the velocity model (once the network has been trained), the processing takes equal time for both the 1D and 3D trained models. Example inversions for three events are shown in Figure 5.

To understand the validity of our location technique we compare our earthquake catalogue, with a catalogue determined using the conventional earthquake location software, NonLinLoc. NonLinLoc is a non linear earthquake technique leveraging finite-difference travel-time solutions; Gaussian or equal-differential likelihood functions; and, likelihood estimations schemes using oct-tree, grid-search or Markov Chain Monte Carlo (MCMC). Travel-times are computed by solving the eikonal using a finite-difference approach outlined in Podvin & Lecomte (1991). For a 1D velocity structure, only varying in depth, the package computes the travel-times as an radial 2D finite-difference travel-time model that depends on the radial distance from the observation point and the depth, saving these as independent travel-time look-up tables. In contrast, for complex three-dimensional velocity structures the travel-times are computed for a user defined gridded series of receiver locations, with each observation saved as a separate travel-time look-up table. Since the storage and computational requirements for a NonLinLoc using the complex 3D velocity for a very high resolution location grid, this method was intractable as it would return large gridding artifacts to the recovered earthquake locations and predicted location uncertainty, which are not directly comparable to the non-gridded solutions of the HypoSVI. Instead we compare the HypoSVI and NonLinLoc locations using the one-dimensional velocity structure, with the NonLinLoc travel-time and initial location grids resolved to 1km and 2km receptively. The location is determined using a Equal-Differential Travel-Time (EDT) likelihood function and octree sampling technique. The location uncertainty of the recovered NonLinLoc catalogue is determined as the standard error in X,Y,Z to 2-std using the diagonal of the covariance matrix. The remaining NonLinLoc user parameters are given in the full control file in the Supplementary Material. The HypoSVI earthquake catalogues for the 1D and 3D velocity structures are given in Figure 6a-b and 6c-d respectively.

For comparison we derive a NonLinLoc catalogue for subregion of $[117^oW, 33^oN]$ to $[116^oW, 33^o45'N]$. This region comprises a total of 6307 events in the HypoSVI 1D catalogue (Figure 7a-b), with the NonLinLoc comprising 6383 events (Figure 7c-d). Manual inspection showed that the events present in the NonLinLoc catalogue but not HypoSVI catalogue, are events that are locate external to the subregion in the HypoSVI catalogue but are projected to the edge of NonLonLoc search grid, having large location uncertainties. For the remaining events we determine the relative location differences between the two catalogues by projecting both catalogues to a local universal transverse mercato (UTM) coordinate system and determining the distance between the events in $km$ in a local XYZ coordinates. The relative distance of the NonLinLoc locations minus the HypoSVI 1D locations are given in Figure 8a-c. The relative locations demonstrate no consistent spatial bias, with the mean location difference given by $[X, Y, Z] = [+0.07\text{km}, +0.19\text{km}, -0.41\text{km}]$, as shown by the red dot in Figure 8a-c. In addition, we normalise the location difference by the location uncertainty from the NonLinLoc catalogue. Figure 8d-f gives the normalized location distances, with 83.29% of the events having a relative distance less than that of the NonLinLoc location uncertainty, as shown by the points within the dashed box.

## 6    Discussion & Conclusions

In this paper, we developed a new approach to performing Bayesian inference on earthquake hypocenters that combines a differentiable forward model (physics-informed neural network) with Stein variational inference.
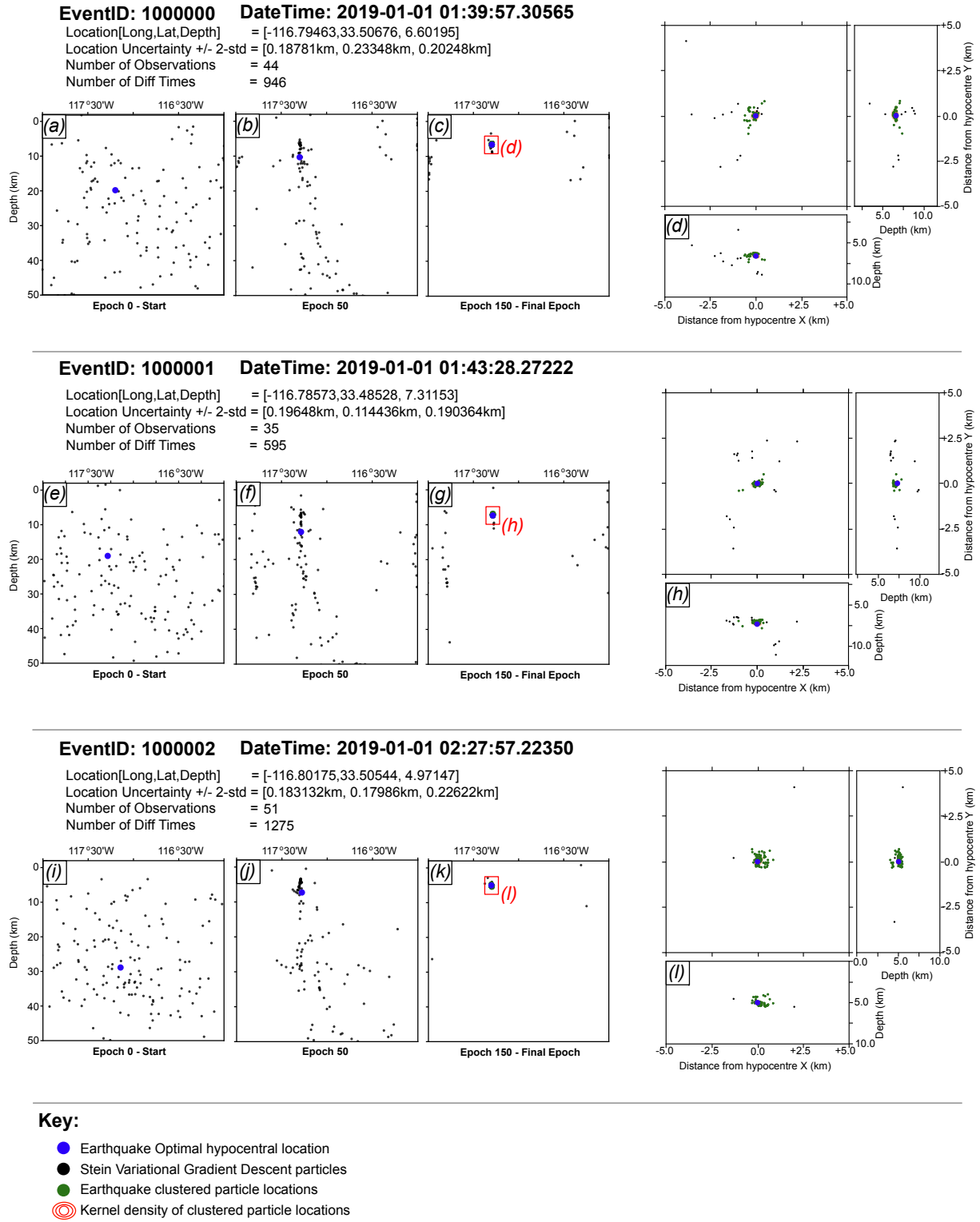
**EventID: 1000000**   **DateTime: 2019-01-01 01:39:57.30565**

Location[Long,Lat,Depth]    = [-116.79463,33.50676, 6.60195]
Location Uncertainty +/- 2-std = [0.18781km, 0.23348km, 0.20248km]
Number of Observations    = 44
Number of Diff Times    = 946

**EventID: 1000001**   **DateTime: 2019-01-01 01:43:28.27222**

Location[Long,Lat,Depth]    = [-116.78573,33.48528, 7.31153]
Location Uncertainty +/- 2-std = [0.19648km, 0.114436km, 0.190364km]
Number of Observations    = 35
Number of Diff Times    = 595

**EventID: 1000002**   **DateTime: 2019-01-01 02:27:57.22350**

Location[Long,Lat,Depth]    = [-116.80175,33.50544, 4.97147]
Location Uncertainty +/- 2-std = [0.183132km, 0.17986km, 0.22622km]
Number of Observations    = 51
Number of Diff Times    = 1275

**Key:**

- ● Earthquake Optimal hypocentral location
- ● Stein Variational Gradient Descent particles
- ● Earthquake clustered particle locations
- ◎ Kernel density of clustered particle locations

Figure 5: Example earthquake locations for three earthquakes in the Catalogues using travel-times derived from the three-dimensional regional velocity model. Left panels represent the particle locations changing at different epochs in the Stein Variational Gradient Descent. Right panels represent a zoom in of the final event locations, with the particle locations shown relative the recovered optimal hypocentral location. Kernel density contours are shown in red for the clustered particles.
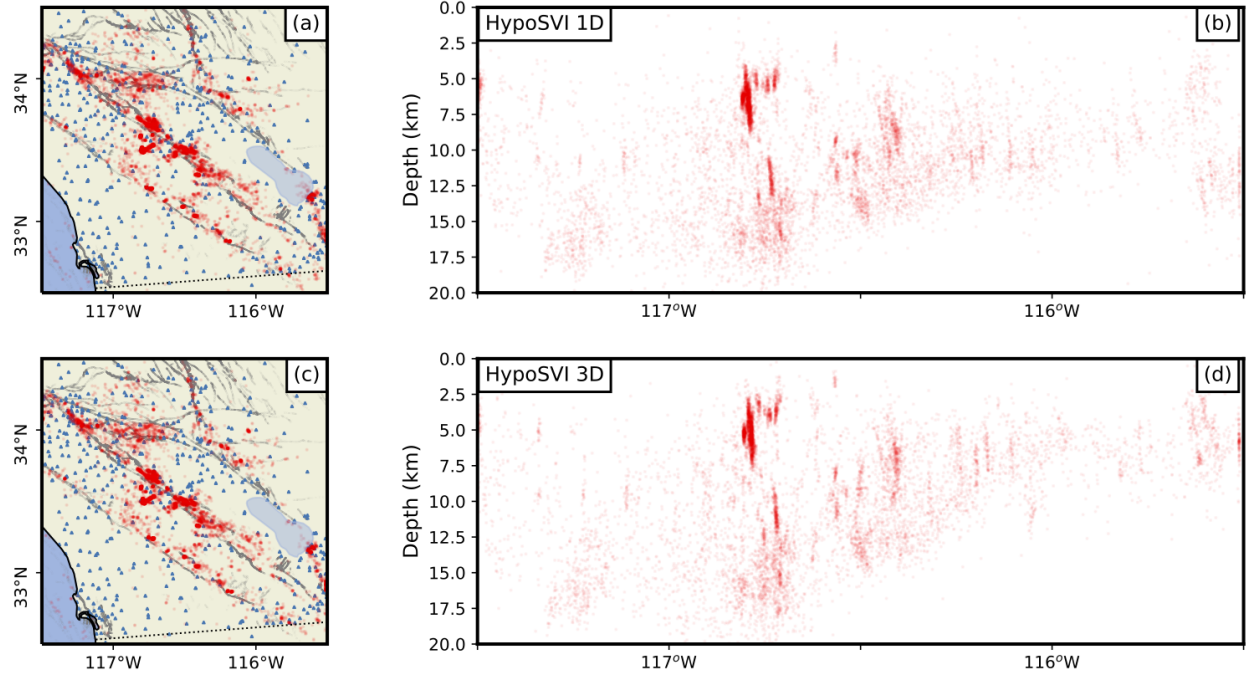
15

Figure 6: Comparison of earthquake locations between the HypoSVI and NLLoc. Left column represents the Latitude/Longitude map of the detected earthquakes given by red dots, observational station locations given by blue triangles and mapped faults by gray lines. Right column represents a Longitude vs Depth cross-sections of earthquakes. (a) and (b) are the locations determined from HypoSVI with a EikoNet model trained on a regional 1D velocity. (c) and (d) are the locations determined from HypoSVI with a EikoNet model trained on a regional SCEC-CVM-H 3D velocity structure.

Figure 7: Zoom in earthquake location comparison for the region for subregion of $[117^oW, 33^oN]$ to $[116^oW, 33^o45'N]$. (a)-(b) are the locations determined from HypoSVI with a EikoNet model trained on a regional 1D velocity. (c)-(d) are the locations determined from the NonLinLoc inversion procedure.

Figure 8: Earthquake distance comparison for the NonLinLoc and HypoSVI 1D catalogue for the region $[117^oW, 33^oN]$ to $[116^oW, 33^o45'N]$, projected to the local $X,Y,Z$ UTM coordinate system. (a)-(c) black dots represent the relative distance between catalogue event locations in X,Y,Z; with red dot representing the mean location. (d)-(f) black points relative distance between catalogue event locations normalized by the NonLinLoc 2-std location uncertainty. Red-dashed region represents the catalogue events with a relative distance less than the location uncertainty.

18

Unlike with MCMC sampling methods, SVI approximates a posterior with a collection of particles, with the set of particle locations jointly optimized. In this paper we use an EikoNet forward model, but this could be replaced with any other differentiable forward model. Thus, HypoSVI is a general variational approach to hypocenter inversion. We validated the method with synthetic tests and compared the locations for $\sim 10000$ events in Southern California with those produced by the Southern California Seismic Network. In particular, we focused on demonstrating the reliability of the method in the presence of non-convex posterior distributions, which SVI is well suited for handling. This is all possible because of the differentiable forward model.

Another advantage of our approach is that it is computationally efficient and can make use of state of the art GPU architectures and modern deep learning APIs like PyTorch. This allows for rapid calculation of the gradients with automatic differentiation. As GPU hardware improves, such as increased memory, these performance gains will be passed on to the algorithm which will allow for even larger datasets to be worked with than currently possible. By combining SVI with EikoNet, we are able to evaluate observations at any point within the 3D volume without retraining, i.e. the forward model is valid for any array geometry. Due to the highly-parallelized nature of calculations with neural networks, our method scales well to very large networks, which may be important for emerging technologies like Distributed Acoustic Sensing (DAS). This was demonstrated herein by the ability to locate an earthquake with 2048 phase picks in 439 seconds. Thus, our HypoSVI approach is ideal for handling the enormous data volumes that are starting to emerge in seismology.

# Acknowledgments

# References

Allam A. & Ben-Zion Y., 2012, Seismic velocity structures in the Southern California plate-boundary environment from double-difference tomography, *Geophysical Journal International*, **190**, 1181-1196

Geiger L., 1912, Probability method for the determination of earthquake epicenters from the arrival time only, *St. Louis Univ. Bull.*,**8**,60-71

Hahsler M., Piekenbrock M., & Doran D., 2019, dbscan: Fast Density-Based Clustering with R, *Journal of Statistical Software*, **91**, 1-30

Hutton, K., Woessner, J., & Hauksson, E., 2010, Earthquake monitoring in southern California for seventy-seven years (1932–2008), *Bulletin of the Seismological Society of America*, **100**, 423-446

Kingma, D.P., & Ba, J., 2014, Adam: A method for stochastic optimization, *arXiv preprint arXiv:1412.6980*

Lee, E., Chen, P., Jordan, T.H., Maechling, P.B., Denolle, M.A.M., & Beroza, G.C., 2014, Full-3-D tomography for crustal structure in Southern California based on the scattering-integral and the adjoint-wavefield methods, *Journal of Geophysical Research: Solid Earth*, **119**, 6421-6451

Lomax, A., Virieux, J., Volant, P., & Catherine, B., 2000, Probabilistic earthquake location in 3D and layered models, *Advances in seismic event location*,101–134

Mousavi, S. M., Ellsworth, W. L., Zhu, W., Chuang, L. Y., & Beroza, G. C., 2020, Earthquake transformer—an attentive deep-learning model for simultaneous earthquake detection and phase picking, *Nature communications*,**11**,1-12

Podvin, P., & Lecomte, I., 1991, Finite difference computation of traveltimes in very contrasted velocity models: a massively parallel approach and its associated tools,*Geophysical Journal International*, **105**, 271-284

Qiang, L. & Dilin, W., 2016, Stein variational gradient descent: A general purpose bayesian inference algorithm, *Advances in neural information processing systems*, 2378-2386

Rawlinson, N., & Sambridge, M., 2005, The fast marching method: an effective tool for tomographic imaging and tracking multiple phases in complex layered media,*Exploration Geophysics*,**36**, 341-350

Ross, Z. E., Meier, M., Hauksson, E., & Heaton, T. H., 2004, Generalized seismic phase detection with deep learning,*Bulletin of the Seismological Society of America*,**108**, 2894-2901

Shaw, J.H., Plesch, A., Tape, C., Suess, M.P., Jordan, T.H, Ely, G., Hauksson, E., Tromp, J., Tanimoto, T. & Graves, R., 2015, Unified structural representation of the southern California crust and upper mantle,*Earth and Planetary Science Letters*,**415**, 1-15

Smith, J. D., Azizzadenesheli, K., & Ross, Z. E., 2020, EikoNet: Solving the Eikonal Equation With Deep Neural Networks,*IEEE Transactions on Geoscience and Remote Sensing*,1–12, 10.1109/TGRS.2020.3039165

Süss, M.P, & Shaw, J.H., 2003, P wave seismic velocity structure derived from sonic logs and industry reflection data in the Los Angeles basin, California,*Journal of Geophysical Research: Solid Earth*,**108**

Tarantol, A., 2004, Inverse Problem Theory and Methods for Model Parameter Estimation,*Society for Industrial and Applied Mathematics*

Thurber, C.H., 1985, Nonlinear earthquake location: theory and examples,*Bulletin of the Seismological Society of America*,**75**,779-790

Treister, E. & Haber, E., 2016, A fast marching algorithm for the factored eikonal equation, *Journal of Computational physics*,**324**, 210-225
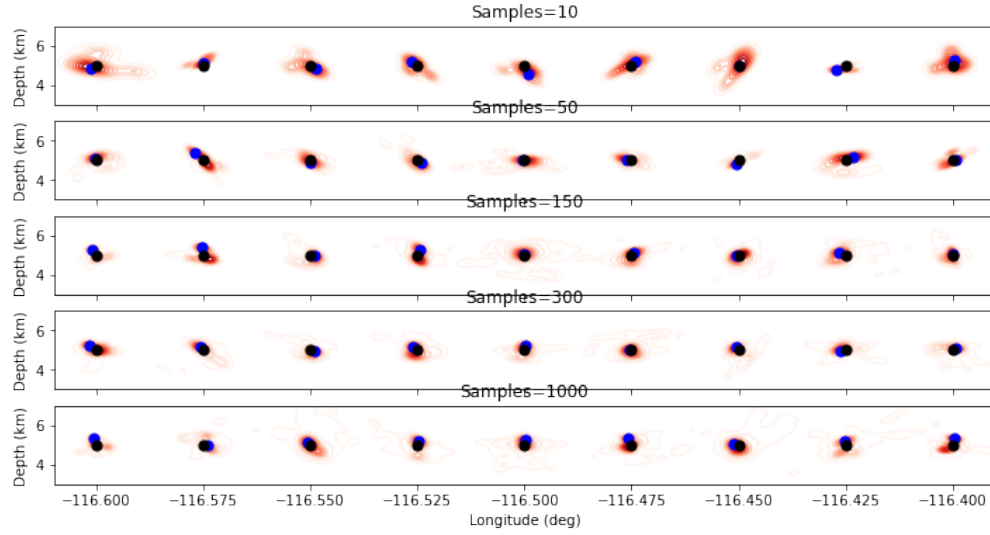
Figure S9: Synthetics earthquake location recovery for changing number of particles. An outline of observation and synthetic locations distributions is given in Section 4. Black points represent the imposed synthetic earthquake location, blue dots the recovered optimal location, red contours present the recovered posterior determined by the particle density.

# A   Supplementary Figures
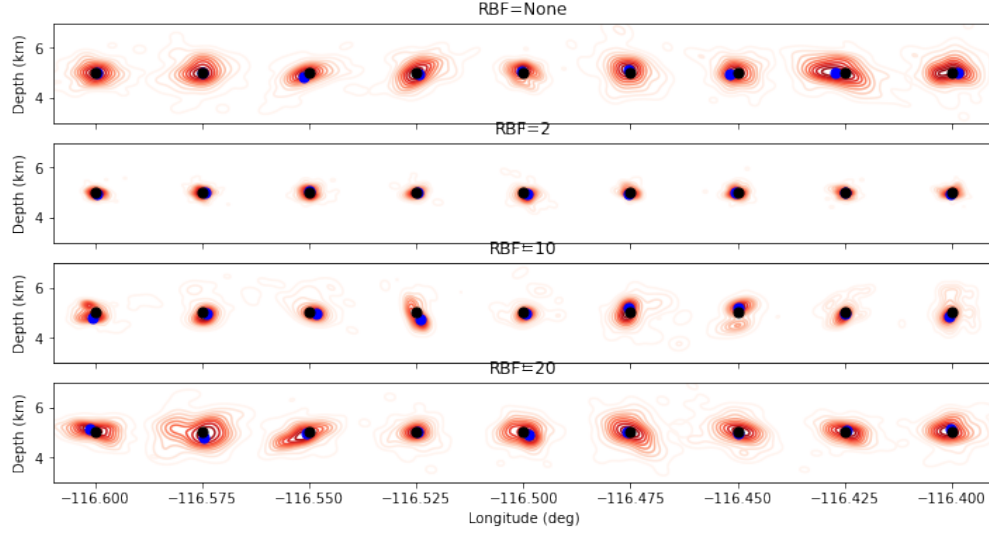
Figure S10: Synthetics earthquake location recovery for changing values for the Radial Basis Function shape parameter value. An outline of observation and synthetic locations distributions is given in Section 4. Black points represent the imposed synthetic earthquake location, blue dots the recovered optimal location, red contours present the recovered posterior determined by the particle density.
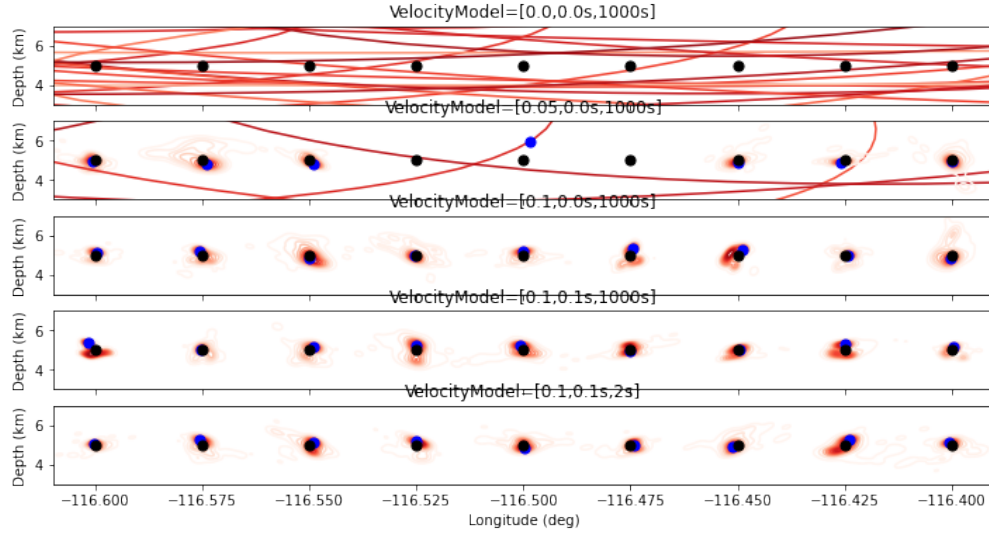


Figure S11: Synthetics earthquake location recovery for changing values for the forward model uncertainty in form $[\sigma_f, \sigma_{min}, \sigma_{max}]$. An outline of observation and synthetic locations distributions is given in Section 4. Black points represent the imposed synthetic earthquake location, blue dots the recovered optimal location, red contours present the recovered posterior determined by the particle density.
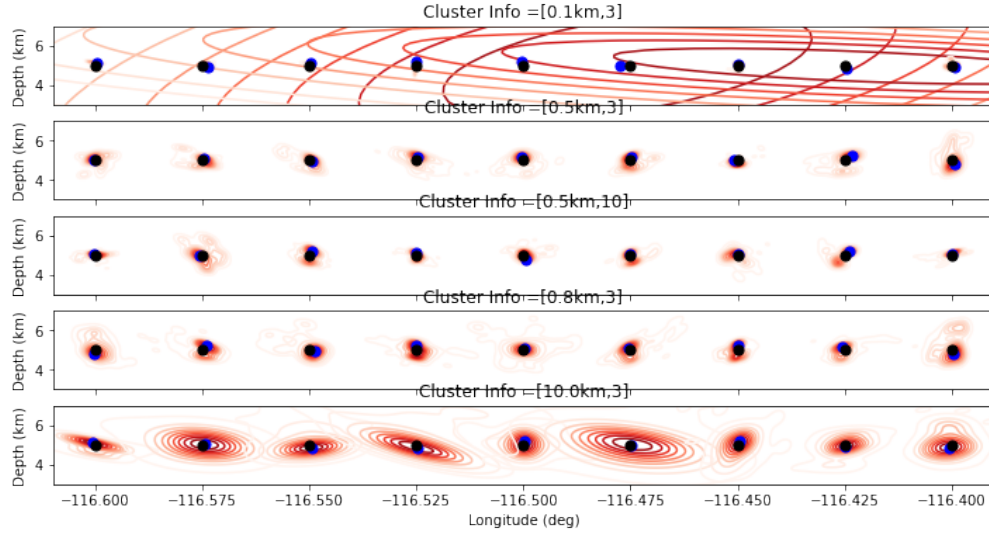
Figure S12: Synthetics earthquake location recovery for changing values for the location uncertainty clustering for parameters: Distance between the particles to define a cluster and minimum number of particles to represent a cluster. An outline of observation and synthetic locations distributions is given in Section 4. Black points represent the imposed synthetic earthquake location, blue dots the recovered optimal location, red contours present the recovered posterior determined by the particle density.