

# A physics-informed Bayesian framework for characterizing ground motion process in the presence of missing data

Yu Chen<sup>1</sup>  | Edoardo Patelli<sup>2</sup>  | Benjamin Edwards<sup>1</sup> | Michael Beer<sup>1,3,4</sup>

<sup>1</sup>Institute for Risk and Uncertainty, University of Liverpool, Liverpool, UK

<sup>2</sup>Department of Civil and Environmental Engineering, University of Strathclyde, Glasgow, UK

<sup>3</sup>Institute for Risk and Reliability, Leibniz Universität Hannover, Hannover, Germany

<sup>4</sup>International Joint Research Center for Resilient Infrastructure & International Joint Research Center for Engineering Reliability and Stochastic Mechanics, Tongji University, Shanghai, China

## Correspondence

Edoardo Patelli.

Email: [edoardo.patelli@strath.ac.uk](mailto:edoardo.patelli@strath.ac.uk)

## Funding information

H2020 Marie Skłodowska-Curie Actions; EU Horizon 2020 – MSCA Actions project URBASIS, Grant/Award Number: 813137

## Abstract

A Bayesian framework to characterize ground motions even in the presence of missing data is developed. This approach features the combination of seismological knowledge (*a priori knowledge*) with empirical observations (even incomplete) via Bayesian inference. At its core is a Bayesian neural network model that probabilistically learns temporal patterns from ground motion data. Uncertainties are accounted for throughout the framework. Performance of the approach has been quantitatively demonstrated via various missing data scenarios. This framework provides a general solution to dealing with missing data in ground motion records by providing various forms of representation of ground motions in a probabilistic manner, allowing it to be adopted for numerous engineering and seismological applications. Notably, it is compatible with the versatile Monte Carlo simulation scheme, such that stochastic dynamic analyses are still achievable even with missing data. Furthermore, it serves as a complementary approach to current stochastic ground-motion models in data-scarce regions under the growing interests of PBEE (performance-based earthquake engineering), mitigating the data-model dependence dilemma due to the paucity of data, and ultimately, as a fundamental solution to the limited data problem in data scarce regions.

## KEYWORDS

Bayesian model updating, earthquake ground motion, evolutionary power spectra, missing data, stochastic variational inference, uncertainty quantification

## 1 | INTRODUCTION

The random nature of earthquake ground motions is well appreciated. Various research efforts and progress, based on stochastic process formulation, have been made towards the problem of characterization, simulation and response evaluation.<sup>1–3</sup> In recent years, the growing interest in performance-based earthquake engineering (PBEE), which requires ground motions of various hazard levels to consider the entire range of structural response, including non-linear behaviour and even collapse,<sup>3</sup> has driven the need for simulating ground motions of various earthquake scenarios. Stochastic

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. *Earthquake Engineering & Structural Dynamics* published by John Wiley & Sons Ltd.

simulations are further utilised for evaluation of future seismic demand and seismic reliability assessment,<sup>4</sup> non-linear stochastic dynamic analyses,<sup>5</sup> developing ground motion prediction equations (GMPEs)<sup>6</sup> or seismic hazard characterization and simulation-based seismic risk assessment.<sup>7,8</sup>

However, their applicability is not without questioning. Empirical ground motions are responsible for developing and calibrating stochastic ground motion models. However, the paucity of recordings (especially strong motions) in data scarce regions leads to a bottleneck that observational data are lacking in the first place to justify modelling and calibration. For instance, in characterizing seismic hazard, a category of predictive-relation-based stochastic ground-motion models (see, e.g., Rezaeian and Der Kiureghian,<sup>9</sup> Laurendeau et al.,<sup>10</sup> Vlachos et al.<sup>11</sup>) is gaining increasing attention for its ability to generate a suite of non-stationary time-histories, given specific earthquake scenarios. The core component of these models is an underlying empirical regression between model parameters and earthquake characteristics over a selected (sometimes limited) subset of records. However these empirical relations are largely bounded by the scope of data being regressed. Significant epistemic uncertainties are expected on further uses of these underlying empirical regressions as *extrapolation* than *interpolation*. Similarly, such uncertainty also applies to those empirical GMPEs developed using stochastic simulations calibrated from small-to-moderate earthquakes often due to a lack of strong motions.<sup>6,12</sup> Concerns have been raised over the subsequent stochastic simulations from these biased models, as the underlying regression are typically not well-constrained by empirical data and their extrapolation may, therefore, not even be physically realistic.<sup>13</sup>

Therefore, for data-scarce regions, where there are stronger needs of synthetic ground motions for abundant earthquake scenarios, however, the paucity of data poses a causality dilemma concerning the dependence between observations and the extracted knowledge/information for the development of models. This raises difficulties, in data scarce regions, in the characterization of ground motions for the seismic risk assessment as well as researches of regional seismicity and Earth regional structures.

As such, a method to make the most of existing data (even where incomplete), robustly characterizing the underlying physical processes from bad measurements (e.g., incomplete), could enrich the observational database, whereby one is able to progressively update the development and calibration of ground motion models, producing more realistic stochastic simulations in the otherwise data scarce regions, for hazard characterization and risk assessment. It serves as a complementary approach to stochastic ground-motion models under the growing interests of PBEE, and ultimately a fundamental solution to the limited data problem. This may be of particular interest to studies of historical earthquakes, which may potentially provide strong-motion records but many of them are discarded due to the presence of data gaps.<sup>14</sup> Furthermore, missing data exist in both historical and modern earthquake time histories due to intermittent instrumentation or data-transmission failure. For instance, old mechanical, short-period high-sensitivity or broadband seismometers are vulnerable to clipping during local strong motions. In addition, sensor malfunctions, instrument tilt, or data contamination may lead to missing or incorrect values, or waveform clipping around the peak motion.<sup>14–16</sup> With the recent use of low-cost temporary instruments, deployed at scale, sometimes in harsh conditions, the fidelity and continuity of recording is also not as reliable as traditional permanent seismological stations, which itself can be understood as a bad- or missing-data problem.

The characterization of ground motions and accounting for their random nature is challenging when only limited and partial recordings are available.<sup>16–18</sup> Pioneering works for analysis in the presence of missing data, such as the Lomb–Scargle periodogram,<sup>19</sup> iterative deconvolution CLEAN,<sup>20</sup> are acknowledgedly to have deficiencies such as bias issue and periodic content limitation.<sup>15,21–23</sup> With different assumptions (hence limitations), many other methods have been proposed in recent years. Notably, a compressive sensing approach is exploited with the sparsity assumption of the underlying spectral representation.<sup>17</sup> By assuming the same frequency contents between the missing portion and the observations, a projection onto convex sets (POCS) method can be used to reconstruct clipped waveforms.<sup>16</sup> Parametric models are also developed based on various formulations, such as autoregressive modelling methods,<sup>21,24,25</sup> with parameterized assumptions on the structure of the underlying stochastic processes. Similarly, Maranò et al.<sup>14</sup> proposed a method to fit a parametric seismological model to earthquake recordings with missing gaps.

Alternatively, a variety of methods are available that explicitly or implicitly transform spectral analysis with missing data into the imputation of missing values, followed by standard full-data spectral analysis.<sup>26–30</sup> This strain of methods provides reconstructed waveforms in a straightforward manner, whereby extensive established spectral analyses, developed on equidistant data, whether stationary or nonstationary, can still be universally harnessed.

Two main challenges are identified in dealing with missing data. First, most current approaches fail to address the uncertainties related to the missing data properly.<sup>18,31</sup> For reconstruction based methods, inaccuracies of the imperfect reconstruction will be propagated to spectral estimates owing to the convolutional nature of Fourier transform. Similarly, for parametric modelling methods that results in a parametric form of spectrum, parameter uncertainties due to the incomplete data are not well captured. More importantly, despite existing approaches that handle uncertainties

(notably Bayesian spectral analyses<sup>32,33</sup>), they are still constrained by the significantly limited information from the very incomplete signal.

Therefore, to exploit additional information besides the incomplete recording and to appropriately quantify the uncertainties brought by the missing data, we propose a novel Bayesian framework that aims to robustly combine prior seismological knowledge with empirical observations (even incomplete). A Bayesian neural network (BNN) model that probabilistically learns the temporal dynamics from earthquake time histories forms the key component of the framework. In particular, it is initially trained from physics-informed simulated ground motions given the event metadata (e.g., magnitude, epicentral distance,  $V_{s30}$ , etc.), as geological *a-priori*, and subsequently updated via Bayesian inference utilising the partial empirical observations. Importantly, uncertainty has been accounted for throughout the framework. Variability of the physics-informed simulations is considered. Epistemic uncertainties on model parameters of the BNN are learnt through stochastic variational inference, whereby an ensemble of reconstructed time histories is obtained by marginalizing over the posterior distribution of model parameters. Furthermore, uncertainties of the spectral representations (e.g., evolutionary power spectral density [EPSD]) of the underlying stochastic process are quantified, with the spectral density values represented by probability distributions. As a result, sample realizations associated with the stochastic process can be further simulated for stochastic dynamic analysis through the spectral representation method, even with incomplete recordings.

Details of the framework are discussed first, then the performance of the proposed method is demonstrated with various missing data scenarios based on an earthquake strong motion recording.

## 2 | A BAYESIAN FRAMEWORK FOR CHARACTERIZATION OF GROUND MOTION WITH MISSING DATA

We build on the premise that a priori seismological knowledge can provide a general, yet insightful, prior expectation of the ground motions of the certain earthquake scenario, which can be combined with the information extracted from empirical observations (even when incomplete).

### 2.1 | Physics-informed stochastic simulations as a geological prior

A stochastic representation that encapsulates the physics of the earthquake process and wave propagation plays the central role, from the seismological perspective, in characterizing the ground motions (see, e.g., Zeng et al.<sup>34</sup>, Boore<sup>35</sup>). One of the most desired advantage is that such representations explicitly distil the knowledge of various factors affecting ground motions (e.g., source, path and site effects) into a parametric formulation. In this study, we have adopted a well-validated stochastic seismological model,<sup>35</sup> as given below, whereby source process, attenuation and site effects are encapsulated in a parameterized form of the Fourier amplitude spectrum. A finite fault strategy is particularly employed to represent the geometry of larger ruptures for large earthquakes.<sup>6,36</sup>

$$A(f; \Theta) = \frac{CM_0}{1 + (f/f_0)^2} Z(R) \exp[-\pi f R/Q(f)\beta] G(f) \quad (1)$$

where  $\Theta = (\Theta_e, \Theta_g)$  represents the event parameters ( $\Theta_e$ ) that are still accessible from the metadata of an incomplete recording, such as seismic moment  $M_0$  and hypocentral distance  $R$ , and region-specific seismological parameters ( $\Theta_g$ ) that embody the source, path and site effects. Specifically,  $f_0$  is the earthquake's source corner frequency given by  $f_0 = 0.4906\beta(\Delta\sigma/M_0)^{1/3}$  red (in SI units);  $R = \sqrt{r^2 + d^2}$  where  $r$  and  $d$  are the epicentral distance and depth to a given sub-fault;  $\Delta\sigma$  is referred to as the stress drop and  $\beta$  represents the shear wave velocity in the vicinity of the source. The constant  $C$  is given by:  $C = R_{\theta\Phi} V F / (4\pi\rho_s\beta^3 R_0)$ , where  $R_{\theta\Phi}$  is the radiation pattern;  $V$  represents the partition of total shear-wave energy into horizontal components;  $F$  accounts for the free-surface effect;  $R_0$  is the a reference distance and  $\rho$  is the density in the vicinity of the source.  $Z(R)$  is the geometrical spreading function defined by a piece-wise series of segments in the form of  $R^{b_n}$ , where  $b_n$  defines the geometrical-spreading coefficient in the  $n$ th segment. The quality factor  $Q(f)$  is an inverse measure of anelastic attenuation. The site effect  $G(f) = \exp(-\pi f\kappa_0)10^v$  is given by the counteraction of a high-cut filter,  $\exp(-\pi f\kappa_0)$ , accounting for the diminution of the high-frequency motions and an amplification factor  $v$  in log units. The specific values for each of the model terms used in this model can be taken from the existing literature, or directly through spectral modelling of waveform data (e.g., Edwards and Fäh<sup>12</sup>).

In particular, the variability of model parameters in the spectral formulation, and hence the uncertainty in stochastic simulations are represented by probability distribution over the input parameters  $\Theta_g$  as proposed by Atkinson and Boore<sup>6</sup> and Vetter and Taflanidis.<sup>37</sup> Note that the above stochastic simulation procedures are distinct from those comprehensive deterministic numerical models that solve the complex 3D equations governing seismic wave propagation. Those models are typically referred to as physics-based numerical models in the literature, see, for example, McCallen et al.<sup>38,39</sup>, Paolucci et al.<sup>40</sup> among others.

## 2.2 | Sequential modelling

In recent years, neural network models have become established in learning complex and non-linear relations. Most recently, successes have been seen for neural networks to learn the temporal dynamics in sequential data (e.g., time series) under an autoregressive setting.<sup>29,41–43</sup> They model the data generating process by formulating the conditional distribution,  $p(y_t|\mathbf{x}_t, \mathbf{w})$ , of the value  $y_t$  based on a window of past lagged values ( $[y_{t-1}, \dots, y_{t-p}]$ ), as given by

$$y_t = f(\mathbf{x}_t; \mathbf{w}) + \epsilon, \text{ with } \mathbf{x}_t = [y_{t-1}, \dots, y_{t-p}] \quad (2)$$

where  $\epsilon$  denotes the noise term;  $f(\cdot)$  represents the neural network model, parameterized by  $\mathbf{w}$ , which learns complex nonlinear temporal dependence in the time series, as opposed to a linear combination of fixed coefficients in a classic autoregressive AR( $p$ ) model.  $y_t$  and  $\mathbf{x}_t$  represent the prediction and the lagged window pair. In practice, training with maximum likelihood estimation (MLE) gives rise to a probabilistic interpretation of the data generating process. The likelihood function, assuming Gaussian noise with variance  $\sigma^2$ , is given by<sup>44</sup>

$$p(y_t|\mathbf{x}_t, \mathbf{w}) = \mathcal{N}(y_t|f(\mathbf{x}_t, \mathbf{w}), \sigma^2) \quad (3)$$

Model parameters  $\mathbf{w}$ , collectively the weights and biases of the neural network model (referred as weights hereafter), are estimated during training by optimizing with the likelihood as the objective as follows:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \sum_t \log p(y_t|\mathbf{x}_t, \mathbf{w}) \quad (4)$$

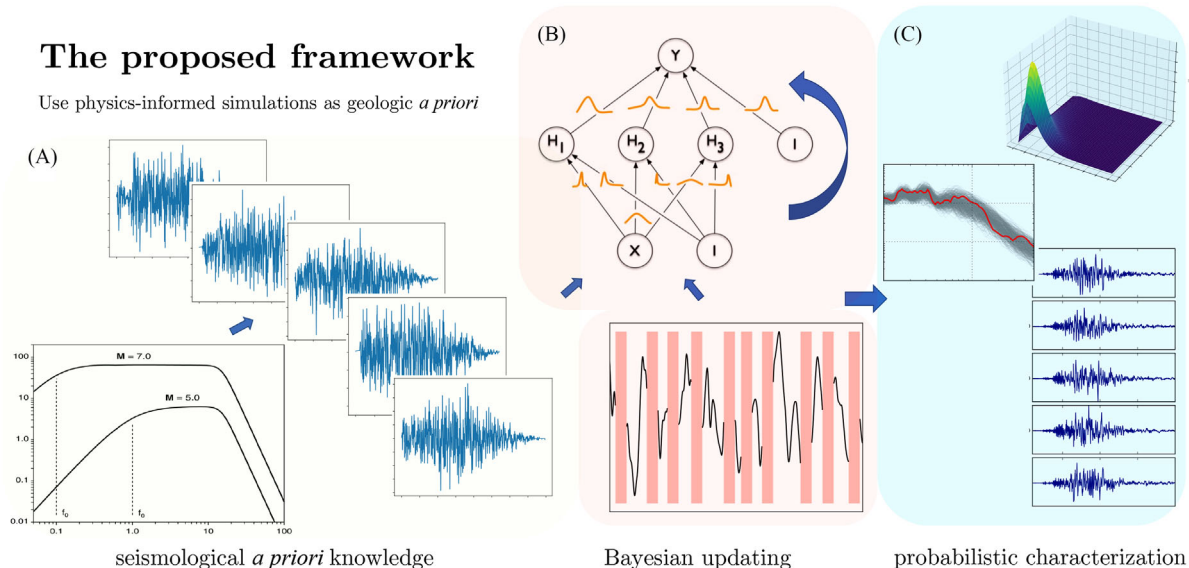
Once trained, its generative power could be employed to generate sequences,<sup>45</sup> forecast time series future values<sup>41</sup> and impute missing values.<sup>29</sup> However, despite accounting for the aleatoric uncertainty using Gaussian noise, the above MLE strategy ignores the uncertainties of the model parameters (i.e., epistemic uncertainties) that can explain the observed data (especially in the context of limited data and missing data) as well as the resulting predictive uncertainties regarding the imputation. Significant uncertainties exist on the model configurations that may have explained the limited data. Consequently, such uncertainties further compromise the generalization power of learned models in that predictions from uncertain/unrepresentative models can still be unreliable and over confident.<sup>46,47</sup>

## 2.3 | Bayesian updating on partial observations

In order to capture the model uncertainty, probability distributions are applied to the neural net model parameters (see Figure 1). Bayesian inference hence formulates the update of the neural network modelling the underlying generating process, when new observations (even incomplete) become available, as given below:

$$p(\mathbf{w}|D) = p(D|\mathbf{w})p(\mathbf{w})/p(D) \quad (5)$$

where  $p(\mathbf{w})$  represents the prior probability distribution of weights learnt from the physics-informed simulations;  $p(D|\mathbf{w})$  stands for the likelihood and  $D$  specifically refers to the partial and incomplete observations.  $p(\mathbf{w}|D)$  is the posterior distribution, in which both the prior seismological knowledge and the real-world empirical observations are collectively considered. The posterior predictive distribution for the prediction of the missing value  $y_t^*$ , based on the lagged window,



**FIGURE 1** A stochastic framework characterizing ground motion process in the presence of missing data. Three components are presented: (A) a seismological model generating physics-informed stochastic simulations with *a priori* seismological knowledge; (B) a Bayesian neural network model initially trained from physics-informed stochastic simulations and later updated by empirical partial observations; (C) a host of model-based probabilistic representations of ground motions (e.g., evolutionary power spectral density EPSD, elastic response spectra, ensemble reconstructed time histories etc.).

can be made for each possible configuration of the weights, by marginalizing over the posterior distribution, as shown below:

$$\begin{aligned}
 p(y_t^* | \mathbf{x}_t, D) &= \int p(\mathbf{w} | D) p(y_t^* | \mathbf{x}_t, \mathbf{w}) d\mathbf{w} \\
 &= \mathbb{E}_{p(\mathbf{w} | D)} [p(y_t^* | \mathbf{x}_t, \mathbf{w})]
 \end{aligned} \tag{6}$$

As a result of considering uncertainties within the neural network, an ensemble of reconstructed time-histories, based on Monte Carlo sampling of the posterior distributions of weights, can be obtained. Subsequently, an ensemble of spectral estimates (e.g., EPSD, response spectra etc.) can be computed from the ensemble reconstructions using established spectral analysis methods. Performing such analyses for many incomplete recordings in the otherwise data scarce region produces an enriched database, which could be further adopted to update the development or calibration of ground motion models (including both stochastic ground-motion models and empirical GMPEs). This scheme is interpreted as an escape from the model-data dependence dilemma, as highlighted earlier, by making the most of the observed data (even when incomplete).

## 2.4 | Stochastic variational inference

A key challenge in Equation (5) is the approximation of the posterior distribution. Analytic Bayesian inference to the true posterior  $p(\mathbf{w} | D)$  is intractable and Markov Chain Monte Carlo (MCMC) based sampling approaches generally have difficulties in scaling to the huge dimensions of neural networks.<sup>47,48</sup> Alternatively, stochastic variational inference (see, e.g., Graves<sup>49</sup>, Kingma and Welling<sup>50</sup>, Blei et al.<sup>51</sup>) approximates the posterior distribution  $p(\mathbf{w} | D)$  efficiently, by turning such inference problem into an optimization problem. It optimizes the parameters of a proposed variational distribution, such that the Kullback–Leibler (KL) divergence between the approximate distribution and the true posterior distribution is minimised:  $\theta^* = \arg \min_{\theta} \text{KL}[q(\mathbf{w} | \theta) \| p(\mathbf{w} | D)]$ . This minimization objective is indeed equivalent to the following cost function<sup>49</sup>:

$$\mathcal{J}(D, \theta) = \text{KL}[q(\mathbf{w} | \theta) \| p(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w} | \theta)} \log p(D | \mathbf{w}) \tag{7}$$



Equation (7), hence, represents the new cost function to which optimization on  $\theta$  is taken. Directly taking derivatives is computationally prohibitive. However, it could be further re-arranged into the form of an expectation, lending itself to known approximate solutions such as Monte Carlo estimator of expectation on samples (see Appendix B). Specifically, prior to rearranging into an expectation, if assuming the variational posteriors have diagonal Gaussian distributions, the KL divergence term of Equation (7) can be further analytically integrated,<sup>50</sup> as given below, leaving only the likelihood-dependent part to be computed by a Monte Carlo estimator:

$$\text{KL}[q(\mathbf{w}|\theta) \parallel p(\mathbf{w})] = \frac{1}{2} \sum_j \left( \sigma_j^2 + \mu_j^2 - \log \sigma_j^2 - 1 \right) \quad (8)$$

where  $\mu_j, \sigma_j$  denote the  $j$ th element of the vectors that represent the variational distribution of weights,  $\theta = (\boldsymbol{\mu}, \boldsymbol{\sigma})$ . Subsequently, a reparameterization operation (see, e.g., Kingma and Welling<sup>50</sup>) is used to remove the dependence on the distribution to which the expectation is taken (i.e.  $q(\mathbf{w}|\theta)$ ) in the likelihood-dependent part, whereby unbiased Monte Carlo gradients can be obtained, as given below:

$$\mathbb{E}_{q(\mathbf{w}|\theta)} \log p(\mathcal{D}|\mathbf{w}) = \mathbb{E}_{\boldsymbol{\epsilon} \sim r(\boldsymbol{\epsilon})} [f(g(\boldsymbol{\epsilon}, \theta))] \simeq \frac{1}{L} \sum_{l=1}^L f(g(\boldsymbol{\epsilon}^{(l)}, \theta)) \quad (9)$$

where  $f(\mathbf{w}, \theta) = \log p(\mathcal{D}|\mathbf{w})$ ;  $L$  is the number of samples drawn for the Monte Carlo estimator;  $g(\cdot)$  is a differentiable function that transforms a parameter free noise sample,  $\boldsymbol{\epsilon}^{(l)} \sim r(\boldsymbol{\epsilon})$ , into a sample of the variational posterior:  $\mathbf{w}^{(l)} = g(\boldsymbol{\epsilon}^{(l)}, \theta) = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}^{(l)}$ , where  $r(\boldsymbol{\epsilon})$  is often modelled as standard Gaussian distribution. Otherwise, when the KL divergence term in Equation (8) is not analytically solvable, the reparameterization operation will then instead be applied to the full expectation from the cost function Equation (7), given as:  $\mathcal{J}(\mathcal{D}, \hat{\boldsymbol{\theta}}) = \mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}|\theta)} [\log q(\mathbf{w}|\theta) - \log p(\mathbf{w}) - \log p(\mathcal{D}|\mathbf{w})]$ .

In practice, when training in mini-batches (i.e., mini-batch optimization), the above implementation should be re-scaled before derivation is taken:

$$\mathcal{J}^M(\mathcal{D}_M, \theta) = \frac{1}{N} \text{KL}[q(\mathbf{w}|\theta) \parallel p(\mathbf{w})] - \frac{1}{M} \mathbb{E}_{r(\boldsymbol{\epsilon})} \log p(\mathcal{D}_M | g(\boldsymbol{\epsilon}, \theta)) \quad (10)$$

where  $M$  and  $N$  are the size of the mini batch and whole training data, respectively. Reparameterization enables the cost function to be differentiated with respect to  $\theta$ , whereby the resulting gradients can still be employed using standard stochastic optimization pipelines (e.g., stochastic gradient descent<sup>52</sup>):

$$\boldsymbol{\theta}^{\tau+1} = \boldsymbol{\theta}^\tau - \eta \nabla_{\boldsymbol{\theta}} \mathcal{J}^M(\mathcal{D}_M, \theta) \quad (11)$$

where the variational parameters are sequentially updated by mini-batches during training;  $\eta$  represents the learning rate.

## 2.5 | Stochastic process representation

For stochastic dynamic response analyses and reliability assessment, in which ground motions are represented as stochastic excitation inputs to engineering structural systems, a Monte Carlo simulation scheme plays a central part (see, e.g., Shinozuka and Deodatis<sup>2,53</sup>, Spanos and Kougoumtzoglou<sup>54</sup>, Jalayer and Beck<sup>55</sup>, Kiureghian and Fujimura<sup>3</sup>, Rezaeian and Luco<sup>56</sup>, Vlachos et al.<sup>5</sup>). Sample realizations are generated, provided the EPSD of the underlying stochastic process, whose estimation is challenging in the presence of missing data.<sup>4,18</sup> Our framework is dedicated to solving this problem. Particularly, the EPSD of the process is estimated from the ensemble average over reconstructions imputed by Equation (6) and the uncertainty on the spectral density estimates is represented by probability distributions.

Established spectral density estimation approaches, either for stationary cases or non-stationary cases, can be employed in this regard (see, e.g., Spanos and Failla<sup>57</sup>, Liang et al.<sup>58</sup>, Spanos and Kougoumtzoglou<sup>54</sup> for a review). Given the EPSD, sample realizations can, hence, be generated via a spectral representation method SRM<sup>58</sup>:

$$m(t) = \sqrt{2} \sum_{n=0}^{N-1} \sqrt{2S_Y(t, \omega_n)} \Delta\omega \cos(\omega_n t + \Phi_n) \quad (12)$$

**TABLE 1** Statistical parameters of the stochastic finite fault model.

Parameter	Distribution	Mean	s.t.d	Min	Max
$\log \Delta\sigma$	Gaussian	1.96	0.31		
$\kappa_0$	Uniform			0.002	0.008
$d$	Gaussian	9.2	10	2	30
$b_1$ (0–70 km)	Gaussian	−1.35	0.1		
$b_2$ (70–140 km)	Gaussian	−0.57	0.5		
$\nu$	Uniform			−0.15	0.15

where  $S_Y(t, \omega)$  is the two-sided EPSD of the underlying stochastic process  $\{Y(t)\}$ ;  $m(t)$  is the simulation,  $\phi_n$  is the independent random phase angle distributed uniformly over the interval  $[0, 2\pi]$ ;  $N$  and  $\Delta\omega$  relate to the discretization of the frequency domain. This enables the proposed approach to still provide a realistic representation of the non-stationary characteristics of earthquake ground motions given incomplete recordings, which is important when the associated earthquake scenarios are of interests to the seismic assessment of engineering structures, under the PBEE practice.

### 3 | APPLICATION EXAMPLES

In this section, we demonstrate the performance of the proposed framework using an accelerogram from the ESM (Engineering Strong Motion) database.<sup>59</sup> Note that when working with recorded time-histories, one can generally have a single observed seismic recording as a realization of a stochastic process, where the true power spectrum of the underlying process is typically unknown.<sup>1</sup> Therefore, the spectral estimates from the otherwise complete recording could then serve as the reference for comparison. Given a ground motion time-history record, power spectral density (PSD) estimates are derived using the Welch method<sup>60</sup> (stationary case), and the evolutionary power spectra (EPSP) are estimated from short time Fourier transform<sup>58</sup> (non-stationary case).

Region-specific parameters to the seismological model (see Equation 1) are inferred from seismographic studies of the region,<sup>61,62</sup> coupled with the event information associated with the target recording (i.e.,  $M_w = 6.5$ , normal faulting,  $R = 18.6$  km, recorded at a class A site in Italy). To consider the variability of ground motions, some key input parameters of significance are modelled as probability distributions, as shown in Table 1, while other deterministic ones are listed in the Appendix in Table A.2. In generating ground motions, the slip distribution and hypocenter location are modelled as random. Specifically, 100 physics-informed simulations with parameter variability are obtained, from which we have trained a BNN model with two hidden layers. Under the autoregressive modelling scheme, as suggested by Equation (2), the input layer is specified by the lagged width  $p$  while the output layer has one output node. Each hidden layer is composed of 16 hidden units, activated by the rectified linear function. This architecture is the result of comprehensive hyperparameter tuning (including the learning rate  $\eta$ ) based on a 20% hold-out validation set from these simulations.

#### 3.1 | Missing gaps at random locations

In this study, we focus on the effect of missing gaps, which suggest a variable length of unknown samples consecutively grouped together from an otherwise continuous set of measurements, significantly decreasing the number of usable empirical records. This situation is of particular interest to studies of historical earthquakes, which may potentially provide strong-motion records but many of them are discarded due to the presence of missing gaps.<sup>63,64</sup> For example, in a study of an Italian earthquake in 1930,<sup>65</sup> only 11 out of the 113 seismograms recovered from seismological observatories across Europe were employed mostly due to the inability to analyse incomplete seismograms.<sup>14</sup> Moreover, the presence of gaps is also common in modern seismograms subject to serious clipping in which consecutive points are clipped during peak motions.<sup>16,66</sup> Instrumentation malfunction or incompetence, or loss of communications may also lead to missing data. Other examples include instrument bandwidth limitations, low-cost temporary instruments in harsh conditions or data contamination and so forth.<sup>15,17,18,29</sup> To comprehensively investigate the effects of data gaps, various scenarios where different combinations of gap sizes (i.e., the number of missing samples) and gap number (i.e., the number of gaps) are randomly removed in the strong motion phase, are conducted in this analysis, as listed in Table A.1.

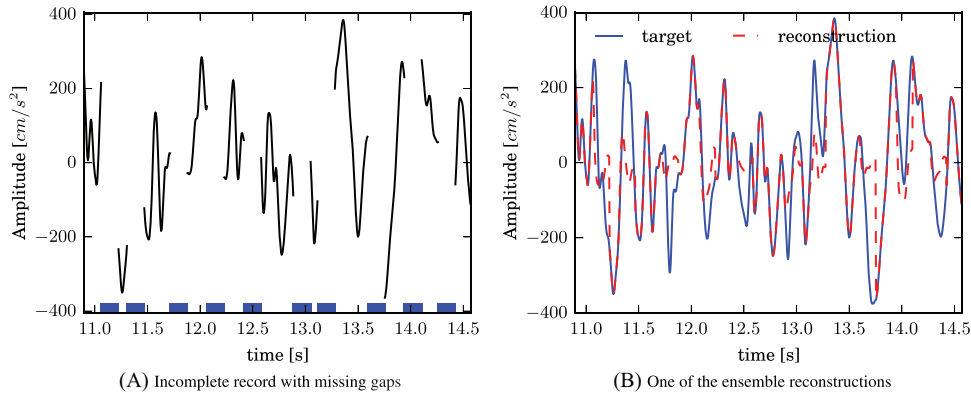


FIGURE 2 Gapped type of missing data and one reconstruction from the ensemble. Missing percentage 44%.

### 3.2 | Quantitative metrics to compare the performance

To evaluate uncertainties and accuracy under different configurations of missing data, three quantitative metrics are designed. These metrics are reported on the PSDs for characterizing the input stochastic process and on pseudo spectral accelerations (5% damped) for characterizing responses of engineering systems.  $P_{95}$  corresponds to an interval coverage probability measure that reflects the percentage of target PSD values being captured by the estimated credible intervals,<sup>67</sup> given as

$$P_{95} = \frac{c_f}{n_f} \quad (13)$$

where  $c_f$  represents the number of frequencies in which the target spectral density is captured within the 95% credible interval. Upon denoting the predicted lower and upper bound as  $y_L$  and  $y_U$ ,  $c_f$  is defined by a variable  $k_i$  of length  $n_f$  (total number of frequency bins) that indexes a frequency value captured by the estimated credible interval:

$$c_f = \sum_{i=1}^n k_i \quad (14)$$

$$k_i = \begin{cases} 1 & y_{Li} \leq y_i \leq y_{Ui} \\ 0 & \text{else} \end{cases} \quad (15)$$

In addition,  $A_{LU}$  represents the area between the lower  $y_U$  and upper bounds  $y_L$  across the frequency range, which illustrates the magnitude of uncertainty levels.  $e$  denotes the mean absolute error of the PSD estimates, which evaluates the accuracy of the mean estimation:

$$e = \frac{1}{n} \sum_{i=1}^n y_i - \hat{y}_i \quad (16)$$

### 3.3 | A detailed scenario case

Of all the scenarios considered (see Table A.1), one serious scenario case corresponding to 10 gaps of size 32, in total equivalent to 44% missing data within the strong motion phase, is specifically demonstrated herein in details for conciseness (see Figures 2–7). Figure 2A shows such incomplete recording with gaps indicated by the blue bar at the bottom. Figure 2B then shows one reconstructed time-history from the ensemble collection of 500 reconstructions by the updated BNN model, which largely resemble the waveform of the original recording. Past studies have suggested the difficulty in restoring the waveform in the time domain with missing values consecutively grouped (as in gaps), compared to missing



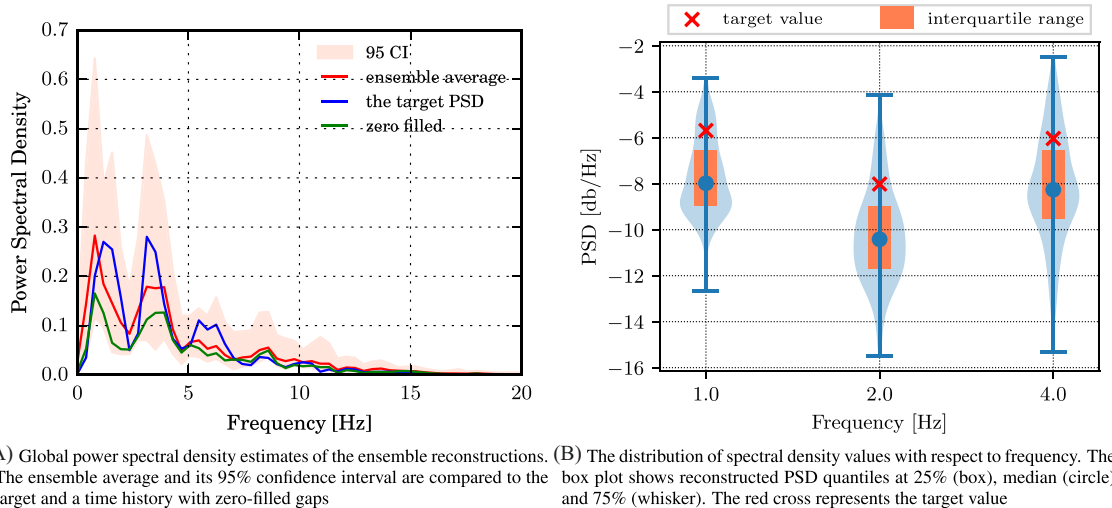


FIGURE 3 Uncertainties in the power spectral density estimates. Missing percentage 44%.

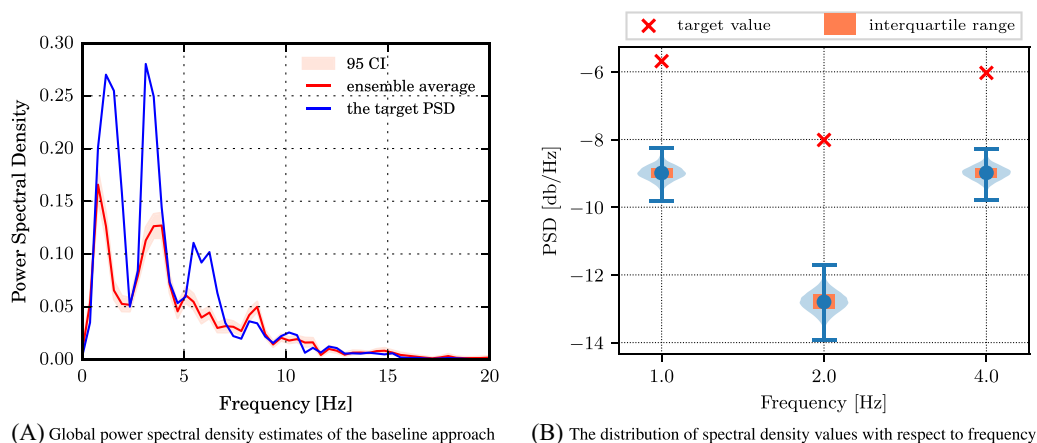


FIGURE 4 An baseline approach for comparison with the proposed approach.

values scattered across the signal.<sup>4,14,33</sup> In fact, this difficulty further justifies the importance of uncertainty quantification due to the propagation of imperfect reconstruction error.

Based on the ensemble reconstructions, the uncertainties over the power spectrum can further be seen in Figure 3A. Despite a significant portion of data missing (44%), the ensemble-averaged PSD agrees well with the target PSD from the otherwise complete recording, whose target spectral values across the whole frequency range are generally captured in the 95% credible interval bounds. The heteroscedasticity of variances with respect to frequencies is observed. As a comparison, significant power loss is seen from the result by a simple zero-padded approach. In more details, Figure 3B illustratively displays the probability distribution shape of spectral density estimates with respect to frequency. In addition, descriptive statistics regarding the ensemble-averaged PSD estimates are also depicted. The box within represents the regular box plot showing the statistics corresponding to quantiles such as 25%, median and 75%. The blue circle represents the median value while the red cross represents the target, that is, the PSD value from the full recording.

In addition, results from another baseline method, in which missing values are filled with samples from standard Gaussian distribution,<sup>31</sup> are shown in Figure 4. By contrast, our ensemble-average estimate has better approximated the target result and our interval bounds have better covered the target, as clearly seen in Figures 3B and 4B. This superior performance could be attributed to our updated BNN's ability to learn the temporal dependence of the underlying process. While the “white noise” imputation approach respects the basic property of a stochastic process, it can hardly know the variance with respect to the random variable at each time stamp and also the covariance structure.

It should be noted that the stationary (global) PSD estimates provide the spectral distribution in an average sense, without time information. But engineering interests, driven by PBEE, are increasingly focused on the time-varying spectral

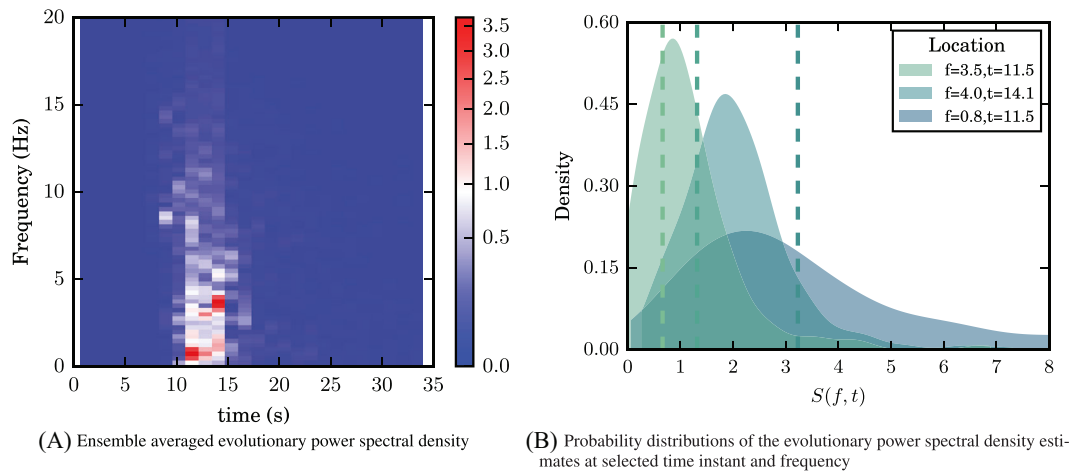


FIGURE 5 Evolutionary power spectral density estimate and its uncertainty.

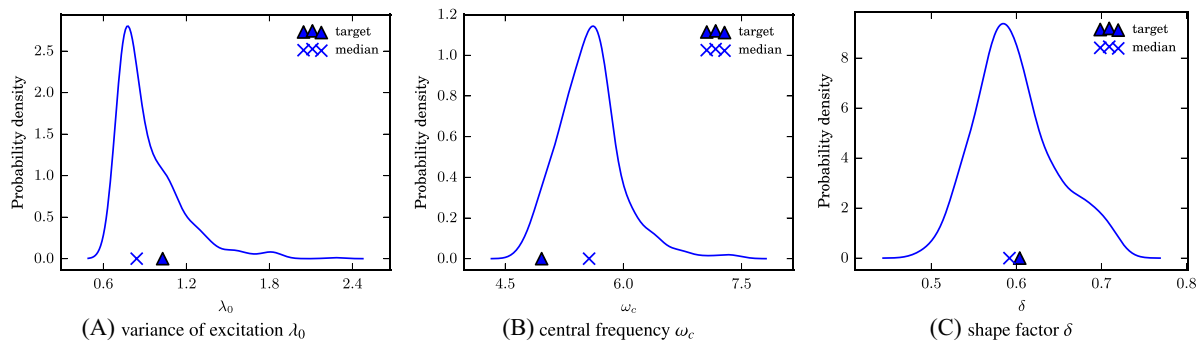


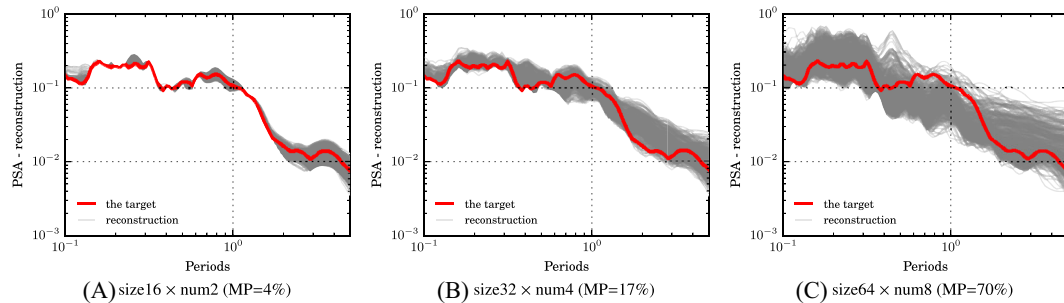
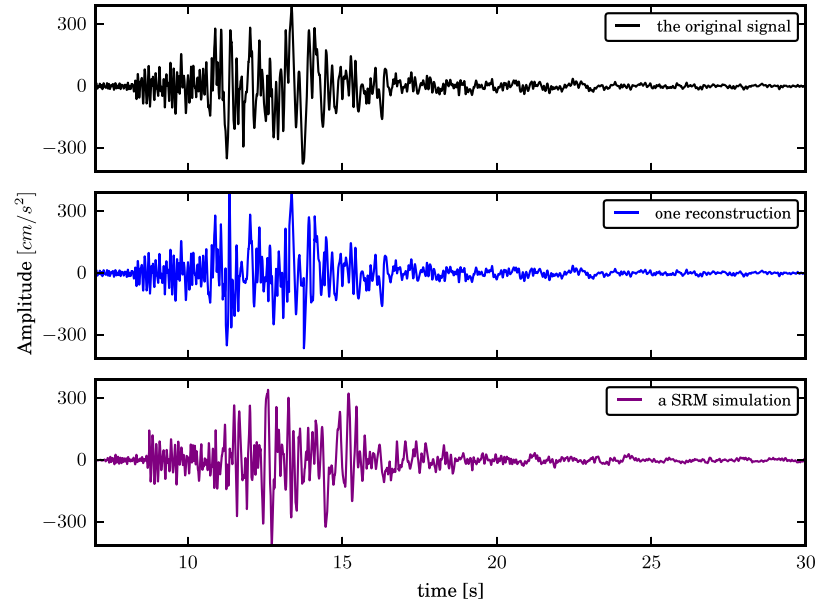
FIGURE 6 The distribution of spectral moments due to incomplete data.

representation due to the “moving resonance” effect of nonlinear structural analysis. As such, an ensemble of estimates of the evolutionary power spectrum is computed, with the averaged EPSP shown in Figure 5A; more importantly, the distribution of spectral density values,  $S(f, t)$ , at selected time instants and frequency bins are displayed in Figure 5B for illustration. Several representative combinations of time instants and frequency bins are selected to show the variance of spectral estimates. The corresponding target values are shown by the vertical lines, which are well captured by the estimated probability distributions.

Figure 6 further displays the distribution of spectral moments (see definition in D), the key parameters of spectral representation of stochastic seismic inputs.<sup>18,68</sup> Uncertainties due to the incomplete data are shown, indicating that the target values from the full recording are well captured even with a missing percentage (MP) of 44%. Spectral moments can be used to calibrate parameterized stochastic process models, for example, the established Kanai Tajimi model via a spectral moment method (see, e.g., Lai<sup>68</sup> for details). Indeed more complex models (e.g., Conte and Peng<sup>69</sup>, Vlachos et al.<sup>11</sup>) that reflect the nonstationary characteristics of ground motions could also similarly be calibrated with the ensemble reconstructions through, for example, spectral fitting. Importantly, it suggests that parameter uncertainties could thus be accounted for when characterising ground motions using parameterized models.

Relying on the Monte Carlo simulation approach,<sup>2</sup> powered by the spectral representation method SRM (Equation 12), sample realizations compatible with the given stochastic process can be simulated for stochastic nonlinear dynamic analyses (see, e.g., Jalayer and Beck<sup>55</sup>, Kiureghian and Fujimura<sup>3</sup>, Rezaeian and Luco<sup>56</sup>, Vlachos et al.<sup>5</sup>). As a result, Figure 7 illustrates, side by side, the sample generation based on the ensemble averaged EPSP estimates, along with the reconstruction directly from our updated BNN model. It suggests that, even in the presence of a significant number of data gaps, both the reconstruction and the generation resemble the target recording very well.

**FIGURE 7** Target recording (top) compared with a direct reconstruction from the updated Bayesian neural network model (middle) and a sample generation of the underlying stochastic process by the stochastic representation method (SRM) from the ensemble-averaged EPSD (bottom).



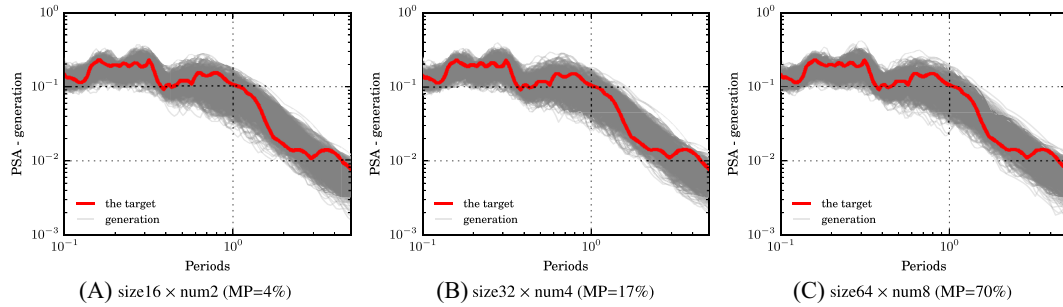
**FIGURE 8** Response spectrum of reconstructions from the updated BNN: three representative missing gap scenarios with increasing missing percentages. The target response spectrum is shown by the thick line, together with response spectra of 500 reconstructions from the ensemble. BNN, Bayesian neural network.

### 3.4 | Performance comparison of many scenarios

In earthquake engineering, accelerograms are also frequently characterized by the pseudo-acceleration (5% damped) elastic response spectra. Figure 8 illustratively shows the variability of spectral amplitudes of the reconstructions associated with three representative levels of missing gaps. The target response spectrum is shown in thick line, together with response spectra of 500 reconstructions from the ensemble. While larger uncertainty is found with increasing levels of missing data, the extreme case with roughly 70% of missing gaps still captures the target spectra to a large extent. For less extreme cases, the target response spectra are well contained within the suite of reconstructed response spectra across the full range of spectral periods. This reflects the ability of the proposed approach to quantify uncertainty in our reconstructions in response to the missing data and suggests the validity for the reconstructions to be used for seismic structural analyses.

On the other hand, the response spectra of our sample generations from the EPSD, along with the target response spectra, are displayed in Figure 9. All the sample realizations have captured the target spectra quite well. Little differences can be seen between the three data-loss scenarios, suggesting the robustness of the ensemble-averaged EPSD even under serious missing data (of up to 70%). This, therefore, validates the representation of the ground motion using estimated evolutionary power spectra by the presented approach and demonstrates its ability to make stochastic dynamic analyses still achievable in the presence of serious missing data. This result furthermore highlights the usefulness of the proposed method within a Monte Carlo simulation scheme.

For completeness, quantitative performance evaluation of the reconstructions in respect to various missing gap scenarios is tabulated in Table 2 (reported in terms of the power spectrum) and Table 3 (reported in terms of the response



**FIGURE 9** Response spectrum of sample generations from the ensemble-averaged EPSD: three representative missing gap scenarios with increasing missing ratio. EPSD, evolutionary power spectral density.

**TABLE 2** Performance comparison on power spectral density of reconstructions under various configurations of missing gaps (averaged over 10 runs).

PSD	Gap size	Number of gaps				
		2	4	6	8	10
$e$ (e-3)	16	0.958	1.181	1.935	2.282	2.879
	32	1.703	2.389	3.202	3.846	4.336
	64	2.806	4.232	5.343	7.986	-
$A_{LU}$	16	0.524	0.630	0.848	1.006	1.205
	32	0.830	1.274	1.618	2.262	2.418
	64	1.707	2.920	3.528	5.301	-
$P_{95}$ (%)	16	86.095	86.243	79.734	74.556	73.077
	32	83.876	83.432	76.479	78.107	80.030
	64	83.136	86.686	81.065	81.361	-

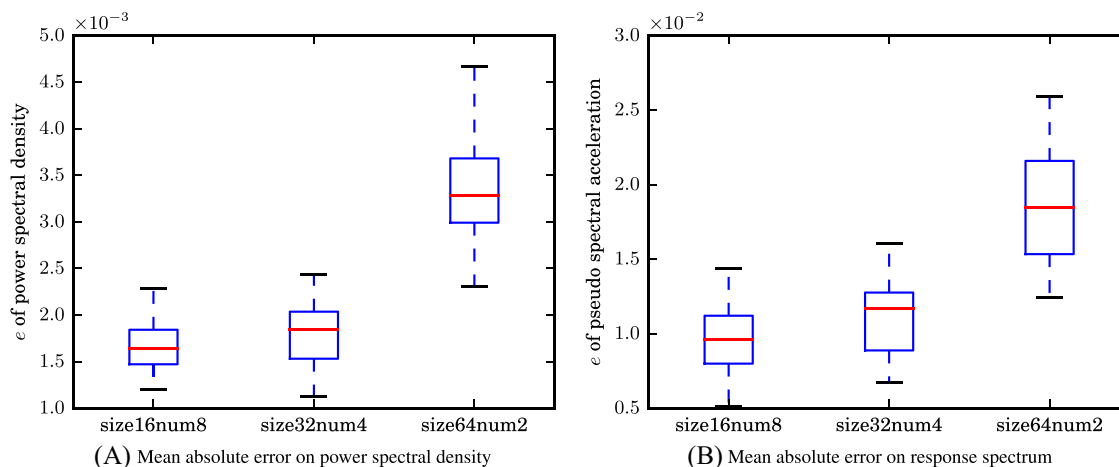
$e$  denotes the mean absolute error;  $A_{LU}$  the area metric;  $P_{95}$  prediction interval coverage probability

**TABLE 3** Performance comparison on response spectrum of reconstructions under various configurations of missing gaps (averaged over 10 runs).

PSA	Gap size	Number of gaps				
		2	4	6	8	10
$e$ (e-2)	16	0.538	0.667	1.134	1.313	1.600
	32	0.925	1.328	1.785	2.039	2.229
	64	1.621	2.029	2.658	3.157	-
$A_{LU}$	16	0.013	0.015	0.020	0.023	0.026
	32	0.020	0.029	0.035	0.043	0.045
	64	0.037	0.049	0.060	0.070	-
$P_{95}$ (%)	16	81.615	89.769	89.231	88.308	83.462
	32	80.385	84.000	86.077	82.923	88.077
	64	85.154	87.615	82.385	85.308	-

$e$  denotes the mean absolute error,  $A_{LU}$  the area metric and  $P_{95}$  prediction interval coverage probability.

spectrum), in which all the metrics are computed and averaged over 10 runs to obtain representative results against randomness. The total MP of various combinations of gap numbers and sizes are listed as a reference in a colour-coded way in Table A.1. For both spectra, larger deviations and higher uncertainties are found as with the increase of MP, which is intuitively understandable as a result of the iterative nature of the approach. Particularly, the error of PSD roughly increases by 60% when doubling the gap length (under the same gap numbers), which suggests the accumulation of errors propagated from the reconstructions. Generally, the estimated credible intervals covered both target spectrum quite well, with  $P_{95}$  higher than 80% for most scenarios. However, it should be noted that the high coverage probability of scenarios with MP



**FIGURE 10** Comparison of mean absolute error for investigating the effects of three different missing gap scenarios with same missing level.

are at the cost of wider interval bounds, as suggested by  $A_{LU}$ . The detailed scenario case in Section 3.3, along with three more scenarios shown in Figures 8 and 9, exemplify the scale of results and demonstrates the performance.

Note that, while included for completeness, the scenario with 10 gaps of size 64 is not compatible with our Bayesian updating setting, since too much of the empirical observations are missing (i.e., 87%), indicating that only very sparse samples of data are left. It is suggested by Equation (2) that the partial chunks adopted for updating should be at least the size of  $p$ .

### 3.5 | Impact of different data-loss scenarios

In addition to exploring the impacts of missing levels, this analysis further investigates more complicated patterns, since a certain missing data percentage could be associated with different scenarios, for example, a 17.41% data loss in the strong motion phase may be attributed to three combinations: eight gaps of size 16, four gaps of size 32 or two gaps of size 64. As a result, Figure 10 shows the comparison of errors on both PSD and response spectral acceleration amplitudes, over 10 runs, in box plots. For power spectral estimates, under the same missing level, the first two scenarios (namely, eight gaps of size 16 and four gaps of size 32) achieve comparable accuracy on average, though the second has slight higher error and slightly larger variability. But more significantly, the third scenario with the longest gap and least number of gaps (i.e., two gaps of size 64) has much higher error and much higher variability. For response spectral acceleration amplitudes, differences manifest a similar trend as the results in terms of power spectra. As with longer gaps, in spite of fewer gaps, the average error increases. Still, the third scenario (two gaps of size 64) results in the worst performance, with largest error and variability. This may suggest that the performance is more sensitive to the gap length (especially quite long gaps) than the quantity of gaps.

## 4 | CONCLUSION

In this paper, a Bayesian framework to characterize ground motions in the presence of missing data is presented. This framework features the setting of Bayesian model updating that allows the combination of seismological a priori knowledge, related to the physical phenomena, with the empirical yet incomplete observations. Uncertainties are accounted for throughout the framework. The effect of missing gaps has been comprehensively studied via various missing scenarios, based on which the performance of the proposed method has been quantitatively demonstrated. Results show that the proposed method is highly effective even in serious cases of data-loss with about half of data missing in the strong motion phase, being capable of providing imputed waveforms, spectral estimates and stochastic synthetic generations that agree well with the target recording. A host of representations of ground motion, consistent with an underlying stochastic process, is provided in a probabilistic manner, suggesting the versatility of the proposed approach as a general solution to



dealing with missing data for various engineering and seismological applications, whether waveform-based or spectrum-based. The proposed approach helps in recovering the information conveyed from faulty or incomplete observations, for example, from low-cost temporary instruments deployed at scale. The Bayesian framework provides a building block on which it could be developed to enrich the database of ground motions in data scarce areas (e.g., near-field strong motions), facilitating stochastic dynamic analyses of engineering structures and boosting the understanding of earth structures. Of particular note is its mechanism that combines a priori information with empirical observations, remedying the causality dilemma concerning the dependence of observations and the extracted knowledge/information. Finally, we consider that such Bayesian framework could serve as a complementary approach to current stochastic ground-motion models under the growing interests of PBEE (performance-based earthquake engineering), and ultimately a fundamental solution to the limited data problem in data scarce regions.

## ACKNOWLEDGEMENTS

This work was supported by the European Union Horizon 2020 Marie Skłodowska-Curie Actions project URBASIS [Project no. 813137].

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are available from the corresponding author upon reasonable request.

## ORCID

Yu Chen  <https://orcid.org/0000-0001-6617-2946>

Edoardo Patelli  <https://orcid.org/0000-0002-5007-7247>

## REFERENCES

- Narayana Iyengar R, Sundara Raja Iyengar K. A nonstationary random process model for earthquake accelerograms. *Bull Seismol Soc Am*. 1969;59:1163-1188.
- Shinozuka M, Deodatis G. Stochastic process models for earthquake ground motion. *Probab Eng Mech*. 1988;3:114-123.
- Kiureghian AD, Fujimura K. Nonlinear stochastic dynamic analysis for performance-based earthquake engineering. *Earthq Eng Struct Dyn*. 2009;38:719-738.
- Comerford L, Jensen H, Mayorga F, Beer M, Kougioumtzoglou I. Compressive sensing with an adaptive wavelet basis for structural system response and reliability analysis under missing data. *Comput Struct*. 2017;182:26-40.
- Vlachos C, Papakonstantinou KG, Deodatis G. Structural applications of a predictive stochastic ground motion model: assessment and use. *ASME J Risk Uncertain Eng Syst A Civ*. 2018;4:04018006.
- Atkinson GM, Boore DM. Earthquake ground-motion prediction equations for eastern North America. *Bull Seismol Soc Am*. 2006;96:2181-2205.
- Vetter C, Taflanidis AA. Comparison of alternative stochastic ground motion models for seismic risk characterization. *Soil Dyn Earthq Eng*. 2014;58:48-65.
- Tsioulou A, Taflanidis AA, Galasso C. Modification of stochastic ground motion models for matching target intensity measures. *Earthq Eng Struct Dyn*. 2018;47:3-24.
- Rezaeian S, Der Kiureghian A. Simulation of synthetic ground motions for specified earthquake and site characteristics. *Earthq Eng Struct Dyn*. 2010;39:1155-1180.
- Pousse G, Bonilla LF, Cotton F, Margerin L. Nonstationary stochastic simulation of strong ground-motion time histories including natural variability: application to the K-Net the Japanese database. *Bull. Seismol. Soc. Am*. 2006;96(6):2103-2117. <https://doi.org/10.1785/0120050134>
- Vlachos C, Papakonstantinou KG, Deodatis G. Predictive model for site specific simulation of ground motions based on earthquake scenarios. *Earthq Eng Struct Dyn*. 2018;47:195-218.
- Edwards B, Fäh D. A stochastic ground-motion model for Switzerland. *Bull Seismol Soc Am*. 2013;103:78-98.
- Baker J, Bradley B, Stafford P. *Physics-based Ground-Motion Characterization*. Cambridge University Press. 2021:196-246. <https://doi.org/10.1017/9781108425056.007>
- Maranò S, Edwards B, Ferrari G, Fäh D. Fitting earthquake spectra: colored noise and incomplete data. *Bull Seismol Soc Am*. 2017;107:276-291.
- Smith-Boughner L, Constable C. Spectral estimation for geophysical time-series with inconvenient gaps. *Geophys J Int*. 2012;190:1404-1422.
- Zhang J, Hao J, Zhao, et al. Restoration of clipped seismic waveforms using projection onto convex sets method. *Sci Rep*. 2016;6:1-10.
- Comerford L, Kougioumtzoglou IA, Beer M. Compressive sensing based stochastic process power spectrum estimation subject to missing data. *Probab Eng Mech*. 2016;44:66-76.
- Zhang Y, Comerford L, Kougioumtzoglou IA, Patelli E, Beer M. Uncertainty quantification of power spectrum and spectral moments estimates subject to missing data. *ASCE ASME J Risk Uncertain Eng Syst A Civ*. 2017;3:04017020.
- Scargle JD. Studies in astronomical time series analysis. II-statistical aspects of spectral analysis of unevenly spaced data. *Astrophys J*. 1982;263:835-853.

20. Roberts DH, Lehar J, Dreher JW. Time series analysis with clean-part one-derivation of a spectrum. *Astron J.* 1987;93:968.
21. Bos R, De Waele S, Broersen PM. Autoregressive spectral estimation by application of the burg algorithm to irregularly sampled data. *IEEE Trans Instrum Meas.* 2002;51:1289-1294.
22. Wang Y, Stoica P, Li J, Marzetta TL. Nonparametric spectral analysis with missing data via the em algorithm. *Digital Signal Process.* 2005;15:191-206.
23. Babu P, Stoica P. Spectral analysis of nonuniformly sampled data – a review. *Digital Signal Process.* 2010;20:359-378.
24. Broersen PM, De Waele S, Bos R. Autoregressive spectral analysis when observations are missing. *Automatica.* 2004;40:1495-1504.
25. Hung JC. A genetic algorithm approach to the spectral estimation of time series with noise and missed observations. *Inf Sci.* 2008;178:4632-4643.
26. Stoica P, Larsson EG, Li J. Adaptive filter-bank approach to restoration and spectral analysis of gapped data. *Astron J.* 2000;120:2163.
27. Kondrashov D, Ghil M. Spatio-temporal filling of missing points in geophysical data sets. *Nonlinear Processes Geophys.* 2006;13:151-159.
28. Kondrashov D, Denton R, Shprits Y, Singer H. Reconstruction of gaps in the past history of solar wind parameters. *Geophys Res Lett.* 2014;41:2702-2707.
29. Comerford L, Kougoumtzoglou IA, Beer M. An artificial neural network approach for stochastic process power spectrum estimation subject to missing data. *Struct Saf.* 2015;52:150-160.
30. Musial JP, Verstraete MM, Gobron N. Comparing the effectiveness of recent algorithms to fill and smooth incomplete and noisy time series. *Atmos Chem Phys.* 2011;11:7905-7923.
31. Comerford L, Kougoumtzoglou IA, Beer M. On quantifying the uncertainty of stochastic process power spectrum estimates subject to missing data. *Int J Sustain Mater Struct Syst.* 2015;2:185-206.
32. Tobar F. Bayesian nonparametric spectral estimation. In: Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R, eds. *Advances in Neural Information Processing Systems.* Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2018/file/abd1c782880cc59759f4112fda0b8f98-Paper.pdf>
33. Christmas J. The effect of missing data on robust Bayesian spectral analysis. In: 2013 IEEE International Workshop on Machine Learning for Signal Processing (MLSP). IEEE; 2013:1-6.
34. Zeng Y, Anderson JG, Yu G. A composite source model for computing realistic synthetic strong ground motions. *Geophys Res Lett.* 1994;21:725-728.
35. Boore DM. Simulation of ground motion using the stochastic method. *Pure Appl Geophys.* 2003;160:635-676.
36. Edwards B, Zurek B, Van Dedem E, et al. Simulations for the development of a ground motion model for induced seismicity in the Groningen gas field, the Netherlands. *Bull Earthq Eng.* 2019;17:4441-4456.
37. Vetter C, Taflanidis AA. Global sensitivity analysis for stochastic ground motion modeling in seismic-risk assessment. *Soil Dyn Earthq Eng.* 2012;38:128-143.
38. McCallen D, Petersson A, Rodgers A, et al. Eqsim – a multidisciplinary framework for fault-to-structure earthquake simulations on exascale computers part I: computational models and workflow. *Earthq Spectra.* 2021;37:707-735.
39. McCallen D, Petrone F, Miah M, Pitarka A, Rodgers A, Abrahamson N. Eqsim – a multidisciplinary framework for fault-to-structure earthquake simulations on exascale computers, part II: regional simulations of building response. *Earthq Spectra.* 2021;37:736-761.
40. Paolucci R, Smerzini C, Vanini M. BB-SPEEDset: a validated dataset of broadband near-source earthquake ground motions from 3D physics-based numerical simulations. *Bull Seismol Soc Am.* 2021;111:2527-2545.
41. Salinas D, Flunkert V, Gasthaus J, Januschowski T. DeepAR: probabilistic forecasting with autoregressive recurrent networks. *Int J Forecast.* 2020;36:1181-1191.
42. Beer M, Spanos PD. A neural network approach for simulating stationary stochastic processes. *Struct Eng Mech.* 2009;32:71-94.
43. Gatti F, Clouteau D. Towards blending physics-based numerical simulations and seismic databases using generative adversarial network. *Comput Methods Appl Mech Eng.* 2020;372:113421.
44. Williams CK, Rasmussen CE. *Gaussian Processes for Machine Learning.* Vol 2. MIT Press Cambridge; 2006.
45. Graves A. Generating sequences with recurrent neural networks. arXiv preprint arXiv:1308.0850.2013.
46. Blundell C, Cornebise J, Kavukcuoglu K, Wierstra D. Weight uncertainty in neural network. In: International Conference on Machine Learning. PMLR; 2015:1613-1622.
47. Gal Y, Ghahramani Z. Dropout as a Bayesian approximation: representing model uncertainty in deep learning. In: International Conference on Machine Learning. PMLR; 2016:1050-1059.
48. Hernández-Lobato JM, Adams R. Probabilistic backpropagation for scalable learning of Bayesian neural networks. In: International Conference on Machine Learning. PMLR; 2015:1861-1869.
49. Graves A. Practical variational inference for neural networks. *Adv Neural Inf Process Syst.* 2011;24:2348-2356.
50. Kingma DP, Welling M. Auto-encoding variational Bayes. arXiv preprint arXiv:1312.6114. 2013.
51. Blei DM, Kucukelbir A, McAuliffe JD. Variational inference: a review for statisticians. *J Am Stat Assoc.* 2017;112:859-877. <https://doi.org/10.1080/01621459.2017.1285773>
52. Bottou L. Stochastic gradient descent tricks. In: *Neural Networks: Tricks of the Trade.* Springer; 2012:421-436.
53. Shinozuka M, Deodatis G. Simulation of stochastic processes by spectral representation. *ASME Appl Mech Rev.* 1991;44(4):191-204. <https://doi.org/10.1115/1.3119501>
54. Spanos P, Kougoumtzoglou I. Harmonic wavelets based statistical linearization for response evolutionary power spectrum determination. *Probab Eng Mech.* 2012;27:57-68.

55. Jalayer F, Beck J. Effects of two alternative representations of ground-motion uncertainty on probabilistic seismic demand assessment of structures. *Earthq Eng Struct Dyn*. 2008;37:61-79.
56. Rezaeian S, Luco N. Example applications of a stochastic ground motion simulation methodology in structural engineering. In: 15th World Conf. Earthquake Engineering (WCEE). 2012.
57. Spanos PD, Failla G. Evolutionary spectra estimation using wavelets. *J Eng Mech*. 2004;130:952-960.
58. Liang J, Chaudhuri SR, Shinozuka M. Simulation of nonstationary stochastic processes by spectral representation. *J Eng Mech*. 2007;133:616-627.
59. Lanzano G, Luzi L, Cauzzi C, et al. Accessing european strong-motion data: an update on ORFEUS coordinated services. *Seismol Res Lett*. 2021;92:1642-1658.
60. Welch P. The use of fast Fourier transform for the estimation of power spectra: a method based on time averaging over short, modified periodograms. *IEEE Trans Audio Electroacoust*. 1967;15:70-73.
61. Bindi D, Kotha S. Spectral decomposition of the engineering strong motion (ESM) flat file: regional attenuation, source scaling and Arias stress drop. *Bull Earthq Eng*. 2020;18:2581-2606.
62. Razafindrakoto HN, Cotton F, Bindi D, Pilz M, Graves RW, Bora S. Regional calibration of hybrid ground-motion simulations in moderate seismicity areas: application to the Upper Rhine Graben. *Bull Seismol Soc Am*. 2021;111:1422-1444.
63. Church ED, Bartlett AH, Jourabchi MA. Raster-to-vector image analysis for fast digitization of historic seismograms. *Seismol Res Lett*. 2013;84:489-494.
64. Palombo B, Pino NA. On the recovery and analysis of historical seismograms. *Ann Geophys*. 2013;56:1-14.
65. Vannoli P, Vannucci G, Bernardi F, Palombo B, Ferrari G. The source of the 30 October 1930 Mw 5.8 Senigallia (Central Italy) earthquake: a convergent solution from instrumental, macroseismic, and geological data. *Bull Seismol Soc Am*. 2015;105:1548-1561. <https://doi.org/10.1785/0120140263>
66. Yang W, Ben-Zion Y. An algorithm for detecting clipped waveforms and suggested correction procedures. *Seismol Res Lett*. 2010;81:53-62.
67. Pearce T, Brintrup A, Zaki M, Neely A. High-quality prediction intervals for deep learning: a distribution-free, ensembled approach. In: International Conference on Machine Learning. PMLR; 2018:4075-4084.
68. Lai SSP. Statistical characterization of strong ground motions using power spectral density function. *Bull Seismol Soc Am*. 1982;72:259-274.
69. Conte J, Peng B. Fully nonstationary analytical earthquake ground-motion model. *J Eng Mech-ASCE*. 1997;123:15-24.
70. Mohamed S, Rosca M, Figurnov M, Mnih A. Monte Carlo gradient estimation in machine learning. *J Mach Learn Res*. 2020;21:1-62.

**How to cite this article:** Chen Y, Patelli E, Edwards B, Beer M. A physics-informed Bayesian framework for characterizing ground motion process in the presence of missing data. *Earthquake Engng Struct Dyn*. 2023;1-17. <https://doi.org/10.1002/eqe.3877>

## APPENDIX A: MISSING PERCENTAGES FOR VARIOUS SCENARIOS

TABLE A.1 The total missing percentage (MP) for various missing scenarios.

Gap size	Number of gaps				
	2	4	6	8	10
16	4.35	8.71	13.06	17.41	21.77
32	8.71	17.41	26.12	34.83	43.54
64	17.41	34.83	52.24	69.66	87.07

## APPENDIX B: MONTE CARLO ESTIMATOR

Consider a general probabilistic objective function of the form:

$$\mathcal{F}(\theta) = \int p(\mathbf{x}; \theta) f(\mathbf{x}; \phi) d\mathbf{x} = \mathbb{E}_{p(\mathbf{x}; \theta)} [f(\mathbf{x}; \phi)] \quad (\text{B.1})$$

where  $f(\mathbf{x}; \phi)$  denotes a general function of an input variable  $\mathbf{x}$  with structural parameters  $\phi$ ;  $p(\mathbf{x}; \theta)$  represents a probability distribution parameterized by  $\theta$ .

The usual Monte Carlo estimator for expectation is given by

$$\mathbb{E}_{p(\mathbf{x};\theta)}[f(\mathbf{x};\boldsymbol{\phi})] \simeq \frac{1}{N} \sum_1^N f(\hat{\mathbf{x}}^{(n)}), \text{ where } \hat{\mathbf{x}}^{(n)} \sim p(\mathbf{x};\theta) \quad (\text{B.2})$$

It suggests that a complex integral in Equation (A.1) can be numerically evaluated by drawing samples from the probability distribution  $p(\mathbf{x};\theta)$  and then computing the average of the function evaluated at these samples. Furthermore, as many problems in Machine Learning generally focused on the computation of gradients, such as  $\nabla_{\theta} \mathbb{E}_{p(\mathbf{x};\theta)}[f(\mathbf{x};\boldsymbol{\phi})]$ . Several techniques exist to do further approximation, see additional details in Mohamed et al.<sup>70</sup> As an example, a Monte Carlo gradient estimator by the score function is given as

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p(\mathbf{x};\theta)}[f(\mathbf{x};\boldsymbol{\phi})] &= \mathbb{E}_{p(\mathbf{x};\theta)}[f(\mathbf{x};\boldsymbol{\phi}) \nabla_{\theta} \log p(\mathbf{x};\theta)] \\ &= \frac{1}{N} \sum_1^N f(\hat{\mathbf{x}}^{(n)}) \nabla_{\theta} \log p(\hat{\mathbf{x}}^{(n)}; \theta) \end{aligned}$$

where  $\hat{\mathbf{x}}^{(n)} \sim p(\mathbf{x};\theta)$

## APPENDIX C: SEISMOLOGICAL PARAMETERS OF THE FINITE-FAULT MODEL

TABLE C.1 Source and path parameters of the stochastic finite fault model (sourced from Refs. 61, 62).

Parameter	Description	Value
$\rho_s$	Density of the medium	2.7
$\beta$	Shear wave velocity	3.2
$V$	Horizontal partition	$1/\sqrt{2}$
$R_{\theta\Phi}$	Radiation pattern	0.55
$F$	Free-surface factor	2
$R_0$	Reference distance	10
$Q$	Quality factor	$Q = 250.4 f^{0.29}$

## APPENDIX D: SPECTRAL MOMENTS

The spectral moments are key statistical parameters in frequency domain analyses, which are of particular importance in evaluating survival probability or reliability assessment for structural systems. Consider stationary random processes, the  $j$ th spectral moment  $\lambda_j$  are given as<sup>18,68</sup>

$$\lambda_j = \int_{-\infty}^{+\infty} \omega^j S(\omega) d\omega \quad (\text{D.1})$$

where  $S(\omega)$  denotes the two-sided PSD. Specifically, the zero spectral moment  $\lambda_0$ , which is also the variance of the excitation, is given as

$$\lambda_0 = \int_{-\infty}^{+\infty} S(\omega) d\omega \quad (\text{D.2})$$

then the central frequency  $\omega_c$ , and the shape factor  $\delta$  (also known as bandwidth measure) of the stochastic process can be computed from the first few spectra moments:

$$\begin{aligned} \omega_c &= [\lambda_1/\lambda_2]^{1/2} \\ \delta &= [1 - (\lambda_1^2/\lambda_0\lambda_2)]^{1/2} \end{aligned}$$