

## 1. Introduction

Probabilistic seismic hazard analysis (PSHA) based on physics-based simulations offers many advantages compared to more traditional empirical Ground Motion Model (GMM) based PSHA. However when a large number of seismic sources is considered, such as distributed seismicity, simulation based PSHA becomes unfeasible due to the high computational cost associated with simulations. Utilising a surrogate model, such as an artificial neural network, gives the physics-based simulation advantages without the associated computational cost. Additionally it provides an extra avenue for investigating simulation results.

## 2. Dataset

The labelled dataset used for training and validation is made up of New Zealand simulation data from both validation of historical events (Lee et al., 2020;  $\approx 600$  simulations), as well as future potential events (Motha et al., 2020;  $\approx 11,000$  simulations).

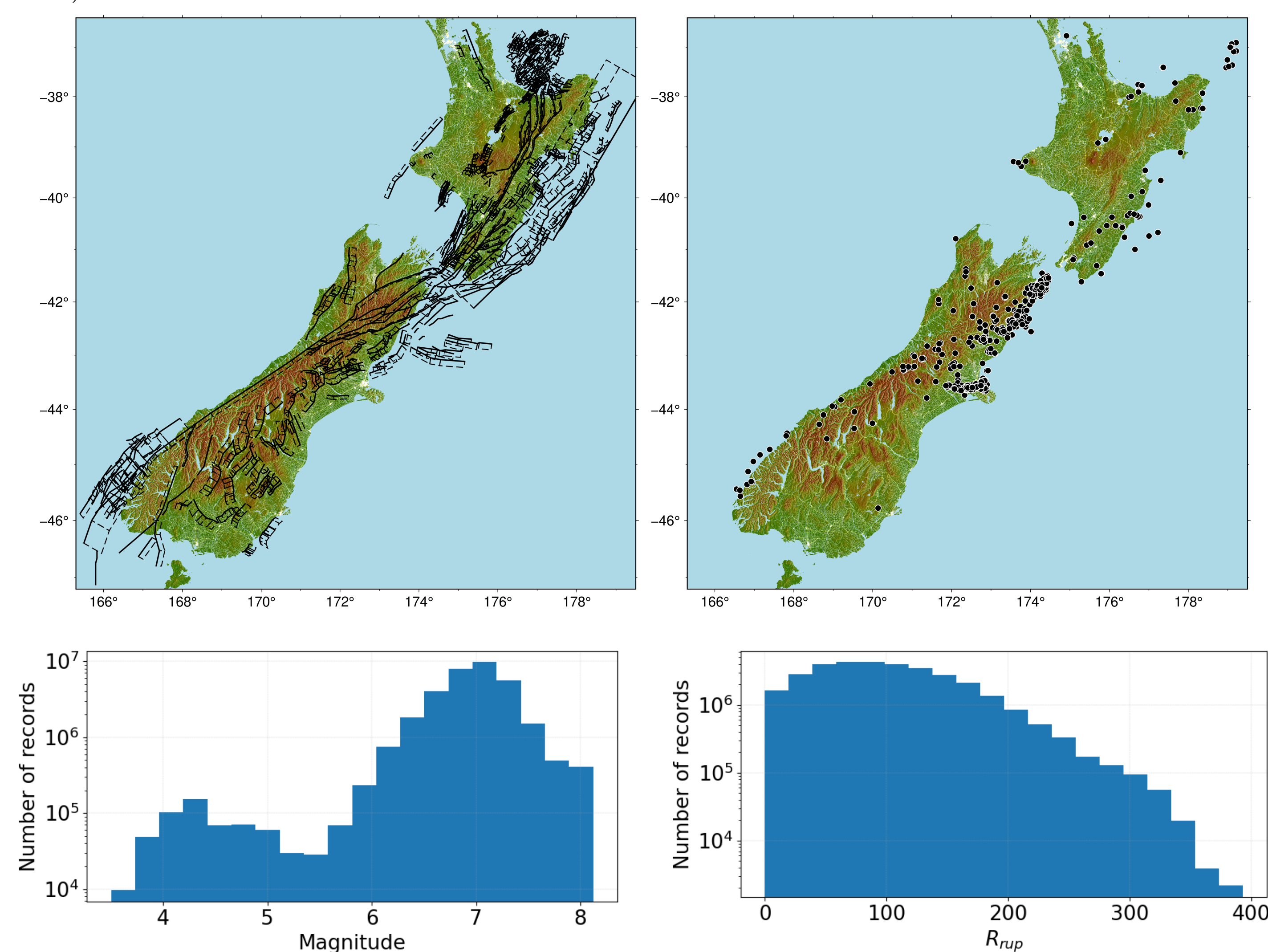


Figure 1: *Top-left*: Simulated faults (potential future events), *Top-right*: Simulated historical events, *Bottom*: Distribution of training data with respect to Magnitude and  $R_{rup}$

The resulting dataset spans magnitudes from 3.5 to 8, allowing training of a purely data-based model.

## 3. Model

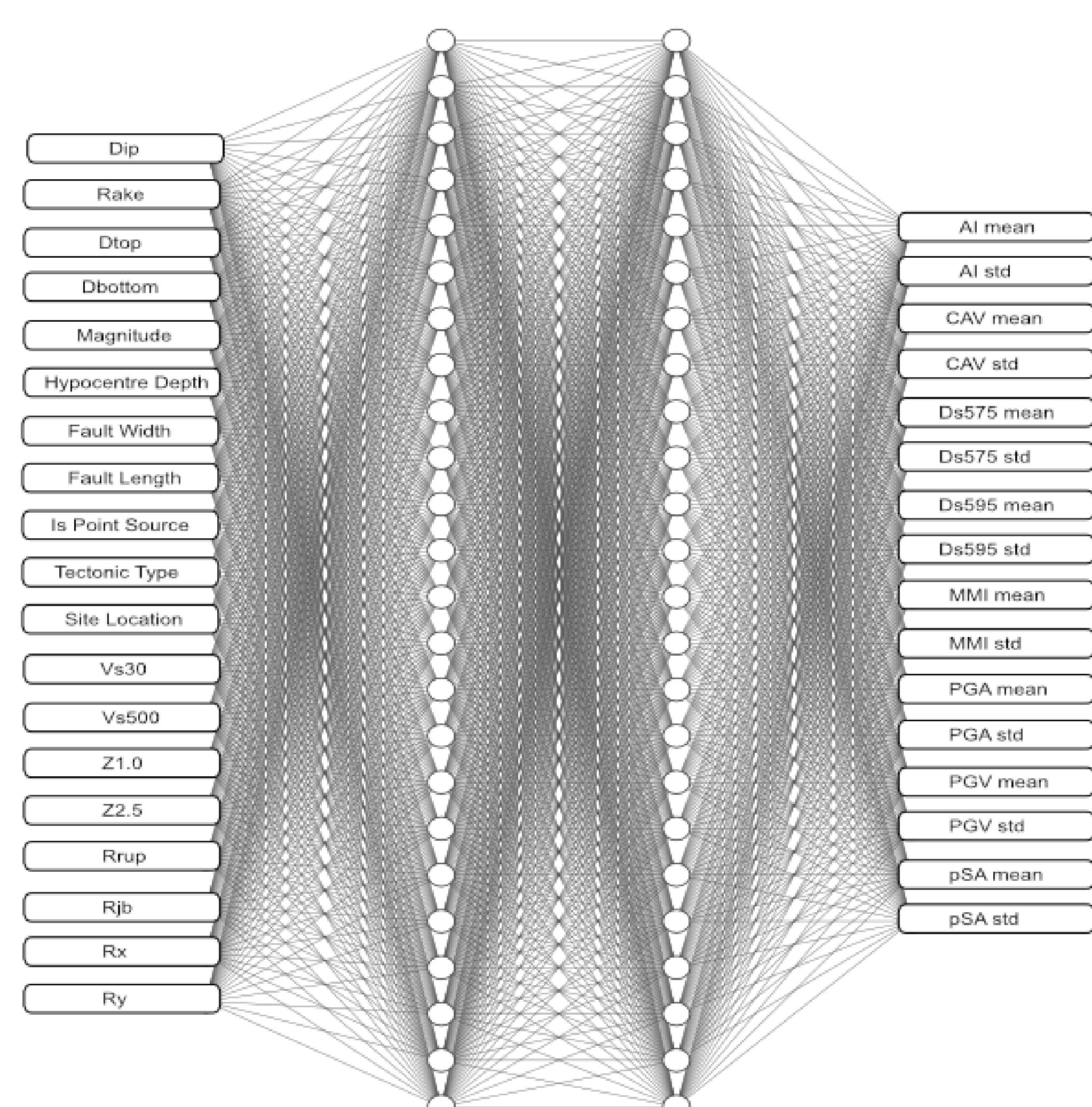


Figure 2: Example model schematics with the inputs on the left connected to outputs on the right via 2 fully connected hidden layers with 24 units each

A fully connected multi-layer perceptron (MLP) neural network is used as the surrogate model. It is trained using the back-propagation algorithm with negative log-likelihood as the cost function. The model has 5 hidden layers each with 256 units using the Rectified Linear Unit (ReLU) activation function; with the output layer using no activation function. Additionally, dropout layers are added following each hidden layer to prevent overfitting and estimate model uncertainty (Gal et al., 2016).

## 4. Results

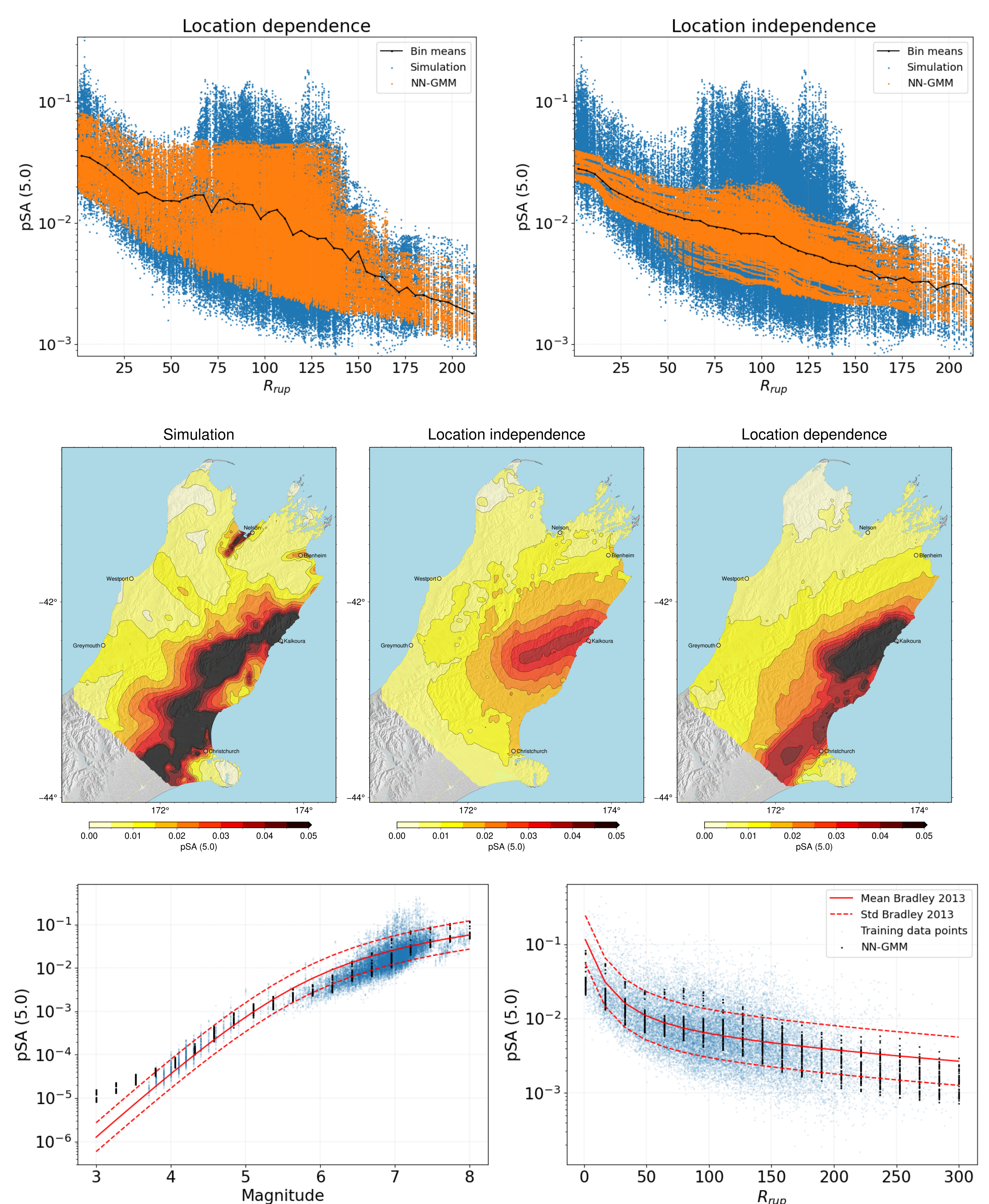


Figure 3: *Top and middle*: Comparison of with and without location data models for the Hope Conway fault data (validation dataset), *Bottom*: Model's general behaviour compared to the Bradley 2013 empirical GMM, with the training data shown in the background

The resulting model is able to learn the general trends from the simulation data, and predict IM values to a good degree of accuracy compared to traditional empirical GMMs. Using site location as a model input additionally allows the model to learn site-specific effects such as basin amplification. However, the model currently still generally underpredicts these regions or even misses them completely (e.g. the Nelson basin shown in the middle plots of Figure 3).

## 5. Next steps

In addition, to continue tuning the model architecture and hyperparameters, improvements to the estimation of site-specific effects based on location is a priority. Additionally, the imbalance in the records with respect to magnitude has to be investigated to ensure correct model behaviour at small magnitudes.