



# CHORUS

This is the accepted manuscript made available via CHORUS. The article has been published as:

## Quantification of uncertainties in thermoelectric properties of materials from a first-principles prediction method: An approach based on Gaussian process regression

Daehyun Wee, Jeeyoung Kim, Semi Bang, Georgy Samsonidze, and Boris Kozinsky

Phys. Rev. Materials **3**, 033803 — Published 25 March 2019

DOI: [10.1103/PhysRevMaterials.3.033803](https://doi.org/10.1103/PhysRevMaterials.3.033803)

1 **Quantification of Uncertainties in Thermoelectric Properties of**  
2 **Materials from a First-Principles Prediction Method: An**  
3 **Approach Based on Gaussian Process Regression**

4 Daehyun Wee,\* Jeeyoung Kim, and Semi Bang

5 *Department of Environmental Science and Engineering,*

6 *Ewha Womans University, Seoul, 03760, Republic of Korea*

7 Georgy Samsonidze

8 *Research and Technology Center, Robert Bosch LLC, Cambridge, MA 02142, U.S.A.*

9 Boris Kozinsky

10 *Harvard John A. Paulson School of Engineering and Applied Sciences,*

11 *Harvard University, Cambridge, MA 02138, U.S.A.*

12 (Dated: February 27, 2019)

## Abstract

We present the electron-phonon averaged via Gaussian process regression (EPA-GPR) method, in which the electron-phonon coupling matrix is represented as a function of two energies and is in turn modeled as a Gaussian process. The EPA-GPR method can be used as an efficient method to estimate thermoelectric properties of materials for fast-screening applications, comparable to the original electron-phonon averaged (EPA) method and the electron-phonon averaged via moving-least-squares (EPA-MLS) method. The EPA-GPR method does not require specification of any open parameter, unlike the other EPA-related methods, since all the hyperparameters in the model can be unambiguously estimated within the type II maximum likelihood (ML-II) approximation. Thus, the EPA-GPR method is a parameter-free estimation method. Additionally, the concept of Gaussian processes in the EPA-GPR method allows us to quantify the uncertainty in estimated properties of thermoelectric materials. One can randomly realize the electron-phonon coupling coefficients from the identified Gaussian process, and those realized samples can be further analyzed in the solution process of the semiclassical Boltzmann transport equation for charge carriers. The results of the semiclassical Boltzmann transport equation provide the statistical properties of the thermoelectric properties of interest. The means, standard deviations, histograms, and confidence intervals of the Seebeck coefficient, the electrical conductivity, and the power factor can be constructed and analyzed. The proposed EPA-GPR method is applied to a  $p$ -type half-Heusler compound, i.e., HfCoSb, as a case example, the results of which clearly present the advantages of the method.

## 13 I. INTRODUCTION

14 The performance of thermoelectric (TE) energy conversion, where thermal energy is di-  
15 rectly converted into electrical energy or vice versa by using thermoelectric materials, de-  
16 pends on the thermoelectric figure of merit, i.e.,  $ZT$  of the material<sup>1</sup>.  $ZT$  is defined as a  
17 combination of thermal and electrical properties of the material, as follows:

$$ZT = \frac{S^2 \sigma T}{k}, \quad (1)$$

18 where  $S$  is the Seebeck coefficient,  $\sigma$  is the electrical conductivity,  $k$  is the thermal con-  
19 ductivity of the material, and  $T$  is the operating temperature. Researchers in the compu-  
20 tational materials science community interested in thermoelectric energy conversion have  
21 exerted tremendous efforts to develop methods that can be used for estimating these prop-  
22 erties from the first principles to discover better thermoelectric materials. For example,  
23 many computational studies estimating  $k$  and other related thermal properties have been  
24 reported<sup>2-4</sup>.

25 On the other hand, the electrical properties of inorganic materials, including  $S$  and  $\sigma$ ,  
26 can be obtained by solving the semiclassical Boltzmann transport equation within the re-  
27 laxation time approximation<sup>5</sup>. The simplest approach that can be used for estimation of  
28 these electronic transport properties is the constant relaxation time (CRT) approximation,  
29 in which one single value for relaxation time  $\tau$  is arbitrarily assumed<sup>6</sup>. However, such an ap-  
30 proach naturally introduces an arbitrary constant, i.e., relaxation time, and does not possess  
31 any predictive capacity, rendering the approach unsatisfactory for screening thermoelectric  
32 materials from the first principles. It is, therefore, necessary to develop a more predictive  
33 method for dealing with the relaxation time.

34 Matthiessen's rule states that the total scattering rate  $\tau^{-1}$  of electrons is the sum of the  
35 rates associated with intrinsic (electron-electron, electron-phonon) and extrinsic (impurities,  
36 grain boundaries, alloy disorder) scattering mechanisms. To screen potentially promising  
37 candidates for thermoelectric applications, one must first identify the intrinsic properties  
38 of the material, since the extrinsic properties are tuned during the synthesis process. In  
39 automotive TE power generation, the relevant temperature is around 400°C at the hot side  
40 of the device, at which electron-phonon (el-ph) interaction becomes the dominant scattering  
41 mechanism<sup>7,8</sup>. The first-principles estimation of the el-ph interaction has been pursued

42 by several different approaches with various levels of computational complication. The  
43 deformation potential (DP) approximation is one of the simplest approaches<sup>9</sup>, but such  
44 simplification often cannot be justified for complex TE materials. The other extreme is  
45 the electron-phonon Wannier (EPW) method<sup>10</sup>, which fully describes the el-ph scattering.  
46 However, the EPW method is not appropriate for fast-screening applications due to its high  
47 computational cost.

48 More recently, a new approach, i.e., the electron-phonon averaged (EPA) method, which  
49 combines simplicity and speed with a fully first-principles treatment of the el-ph interaction,  
50 has been introduced<sup>11</sup>. By turning the complex momentum-space integration into an inte-  
51 gration over energies and simultaneously replacing several terms with their averages within  
52 bins over an energy range, the EPA method allows for automated rapid calculations for op-  
53 timization of electronic transport quantities, while being more predictive than the CRT and  
54 DP approximations. The method has been successfully used for screening potential TE ma-  
55 terials from a group of half-Heusler (HH) compounds<sup>11</sup>. Later, it was proposed to modify the  
56 standard EPA method through combination with a moving-least-squares (MLS) averaging  
57 strategy<sup>12</sup>. It was demonstrated that the electron-phonon averaged via moving-least-squares  
58 (EPA-MLS) method could make a similar prediction of thermoelectric properties of materials  
59 with a much coarser momentum grid than was required for the standard EPA method<sup>12</sup>.

60 However, several limitations remain in the EPA and EPA-MLS methods. First, these  
61 methods require specification of an open parameter, i.e., either the bin size or the length  
62 scale of the smoothing kernel. Second, although the sample variance during the averaging  
63 process can roughly provide the amount of uncertainty in the estimated electron-phonon  
64 coupling effects, rigorous analysis of uncertainty and sensitivity can be difficult within the  
65 EPA and EPA-MLS methods. The first problem is of minor importance, especially because  
66 the result of the EPA-MLS method seems rather insensitive to particular choices of the open  
67 parameter<sup>12</sup>. The lack of a rigorous strategy for uncertainty quantification in the numerical  
68 procedure is a more serious issue that requires immediate attention. One should not place  
69 blind confidence in his or her prediction without describing the underlying uncertainty. The  
70 same issue is essentially shared by most of the first-principles methods used in the study of  
71 thermoelectric properties. None of the methods we have mentioned so far, i.e., the CRT,  
72 DP, and EPW methods, currently has a rigorous quantification strategy of uncertainty in  
73 its numerical procedures.

74 In this paper, we investigate the possibility of using a mathematically rigorous alternative  
75 method. Here, the electron-phonon coupling matrix is modeled as a Gaussian process, which  
76 is widely used in the context of machine learning<sup>13</sup>. During regression, the characteristic  
77 length-scale of the covariance function of the Gaussian process, which serves a similar pur-  
78 pose to the smoothing scale in the EPA-MLS method, can be estimated within the type II  
79 maximum likelihood (ML-II) approximation without any ambiguity. At the same time, all  
80 the statistical tools that can be used for Gaussian processes are readily available for further  
81 analysis of uncertainty and sensitivity of the results. The resulting formulation, i.e., the  
82 electron-phonon averaged via Gaussian process regression (EPA-GPR) method, may resolve  
83 the above issues.

84 The paper is organized as follows. We first describe the basic theory of the EPA-GPR  
85 method in Section II. We continue to test the method on a *p*-type HH compound in Sec-  
86 tion III. The values of the thermoelectric properties, i.e.,  $S$  and  $\sigma$ , estimated by the EPA-  
87 GPR method are compared to those using other related methods and experiments. The  
88 uncertainties in the thermoelectric properties are also quantified by the method described  
89 in Section II. A brief summary follows in Section IV.

## 90 II. THEORY

### 91 A. The EPA and EPA-MLS Methods

92 We first briefly review the main features of the EPA and EPA-MLS methods. Details  
93 may be found in<sup>12</sup>. The main task of predicting the electronic transport coefficients for  
94 electrons within the relaxation time approximation is evaluation of the inverse of the electron  
95 energy relaxation time induced by the electron-phonon (el-ph) interaction, which is given as  
96 follows<sup>14,15</sup>:

$$\begin{aligned} \tau_{n\mathbf{k}}^{-1}(\mu, T) = & \frac{\Omega}{(2\pi)^2 \hbar} \sum_{m\nu} \int_{BZ} d\mathbf{q} |g_{m\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2 \\ & \times \left\{ \left[ n(\omega_{\nu\mathbf{q}}, T) + f(\epsilon_{m\mathbf{k}+\mathbf{q}}, \mu, T) \right] \delta(\epsilon_{n\mathbf{k}} + \omega_{\nu\mathbf{q}} - \epsilon_{m\mathbf{k}+\mathbf{q}}) \right. \\ & \left. + \left[ n(\omega_{\nu\mathbf{q}}, T) + 1 - f(\epsilon_{m\mathbf{k}+\mathbf{q}}, \mu, T) \right] \delta(\epsilon_{n\mathbf{k}} - \omega_{\nu\mathbf{q}} - \epsilon_{m\mathbf{k}+\mathbf{q}}) \right\}, \quad (2) \end{aligned}$$

97 where  $\Omega$  is the volume of the primitive cell,  $m$  and  $n$  are the electron band indices,  $\nu$  is the  
 98 phonon mode index,  $\mathbf{k}$  is the electron wavevector,  $\mathbf{q}$  is the phonon wavevector,  $\epsilon_{n\mathbf{k}}$  is the  
 99 electron energy,  $\omega_{\nu\mathbf{q}}$  is the phonon energy,  $g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})$  is the el-ph coupling matrix element,  
 100  $n(\omega, T)$  is the Bose-Einstein distribution function,  $f(\epsilon, \mu, T)$  is the Fermi-Dirac distribution  
 101 function,  $\delta$  is the Dirac delta function,  $\mu$  is the chemical potential of electrons,  $k_B$  is the  
 102 Boltzmann constant, and  $\hbar$  is the reduced Planck constant.

103 The el-ph coupling matrix elements, i.e.,  $g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})$ , can be obtained from the DFPT  
 104 calculations<sup>16</sup>, which are relatively costly for materials with a large unit cell. The main ele-  
 105 ment of the EPA approximation is to replace the momentum-dependent quantities in Eq. (2)  
 106 with their energy-dependent averages. Accordingly, the el-ph coupling matrix elements are  
 107 averaged over the directions of  $\mathbf{k}$  and  $\mathbf{k} + \mathbf{q}$  wavevectors:

$$|g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2 \mapsto g_\nu^2(\epsilon_{n\mathbf{k}}, \epsilon_{m\mathbf{k}+\mathbf{q}}). \quad (3)$$

108 As a result,  $g_\nu^2$  becomes a function of two energies,  $\epsilon_1$  and  $\epsilon_2$ , which represent the energy of the  
 109 incoming electron state and that of the outgoing electron state, respectively. Additionally,  
 110  $\omega_{\nu\mathbf{q}}$  is also replaced with its average:

$$\omega_{\nu\mathbf{q}} \mapsto \bar{\omega}_\nu. \quad (4)$$

111 With these substitutions, the integration over  $\mathbf{q}$  and the summation over  $m$  in Eq. (2) can  
 112 be evaluated analytically, yielding

$$\begin{aligned} \tau^{-1}(\epsilon, \mu, T) = \frac{2\pi\Omega}{g_s \hbar} \sum_\nu & \left\{ g_\nu^2(\epsilon, \epsilon + \bar{\omega}_\nu) \left[ n(\bar{\omega}_\nu, T) + f(\epsilon + \bar{\omega}_\nu, \mu, T) \right] \rho(\epsilon + \bar{\omega}_\nu) \right. \\ & \left. + g_\nu^2(\epsilon, \epsilon - \bar{\omega}_\nu) \left[ n(\bar{\omega}_\nu, T) + 1 - f(\epsilon - \bar{\omega}_\nu, \mu, T) \right] \rho(\epsilon - \bar{\omega}_\nu) \right\}. \quad (5) \end{aligned}$$

113 Here,  $\rho(\epsilon)$  is the electron density of states defined as the number of electronic states per unit  
 114 energy and unit volume, and  $g_s = 2$  is the spin degeneracy.

115 Various methods can be used to achieve the mapping of Eq. (3). In the original EPA  
 116 method<sup>11</sup>, a bin-based averaging strategy was employed with a predefined bin size  $\delta_{\text{Bin}}$ . On

117 the other hand, in<sup>12</sup>, we proposed the use of an MLS averaging strategy<sup>17</sup>, where  $g_\nu^2(\epsilon_1, \epsilon_2)$   
 118 for each pair of  $\epsilon_1$  and  $\epsilon_2$  is obtained by minimizing

$$\sum_{mn} \iint_{BZ} d\mathbf{k}d\mathbf{q} \left( g_\nu^2(\epsilon_1, \epsilon_2) - |g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2 \right)^2 \times \exp \left( -\frac{(\epsilon_{n\mathbf{k}} - \epsilon_1)^2 + (\epsilon_{m\mathbf{k}+\mathbf{q}} - \epsilon_2)^2}{2\sigma_{\text{Gauss}}^2} \right), \quad (6)$$

119 in which  $\sigma_{\text{Gauss}}$  represents the smoothing scale of the Gaussian function. Since BZ inte-  
 120 grations are typically performed by sampling over  $\mathbf{k}$  and  $\mathbf{q}$ -point grids, the expression for  
 121  $g_\nu^2(\epsilon_1, \epsilon_2)$  is given by the weighted sample mean of  $|g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$ . Setting  $V_1 = \sum_{mn\mathbf{k}\mathbf{q}} w_{mn\mathbf{k}\mathbf{q}}$ ,  
 122 where  $w_{mn\mathbf{k}\mathbf{q}}$  is the weight of each sample, including both the degeneracy of the sample point  
 123 in the Brillouin zone and the Gaussian factor shown in Eq. (6), we get

$$g_\nu^2(\epsilon_1, \epsilon_2) = \frac{1}{V_1} \sum_{mn\mathbf{k}\mathbf{q}} w_{mn\mathbf{k}\mathbf{q}} |g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2. \quad (7)$$

124 Since the phonon calculations typically dominate the computational cost during the en-  
 125 tire calculation process during the evaluation of electron-phonon coupling matrix, the use  
 126 of a coarser  $\mathbf{q}$ -point grid directly leads to an almost proportional reduction of the overall  
 127 computational cost. It was reported in<sup>11</sup> that the el-ph calculation, i.e., the phonon calcu-  
 128 lation, took about 100 core-hours on  $4 \times 4 \times 4$   $\mathbf{q}$ -point grids and 4600 core-hours on  $8 \times 8 \times 8$   
 129 grids for a single HH compound. It was also reported that, for a given chemical potential  
 130 of electrons and temperature, the CRT and EPA calculations took about 0.15 core-hours  
 131 each when using  $8 \times 8 \times 8$   $\mathbf{q}$ -point grids for phonon calculations and  $48 \times 48 \times 48$   $\mathbf{k}$ -point grids  
 132 for the band structure calculations, while a comparable EPW calculation took about 2600  
 133 core-hours when using  $4 \times 4 \times 4$  and  $32 \times 32 \times 32$  grids<sup>11</sup>. On the other hand, it was shown in<sup>12</sup>  
 134 that the EPA-MLS method could allow the use of a much coarser grid, i.e.,  $2 \times 2 \times 2$   $\mathbf{q}$ -point  
 135 grid, for the phonon calculation with an acceptable result for fast-screening purposes.

136 While the use of the EPA-MLS method achieved a reasonable balance between perfor-  
 137 mance and accuracy<sup>12</sup>, there are still problems. First, the method still requires specification  
 138 of an open parameter, the smoothing scale of the smoothing kernel. Although the computed  
 139 results were not very sensitive to this parameter, it is an annoying nuisance. Second, al-  
 140 though the sample variance can be used for a rough estimate of uncertainty in the estimated  
 141 el-ph coupling effects<sup>12</sup>, a rigorous analysis of uncertainty and sensitivity is rather difficult  
 142 within the EPA-MLS method.



143 **B. Gaussian Process Regression of Electron-Phonon Coupling**

144 A more rigorous alternative method to achieve the transformation of Eq. (3) is to model  
 145  $g_\nu^2(\epsilon_1, \epsilon_2)$  as a Gaussian process<sup>13</sup> and to perform regression based on the observed elements of  
 146 the electron-phonon coupling matrix. Then, during Gaussian process regression (GPR), the  
 147 characteristic length-scale of the covariance function of the Gaussian process, which serves  
 148 the same purpose as the smoothing scale in the EPA-MLS method, can be estimated using  
 149 the type II maximum likelihood (ML-II) approximation. At the same time, the analysis of  
 150 uncertainty and sensitivity can become theoretically more straightforward.

151 Formally, a Gaussian process is a collection of random variables, any finite number of  
 152 which have a joint Gaussian distribution<sup>13</sup>. In this paper,  $g_\nu^2$  is modeled as a Gaussian  
 153 process. Thus,

$$g_\nu^2(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')), \quad (8)$$

154 where  $\mathbf{x}$  represents the two-dimensional vector coordinate  $(\epsilon_1, \epsilon_2)$ , and  $m(\mathbf{x})$  and  $k(\mathbf{x}, \mathbf{x}')$  are  
 155 the mean and covariance functions of  $g_\nu^2(\mathbf{x})$ , respectively. We consider 0 as the mean, since  
 156 virtually no prior knowledge is available. Many different covariance functions can be used,  
 157 but a simple square exponential kernel is employed here as the covariance function of the  
 158 choice:

$$k(\mathbf{x}, \mathbf{x}') = \sigma_{\text{SEK}}^2 \exp\left(-\frac{(\epsilon_1 - \epsilon'_1)^2 + (\epsilon_2 - \epsilon'_2)^2}{2\ell_{\text{SEK}}^2}\right). \quad (9)$$

159  $\ell_{\text{SEK}}$  is the correlation length scale of the Gaussian process, essentially playing the same role  
 160 as that of  $\sigma_{\text{Gauss}}$  in the EPA-MLS method.

161 Then, we make noisy observations of  $g_\nu^2$  at various training points in two-dimensional  
 162 energy space. The set of training points is denoted as  $X$ , and the DFPT calculations of the  
 163 values of  $|g_{m\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$  on these training points are considered to be such observations. That  
 164 is,

$$|g_{m\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2 = g_\nu^2(\epsilon_{n\mathbf{k}}, \epsilon_{m\mathbf{k}+\mathbf{q}}) + \epsilon_{\text{noise}}, \quad (10)$$

165 where  $\epsilon_{\text{noise}}$  is additive, independent, identically distributed Gaussian noise with variance  
 166  $\sigma_{\text{noise}}^2$ . The total number of observed  $|g_{m\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$  is  $N_S$ .

167 The objective of Gaussian process regression is to predict the values of  $g_\nu^2$  at  $N_T$  test  
 168 points  $\mathbf{x}_{*,j}$  ( $1 \leq j \leq N_T$ ). The training vector is constructed by combining the noisy  
 169 observations, i.e.,  $\mathbf{y} = \left[ |g_{mn\nu}^{\text{SE}}|_1^2 \ |g_{mn\nu}^{\text{SE}}|_2^2 \ |g_{mn\nu}^{\text{SE}}|_3^2 \ \cdots \ |g_{mn\nu}^{\text{SE}}|_{N_S}^2 \right]^\top$ . We define the test output  
 170 vector as  $\mathbf{f}_* = \left[ g_{\nu*,1}^2 \ g_{\nu*,2}^2 \ g_{\nu*,3}^2 \ \cdots \ g_{\nu*,N_T}^2 \right]^\top$ , where  $g_{\nu*,j}^2$  is the estimated value of  $g_\nu^2(\mathbf{x}_{*,j})$  plus  
 171 an additive noise:

$$g_{\nu*,j}^2 = g_\nu^2(\mathbf{x}_{*,j}) + \epsilon_{\text{noise}}. \quad (11)$$

172 Here,  $\epsilon_{\text{noise}}$  has the same variance  $\sigma_{\text{noise}}^2$  as in Eq. (10). The definition of the test output  
 173 vector in this paper is slightly different from that of typical GPR cases. Typically, the  
 174 test output is specified as the estimated value of the Gaussian process only, excluding any  
 175 additive noise. However, in our case, the Gaussian process, i.e.,  $g_\nu^2$ , is only an approximate  
 176 representation of the physical quantity of interest, i.e.,  $|g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$ . Since  $|g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$   
 177 is the sum of  $g_\nu^2$  and  $\epsilon_{\text{noise}}$  as represented in Eq. (10), it is more appropriate to include  
 178  $\epsilon_{\text{noise}}$  during realization of the random Gaussian process, which must reproduce not  $g_\nu^2$  but  
 179  $|g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$ .

180 According to the prior, the joint distribution of the training vector and the test output  
 181 vector is given as follows:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \right), \quad (12)$$

182 where  $\mathbf{A} = \mathbf{K}_{X,X} + \sigma_{\text{noise}}^2 \mathbf{I}$ ,  $\mathbf{B} = \mathbf{K}_{X_*,X_*} + \sigma_{\text{noise}}^2 \mathbf{I}$ , and  $\mathbf{C} = \mathbf{K}_{X,X_*}$ .  $\mathbf{K}_{X,X_*}$  denotes the  
 183  $N_S \times N_T$  matrix of the covariances evaluated at all pairs of training and test points, and the  
 184 other entries are defined in a similar way.  $\mathbf{I}$  represents an identity matrix of an appropriate  
 185 size.

186 Applying a standard argument for multivariate Gaussian distributions to this distribution<sup>13</sup>,  
 187 we construct the conditional distribution, which provides the key predictive equations for  
 188 Gaussian process regression:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \quad (13)$$

189 where

$$\bar{\mathbf{f}}_* = \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = \mathbf{C}^\top \mathbf{A}^{-1} \mathbf{y}, \quad (14)$$

190 and

$$\text{cov}(\mathbf{f}_*) = \mathbf{B} - \mathbf{C}^\top \mathbf{A}^{-1} \mathbf{C}. \quad (15)$$

191 The log marginal likelihood is given as follows:

$$\log p(\mathbf{y} | X) = -\frac{1}{2} \mathbf{y}^\top \mathbf{A}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{A}| - \frac{N_S}{2} \log 2. \quad (16)$$

192 To complete the specification of the model, we need to determine hyperparameters. There  
 193 are three hyperparameters in the model:  $\ell_{\text{SEK}}$ ,  $\sigma_{\text{SEK}}$ , and  $\sigma_{\text{noise}}$ . One of the most widely used  
 194 methods for identification of hyperparameters is the type II maximum likelihood (ML-II)  
 195 approximation, in which the marginal likelihood of the available observations, i.e., Eq. (16),  
 196 under the model is maximized with respect to the hyperparameters<sup>13</sup>. In this fashion, all  
 197 the hyperparameters in the model, i.e.,  $\ell_{\text{SEK}}$ ,  $\sigma_{\text{SEK}}$ , and  $\sigma_{\text{noise}}$ , are estimated.

198 The actual GPR procedure is performed using the Gaussian Processes for Machine Learn-  
 199 ing (GPML) Toolbox<sup>18</sup>. To reduce the computational cost, we use the KISS-GP (Kernel  
 200 Interpolation for Scalable Structured Gaussian Processes) method<sup>19</sup>, in which evaluation of  
 201 the covariance function is replaced with interpolation from a well-defined grid. The cal-  
 202 culation on the grid can exploit its underlying Kronecker-Toeplitz structure to boost the  
 203 calculation speed, which renders the entire method practically feasible. The current imple-  
 204 mentation of the KISS-GP method requires us to use two separate correlation length scales,  
 205 i.e., one for  $\epsilon_1$  and the other for  $\epsilon_2$ . To recover one single length scale, the mean of these  
 206 two lengths is calculated as  $\ell_{\text{SEK}}$  after application of the ML-II approximation, which is the  
 207 method that we use in this study.

### 208 C. Uncertainty Quantification of Thermoelectric Properties

209 One of the main advantages of the proposed GPR procedure is that it provides information  
 210 on uncertainty. For example, the variance from Eq. (15) can be used as an indicator of the  
 211 confidence interval for the estimated values of the el-ph coupling matrix elements at the test  
 212 points. However, our main interest is not to quantify the uncertainty in the el-ph coupling

213 matrix elements but to quantify the uncertainty in the thermoelectric properties of the  
 214 material itself, i.e.,  $S$  and  $\sigma$ . For this purpose, we propose a simple Monte Carlo approach.  
 215 Statistically,  $\mathbf{f}_*$  follows a multivariate Gaussian distribution, whose mean and covariance are  
 216 given by Eqs. (14-15). Thus, we can randomly realize  $\mathbf{f}_*$  using the multivariate Gaussian  
 217 statistic. Since  $\mathbf{f}_*$  represents the estimated values of the el-ph coupling matrix elements,  
 218 the semiclassical Boltzmann transport equation can be solved for each realization of  $\mathbf{f}_*$  to  
 219 create a sample of  $S$  and  $\sigma$ . By repeating this realization, one can create a large-sized set  
 220 of realized samples, which will be further diagnosed to obtain the statistics of  $S$  and  $\sigma$ .

221 To reduce the computational cost further, we employ the following approximation:

$$g_\nu^2(\epsilon_1, \epsilon_2) \approx g_\nu^2\left(\frac{\epsilon_1 + \epsilon_2}{2}, \frac{\epsilon_1 + \epsilon_2}{2}\right). \quad (17)$$

222 This is a valid approximation, since we only need  $g_\nu^2(\epsilon, \epsilon \pm \bar{\omega}_\nu)$  for evaluating Eq. (5). The  
 223 values of  $\bar{\omega}_\nu$  are typically smaller than 0.1 eV, while the values of  $\ell_{\text{SEK}}$  are about 1 eV. Since  
 224  $\bar{\omega}_\nu \ll \ell_{\text{SEK}}$ ,

$$g_\nu^2(\epsilon, \epsilon \pm \bar{\omega}_\nu) \approx g_\nu^2\left(\epsilon \pm \frac{\bar{\omega}_\nu}{2}, \epsilon \pm \frac{\bar{\omega}_\nu}{2}\right), \quad (18)$$

225 because  $g_\nu^2$  will not vary much within the smoothing length scale of the Gaussian process,  
 226 i.e.,  $\ell_{\text{SEK}}$ . The approximation allows us to use test points on the diagonal line ( $\epsilon_1 = \epsilon_2$ ) only,  
 227 which can be later extrapolated onto the two-dimensional energy space, using Eq. (17).

228 The overall procedure of uncertainty quantification within the electron-phonon averaged  
 229 via Gaussian process regression (EPA-GPR) method can be summarized as follows.

- 230 1. Generate the training data set, i.e.,  $g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})$  on a coarse  $\mathbf{q}$  mesh, from the DFPT  
 231 calculations. For this purpose, we use the QUANTUM ESPRESSO package<sup>20</sup>.
- 232 2. Perform the GPR procedure using the training data set. Fix the hyperparameters, i.e.,  
 233  $\ell_{\text{SEK}}$ ,  $\sigma_{\text{SEK}}$ , and  $\sigma_{\text{noise}}$  by applying the ML-II approximation. All the operations in the  
 234 GPR procedure are performed using the KISS-GP method in the GPML Toolbox<sup>18</sup>.
- 235 3. Construct  $\bar{\mathbf{f}}_*$  and  $\text{cov}(\mathbf{f}_*)$  for the test points on the diagonal line ( $\epsilon_1 = \epsilon_2$ ) using  
 236 Eqs. (14-15).
- 237 4. Randomly create  $N_R$  realized samples on the diagonal test points, using `mvnrnd`, which  
 238 is a MATLAB function for random realization of the multivariate Gaussian statistic<sup>21</sup>.

239 Extrapolate onto the two-dimensional energy space using Eq. (17). If an unrealistic  
240 negative value of  $g_\nu^2$  does occur, we put zero instead.

241 5. For each realized sample, construct the input files for a run of BoltzTraP<sup>6</sup>. BoltzTraP  
242 is a standard program for solving the semiclassical Boltzmann equation for inorganic  
243 semiconductors. The version used in this work is slightly modified from the original  
244 BoltzTraP program to incorporate an energy-dependent relaxation time.

245 6. Run BoltzTraP for  $N_R$  sample cases using the input files. Each BoltzTraP run is  
246 independent from the others. A large number of runs can be carried out in a parallel  
247 fashion within a relatively short time if enough computing power is available.

248 7. Statistically analyze the results of the BoltzTraP runs to quantify the uncertainty in  
249  $S$  and  $\sigma$ .

### 250 III. NUMERICAL RESULTS

251 In this section, electronic transport properties for a TE material from the family of HH  
252 compounds, the  $p$ -type HfCoSb<sup>22-24</sup>, are estimated to demonstrate the procedure explained  
253 in Section II. The HH compound has a MgAgAs structure type, whose space group is  
254  $F\bar{4}3m$ <sup>25,26</sup>. After structural relaxation, the lattice parameter of the conventional cubic unit  
255 cell of the HfCoSb compound has a value around 6.0471 Å. The carrier concentration is  
256 fixed first at the value obtained from a Hall measurement at room temperature:  $p = 0.06$   
257 per formula unit ( $1.1 \times 10^{21} \text{ cm}^{-3}$ ) for Hf<sub>0.5</sub>Zr<sub>0.5</sub>CoSb<sub>0.8</sub>Sn<sub>0.2</sub><sup>22</sup>.

258 The electron energy relaxation times and the electronic transport coefficients are calcu-  
259 lated with the original EPA, EPA-MLS, and EPA-GPR methods. DFT and DFPT calcu-  
260 lations are performed using the generalized gradient approximation in the PBE form<sup>27</sup> for  
261 exchange-correlation functional, ultrasoft pseudopotentials<sup>28,29</sup> for core-valence interaction  
262 and a plane wave basis set with 80 and 700 Ry kinetic energy cutoffs for wavefunctions and  
263 charge density. A uniform  $8 \times 8 \times 8$   $\Gamma$ -centered  $\mathbf{k}$ -point grid is used for self-consistent cal-  
264 culation of charge density, and  $48 \times 48 \times 48$  grids are used for band structure and transport  
265 calculations.

266 For the EPA method, a uniform  $8 \times 8 \times 8$   $\Gamma$ -centered  $\mathbf{q}$ -point grid is used for sampling  
267  $|g_{mn\nu}^{\text{SE}}(\mathbf{k}, \mathbf{q})|^2$  by direct el-ph calculations, which was the resolution used in a previous screen-

ing study<sup>11</sup>. For the EPA-MLS method, the same  $8 \times 8 \times 8$   $\Gamma$ -centered  $\mathbf{q}$ -point grid and a uniform  $2 \times 2 \times 2$   $\Gamma$ -centered  $\mathbf{q}$ -point grid are employed. For the EPA-GPR method, only the uniform  $2 \times 2 \times 2$   $\Gamma$ -centered  $\mathbf{q}$ -point grid is employed. Averaging in the EPA calculation is performed over the bins with  $\delta_{\text{Bin}} = 0.2$  eV—the smallest bin size at which all cells in the energy grid are filled with  $\mathbf{k}$ -points. For the EPA-MLS method,  $\sigma_{\text{Gauss}} = 0.2$  eV. The hyperparameters in the EPA-GPR method are all identified within the ML-II approximation as described in Section II. To quantify the uncertainty in the EPA-GPR method, one thousand realized samples of  $g_{\nu}^2$  were generated and statistically analyzed.

First, the identified hyperparameters are presented in Table I. The values of  $\ell_{\text{SEK}}$  identified from the process are slightly less than 1 eV. Hence, the approximation of Eq. (17) can be considered valid.  $\sigma_{\text{SEK}}$  and  $\sigma_{\text{noise}}$  exhibit nontrivial sizes, suggesting that the averaging process of Eq. (3) involves significant uncertainty. The Seebeck coefficient and the electrical conductivity, computed from the original EPA, EPA-MLS, and EPA-GPR methods, are shown in Figure 1. We also plot the experimental data at similar conditions<sup>22–24</sup>. As mentioned earlier, the Hall measurement in<sup>22</sup> reported that the carrier concentration was  $1.1 \times 10^{21} \text{ cm}^{-3}$  for  $\text{Hf}_{0.5}\text{Zr}_{0.5}\text{CoSb}_{0.8}\text{Sn}_{0.2}$ , which corresponds to  $p = 0.06$  per formula unit. The doping concentration of the main dopant, i.e., Sn, for the nanostructured sample of  $\text{Hf}_{0.8}\text{Zr}_{0.2}\text{CoSb}_{0.8}\text{Sn}_{0.2}$  in<sup>23,24</sup> was essentially the same to that for  $\text{Hf}_{0.5}\text{Zr}_{0.5}\text{CoSb}_{0.8}\text{Sn}_{0.2}$  in<sup>22</sup>, and hence it is expected that the nanostructured sample would have a similar carrier concentration, allowing us to compare our numerical results against those experimental data. It is clear that the EPA-GPR method shows good agreement to the other types of EPA method. In particular, the Seebeck coefficient is affected very little by the choice of a particular type of EPA method. On the other hand, the electrical conductivity shows slightly greater sensitivity to the choice. Also, our predicted values clearly exhibit reasonable correspondence to the experimental observations. Although the correspondence is not perfect, our predicted values maintain a level of accuracy that can be used in fast-screening applications.

Statistical properties of three different thermoelectric properties, i.e., the Seebeck coefficient ( $S$ ), the electrical conductivity ( $\sigma$ ), and the power factor ( $\text{PF} = S^2\sigma$ ), at two different temperatures (300K and 700K) are summarized in Table II. The reference value ( $a_{\text{ref}}$ ) obtained from  $\bar{\mathbf{f}}_*$  of Eq. (14), the mean ( $\bar{a}$ ) and standard deviation ( $S_a$ ) of the realized samples, and the ratio between the mean and the standard deviation ( $S_a/\bar{a}$ ) are presented. There exists discrepancy between the sample mean and the reference value. The discrepancy may be

$\nu$	$\ell_{\text{SEK}}$ (eV)	$\sigma_{\text{SEK}}$ (eV <sup>2</sup> )	$\sigma_{\text{noise}}$ (eV <sup>2</sup> )
1	0.815	$1.47 \times 10^{-4}$	$2.11 \times 10^{-4}$
2	0.655	$1.07 \times 10^{-4}$	$1.11 \times 10^{-4}$
3	0.668	$3.08 \times 10^{-4}$	$3.00 \times 10^{-4}$
4	0.469	$2.39 \times 10^{-4}$	$2.22 \times 10^{-4}$
5	0.470	$2.38 \times 10^{-4}$	$2.22 \times 10^{-4}$
6	0.578	$3.76 \times 10^{-4}$	$4.80 \times 10^{-4}$
7	0.599	$2.20 \times 10^{-4}$	$2.06 \times 10^{-4}$
8	0.594	$2.18 \times 10^{-4}$	$2.06 \times 10^{-4}$
9	0.813	$1.52 \times 10^{-4}$	$2.73 \times 10^{-4}$

TABLE I. The hyperparameters, i.e.,  $\ell_{\text{SEK}}$ ,  $\sigma_{\text{SEK}}$ , and  $\sigma_{\text{noise}}$ , for the valence bands of HfCoSb, identified within the ML-II approximation.  $\nu$  is the index of the corresponding phonon branch.

300 attributed to two factors. One obvious reason that can be considered is the limited sample  
301 size, although this is not the most decisive factor in this case. Rather, the central reason  
302 for the discrepancy is that the statistical distributions of these thermoelectric properties are  
303 not normal, which will be discussed in more detail later. As shown in Figure 1, the Seebeck  
304 coefficients exhibit relatively little dependency on the changes in  $g_{\nu}^2$  values. Similarly, the  
305 Seebeck coefficients show small standard deviations in Table II, which are only 3-8% of the  
306 corresponding mean values. On the other hand, the electrical conductivity and the power  
307 factor exhibit much larger standard deviations, amounting to about 20% of the correspond-  
308 ing mean values. This is probably a natural behavior, since the electrical conductivity is  
309 directly proportional to the relaxation time, which is directly affected by the uncertainty in  
310  $g_{\nu}^2$ . The power factor is again proportional to the electrical conductivity, and hence experi-  
311 ences a similar level of uncertainty. Overall, the result clearly indicates that we can place  
312 more confidence in our predicted values of the Seebeck coefficient than in those of the other  
313 properties.

314 Figure 2 shows the histograms of the thermoelectric properties at two different temper-  
315 atures (300K and 700K). We additionally present the histograms of the resistivity ( $1/\sigma$ ),  
316 which is the inverse of the electrical conductivity. As previously mentioned, the statistical  
317 distributions of the Seebeck coefficient, the electrical conductivity, and the power factor are

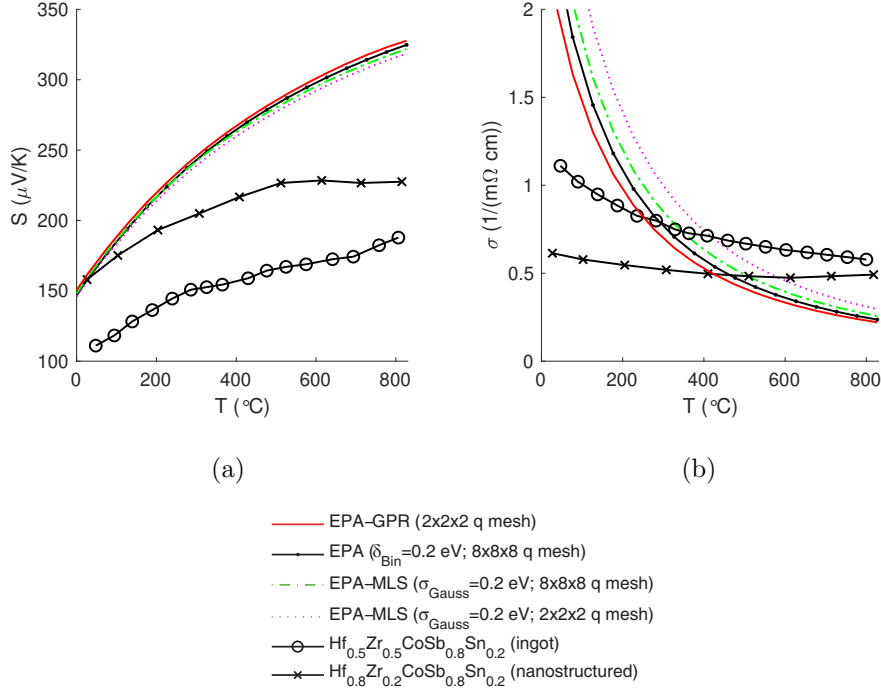


FIG. 1. (Color online.) The Seebeck coefficient  $S$  and the electrical conductivity  $\sigma$  for  $p$ -type HfCoSb as a function of temperature  $T$  calculated with the EPA method and the EPA-MLS method. Consult the legend for the condition represented by each curve. Calculations are performed at the carrier concentration  $p = 0.06$  per formula unit. The open circles and the crosses show the experimental data for the ingot sample of  $p$ -type  $\text{Hf}_{0.5}\text{Zr}_{0.5}\text{CoSb}_{0.8}\text{Sn}_{0.2}$ <sup>22</sup> and the experimental data for the nanostructured sample of  $p$ -type  $\text{Hf}_{0.8}\text{Zr}_{0.2}\text{CoSb}_{0.8}\text{Sn}_{0.2}$ <sup>23,24</sup>, respectively.

318 not normal, as clearly seen in Figures 2 (a), (b), and (d). We have also quantitatively tested  
 319 the normality of these distributions by applying the Jarque-Bera test<sup>30</sup> to each set of realized  
 320 samples. The Jarque-Bera test checks the null hypothesis that each data set comes from  
 321 a normal distribution with an unknown mean and variance. The  $p$ -value of the test is the  
 322 probability of observing a test statistic as extreme as, or more extreme than, the observed  
 323 sample under the null hypothesis. The  $p$ -values of our test for the Seebeck coefficient, the  
 324 electrical conductivity, and the power factor at two different temperatures had values much  
 325 less than 1%, clearly rejecting the null hypothesis for these thermoelectric properties. On the  
 326 other hand, we visually recognize that the resistivity in Figure 2 (c) exhibits distributions  
 327 very close to normal. Indeed, the  $p$ -value of the Jarque-Bera test for the resistivity is 28.2%  
 328 at 300K and 15.8% at 700K. The resistivity is the inverse of the electrical conductivity and



$a$	$S$ ( $\mu\text{V}/\text{K}$ )		$\sigma$ ( $1/(\text{m}\Omega \text{ cm})$ )		$\text{PF} = S^2\sigma$ ( $\text{mW}/(\text{m K}^2)$ )	
	300K	700K	300K	700K	300K	700K
$a_{\text{ref}}$	161.7	272.3	2.112	0.4916	5.520	3.646
$\bar{a}$	162.4	272.5	2.115	0.4898	5.537	3.625
$S_a$	12.95	10.45	0.4289	$6.910 \times 10^{-2}$	1.092	0.4612
$S_a/\bar{a}$	$7.977 \times 10^{-2}$	$3.835 \times 10^{-2}$	0.2028	0.1411	0.1973	0.1272

TABLE II. Statistical properties of the Gaussian process  $g_\nu^2$  and its realized samples at two different temperatures (300K and 700K). The statistical properties of three different thermoelectric properties, i.e., the Seebeck coefficient ( $S$ ), the electrical conductivity ( $\sigma$ ), and the power factor ( $\text{PF} = S^2\sigma$ ), are provided. For each thermoelectric property ( $a$ ), the reference value ( $a_{\text{ref}}$ ) obtained from  $\bar{\mathbf{f}}_*$  of Eq. (14), the mean ( $\bar{a}$ ) and standard deviation ( $S_a$ ) of the realized samples, and the ratio between the mean and the standard deviation ( $S_a/\bar{a}$ ) are presented. All the values except those of  $S_a/\bar{a}$ , which are dimensionless, are reported in the unit corresponding to each quantity.

329 hence can be considered to be roughly proportional to the scattering rate, i.e.,  $\tau^{-1}$ , which is  
330 in turn proportional to the value of  $g_\nu^2$ . Since  $g_\nu^2$  follows a multivariate Gaussian statistic in  
331 our numerical model, it is more natural for the resistivity to follow a normal statistic, which  
332 is indeed the case in our numerical test.

333 One of the most important statistical properties that are relevant to the fast-screening  
334 procedure of thermoelectric materials is the confidence interval of the estimation. The  
335 sample statistics can be utilized to provide such information. In Figure 3, the 5%, 50%,  
336 and 95% percentiles for the thermoelectric properties of interest are provided, along with  
337 the reference curve directly computed from  $\bar{\mathbf{f}}_*$  of Eq. (14). The colored range between the  
338 5% and 95% percentiles indicates a confidence interval of 90%. As discussed already, the  
339 confidence intervals of the electrical conductivity and the power factor are relatively large.  
340 For example, at 300K, the 95% percentile value, i.e.,  $7.469 \text{ mW}/(\text{m K}^2)$  and the 5% percentile  
341 value, i.e.,  $4.113 \text{ mW}/(\text{m K}^2)$ , of the power factor deviate by 35% and 25%, respectively,  
342 from the reference value, i.e.,  $5.520 \text{ mW}/(\text{m K}^2)$ . Clearly, the range is still acceptable for  
343 fast-screening applications, but one must remain cautious not to place blind faith on the  
344 values from computational estimations.

345 So far, we have considered the uncertainty in the electronic transport properties of the

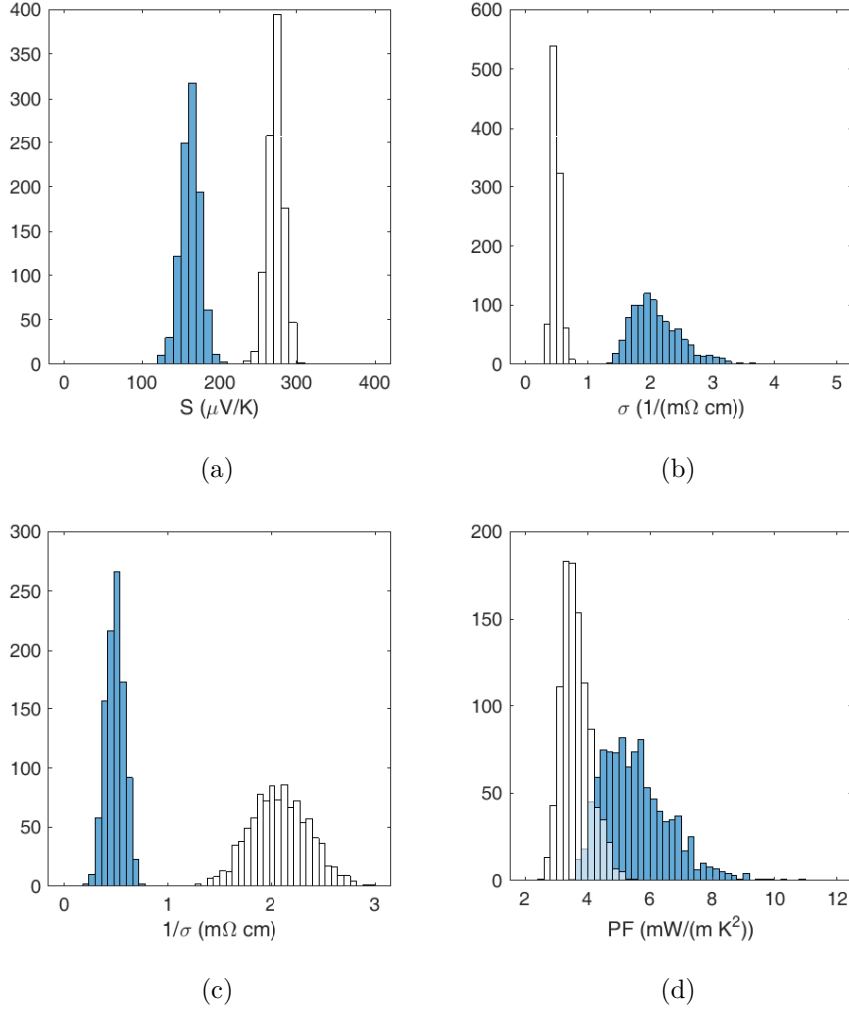


FIG. 2. (Color online.) The histograms of thermoelectric properties at two different temperatures (300K and 700K): (a) the Seebeck coefficient ( $S$ ); (b) the electrical conductivity ( $\sigma$ ); (c) the resistivity ( $1/\sigma$ ); and (d) the power factor (PF). The blue-faced bars and the semi-transparent bars represent data at 300K and at 700K, respectively.

346 material at a fixed carrier concentration, i.e.,  $p = 0.06$  per formula unit, but the electronic  
 347 transport coefficients of materials depend strongly on the carrier concentration<sup>11</sup>. One of  
 348 the most important objectives of computational prediction is to suggest an optimal carrier  
 349 concentration for a given composition. The values of  $g_v^2$  bear certain uncertainty, and hence  
 350 the predicted optimal carrier concentration will also involve uncertainty. In a previous  
 351 study<sup>11</sup>, it was reported that the values of the optimal carrier concentration maximizing  $ZT$   
 352 were only about 10% different in average from the values of carrier concentration maximizing  
 353 PF. Therefore, the PF values from several realized samples at 700K are plotted versus the

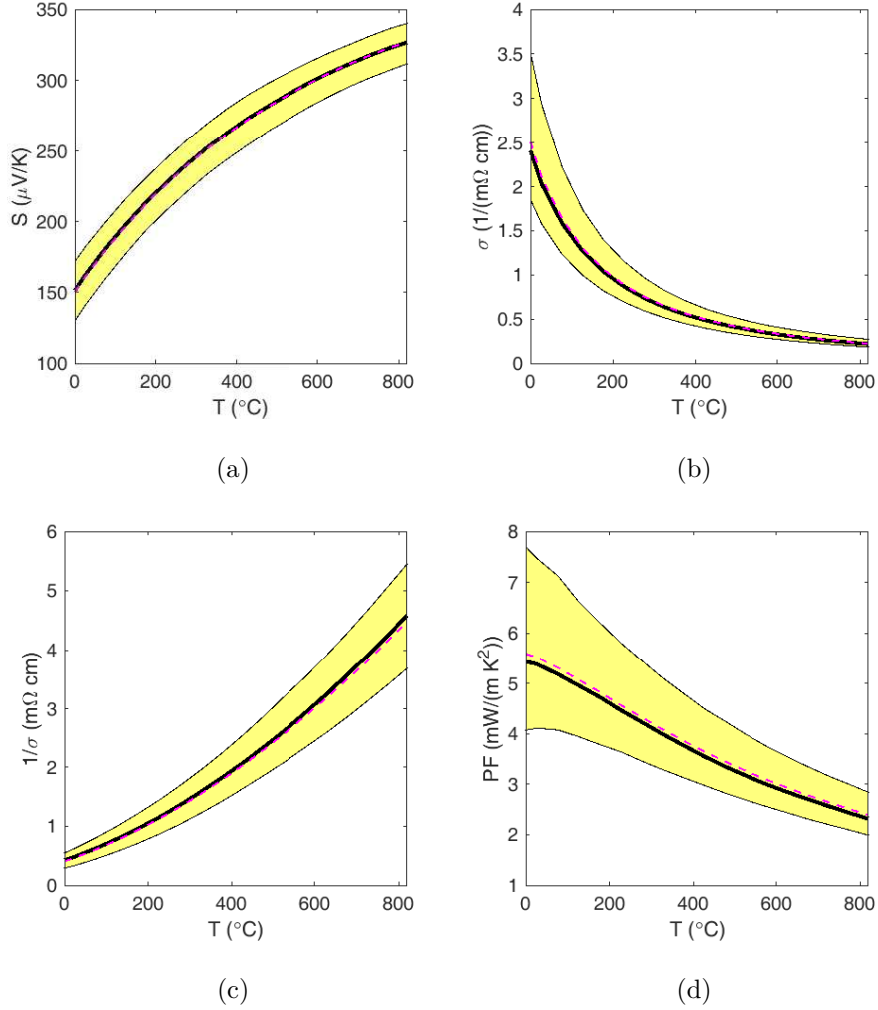


FIG. 3. (Color online.) The 5%, 50%, and 95% percentile curves of thermoelectric properties plotted versus  $T$ , presented along with the reference curve directly computed from  $\bar{\mathbf{f}}_*$  of Eq. (14): (a) the Seebeck coefficient ( $S$ ); (b) the electrical conductivity ( $\sigma$ ); (c) the resistivity ( $1/\sigma$ ); and (d) the power factor (PF). The thick solid curves represent the 50% percentiles, while the thin solid curves represent the 5% and 95% percentiles. The dashed curves represent the corresponding reference curves. The colored range between the 5% and 95% percentiles represents the confidence interval of 90%.

354 hole concentration  $p$  in Figure 4 (a). The result shows that there exist large variations in the  
 355 maximum values of PF and the optimal values of  $p$  associated, among the chosen samples.  
 356 Such uncertainty in the values of the optimal carrier concentration and that in the associated  
 357 PF values should be carefully quantified, and our method can be utilized for serving such  
 358 a purpose. In Figure 4 (b), a scatter plot showing the maximum value of PF ( $\text{PF}_{\max}$ ) and

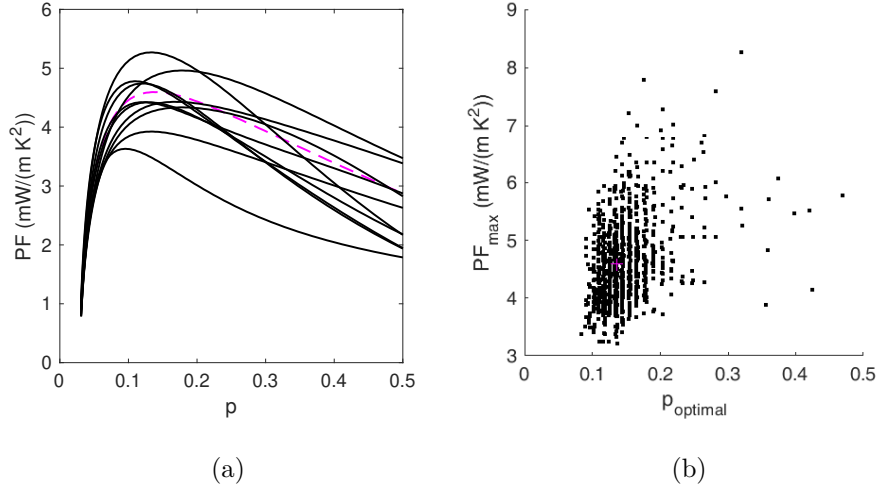


FIG. 4. (Color online.) Variations in the maximum power factor ( $PF_{\max}$ ) and the optimal value of hole concentration ( $p_{\text{optimal}}$ ) at 700K: (a) PF versus  $p$  for 10 realized samples (solid black) and the reference (dashed magenta) from  $\bar{f}_*$  of Eq. (14); and (b) a scatter plot showing  $p_{\text{optimal}}$  and the associated maximum PF ( $PF_{\max}$ ) for each realized sample (black dots) and for the reference (a magenta cross). In both plots, data from too small hole concentrations ( $p < 0.03$ ) have been excluded, since the character of the material changes from a  $p$ -type to an  $n$ -type there.

359 associated optimal  $p$  value ( $p_{\text{optimal}}$ ) of each realized samples. A fairly large variation is  
 360 observed, but there is an underlying trend. The samples are scattered around the reference  
 361 point, where  $p_{\text{optimal}} = 0.136$  per formula unit, with an area with high density in the range  
 362 of  $0.1 < p_{\text{optimal}} < 0.2$ . The value of  $PF_{\max}$  at this new reference point is  $4.595 \text{ mW}/(\text{m K}^2)$ ,  
 363 which is higher than that reported for  $p = 0.06$  (Table II).

364 In Figure 5 (a), we present the histograms of the optimal value of hole concentration  
 365 ( $p_{\text{optimal}}$ ) at 300K and 700K. At 300K, the most probable value of  $p_{\text{optimal}}$  turns out to be  
 366 around 0.06, which was the value employed for our study mentioned above, i.e., Tables I-II  
 367 and Figures 1-3. On the other hand, the most probable value of  $p_{\text{optimal}}$  at 700K occurs  
 368 in between 0.135 and 0.145, which is larger than 0.06. Increase in temperature activates  
 369 carriers in a wider energy range, and a too low  $p$  value may result in a conflict between  
 370 two different charge carriers, i.e., electrons and holes, resulting in a very low value or even  
 371 a sign reversal of the Seebeck coefficient at high temperature. Thus, it is natural to find  
 372 that the most probable  $p_{\text{optimal}}$  value at 700K is larger than that at 300K. The range of  
 373 the most probable  $p_{\text{optimal}}$  at 700K, observed from the histogram, includes  $p_{\text{optimal}}$  of the

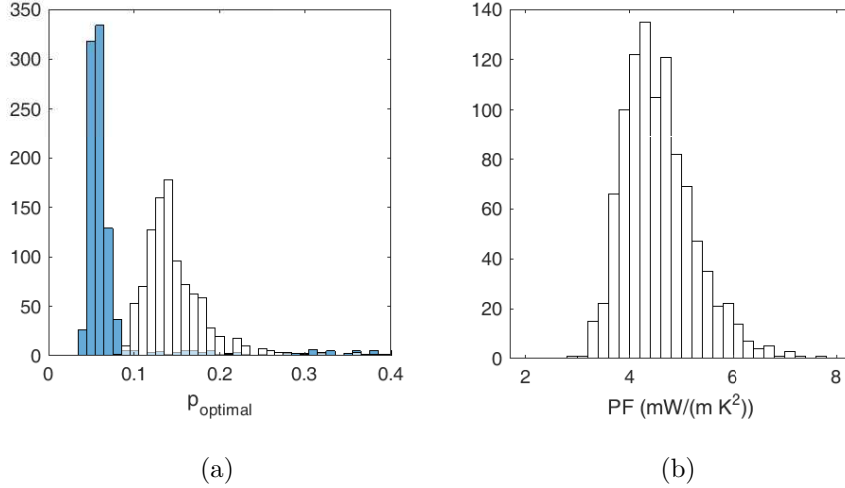


FIG. 5. (Color online.) The histograms of the optimal value of hole concentration ( $p_{\text{optimal}}$ ) and the PF values associated: (a)  $p_{\text{optimal}}$  at 300K and 700K; and (b) PF at 700K with  $p = 0.14$  per formula unit. The blue-faced bars and the semi-transparent bars represent data at 300K and at 700K, respectively, if there are histograms for both temperatures.

374 reference point, which is 0.136, observed in Figure 4 (b). A large number of realized samples  
 375 exhibit maximum values in PF within the range of  $0.1 < p_{\text{optimal}} < 0.2$ , which confirms our  
 376 observation made from Figure 4 (b).

377 In Figure 5 (b), we additionally present the histogram of the PF values with  $p = 0.14$  per  
 378 formula unit, which is the value lying at the center of the range exhibiting the most probable  
 379  $p_{\text{optimal}}$  value ( $0.135 < p < 0.145$ ), at 700K. This plot shows the distribution of possible PF  
 380 values, which are expected to be achieved if the carrier concentration is optimized at the  
 381 computationally predicted  $p$  value during the synthesis process. The expected PF values  
 382 are lying mostly in between  $3 \text{ mW}/(\text{m K}^2)$  and  $6 \text{ mW}/(\text{m K}^2)$ , exhibiting a significant  
 383 variation. Such information on the potential uncertainty in the predicted values can be  
 384 useful for assessing the feasibility of a candidate material.

#### 385 IV. CONCLUSIONS

386 We presented the EPA-GPR method where the el-ph coupling matrix is represented as a  
 387 function of two energies, which is in turn modeled as a Gaussian process. Unlike the other  
 388 EPA-related methods, the EPA-GPR method is a truly parameter-free estimation method,

389 since all the hyperparameters in the model can be unambiguously determined within the  
390 ML-II approximation. On top of that, the use of a Gaussian process allows us to quantify  
391 the uncertainty in the estimated thermoelectric properties.

392 To demonstrate the effectiveness of the EPA-GPR method, we applied it to a  $p$ -type half-  
393 Heusler compound, i.e., HfCoSb. Our numerical results clearly exhibit the advantages of the  
394 method. In particular, we note that the estimated power factor can vary up to about 35%  
395 at room temperature within a confidence level of 90%, which is acceptable for fast-screening  
396 applications but still requires a certain level of caution in fast-screening applications. Overall,  
397 the information on the potential uncertainty in computational prediction can be valuable in  
398 future decision-making processes of the research and development of new TE materials.

## 399 ACKNOWLEDGMENTS

400 The first, second, and third authors were mainly supported by Solvay SA, an advanced  
401 materials and specialty chemicals company, through an Ewha-Solvay collaboration agree-  
402 ment, during the course of the study reported in this article. The fourth and fifth authors  
403 were partly supported by the U.S. Department of Energy under Award DE-EE0004840.

---

404 \* dhwee@ewha.ac.kr; corresponding author.

405 <sup>1</sup> S. Twaha, J. Zhu, Y. Yan, and B. Li, *Renewable and Sustainable Energy Reviews* **65**, 698  
406 (2016).

407 <sup>2</sup> D. A. Broido, M. Malorny, G. Birner, N. Mingo, and D. A. Stewart, *Applied Physics Letters*  
408 **91**, 231922 (2007).

409 <sup>3</sup> J. Carrete, W. Li, N. Mingo, S. Wang, and S. Curtarolo, *Phys. Rev. X* **4**, 011019 (2014).

410 <sup>4</sup> A. van Roekeghem, J. Carrete, C. Oses, S. Curtarolo, and N. Mingo, *Phys. Rev. X* **6**, 041061  
411 (2016).

412 <sup>5</sup> J. M. Ziman, *Principles of the Theory of Solids*, 2nd ed. (Cambridge University Press, 1972).

413 <sup>6</sup> G. K. Madsen and D. J. Singh, *Computer Physics Communications* **175**, 67 (2006).

414 <sup>7</sup> M. Bernardi, D. Vigil-Fowler, J. Lischner, J. B. Neaton, and S. G. Louie, *Phys. Rev. Lett.* **112**,  
415 257402 (2014).

- 416 <sup>8</sup> J. Yan, P. Gorai, B. Ortiz, S. Miller, S. A. Barnett, T. Mason, V. Stevanovic, and E. S. Toberer,  
417 Energy Environ. Sci. **8**, 983 (2015).
- 418 <sup>9</sup> A. J. Hong, L. Li, R. He, J. J. Gong, Z. B. Yan, K. F. Wang, J.-M. Liu, and Z. F. Ren, Scientific  
419 Reports **6**, 22778 (2016).
- 420 <sup>10</sup> J. Noffsinger, F. Giustino, B. D. Malone, C.-H. Park, S. G. Louie, and M. L. Cohen, Computer  
421 Physics Communications **181**, 2140 (2010).
- 422 <sup>11</sup> G. Samsonidze and B. Kozinsky, Advanced Energy Materials **8**, 1800246 (2018).
- 423 <sup>12</sup> S. Bang, J. Kim, D. Wee, G. Samsonidze, and B. Kozinsky, Materials Today Physics **6**, 22  
424 (2018).
- 425 <sup>13</sup> C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning* (The MIT  
426 Press, 2006).
- 427 <sup>14</sup> F. Giustino, M. L. Cohen, and S. G. Louie, Phys. Rev. B **76**, 165108 (2007).
- 428 <sup>15</sup> W. Li, Phys. Rev. B **92**, 075405 (2015).
- 429 <sup>16</sup> S. Baroni, S. de Gironcoli, A. Dal Corso, and P. Giannozzi, Rev. Mod. Phys. **73**, 515 (2001).
- 430 <sup>17</sup> D. Levin, Math. Comput. **67**, 1517 (1998).
- 431 <sup>18</sup> C. E. Rasmussen and H. Nickisch, J. Mach. Learn. Res. **11**, 3011 (2010).
- 432 <sup>19</sup> A. G. Wilson and H. Nickisch, in *Proceedings of the 32nd International Conference on Interna-*  
433 *tional Conference on Machine Learning - Volume 37, ICML'15 (JMLR.org, 2015)* pp. 1775–1784.
- 434 <sup>20</sup> P. Giannozzi, S. Baroni, N. Bonini, M. Calandra, R. Car, C. Cavazzoni, D. Ceresoli, G. L.  
435 Chiarotti, M. Cococcioni, I. Dabo, A. Dal Corso, S. de Gironcoli, S. Fabris, G. Fratesi,  
436 R. Gebauer, U. Gerstmann, C. Gougoussis, A. Kokalj, M. Lazzeri, L. Martin-Samos, N. Marzari,  
437 F. Mauri, R. Mazzarello, S. Paolini, A. Pasquarello, L. Paulatto, C. Sbraccia, S. Scandolo,  
438 G. Sciauzero, A. P. Seitsonen, A. Smogunov, P. Umari, and R. M. Wentzcovitch, Journal of  
439 Physics: Condensed Matter **21**, 395502 (19pp) (2009).
- 440 <sup>21</sup> MATLAB, *Statistics and Machine Learning Toolbox™ User's Guide* (The MathWorks, Inc.,  
441 Natick, Massachusetts, 2018).
- 442 <sup>22</sup> S. R. Culp, J. W. Simonson, S. J. Poon, V. Ponnambalam, J. Edwards, and T. M. Tritt,  
443 Applied Physics Letters **93**, 022105 (2008).
- 444 <sup>23</sup> X. Yan, W. Liu, H. Wang, S. Chen, J. Shiomi, K. Esfarjani, H. Wang, D. Wang, G. Chen, and  
445 Z. Ren, Energy & Environmental Science **5**, 7543 (2012).
- 446 <sup>24</sup> X. Yan, G. Joshi, W. Liu, Y. Lan, H. Wang, S. Lee, J. W. Simonson, S. J. Poon, T. M. Tritt,

- 447 G. Chen, and Z. F. Ren, Nano Letters **11**, 556 (2011).
- 448 <sup>25</sup> A. Page, P. Poudeu, and C. Uher, Journal of Materiomics **2**, 104 (2016).
- 449 <sup>26</sup> F. Casper, T. Graf, S. Chadov, B. Balke, and C. Felser, Semiconductor Science and Technology  
450 **27**, 063001 (2012).
- 451 <sup>27</sup> J. P. Perdew, K. Burke, and M. Ernzerhof, Phys. Rev. Lett. **77**, 3865 (1996).
- 452 <sup>28</sup> D. Vanderbilt, Phys. Rev. B **41**, 7892 (1990).
- 453 <sup>29</sup> A. D. Corso, Computational Materials Science **95**, 337 (2014).
- 454 <sup>30</sup> C. M. Jarque and A. K. Bera, Economics Letters **6**, 255 (1980).