

Elsevier Editorial System(tm) for Computer
Methods in Applied Mechanics and Engineering
Manuscript Draft

Manuscript Number: CMAME-D-18-01024R1

Title: A data-driven framework for sparsity-enhanced surrogates with arbitrary mutually dependent randomness

Article Type: Research Paper

Keywords: data driven; arbitrary randomness; mutual dependence; compressed sensing; sparsity enhancement; uncertainty quantification

Corresponding Author: Dr. Nathan Andrew Baker, Ph.D.

Corresponding Author's Institution: Pacific Northwest National Laboratory

First Author: Huan Lei

Order of Authors: Huan Lei; Jing Li; Peiyuan Gao; Panagiotis Stinis; Nathan Baker

Abstract: The challenge of quantifying uncertainty propagation in real-world systems is rooted in the high-dimensionality of the stochastic input and the frequent lack of explicit knowledge of its probability distribution. Traditional approaches show limitations for such problems, especially when the size of the training data is limited. To address these difficulties, we have developed a general framework of constructing surrogate models on spaces of stochastic input with arbitrary probability measure irrespective of the mutual dependencies between individual components of the random inputs and the analytical form. The present Data-driven Sparsity-enhancing Rotation for Arbitrary Randomness (DSRAR) framework includes a data-driven construction of multivariate polynomial basis for arbitrary mutually dependent probability measures and a sparsity enhancement rotation procedure. This sparsity-enhancing rotation method was initially proposed in our previous work for Gaussian density distributions, which may not be feasible for non-Gaussian distributions due to the loss of orthogonality after the rotation. To remedy such difficulties, we developed a new data-driven approach to construct orthonormal polynomials for arbitrary mutually dependent randomness, ensuring the constructed basis maintains the orthogonality/near-orthogonality with respect to the density of the rotated random vector, where directly applying the regular polynomial chaos including arbitrary polynomial chaos shows limitations due to the assumption of the mutual independence between the components of the random inputs. The developed DSRAR framework leads to accurate recovery, with only limited training data, of a sparse representation of the target functions. The effectiveness of our method is demonstrated in challenging problems such as partial differential equations and realistic molecular systems within high-dimensional conformational spaces where the underlying density is implicitly represented by a large collection of sample data, as well as systems with explicitly given non-Gaussian probabilistic measures.

Research Data Related to this Submission

There are no linked research data sets for this submission. The following reason is given:
Data will be made available on request

Response to referee's report for the manuscript
#CMAME-D-18-01024

December 8, 2018

Manuscript: #CMAME-D-18-01024.

Title: A data-driven framework for sparsity-enhanced surrogates with arbitrary mutually dependent randomness

Previous Title: Data-driven approach of quantifying uncertainty in complex systems with arbitrary randomness.

Authors: H. Lei, J. Li, P. Gao, P. Stinis and N. A. Baker.

Dear Editor:

We have addressed *all* of the points that the referees raised and also made all the recommended changes regarding our manuscript. We provide point-by-point answers below. In the revised manuscript, revisions are marked in blue color.

We thank the referees for their thorough review and constructive comments that helped us to improve the quality of our manuscript. We also thank you for your time and consideration of this paper.

Sincerely,
Huan Lei
Jing Li
Peiyuan Gao
Panos Stinis
Nathan Baker

Reviewer 1

Reviewer Comment: *This work is an extension of the authors' previous work on uncertainty quantification in complex systems from observation data by removing the assumption of Gaussianity. The novelty in this work includes the construction of an orthogonal basis in random space using the Gram-Schmidt procedure and application of the authors' work to deal with non-Gaussian processes. This work is an important and interesting exploration toward real applications using the state-of-the-art uncertainty quantification techniques. However, many places in this highly technical paper are unclear to me in both mathematics and descriptions. I would suggest a major revision.*

1. **Reviewer Comment:** *Page 8. Line 45-46. The claim that “Consequently, the orthogonal polynomials basis ...” is not true. Whether the involving random variables are independent or not, we can always represent the considered process with a tensor product of univariate basis functions.*

Author Response: We refer to the case where each component of ξ can be mutually dependent. We have clarified this and revised the sentence by “Consequently, the orthogonal polynomial basis $\psi_\alpha(\xi)$ cannot be directly constructed as a tensor product of univariate orthonormal basis functions in each component of ξ . ”

2. **Reviewer Comment:** *Page 9. The Gram-Schmidt orthogonalization often becomes unstable when the orthogonality is lost after several steps. What do you do when you have this problem? Also, what are your requirements on your samples to assure the success of the orthogonalization.*

Author Response: The numerical instability can be addressed by employing modified Gram-Schmidt orthogonalization. For the current work, we checked the orthonormality after the basis was constructed in all examples we presented in this manuscript, the normalized Gram-Schmidt orthogonalization worked quite well, and a possible reason is that we implemented with the polynomials basis up to relatively low degree. We add the modified Gram-Schmidt orthogonalization as an alternative if there exists an issue with instability. There is no requirements on the samples to assure the success of the orthogonalization.

3. **Reviewer Comment:** *Page 9. Line 30-31. By “an i.i.d random vector”, do you mean a random vector with i.i.d elements? Similar issues appear in several places in the paper.*

Author Response: Yes. We replaced all instances of “i.i.d. random vector” with “random vector with i.i.d. components” to clarify this point.

4. **Reviewer Comment:** *How (3.6) is derived from (3.4)? I don't see any connection between these two formulas.*

Author Response: We do not derive (3.6) from (3.4). We use the orthonormal conditions and the representations of (3.4) to calculate the coefficients f_β^α .

5. **Reviewer Comment:** *Page 11. Lines 50-51. If an ordered monomial basis is used, will there be a problem of ill-conditioning?*

Author Response: The ordered monomial basis is used for representation of our orthonormal basis. When we construct the measurement matrix for compressive sensing, we only use the constructed orthonormal basis or the near-orthonormal basis. There are no conditioning problems with these bases.

6. **Reviewer Comment:** *Page 13. Equation (3.14). Why is it important to consider the gradient matrix? For moderately nonlinear f and Gaussian ξ , the gradient might be a good*

indicator of the covariance matrix. However, for non-Gaussian ξ , this may not be the case. The gradient can be misleading for a non-Gaussian case.

Author Response: We evaluate the gradient matrix $\mathbf{G} = \mathbb{E} [\nabla f(\xi) \nabla f(\xi)^T]$ to find the sorted variability directions with respect to the QoI. Eigendecomposition $\mathbf{G} = \mathbf{QKQ}^T$, $\mathbf{Q} = [\mathbf{q}_1 \mathbf{q}_2 \cdots \mathbf{q}_d]$, $\mathbf{K} = \text{diag}(k_1, \dots, k_d)$ shows that the QoI f exhibits the largest variability along \mathbf{q}_1 and the least along \mathbf{q}_d with respect to $\rho(\xi)$. We refer to Ref. [5, 2] for details. This is true for nonlinear f and non-Gaussian ξ . As a special case, if f is linear, $k_2 = \dots = k_d \equiv 0$ and QoI only depends on the univariate $\mathbf{q}_1^T \xi$.

We note that the gradient matrix \mathbf{G} depends on the QoI f and does not serve as an approximation of the covariance matrix. Alternatively, we use \mathbf{G} to define a new random vector $\chi = \mathbf{Q}^T \xi$ s.t. the QoI $g(\chi) = f(\mathbf{Q}\chi)$ shows sparse coefficient (e.g., see Fig. 5). Accordingly, the surrogate model $\tilde{g}(\chi)$ constructed by compressed sensing is more accurate than $\tilde{f}(\xi)$. Since f is nonlinear in general, here we does not conduct dimension reduction, i.e., the dimension of χ equals to dimension of ξ .

To construct \tilde{g} , we note that for non-Gaussian ξ , the rotated vector $\chi = \mathbf{Q}^T \xi$ does not retain the same density distribution $\rho(\cdot)$. This is different from the case where both ξ and χ are Gaussian. The orthonormal basis $\psi(\xi)$ does not retain orthonormal with respect to χ . This is one of the major motivations for the current work, where we developed a general framework to construct orthonormal bases $\phi(\chi)$ with respect to arbitrary non-Gaussian density. The present method enables us to retain the orthonormal properties after the rotation and construct surrogate $\tilde{g}(\chi)$ with respect to the new random vector χ .

7. **Reviewer Comment:** *Page 13. Line 40. Why the model is $f(\chi)$ instead of $f \cdot Q(\chi)$?*

Author Response: To answer this question, we revised the document and replaced “where the surrogate model $f(\chi)$ can be represented by a sparser coefficient vector c as compared with $f(\xi)$ ” by “ $f((\mathbf{Q}^T)^{-1}\chi) = f(\mathbf{Q}\chi) = f(\xi)$ can be approximated by an expansion of orthonormal polynomial basis of χ with sparser coefficient vector c as compared to $f(\xi)$ being expanded by orthonormal basis of ξ .”

8. **Reviewer Comment:** *Numerical results. · It is assumed that ϵ is from a Gaussian mixture with its components whose mean, and covariance are randomly chosen. I don't understand why the Legendre polynomial chaos is chosen for comparison. · It seems that none of the experiments can be reproduced. No random generators are reported.*

Author Response: In the present study, we consider systems with arbitrary randomness (e.g., samples generated by Gaussian mixture, linear transformation from i.i.d. uniform variables, trajectory data set from Molecular Dynamics simulation). For such systems, standard polynomial bases such as Legendre polynomials and Hermite polynomials are commonly chosen for surrogate construction. However, as the reviewer suggests (see also the reviewer comment 15), directly applying such bases has limitations. Therefore, we proposed the present approach based on the data-driven construction of orthonormal basis. We compared our results with ones regressed with Legendre chaos and Hermite chaos. For the purposes of reproducibility, we have provided the seed numbers of the random number generators and name of generation function called to generate the Gaussian Mixture data set in Appendix D as suggested (all of the simulations are done in Matlab).

9. **Reviewer Comment:** *Page 16. Lines 44-45. Why the loss of orthonormality leads to worse performance?*

Author Response: There are 2 possible reasons. One is that in approximation theory,

the orthogonal basis expansion gives the best l_2 approximation. Another reason is that in compressive sensing, the orthonormal basis gives better recovery.

10. **Reviewer Comment:** *Page 16. Lines 48-49. What is the condition? Can this condition be given briefly here?*

Author Response: The sufficient number stated in [52] ([30] in previous version) is given in Theorem 2.5 and we revised the sentence to “ M suggested by the sufficient condition (Theorem 2.5) originally given” for clarification.

11. **Reviewer Comment:** *Page 17. Line 37. Why is the error evaluated in l_2 -norm? instead of l_1 norm.*

Author Response: Due to the orthonormality of the basis functions in probability measure, the l_2 error of the surrogate model is equal to the l_2 error of the coefficient vector. Accordingly, we examined the l_2 -error of the coefficient vector to measure the accuracy of the surrogates. Also the l_2 -error is conventional in the community of compressive sensing.

12. **Reviewer Comment:** *Page 19. What is the level 4 sparse grid integration rule? What quadrature rule is used? Are they sufficient to give the desired accuracy? Also, I don't see anything plotted in Figure 3 on the integration of \tilde{f} using a sparse grid.*

Author Response: We used Smolyak's sparse grid quadrature rule of level 4. It is exact for integration of polynomial of degree 7 and is sufficient to calculate numerical l_2 error of surrogate model (polynomial degree 3). Figure 3 presents the relative l_2 -errors of different regressed surrogates. The l_2 -errors, defined by Eq. (4.12), are evaluated by Smolyak's sparse grid quadrature rule of level 4.

13. **Reviewer Comment:** *Page 20. Why is the density function of the form (4.10) chosen? PDFs of this type are extremely difficult because of the singularity on the boundary.*

Author Response: We chose the density function of the form (4.10) as the random variables with this density function have been used a lot for benchmark problems in compressive sensing for UQ. And also as the reviewer suggested, PDFs of this type are extremely difficult, we would like to test our method for such difficult problems.

14. **Reviewer Comment:** *Page 21. Line 43-44. What is “each training sample size number M ”?*

Author Response: For clarification, we added a sentence “This small set $\{\mathbf{z}^{(i)}\}_{i=1}^M$ is usually called the training sample set and M is the training sample size.” in Section 3 Methods 3 lines below the section title.

15. **Reviewer Comment:** *Page 23. Section 4.3.2. Again why Legendre basis (and Hermite basis) is chosen? The sparsity here doesn't mean anything since these bases can be useless at all (no capability of representing the current function).*

Author Response: Both Legendre basis and Hermite basis up to degree P can span polynomial spaces of degree P . This means that whenever the target function can be approximated by polynomials of order P , then both Legendre basis and Hermite basis can be used. However, in the approach of compressive sensing, using both Legendre basis and Hermite basis may not give good representations. Currently, Legendre polynomials and Hermite polynomials are commonly chosen for surrogate construction in UQ studies in both benchmark test and real-world application. However, as the reviewer suggested, direct employing such bases may not be appropriate in practice. We totally agree on this point, which motivates us to develop the present method to handle such problems. Here,

we compare the surrogates constructed by our approach with the ones constructed with Legendre basis and/or Hermite basis to demonstrate effectiveness of our approach.

16. **Reviewer Comment:** *Page 42-43. What do you mean by “ ξ ... represents 99.99% of the variance”?*

Author Response: Since there is no page 42-43 of this manuscript, we assume the reviewer meant the Page 24, line 42-43. This is a standard approach for dimensionality reduction in molecular systems. Here we employ PCA (principal component analysis) to construct the random (conformational) space of the molecule, and after the dimension reduction we kept 99.99% of the total variance of the original data set.

17. **Reviewer Comment:** *Not sure whether arbitrary randomness is appropriate in the title and abstract. Observe that the authors have two examples: one is a randomized Gaussian mixture: the other also uses Gaussian mixtures.*

Author Response: In this manuscript, we implemented our approach for various cases. Some of the data sets are generated from Gaussian mixture distributions (example 4.1, 4.3.1, 4.3.2), some are generated linear transformation from other i.i.d. distributions such as uniform distribution and distribution with PDF (4.6) (example 4.2.1) and PDF (4.10) (example 4.2.2), and some are from MD simulations (example 4.3.3). Another reason we employ the Gaussian mixture model is that a wide range of arbitrary distributed data can be approximated by Gaussian mixture model.

Moreover, there is a major difference between the present study and the previous UQ studies on based gPC (and aPC). The DSRAR framework developed in the present study does not assume the mutual independence between the components of random inputs; and therefore can be more realistic and particularly well-suited for uncertainty quantification for complex systems, where the knowledge of the underlying randomness can be implicit. Accordingly, we have changed the title to “A data-driven framework for sparsity-enhanced surrogates with arbitrary mutually dependent randomness” and emphasized such difference in the abstract and Introduction of the revised manuscript.

18. **Reviewer Comment:** *Some details Here is a partial list I have when I read the paper. The authors should use more caution when proofreading the paper.*

Author Response: We thank the reviewer’s careful reading for helping us to improve the quality of this manuscript. We have addressed the comments listed below and also carefully checked the revised manuscript.

19. **Reviewer Comment:** *Page 8. Line 43. What does “the dimension of ξ are independent ” mean?*

Author Response: We replaced “the dimension of ξ are independent” by “each component of ξ is mutually independent”.

20. **Reviewer Comment:** *Page 12. Line 52 Insert “a” before “smaller basis bound.” This also applies to Line 42 on Page 27.*

Author Response: Fixed.

21. **Reviewer Comment:** *Page 13. Lines 51-52 what does it mean by saying “ ξ are i.i.d random variables”?*

Author Response: We rewrite this sentence with “In particular, if ξ is a random vector with i.i.d. Gaussian components, χ is also a random vector with i.i.d. Gaussian components” for clarification.

22. **Reviewer Comment:** *Page 15. Line 42. An article “a” is missing before “uniform distribution”.*

Author Response: Fixed.

23. **Reviewer Comment:** *Appendix B.2. What is the basis bound K I checked throughout the paper and didn’t find where it is defined.*

Author Response: K is defined in (2.9). For clarification, we revised that part by adding “where K is called the basis bound.”

Reviewer 2

Reviewer Comment: *This paper presents an interesting strategy to construct a surrogate model using gPC in order to perform Uncertainty Quantification. The novelty of the research lies in the ability of the method to construct a surrogate model with arbitrary probability density functions for the input random variables. Thanks to their sparsity-enhancing procedure, the method is able to recover the quantity of interest with a limited number of samples. The overall quality of the manuscript is good. The method is presented in detailed and the authors tested their approach using various test cases. I find that while the article is, in essence, worthy of publication, I would request the authors to please respond to my questions and kindly address some specific corrections that I have noted.*

Author Response: We thank the referee for the detailed reading of the manuscript and constructive suggestions that help improve our manuscript. In particular, we appreciate the referee's suggestions on computing the predictive coefficient Q_2 and the Sobol indices on the test systems, which provide complementary information for the tested systems and demonstrate the effectiveness of our approach.

1. Reviewer Comment: *From the title and abstract, an important information is missing. It is unclear that you are using gPC as a surrogate model. This is also the case in your highlights were it should be clear. Also, the first highlight is not very informative as this is the essence of UQ as each parameter is naturally described by a PDF (being known or not)..*

Author Response: We have changed the title to “A data-driven framework for sparsity-enhanced surrogates with arbitrary mutually dependent randomness”. We have also revised the abstract and high-light accordingly to specify the motivation and the approach of the present study.

2. Reviewer Comment: *In the introduction, the authors start by clearly stating their UQ problem. The constraints are: (i) high dimensions with (ii) dependent random variables and (iii) unknown probability density functions. Following, the attention is directly shifted to gPC methods and some strategies to overcome the three constraints are given. But the bibliographical review lack information about other surrogate strategies such as Gaussian Processes that do not intrinsically suffer from these problems. As the main opponent method, this could be mentioned through classical references (Rasmussen, Gaussian processes for machine learning, 2006) and comparative studies (P.Roy, Comparison of polynomial chaos and Gaussian process surrogates for uncertainty quantification and correlation estimation of spatially distributed open-channel steady flows, Stochastic Environmental Research and Risk Assessment, 2017). Also, sparse grid methods are commonly used in high-dimensional problems.*

Author Response: We appreciate the referee's suggestions and have included a brief literature review of the Gaussian Process and sparsity grid method in the introduction of the revised manuscript.

3. Reviewer Comment: *In the explanation of the CS method, it is stated that the method is suited for linear functions. Does it mean that the overall method presented in this work is not suited for non-linear and complex system? This is not clear as the method is presented to be a general framework.*

Author Response: The method presented in the present study works for non-linear and complex systems. The surrogate model is represented as the expansion of a set of multivariate polynomial basis, i.e., $\hat{f}(\xi) = \sum_{i=1}^N c_i \psi_i(\xi)$, where $\psi_i(\cdot)$ is the basis and c_i is the

coefficient. To determine the coefficients with training set $\{\boldsymbol{\xi}^{(i)}, f(\boldsymbol{\xi}^{(i)})\}_{i=1}^M$, we employ CS to solve the linear system $\mathbf{A}\mathbf{c} = \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{M \times N}$ with $[\mathbf{A}]_{ij} = \psi_j(\boldsymbol{\xi}^{(i)})$ and $\mathbf{b} \in \mathbb{R}^M$ with $b_i = f(\boldsymbol{\xi}^{(i)})$. This does not mean the obtained surrogate function is linear. In this study, we considered the scenario when the available training points is limited (i.e., $M < N$), and we employ the CS method to determine \mathbf{c} .

4. Reviewer Comment: *May I suggest the use of classical test functions used in the field? g-function and Ishigami are notably used. g-function is particularly interesting due to the ability to set the dimensionality and the influent parameters (example here: Kucherenko, Exploring multidimensional spaces: a Comparison of Latin Hypercube and Quasi Monte Carlo Sampling Techniques, 2015). I am concerned about the ability of the method to really perform well with complex, non-linear, non-monotonous and high-dimensional functions.*

Author Response: The present work is based on the assumption of compressed sensing, which assumes that there exists a sparse representation of the multivariate basis for the QoI. Accordingly, it is not suited for g-functions or Ishigami functions, which are not in the scope of CS. We have added a clarification in the introduction and summary of revised manuscript.

5. Reviewer Comment: *Complementary to the classical l2 error, the predictivity coefficient Q_2 is most of the time used (Marrel, Calculations of Sobol indices for the Gaussian process meta model. Reliability Engineering & System Safety. 2009).*

Author Response: We have computed the predictivity coefficient Q_2 (following the definition in Ref. [3]) for the test cases of Gaussian Mixture systems (with $d = 25$ and $p = 3$) and the biomolecular systems. The results are listed in Tab. 1, where the surrogate models are constructed by the present data-driven basis approach. We have included the such results in the Appendix of the revised manuscript.

Table 1: The predictivity coefficient Q_2 for polynomial function with Gaussian Mixture measure ($d = 25$ and $p = 3$) and the molecular system for solvation energy and SASA of atom H9.

molecule solvation	M	80	160	240	320	400
	Q_2	0.995715	0.999132	0.999731	0.999864	0.999911
molecule SASA	M	200	300	400	500	600
	Q_2	0.988675	0.996069	0.998272	0.998709	0.999027
Gaussian Mixture	M	200	300	400	500	600
	Q_2	0.998372	0.999347	0.999844	0.999892	0.999941

6. Reviewer Comment: *It would be interesting to add information about the effective number of important variables in each case. $d = 25$ in most cases but what are the Sobol indices? Are all these dimensions as important? Analysis using Type A, B and C functions as proposed in Kucherenko2015 is quite informative and allows to really grasp the complexity of the cases. Especially for the motivation case, it lacks a UQ analysis.*

Author Response: We would like to thank the referee's suggestion on studying the Sobol indices. The present method allows us to construct surrogate models for systems with dependent random variables. Accordingly, we have further computed the Sobol sensitivity indices for the test cases of Gaussian Mixture systems (with $d = 25$ and $p = 3$) and the molecular systems with dependent random variables following Ref. [1]. The results

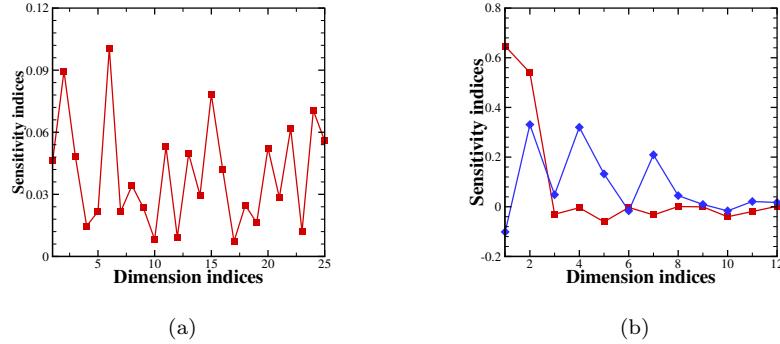


Figure 1: The first-order sobol sensitivity indices for (a) polynomial function with Gaussian Mixture mesure ($d = 25$, $p = 3$) (b) molecular system for solvation energy (“ $\text{---} \blacksquare \text{---}$ ”) and SASA of atom H9(“ $\text{---} \blacklozenge \text{---}$ ”).

are presented in Fig. 1, where the surrogate models are constructed by the present data-driven basis approach using $M = 800$, $M = 240$ and $M = 600$ training points, respectively. Based on the analysis, it is shown the dominant components are on the dimensions $(1, 2, 3, 6, 11, 13, 14, 15, 16, 20, 22, 24, 25)$, $(1, 2, 5)$ and $(1, 2, 4, 5, 7)$ (90% of total variance). We have also included these results in the Appendix of the revised manuscript.

7. Reviewer Comment: *Regarding the results. I found the sizes of the samples quite large. The classical rule of thumb to construct a surrogate model is $10n$ (Jones, Efficient global optimization of expensive black-box functions, J Glob Optim, 1998 and here is a comprehensive review Liu, A survey of adaptive sampling for global meta modeling in support of simulation based complex engineering design, Structural and Multidisciplinary Optimization, 2018). This may be linked to my previous comment. Using classical functions would allow fair comparisons between different surrogate strategies.*

Author Response: In this study, we utilize a relatively large size of sample set S to approximate the underlying probability measure $\rho(\cdot)$ and construct the orthonormal basis. On the other hand, the size of the training set M is still limited. We have examined the error of the surrogate model constructed using a broad range of training sample size similar to Ref. [4]. As shown in Tab. 1, the predictivity coefficient Q_2 is larger than 0.98 when the training size is on order $O(10d)$, which means that our approach still produced good approximations.

8. Reviewer Comment: *In all the test cases, have you used QMC or other sampling technic other than crude MC in order to generate the samples from the sought PDFs? This is known to impact a lot the quality of the surrogates. Furthermore, when constructing a surrogate model, the design of experiments does not have to follow any specific distribution. Once you obtained the surrogate, you can evaluate it using these distributions. This should be discussed even though you are in a data-driven approach as you could choose to filter the data or some time you are able to fetch new data.*

Author Response: In the present study, we focus on the scenario where the underlying measure is implicitly represented sample set and the training set is randomly collected from

the sample set and paired with the values of QoI. Accordingly, the surrogate model is constructed in a non-intrusive way; and we do not have the samples on the pre-selected points. Indeed, to construct the surrogate model, the measure of the sampling points does not need to follow any specific distributions (e.g., uniform, Gaussian, Gamma, etc.) We choose the training sample points and the data-driven orthonormal basis following the underlying measure mainly for two reasons (I) the constructed surrogate model is l_2 optimal with respect to the underlying measure (II) for some systems such as molecular system, it may not be straightforward to sample the QoI following a specific distribution; it would be more natural and convenient to collect samples directly following the underlying distribution (i.e., molecular dynamics trajectories).

9. Reviewer Comment: *Lastly, the introduction lack a concise explanation of the previously introduced method*

Author Response: We have included a brief explanation of the previous method in the revised manuscript.

10. Reviewer Comment: *A reference to classical UQ work can be added for the unfamiliar reader such as: Saltelli, Global Sensitivity Analysis. The Primer, 2007.*

Author Response: We have added the reference.

11. Reviewer Comment: *Two ideas: numerical simulations and experimental simulations. Revise to make it clear that the tools are the same and specify that you will focus on numerical experiments.*

Author Response: We have revised the Introduction accordingly.

12. Reviewer Comment: *Lines 10-11. Trying to reduce the dimensionality is a consequence of the numerical cost. The two sentences are confusing. It is suggested that the objective is to reduce the dimensionality. This should be revised. Also add a reference to the dimensionality issue.*

Author Response: We have revised the Introduction and added a reference.

13. Reviewer Comment: *Line 14 Add a reference of such realistic systems with not independent variables.*

Author Response: We have added a reference with a brief explanation.

14. Reviewer Comment: *Page 2, line 15: You are pointing out 3 challenges as opposed to 2 on line 7. You are focussing on these aspects but there are other challenges. Revise your statement about the number of important challenge by pointing out that these are your challenge of interest.*

Author Response: This refers to the part of the second challenge. We have modified the introduction to clarify this.

15. Reviewer Comment: *Page 2, lines 21-22: you referred to UQ, talk about sensitivity indices.*

Author Response: We have added sensitivity indices.

16. Reviewer Comment: *A better reference can be used in an UQ context such as: Kucherenko, Exploring multidimensional spaces: a Comparison of Latin Hypercube and Quasi Monte Carlo Sampling Techniques, 2015.*

Author Response: We have added this reference.

17. **Reviewer Comment:** *page2, line25 : inefficient is not the word. It is the numerical cost.*

Author Response: We have modified the sentence.

18. **Reviewer Comment:** *page2, line26 : Remove etc..*

Author Response: We have removed the word.

19. **Reviewer Comment:** *Page 2, line 26: Preferably refer to the original papers instead of 6-7: Mckay, A Comparison of Three Methods for Selecting Values of Input Variables in the Analysis of Output from a Computer Code, 1979*

Author Response: We have added the above reference.

20. **Reviewer Comment:** *Page 2, lines 26-28: explain why QMC or LHS are not suited with unknown distributions. This statement is strong and should be backed up.*

Author Response: We have modified this sentence.

21. **Reviewer Comment:** *Page 3, lines 41-50: For someone new to gPC, it will be hard to understand the curse of dimensionality from this.*

Author Response: We have modified this paragraph.

22. **Reviewer Comment:** *Page3, lines 54-61: Missing reference about the biomolecular example. Also, this part is hard to follow without prior knowledge about the physics. For instance, what is a conformational state? This should be revised.*

Author Response: We have modified this paragraph and added a relevant reference.

23. **Reviewer Comment:** *Page 4, line 77: Reference 47 does not seem relevant for copula*

Author Response: Copulas are discussed in Sec. 3.4 of Reference 47.

24. **Reviewer Comment:** *general UQ framework for.*

Author Response: We have modified this sentence.

25. **Reviewer Comment:** *Page 4, line 98: it is a strong statement to make saying that you have access to the analytical form of the PDF in real-world applications*

Author Response: We may not know the analytical form of PDF in real-world applications; and this scenario is discussed as the situation (I) in the paragraph; while there are also systems where the underlying distribution is accessible from physical principle (e.g., Boltzmann distribution of Hamiltonian systems), or the randomness is specified by users as direct input.

26. **Reviewer Comment:** *page 4, line 105: Remove rest.*

Author Response: We have removed the word.

27. **Reviewer Comment:** *Page 5, line 112: biomolecular.*

Author Response: We have changed the word.

28. **Reviewer Comment:** *Page 6, Eq 2.4: c is not defined clearly and it is central in the following.*

Author Response: We have modified this part.

29. **Reviewer Comment:** *Explain what is N, this is important.*

Author Response: We have modified this part.

30. **Reviewer Comment:** You clearly introduce an notation for l_1 norm before. Be consistent. Also explain what are A and b .

Author Response: We have specified a previous definition $|\cdot|$ for l_1 norm on multi-index. We have added the explanation of A and b .

31. **Reviewer Comment:** page 6: Missing equation number after 2.5.

Author Response: We appreciate the suggestion. As this equation is not referred elsewhere in the manuscript, we skipped the equation label.

32. **Reviewer Comment:** Page 7, Def 2.2: $s = 1, 2, \dots$, up to ?

Author Response: We have modified this part.

33. **Reviewer Comment:** Page 7, line 139: correct ϵ .

Author Response: We have corrected this typo.

34. **Reviewer Comment:** Page 7, Th 2.3: Define M and N same as page 6?

Author Response: M refers to the row of matrix \mathbf{B} . This theorem holds under the conditions specified and N is not necessarily any specific value.

35. **Reviewer Comment:** Page 7, Def 2.4: n defined after.

Author Response: We have modified this part.

36. **Reviewer Comment:** page7: Missing two equations number before 2.8 and one after 2.9.

Author Response: We appreciate the suggestion. As this equation is not referred elsewhere in the manuscript, we skipped the equation label.

37. **Reviewer Comment:** Page 6-7: Introduce beforehand what you will express mathematically. The list of definition and theorem as it appears unclear. Also, what is the conclusion of this subsection(CS)?

Author Response: We have modified this part with a few transition sentences.

38. **Reviewer Comment:** Page 8, lines 150-155: The definition of A should appear in the CS subsection. Some elements found here could serve as an introduction to CS subsection.

Author Response: The \mathbf{A} in CS subsection refers to a matrix considered in the linear system $\mathbf{Ac} = \mathbf{b}$. The theory of CS discussed in Section 2 holds in general. As an application, we specify \mathbf{A} as the measurement matrix and \mathbf{b} as the samples of QoI to recover the coefficient \mathbf{c} . Therefore, we specified the definition of \mathbf{A} in Section 3.

39. **Reviewer Comment:** Page 8, line 155: our approach, which is? From the text before you presented previous work and your contribution does not appear.

Author Response: We have modified this sentence.

40. **Reviewer Comment:** page 8, line 159: Clarify notations of the random vector. Here ξ and in page 6 line 128 it is z . Same goes with N sand M .

Author Response: ξ in Section 3 is related to z in Section 2. z in Section 2 represents a random vector as the input of a multivariate polynomial function, and ξ in Section 3 refers to the random vector of the surrogate model we aim to construct. Later on, we constructed the surrogate model with respect to the rotated random vector χ . Under such consideration,

we employed different notations in section 2 and 3. Here we use N_s with a subscript “s” to represent the total number of samples in set S .

41. Reviewer Comment: *Page 10: see previous remark, N is not the same as N_s but this is really confusing as there is no clear explanation of N*

Author Response: N is the number of basis specified in Eq. (2.5) and text following in the revised manuscript.

42. Reviewer Comment: *page 10: Missing equation number before 3.8.*

Author Response: We appreciate the suggestion. As this equation is not referred elsewhere in the manuscript, we skipped the equation label.

43. Reviewer Comment: *Page 13, Eq 3.14: By using gradient information, you are not ensured to have the global interaction of the parameters. Only global sensitivity indices such as Sobol can guarantee it. Especially in high dimensions with a low number of samples, if the function is not smooth enough you can be assured that gradient information will be erroneous.*

Author Response: The gradient information alone is insufficient to characterize the global interaction of the parameters, and the accuracy of the gradient matrix is not guaranteed if QoI is non-smooth. This is the major reason we did not conduct further dimension reduction following the evaluation of the gradient matrix. Instead, we conducted eigendecomposition and establish a linear transform of the random vector ξ into χ (with the same dimension) such that the coefficient vector c with respect to χ is sparser (therefore more favorable for CS), if the gradient matrix can partially capture the variability of QoI along the d directions. There could be alternative approaches to establish such transformation (i.e., via nonlinear transformation), which admit sparse representation of QoI. This direction is worth further exploration in future studies.

44. Reviewer Comment: *Page 15, line 254: Better define ω*

Author Response: We have added a definition to the paragraph.

45. Reviewer Comment: *page 15: Missing equation number after 4.3.*

Author Response: We appreciate the suggestion. As this equation is not referred elsewhere in the manuscript, we skipped the equation label.

46. Reviewer Comment: *page 17, fig1: The mean is interesting but figure would benefit from information about the variance also.*

Author Response: We have added error bars.

47. Reviewer Comment: *Page 24, line 360: The motivation example should be a separated subsection to separate analytical tests from it.*

Author Response: We have separated this example as a new subsection.

48. Reviewer Comment: *Page 27, lines 407-408: To be correct, you have done an uncertainty propagation study. UQ is a more general term encompassing sensitivity analysis as well.* **Author Response:** We have modified this sentence.

49. Reviewer Comment: *Ref 56 missing year.* **Author Response:** We have included the year.

50. Reviewer Comment: *Ref 60 solver.* **Author Response:** We have fixed the typo.

Reviewer 4

Reviewer Comment: *The manuscript deals with an approach designed to enhance the performance of forward Uncertainty Quantification analysis in the presence of high-dimensionality and arbitrary probability density distributions of the inputs. The two main components of the framework are (i) the possibility to construct the optimal (orthogonal/near-orthogonal) polynomial basis for the compressive sensing and (ii) the use of a Active Subspace based rotation to enhance the sparsity of the coefficient vector. The manuscript consists of a theoretical introduction followed by a set of numerical results where the proposed approach is compared against other related strategies.*

I found the manuscript overall carefully written, the idea interesting, the numerical examples relevant and I think it perfectly fits the scope of the Journal. However, I would recommend some modifications before accepting it for publication. In the following I report my suggestions without ranking them in a particular order of importance.

1. Reviewer Comment: *The introduction presents a relevant literature survey. I found also interesting and appropriate that the authors include a reference to Monte Carlo sampling methods as a robust and alternative way to deal with the same problem. My suggestion here is to extend this reference to Monte Carlo by adding some of the most recent directions in this context. I am referring to the large recent literature based on multilevel and multifidelity sampling strategies. The reason for that is that I believe the approach proposed in this work is somehow complementary with respect to sampling based approach because it is more effective when the dimensionality of the input space is not extremely large and the solution is regular. If those conditions are not met, it is still mandatory to rely on sampling based approach. I think that adding a comment along those lines might help the reader.*

Author Response: Recent work on multilevel-MC and multifidelity-MC have been developed to improve sampling efficiency and are relevant within the scope of Monte Carlo sampling. We appreciate the referee pointed this out and we have added this discussion to the revised manuscript.

2. Reviewer Comment: *Please check page 7 line 14*

Author Response: We would like to thank the referee for careful reading of our manuscript. We have fixed this typo in the revised manuscript.

3. Reviewer Comment: *I really found the discussion carried over at pages 11 and 12 about the near-orthonormal basis very interesting. However, only the case of two disjoint sets is considered, whereas in many cases one would like to add new samples to a previous computation. This might happen for instance if new data become available after the first regression is performed. What does it happen in this case? Is the orthonormality retained? Should Algorithm 3 be modified? I think it is fair to mention this issue without fully addressing it and, for instance, add it as one of the future directions.*

Author Response: In this study, we construct the near-orthonormal basis in a heuristic way by directly splitting the sampling set S into two disjoint sets. In practice, if the new available data set shares the same underlying probability $\rho(\cdot)$ with the previous sample set S , then constructed bases will retain proper orthonormality. However, it is worth exploring how to utilize the new data set to design some more sophisticated (cross-validation) strategies to further optimize the orthonormal threshold values and the near-orthonormal basis construction procedure. We have added a brief discussion in the revised manuscript.

4. Reviewer Comment: *I would suggest to be more explicit about the amount of data*

required for the identification of the probability measure. In the summary the authors clearly state that this is possible under the assumption that a large collection of samples is available. Even if this might be the case in some applications, in many other cases the problem is that a limited number of samples are available. I would recommend to include a comment on that already in the abstract/introduction.

Author Response: For the case of implicit distributions, the current study assumes that a large collection of samples is available, which represents the underlying probability measure (while the size of the training set can still be small). We have included this comment in the abstract and introduction of the revised manuscript.

5. Reviewer Comment: *I found the numerical results interesting and relevant. However, I found the introduction of Section 4 (Results) a little bit too concise. In particular the strategy adopted for the numerical tests is not completely clear to me. The data set S is divided (evenly) in two subsets (I assume with the intersection set empty) and the first is used for constructing the basis, while the second is used to construct and test the surrogate. We have here two training sets S_1 and part of $S_{2,1}$ which are used to build the basis and to compute the coefficients, and a set $S_{2,2} = S_2 \setminus S_{2,1}$ which should be used for the test of the surrogate. My main point is that the surrogate should be built on a training data set which need to be disjoint from the test dataset. I think the authors need to clarify this point and make it more explicit in the text (i.e. how many training points? how many test points? etc.).*

Author Response: The numerical tests were conducted exactly as the referee described above. The training set $S_{2,1}$ and the test set $S_{2,2}$ are disjointly chosen from S_2 . The size of the training set is $O(10^2) - O(10^3)$ and size of the test set is $O(10^5)$ in the present study. We have clarified this in the revised manuscript.

6. Reviewer Comment: *Related to the previous point, for instance if one looks at Section 4.1 the presentation of the numerical experiment gets a little bit confusing (at least for me). $S = 2 \times 10^5$, so S_1 and S_2 are 10^5 , however the minimization problem is solved for increasing M (number of measurements) which I think it is the right thing to do. However, it is not clear what happens to S_1 , is the cardinality of S_1 also equal to M ? It would be inconsistent in my opinion to use the full set of samples S_1 to build the base and M samples for the output. I suggest to clarify this point.*

Author Response: In this numerical test, the cardinality of S_1 is 10^5 . $M \sim O(10)$ represents the size of the training set $S_{2,1}$ chosen from S_2 . S_1 and $S_{2,1}$ are disjoint sets (please also refer to the above response.) The motivation is to examine, for various bases, how many training points (e.g., M) is needed to accurately recover the vector \mathbf{c} . We have clarified this in the revised manuscript.

Reviewer 5 Reviewer Comment: *This paper addresses some challenging problems in real-life applications that exhibit randomness with imprecise and non-standard input distributions often represented by a finite set of samples. I think this paper is worth to be published, but please note the following comments and suggestions.*

1. **Reviewer Comment:** *First of all, the authors are commended for a relatively comprehensive survey of the current state-of-the-art methods setting a good stage for subsequent sections.*

Author Response: We appreciate the referee's comment on the literature review of our manuscript. Following the suggestions of Review 4, we have also added the relevant literatures of recent work on multilevel-MC and multifidelity-MC in the scope of Monte Carlo sampling in the revised manuscript.

2. **Reviewer Comment:** *The reference to 'high-dimensionality' in the paper seems to be somewhat misleading. The highest dimension studied in the numerical experiments is about 20. Typically, high-dimensionality refers to hundreds to millions of dimensions such as those corresponding to random fields. Sparsification and active subspace were mentioned as ways to reduce the number of effective basis functions, but they may not be able to handle really 'high' dimensions. Claiming to deal with high dimensionality may invite many questions about why this or that machine learning method is not considered. So it is better to clarify what type of high-dimensional problems you are trying to address.*

Author Response: Indeed we address the problems with arbitrary randomness in dimension of $O(10)$ in the present study. In the original version of the manuscript, we referred to "high-dimensional" in the scope of stochastic modeling of molecular systems (within the conformational space). In this field, the number of collective variables (i.e., the dimension of randomness) is often limited to 5 or less due to the numerical challenges in accurate quantification of free energy landscape (e.g., please refer to molecular system where the underlying density is approximated by Gaussian mixture in Sec. 4.3.3) The data-driven and sparsity enhancement approach developed in the present study provides a work-around approach to accurately quantify the uncertainty of such systems.

We appreciate the referee for pointing out this possible confusion. In the revised manuscript, we have specified the dimensionality of problems we aim to address.

3. **Reviewer Comment:** *The authors should spell out what is new in the orthonormalization procedure, i.e. what is new compared to Dunkl and Xu's approach? I think the orthonormalization procedure is somewhat basic. The potentially interesting and useful idea is near-orthonormalization because in practice the data set may not be complete. I suggest more focus be put on demonstrating the effectiveness and limitations of near-orthogonality. For example, In Remark 3.2 you mentioned that Eqn 3.9 "provides a heuristic approach to construct near-orthonormal basis." While this heuristic currently is not accompanied by any rigorous error analysis but it seems to outperform the exact orthonormal basis, it is desirable to add some numerical experiments to study how performance degrades when the relaxation parameters varies (that is, how —c-c— deteriorates with respect to loss of orthonormality?) The authors do state in the Conclusion section that this is something they plan to investigate, but the inclusion of a small numerical study in this paper will make the paper more complete.*

Author Response: The orthonormal basis is constructed following the Dunkl and Xu's approach. We highlighted this in the revised manuscript. In the current study, we did not include a systematic study of the near-orthonormal basis. In our understanding, the differ-

ent numerical performance between the orthonormal and near-orthonormal bases arises from the different null spaces of the measurement matrices constructed from the two bases. Accordingly, there might be alternative approach to optimize the basis construction; however, we feel that it is best to defer such investigation to a separate study.

4. Reviewer Comment: *It will be more informative if the authors include a cost comparison between Algorithm 4 and Algorithm 4 without the use of different bases (χ) - that is, Step 4-6. Numerically, the examples show better accuracy of the former but for what additional costs*

Author Response: The additional cost of Step 4 - 6 mainly consists of two parts: construction of the basis with respect to χ and construction of the surrogate model $\tilde{f}(\chi)$. For the current numerical examples considered in the present study, this additional cost is less than 0.6 CPU (3.7 GHz Quad-Core Intel Xeon E5) hour in general. For realistic applications, the overhead of Step 4 - 6 could be relatively small if sampling of QoI is expensive or the available training set is limited. We have included this discussion in the revised manuscript.

5. Reviewer Comment: *Minor: on p. 21 line 34: What is Λ_p^d ?*

Author Response: Λ_p^d is a multi-indices set defined in Eq. (2.4). We have clarified this in the revised manuscript.

References

- [1] G. Chastaing, F. Gamboa, and C. Prieur. Generalized sobol sensitivity indices for dependent variables: numerical methods. *Journal of Statistical Computation and Simulation*, 85(7):1306–1333, 2015.
- [2] Paul G Constantine, Eric Dow, and Qiqi Wang. Active subspace methods in theory and practice: Applications to kriging surfaces. *SIAM J. Sci. Comput.*, 36(4):A1500–A1524, 2014.
- [3] Amandine Marrel, Bertrand Iooss, Béatrice Laurent, and Olivier Roustant. Calculations of sobol indices for the gaussian process metamodel. *Reliability Engineering & System Safety*, 94(3):742 – 751, 2009.
- [4] Pamphile T. Roy, Nabil El Moçayd, Sophie Ricci, Jean-Christophe Jouhaud, Nicole Goutal, Matthias De Lozzo, and Mélanie C. Rochoux. Comparison of polynomial chaos and gaussian process surrogates for uncertainty quantification and correlation estimation of spatially distributed open-channel steady flows. *Stochastic Environmental Research and Risk Assessment*, 32(6):1723–1741, Jun 2018.
- [5] Trent Michael Russi. *Uncertainty Quantification with Experimental Data and Complex System Models*. PhD thesis, University of California, Berkeley, 2001.

Highlights

- Surrogate construction for arbitrary mutually dependent probability measures
- A general framework for both explicit non-Gaussian densities and implicit probability measures
- Data-driven orthonormal-basis construction coupled with sparsity-enhancing rotations
- Application to molecular systems with non-Gaussian mutually dependent randomness

A data-driven framework for sparsity-enhanced surrogates with arbitrary mutually dependent randomness

Huan Lei^{a,1,*}, Jing Li^{a,1}, Peiyuan Gao^a, Panagiotis Stinis^{a,b}, Nathan A. Baker^{a,c,*}

^aPacific Northwest National Laboratory, Richland, WA 99352

^bDepartment of Applied Mathematics, University of Washington, Seattle, WA 98195

^c*Division of Applied Mathematics, Brown University, Providence, RI 02912*

Abstract

The challenge of quantifying uncertainty propagation in real-world systems is rooted in the high-dimensionality of the stochastic input and the frequent lack of explicit knowledge of its probability distribution. Traditional approaches show limitations for such problems, especially when the size of the training data is limited. To address these difficulties, we have developed a general framework of constructing surrogate models on spaces of stochastic input with arbitrary probability measure irrespective of **the mutual dependencies between individual components of the random inputs** and the analytical form. The present *Data-driven Sparsity-enhancing Rotation for Arbitrary Randomness* (DSRAR) framework includes a **data-driven construction of multivariate polynomial basis for arbitrary mutually dependent probability measures** and a sparsity enhancement rotation procedure. This sparsity-enhancing rotation method was initially proposed in our previous work [1] for Gaussian density distributions, which may not be feasible for non-Gaussian distributions due to the loss of orthogonality after the rotation. To remedy such difficulties, **we developed a new data-driven approach to construct orthonormal polynomials for arbitrary mutually dependent randomness, ensuring the constructed basis** maintains the orthogonality/near-orthogonality with respect to the density of the rotated random vector, **where directly applying the regular polynomial chaos including arbitrary polynomial chaos (aPC)** [2] shows limitations due to the assumption of the mutual independence between the components of the random inputs. The developed DSRAR framework leads to accurate recovery, with only limited training data, of a sparse representation of the target functions. The effectiveness of our method is demonstrated in challenging problems such as partial differential equations and realistic molecular systems **within high-dimensional ($O(10)$) conformational spaces** where the underlying density is implicitly represented by a **large collection of sample data**, as well as systems with explicitly given non-Gaussian probabilistic measures.

Keywords: data-driven, arbitrary randomness, mutual dependence, compressed sensing, sparsity enhancement, uncertainty quantification

*Corresponding author

Email addresses: huan.lei@pnnl.gov (Huan Lei), nathan.baker@pnnl.gov (Nathan A. Baker)

¹The first two authors contributed equally.

1
2
3
4 **1. Introduction**
5

6 A fundamental problem in uncertainty quantification (UQ) [3] is to calculate the statistical properties of a
7 quantity of interest (QoI) due to various sources of randomness, e.g., numerical simulations subject to uncer-
8 tain parameters, initial conditions and/or boundary conditions, as well as experimental measurements in the
9 presence of material heterogeneity, thermal fluctuations. Such sources of uncertainty are usually character-
10 ized by high-dimensional random variables whose probability measures can be either discrete or continuous.
11 In real-world systems, there are usually two crucial challenges to accurately quantify the propagation of the
12 randomness from the input to the system response. The first challenge comes from the high-dimensionality of
13 the random inputs. For such systems, limited computational resources often motivates further dimensionality
14 reduction [4]. However, it is often non-trivial to accurately transfer the high-dimensional random space into
15 a low-dimensional random space. This results in the numerical intractability of quantifying the uncertainty
16 of the QoI from training data of limited size. The second challenge arises from frequent dependencies and ar-
17bitrary distribution of the random inputs. Typically, random inputs are represented by random vectors with
18 mutually independent components. For realistic systems, the underlying distribution of the inputs can often
19 involve dependencies that cannot be ignored (e.g., see molecule systems in Ref. [5] and Sec. 4.4). Moreover,
20 the input distribution could be even unknown and thus we may only have access to it implicitly through a
21 collection of samples. This creates further numerical obstacles in characterizing the random inputs as well
22 as their effect on the system response. In the current work, we present a *Data-driven Sparsity-enhancing*
23 *Rotation for Arbitrary Randomness* (DSRAR) framework for dealing with all of the aforementioned chal-
24 lenges. While we focused on numerical experiments in the present study, the developed framework can be
25 also applied to UQ in experimental studies.

26 In practice, a straightforward and robust approach is the Monte Carlo (MC) method, which involves
27 collecting a large number of samples of the random inputs from their distribution, evaluating the QoI at each
28 sample point, and then obtaining the statistical properties (mean, variance, sensitivity indices, probability
29 density function, probability of a certain event etc.) of the QoI. Unfortunately, to get an accurate estimate,
30 the MC method requires a large number of simulations due to its slow convergence rate [6, 7]. Furthermore,
31 for large or complex systems, even a single instance of these simulations may require very large computational
32 resources. Under such circumstances, the computational cost of MC method can become extremely large.
33 Several approaches have been developed to alleviate such difficulties. For instance, sampling approaches such
34 as multilevel-MC [8, 9, 10] and multifidelity-MC [11, 12] have been designed to optimize the computational
35 load when samples of the QoI are available at hierarchical levels of accuracy; sampling approaches like quasi-
36 MC [13, 14, 15] and Latin Hypercube sampling [16, 17, 18], have been designed to accelerate convergence.
37 However, when the underlying distribution of the inputs is arbitrary and not explicitly given, the latter two
38 sampling strategies may lose their advantage if it is not straightforward to generate quasi-random sequences
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 35 following the underlying distribution.
5

6 An alternative approach approximates the QoI via constructing the surrogate model of the random inputs
7 and then calculates the statistics of the QoI analytically or numerically. Among such approaches, the most
8 popular are the Gaussian Process [19, 20, 21], and the polynomial chaos expansion originally introduced
9 by Wiener [22], applied to UQ by Ghanem [23] and extended to the generalized polynomial chaos (gPC)
10 expansion by Xiu [24]. The Gaussian Process (GP) is a stochastic process which approximates the values
11 of the QoI at every finite sets of sample point as multivariate Gaussian random vectors. The flexibility of
12 the mean and covariance functions enables GP to characterize a wide range of function behavior with broad
13 applications on UQ [21, 25, 26, 27, 28]. The gPC expansion approximates the QoI by a set of simple basis
14 functions. It is known to be a *mathematically optimal* approximation of the QoI when the basis functions
15 are chosen to be orthogonal with respect to the probability measure of the random inputs. This approach
16 has been demonstrated for diverse applications in UQ [29, 30, 31, 32, 33, 34, 35, 36] due to its spectral
17 convergence under certain situations. In this study, we focus on the approach developed based on gPC and
18 we refer to previous publications [37, 38, 39, 40] (and the references therein) for comparative studies of the
19 two approaches.

20
21 50 In principle, if the orthogonal polynomial type and the corresponding random variables are determined,
22 both intrusive and non intrusive methods can be used to evaluate the coefficients of the expansion. For
23 example, stochastic collocation, based on tensor products of one-dimensional quadrature rules, is often
24 employed when dimensionality is small [41, 42, 43], with the number of basis functions given by $(p+d)!/p!d!$,
25 where p is polynomial order and d is the dimension. However, as the dimension increases, the number
26 of quadrature points needed for the tensor product rule increases exponentially. To mitigate this issue,
27 sparse grid and adaptive collocation methods have been proposed to deal with moderate dimensionality
28 [42, 44, 45, 46, 47, 48, 49]. When the dimension of the random inputs is large, none of the above collocation
29 methods is feasible. In the case of a limited number of available simulations and large dimensionality,
30 compressed sensing (CS) approaches have been used to construct sparse polynomial approximations of the
31 QoI [50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62]. Finally, we note that gPC (including extensions such
32 as arbitrary polynomial chaos [2]), in its current form, can only handle random vector with independent
33 identically distributed (i.i.d.) components in standard types (uniform, Gaussian, gamma, beta, etc.). For
34 other distributions, a pre-processing step is required to transform the original random variables into i.i.d.
35 random variables of standard types. In general, these transformations are highly nonlinear which result in
36 the final QoI function approximation to be a high-degree polynomial in order to maintain accuracy.

37
38 53 The methods discussed above rely on the *explicit* knowledge of the underlying probability measures
39 and/or the assumption of mutual independence between the components of the random inputs. However,
40 such assumptions on the random inputs can be quite restrictive for realistic applications. One such example
41

is the UQ for molecular system properties QoIs due to conformational fluctuations [63]. For such systems, the random inputs are the various conformational states (i.e., the instantaneous structure) of the molecule. The underlying distribution is determined by the free energy function of the system, which is essentially the multi-dimensional marginal density distribution with respect to the (Boltzmann) distribution of the full Hamiltonian system. Unfortunately, numerical evaluation of the free energy function is a well-known challenging problem. Although various sampling strategies have been developed [64, 5, 65], the explicit free energy function is usually unknown for dimensions greater than 4. In practice, the underlying density is only known implicitly through a large collection of the molecule conformational states obtained from experiments or simulated trajectories. Another commonly encountered example arises in our recent work [1] on constructing sparse representations of a QoI based on CS. Inspired by the active subspace method [66], we proposed a method to enhance the sparsity of polynomial expansion in terms of a new random vector via unitary rotation of the original random vector. For i.i.d. Gaussian random inputs, the new random vector retains the same distribution. However, for non-Gaussian random inputs, which are more realistic for applications, the new random vector does not retain the mutual independence even if the original random vector elements are i.i.d.

For problems with non-Gaussian random inputs, the traditional approach is to cast the available statistics into a family of standard distributions and then to apply the gPC techniques discussed above. Gaussian mixture models, due to their flexibility, are broadly employed to approximate the distribution of the data. With the distribution approximated, a gPC expansion of the QoI can be constructed for each Gaussian component. The statistical properties of the QoI are derived by combining the statistical properties of all components [67, 68]. However, there are two drawbacks of the Gaussian mixture approach: (i) it lacks one-to-one correspondence between one instance of random inputs and the approximated function evaluation, (ii) it is difficult to determine an appropriate and accurate probability density approximation when the dimension is larger than one. Copulas have been employed to treat dependent probabilistic models for surrogate construction in [69]. Zabararas [70] has established a graph-based approach to factorize the joint distribution into a set of conditional distributions based on the dependence structure of the variables. Alternatively, several studies have been devoted to constructing orthogonal polynomial bases using the moments of the random variables. Orthogonal polynomial chaos for random vectors with independent components of arbitrary measure was proposed in [2, 71, 72, 73, 74]. Ahlfeld investigated the quadrature rule of this arbitrary polynomial chaos (aPC) and proposed a sparse quadrature rule for the integration which can facilitate the evaluation of the expansion coefficients [75]. However, those quadrature rules of arbitrary polynomial chaos again assume *the components of the random inputs are mutually independent*.

In this paper, we develop a general UQ framework for constructing surrogate models via DSRAR *irrespective of possible mutual dependencies between the random input components*. This approach is different

from the aforementioned studies based on polynomial chaos expansions and, therefore, can be particularly useful for realistic systems where the input distributions can be non-standard or unknown analytically. The key idea is a data-driven approach for basis construction, consisting of multivariate orthonormal *polynomials for arbitrary mutually dependent* (amdP), coupled with the previously developed rotation-based sparsity enhancement approach [1]. This can be viewed a special case of the present method when the random inputs are from a Gaussian distribution. When the size of the training set is limited, the method can recover the expansion coefficients by CS, under the assumption that there exists a sparse representation of surrogate model. As we will show, directly employing a regular polynomial basis and/or the sparsity enhancement rotation on the random input may result in large recovery error due to the violation of orthogonality for non-standard density distributions. The procedure of data-driven basis construction described in the present study retains proper orthogonality with respect to the associated random inputs and therefore ensures more accurate recovery. In this sense, the present method takes advantage of both the orthonormal basis expansion and the enhanced sparsity of the expansion coefficients. The method deals with two situations widely encountered in real-world applications: (I) probability measures that are implicitly represented by a large collection of samples and (II) non-Gaussian probability measures with explicit (analytical) forms. For the first situation, we construct orthonormal polynomial bases with respect to discrete measures on the sample set. Besides the exact orthonormal basis, we also propose a heuristic method to construct a *near-orthonormal* basis, which yields a smaller basis bound than the exact orthonormal basis and results in more accurate recovery of the sparse representation. For the second situation, we construct the orthonormal basis when the quadrature rules for polynomial integration are known. This construction is especially well suited to random variables obtained from sparsity enhancement of non-Gaussian distributions.

The paper is organized as follows. In Section 2, we present the problem setup and briefly review preliminary background on multivariate orthogonal polynomials and compressed sensing. In Section 3, we present the DSRAR framework by first introducing the methods to construct data-driven orthonormal amdP basis. When the underlying density is implicitly represented by a large collection of random input samples, we propose a heuristic approach to construct a near-orthonormal basis along with some heuristics on the advantage over an exactly orthonormal basis. Then we introduce the rotation-based sparsity enhancement method and provide algorithmic details on how to combine the data-driven basis construction and sparsity enhancement rotation. In Section 4, we demonstrate the developed framework in a realistic molecular system fluctuating in a high-dimensional conformational space ($O(10)$) as well as partial differential equations (PDEs) with arbitrary randomness where the underlying distributions are either explicitly known or implicitly represented by a large collection of samples. Concluding remarks and directions for future work are provided in Section 5.

2.1. Approximation with orthogonal polynomials

We begin with a few facts about multivariate

nomials in d variables on \mathbb{R}^d . Polynomials in \mathbb{R}^d are naturally indexed by the multi-indices set \mathbb{N}_0^d . For $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}_0^d$ and $z = (z_1, \dots, z_d)$, a monomial z^α is defined by $z^\alpha = z_1^{\alpha_1} \cdots z_d^{\alpha_d}$ and the degree of z^α is defined by $|\alpha| = \alpha_1 + \cdots + \alpha_d$. From now on, without confusion, $|\cdot|$ operating on a **multi-index** α denotes the ℓ_1 norm of α while $|\cdot|$ operating on a set T denotes the cardinality of T . The degree of a polynomial is defined by the largest degree of its monomial terms. Then the space of polynomials of degree at most p is defined by

$$\Pi_p^d := \text{span}\{z^\alpha : |\alpha| \leq p, \alpha \in \mathbb{N}_0^d\} \text{ and } \dim \Pi_p^d = \binom{p+d}{p}. \quad (2.1)$$

If we equip \mathbb{R}^d with a probability measure ρ , then we can define an inner product on Π^d ,

$$\langle f, g \rangle_\rho = \int_{\mathbb{R}^d} fg \, d\rho \quad f, g \in \Pi^d. \quad (2.2)$$

f and g are said to be orthogonal with respect to ρ if $\langle f, g \rangle_\rho = 0$. Given such an inner product, and an order of the set \mathbb{N}_0^d , we can apply the Gram-Schmidt process on the ordered set $\{z^\alpha : \alpha \in \mathbb{N}_0^d\}$ to generate a sequence of orthogonal polynomials. We will revisit this construction in Section 3.1. When $d > 1$, there is no natural order among monomials. As a result, multivariate orthogonal polynomials are, in general, not uniquely determined. In this paper, we choose the *graded lexicographic order* when applying the Gram-Schmidt process, that is, $z^\alpha \succ z^\beta$ if $|\alpha| > |\beta|$ or if $|\alpha| = |\beta|$ and the first nonzero entry in the difference $\alpha - \beta$ is positive.

When a simulation model is expensive to run, building an approximation of the response of the model output with respect to the variations in the model input can often be an efficient approach to quantify uncertainty propagation. The polynomial approximation of a function (model) $f(\mathbf{z}) : \mathbb{R}^d \rightarrow \mathbb{R}$, $d \geq 1$ where $\mathbf{z} = (z_1, \dots, z_d) : \Omega \rightarrow \mathbb{R}^d$ is a d -dimensional random variable with associated probability measure $\rho(\mathbf{z})$, which is widely used due to its fast convergence when $f(\mathbf{z})$ is analytic. In this paper, we will approximate f using an orthogonal polynomial basis. It is a generalization of the gPC expansion which usually deals with i.i.d. random variables.

Let $\Psi = \{\psi_{\alpha}(z) : \alpha \in \mathbb{N}_0^d\}$ be a set of orthonormal polynomial basis of Π^d associated with the measure $\rho(z)$, that is,

$$\int \psi_{\alpha}(z)\psi_{\beta}(z) \, d\rho(z) = \delta_{\alpha\beta}, \quad \alpha, \beta \in \mathbb{N}_0^d, \quad (2.3)$$

where $\delta_{\alpha\beta} := \prod_{i=1}^d \delta_{\alpha_i, \beta_i}$ to be the multi-index Kronecker delta. Then the p th-degree arbitrary orthogonal polynomial expansion $f_p(\mathbf{z})$ of function $f(\mathbf{z})$ associated with ψ is defined as,

$$f(\mathbf{z}) \approx f_p(\mathbf{z}) := \sum_{\alpha \in \Lambda_p^d} c_\alpha \psi_\alpha(\mathbf{z}), \quad \Lambda_p^d = \{\alpha \in \mathbb{N}_0^d : |\alpha| \leq p\}, \quad (2.4)$$

where c_α is the coefficient to be evaluated. Using an ordering of the orthonormal polynomial basis, we can change (2.4) into the following single index version

$$f_p(\mathbf{z}) = \sum_{\alpha \in \Lambda_p^d} c_\alpha \psi_\alpha(\mathbf{z}) = \sum_{n=1}^N c_n \psi_n(\mathbf{z}), \quad (2.5)$$

where N is the total number of basis and is given by

$$N = \dim \Pi_p^d = |\Lambda_p^d| = \binom{d+p}{p}.$$

2.2. Compressed sensing

Compressed sensing is a well-studied and popular approach to find sparse solutions to linear equations [77, 78, 79, 80]. In this subsection, we briefly review the theory of CS and discuss the conditions which allow accurate recovery of solutions to underdetermined linear system.

Under certain assumptions, the solution—or its approximation—can be found by the well-studied ℓ_1 minimization, i.e., finding the minimizer

$$\min \|\mathbf{c}\|_1 \quad \text{subject to } \mathbf{A}\mathbf{c} = \mathbf{b}, \quad (2.6)$$

where $\mathbf{A} \in \mathbb{R}^{M \times N}$, $\mathbf{b} \in \mathbb{R}^M$ and $\|\mathbf{c}\|_1 = \sum_{i=1}^N |c_i|$ is the ℓ_1 norm of the vector \mathbf{c} .

When the data \mathbf{b} is contaminated by noise, the constraint in (2.6) is relaxed to obtain the basis pursuit denoising problem,

$$\min \|\mathbf{c}\|_1 \quad \text{subject to } \|\mathbf{A}\mathbf{c} - \mathbf{b}\|_2 \leq \sigma, \quad (2.7)$$

where σ is an estimate of the ℓ_2 norm of the noise. The optimization problems (2.6) and (2.7) can be solved with efficient algorithms from convex optimization [81].

Next we discuss the conditions for the sparse recovery of \mathbf{c} .

Definition 2.1. A vector \mathbf{c} is said to be s -sparse if it has at most s nonzero entries, i.e., \mathbf{c} is supported on $T \subset \{1, \dots, N\}$ with $|T| \leq s$.

Definition 2.2 (Restricted isometry constant [82, 83]). For each integer $s = 1, 2, \dots, N$ define the isometry constant δ_s of a matrix \mathbf{A} as the smallest number such that

$$(1 - \delta_s) \|\mathbf{c}\|_2^2 \leq \|\mathbf{A}\mathbf{c}\|_2^2 \leq (1 + \delta_s) \|\mathbf{c}\|_2^2$$

holds for any s -sparse vector $\mathbf{c} \in \mathbb{R}^N$.

The restricted isometry constants (RICs) characterizes matrices that are nearly orthonormal. The sparse recovery is established by the following theorem.

Theorem 2.3 (Sparse Recovery for restricted isometry property (RIP)-Matrices). *Let $\mathbf{A} \in \mathbb{R}^{M \times N}$. Assume that its isometry constant δ_{2s} satisfies $\delta_{2s} < 0.4931$. Let $\mathbf{c} \in \mathbb{R}^N$, and assume noisy measurements $\mathbf{b} = \mathbf{Ac} + \boldsymbol{\eta}$ are given with $\|\boldsymbol{\eta}\|_2 \leq \sigma$, then the minimizer \mathbf{c}^* of*

$$\min \|\mathbf{c}\|_1 \quad \text{subject to } \|\mathbf{Ac} - \mathbf{b}\|_2 \leq \sigma,$$

satisfies

$$\begin{aligned} \|\mathbf{c} - \mathbf{c}^*\|_2 &\leq C_1 \frac{\sigma_s(\mathbf{c})}{\sqrt{s}} + C_2 \sigma, \\ \|\mathbf{c} - \mathbf{c}^*\|_1 &\leq C_3 \sigma_s(\mathbf{c}) + C_4 \sqrt{s} \sigma. \end{aligned} \tag{2.8}$$

where constants C_1, C_2, C_3 and C_4 depend only on δ_{2s} , and $\sigma_s(\mathbf{c}) = \inf_{\mathbf{c}_s: \|\mathbf{c}_s\|_0 \leq s} \|\mathbf{c} - \mathbf{c}_s\|_1$ with $\|\mathbf{c}_s\|_0$ indicates the number of nonzero entries of \mathbf{c}_s . In particular, if \mathbf{c} is s -sparse, then the reconstruction is exact.

Proof. See Rauhut and Ward [52]. \square

A bounded orthonormal system has the following definition.

Definition 2.4. $\{\psi_n\}, n = 1, \dots, N$ is a bounded orthonormal system, if

$$K := \sup_n \|\psi_n\|_\infty = \sup_n \sup_{\mathbf{z}} |\psi_n(\mathbf{z})| < \infty, \tag{2.9}$$

where K is called the basis bound.

These definitions allow us to establish the recoverability of (2.6) based on the RIP.

Theorem 2.5 (RIP for bounded orthonormal systems). *Let $\mathbf{A} \in \mathbb{R}^{M \times N}$ be the interpolation matrix with entries $\{a_{j,n} = \psi_n(\mathbf{z}^{(j)})\}_{1 \leq n \leq N, 1 \leq j \leq M}$ (see (3.2)), where $\{\psi_n\}$ is a bounded orthonormal system satisfying (2.9). Assume that*

$$M \geq C \delta^{-2} K^2 s \log^3(s) \log(N),$$

then with probability at least $1 - N^{-\gamma \log^3(s)}$, the RIC δ_s of $1/\sqrt{M}\mathbf{A}$ satisfies $\delta_s \leq \delta$. Here, $C, \gamma > 0$ are universal constants.

Proof. See Rauhut and Ward [52]. \square

Theorem 2.3 and Theorem 2.5 establish the sparse recoverability of the bounded orthonormal systems.

1
2
3
4 **3. Methods**
5
6

In this section, we introduce the DSRAR framework to construct surrogate model. The goal of this study is to determine, given a small set of $M \ll N$ unstructured realizations $\{\mathbf{z}^{(i)}\}_{i=1}^M$ and the corresponding outputs $b = (f(\mathbf{z}^{(1)}), \dots, f(\mathbf{z}^{(M)}))^T$, the polynomial approximation in (2.4) or (2.5) when $f(\mathbf{z})$ has a sparse representation. This small set $\{\mathbf{z}^{(i)}\}_{i=1}^M$ is usually called *training set* and M is the *training sample size*. There are two quantities we need to compute: (i) an appropriate orthonormal polynomial basis ψ and (ii) an interpolation-type sparse solution $\mathbf{c} = (c_1, \dots, c_N)^T \in \mathbb{R}^N$ such that $f_p(\mathbf{z}^{(i)}) = f(\mathbf{z}^{(i)})$ for $i = 1, \dots, M$ with the smallest possible number nonzero \mathbf{c} . The basis construction, step (i), will be discussed in detail in Section 3.1. We can reformulate the second part as the following constrained optimization problem,

$$\min \|\mathbf{c}\|_0 \quad \text{subject to } \mathbf{A}\mathbf{c} = \mathbf{b}, \quad (3.1)$$

where $\|\mathbf{c}\|_0$ indicates the number of nonzero entries of \mathbf{c} and $\mathbf{A} \in \mathbb{R}^{M \times N}$ (usually called the measurement matrix) is written as

$$\mathbf{A} = (a_{ij})_{1 \leq i \leq M, 1 \leq j \leq N}, \quad a_{ij} = \psi_j(\mathbf{z}^{(i)}). \quad (3.2)$$

It is well known that this ℓ_0 minimization problem (3.1) is NP-hard [84]. As mentioned in Section 2.2, CS is a well-studied and popular approach to find sparse solutions to (3.1) through ℓ_1 -minimization shown in (2.6) (no noise) or (2.7) (with noise). Therefore, the approach introduced below can be viewed as a method for data-driven construction of bases that allow sparse representation and accurate recovery for QoIs in UQ applications.

36 **3.1. Data-driven construction of the amdP basis**
37

Let us start with a set of samples of d -dimensional random vector $\xi \in \mathbb{R}^d$, i.e., $S := \{\xi^{(k)}\}_{k=1}^{N_s}$ with the underlying probability measure $\rho(\xi)$. S is usually called the *sample set*. We aim to construct a set of orthonormal polynomial basis functions $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ with respect to $\rho(\xi)$ in Π_p^d , the space of polynomials up to degree p . Since $\rho(\xi)$ can be non-Gaussian or even unknown, we do not make the assumption that each component of ξ is mutually independent, even under a linear transformation such as those based on principal component analysis (PCA). Consequently, the orthogonal polynomial basis $\psi_\alpha(\xi)$ cannot be directly constructed as a tensor product of univariate orthonormal basis functions in each component of ξ . Below, we introduce a data-driven approach to construct multivariate amdP randomness.

51 **3.1.1. Orthonormal basis**
52

When we have a collection of random samples S , and the underlying probability measure $\rho(\xi)$ can be approximated by the discrete measure $\nu_S(\xi)$

$$\rho(\xi) \approx \nu_S(\xi) := \frac{1}{N_s} \sum_{\xi^{(k)} \in S} \delta_{\xi^{(k)}}(\xi), \quad (3.3)$$

where $\delta_{\xi^{(k)}}$ is the Dirac measure, that is $\delta_{\xi^{(k)}}(\xi)$ is equal to 1 when $\xi = \xi^{(k)}$ and 0 otherwise. Given the inner product defined as in (2.2) with ρ replaced by the discrete measure ν_S , we can construct a set of orthonormal multivariate polynomial basis functions $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ via the Gram-Schmidt orthogonalization process on an ordered monomial basis $\{\hat{\psi}_\alpha(\xi)\}_{|\alpha|=0}^p$. Here, we use the aforementioned graded lexicographic ordering of the multi-index.

Similar to Dunkl and Xu [76], ψ_α can be constructed using the recursive formulation

$$\psi_\alpha(\xi) = f_\alpha^\alpha \hat{\psi}_\alpha(\xi) - \sum_{\beta \prec \alpha} f_\beta^\alpha \psi_\beta(\xi), \quad (3.4)$$

where $\hat{\psi}_\alpha(\xi) := \prod_{k=1}^d \xi_k^{\alpha_k}$ represents the multivariate monomial basis function. The expression $\beta \prec \alpha$ means that the multi-index β comes before α under the chosen ordering. The coefficients f_β^α are determined by imposing an orthonormal condition with respect to the discrete measure ν_S , i.e.,

$$\begin{aligned} \int \psi_\alpha(\xi) \psi_\beta(\xi) d\rho(\xi) &\approx \int \psi_\alpha(\xi) \psi_\beta(\xi) d\nu_S(\xi) \\ &= \frac{1}{N_s} \sum_{k=1}^{N_s} \psi_\alpha(\xi^{(k)}) \psi_\beta(\xi^{(k)}) \\ &\equiv \delta_{\alpha,\beta}, \quad \beta \preceq \alpha. \end{aligned} \quad (3.5)$$

Equations (3.4) and (3.5) generate a set of orthonormal basis functions on the discrete measure ν_S irrespective of the mutual dependence between the components of ξ . We employ $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ as the amdP basis on $\rho(\xi)$.

Numerically, the modified Gram-Schmidt orthogonalization can be used as an alternative approach when the number of basis is too large and there exists instability in the standard Gram-Schmidt orthogonalization.

When $\rho(\xi)$ is known explicitly, orthonormal basis functions can also be constructed by taking the general formulation in Equation (3.4) and imposing the inner product in Equation (2.2) with respect to ρ . Here we will also consider a special case when ξ is a random vector that is linearly transformed from a random vector z with i.i.d. components via $\xi = Qz$. This case is motivated by the sparsity enhancement approach discussed in Sec. 3.2. In particular, we assume that the quadrature rule of polynomial integration with respect to the probability measure of z is explicitly known.

Given these assumptions, f_β^α can be determined by

$$\begin{aligned} \int \psi_\alpha(\xi) \psi_\beta(\xi) d\rho(\xi) &= \sum_{k=1}^{N_Q} \psi_\alpha(Qz_Q^{(k)}) \psi_\beta(Qz_Q^{(k)}) w_k \\ &\equiv \delta_{\alpha,\beta}, \quad \beta \preceq \alpha. \end{aligned} \quad (3.6)$$

where $\{z_Q^{(k)}\}_{k=1}^{N_Q}$ and $\{w_k\}_{k=1}^{N_Q}$ represent the quadrature points and weights constructed to yield an exact integration with probability measure of z for polynomials of degree $|\alpha| + |\beta|$ or less.

Algorithms 1 and 2 summarize the procedure of orthonormal basis construction when $\rho(\xi)$ is implicitly

represented by a sample set S and known explicitly, respectively. There is no unique system of orthogonal

1
2
3
4
5 ALGORITHM 1: Construct the orthonormal **amdP** basis $\{\psi_{\alpha}(\xi)\}_{|\alpha|=0}^p$ on discrete sample set S .

6
7 1: Given sample set $S = \{\xi^{(k)}\}_{k=1}^{N_s}$.
8 2: Given a fixed multi-index order $\{\alpha^{(l)}\}_{l=1}^N$.
9 3: **for** $l = 1$ to N **do**
10 4: Let $\alpha = \alpha^{(l)}$, construct $\psi_{\alpha}(\xi) = f_{\alpha}^{\alpha} \hat{\psi}_{\alpha}(\xi) - \sum_{\beta < \alpha} f_{\beta}^{\alpha} \psi_{\beta}(\xi)$ subject to Equation (3.5).
11 5: **end for**

12
13
14 ALGORITHM 2: Construct the orthonormal **amdP** basis $\{\psi_{\alpha}(\xi)\}_{|\alpha|=0}^p$ with probability measure $\rho(\alpha)$.

15
16 1: Given a multi-index order $\{\alpha^{(l)}\}_{l=1}^N$.
17 2: **for** $l = 1$ to N **do**
18 3: Let $\alpha = \alpha^{(l)}$, construct $\psi_{\alpha}(\xi) = f_{\alpha}^{\alpha} \hat{\psi}_{\alpha}(\xi) - \sum_{\beta < \alpha} f_{\beta}^{\alpha} \psi_{\beta}(\xi)$ by evaluating the basis inner product using
19 existing quadrature rule or Equation (3.6) if ξ can be linearly transformed from **a random vector with**
20 **i.i.d. components** z with an explicitly known quadrature rule.
21 4: **end for**

22
23 polynomial basis functions for both scenarios if $d > 1$; different orderings of α lead to different orthogonal
24 basis [76]. On the other hand, the constructed orthonormal basis is unique up to unitary transformations as
25 we prove in Theorem 3.1.

26
27 **Theorem 3.1.** Let $\{\psi_{\alpha}(\xi)\}_{|\alpha|=0}^p$ be a set of orthonormal polynomial basis with respect to the measure $\rho(\xi)$,
28 $\xi \in \mathbb{R}^d$. Denote by $\Psi(\xi)$ the polynomial basis vector

29
30
31
32
33
34
35
36
37
38
39
40
$$\Psi(\xi) := (\psi_{\alpha^{(1)}}, \dots, \psi_{\alpha^{(N)}})^T, \quad (3.7)$$

41 where $\alpha^{(1)}, \dots, \alpha^{(N)}$ is the arrangement of multi-index α according to a fixed multi-index order. Let $\chi = Q\xi$,
42 where $Q \in \mathbb{R}^{d \times d}$ is invertible. Let $\{\phi_{\beta}(\chi)\}_{|\beta|=0}^p$ be a set of orthonormal polynomial basis functions with
43 respect to a measure $\rho'(\chi)$ constructed with order $\beta^{(1)}, \dots, \beta^{(N)}$, where $\rho'(\chi)$ is induced from $\rho(\xi)$. Then
44 there exists a unitary matrix U such that $\Phi(\chi) = U\Psi(\xi)$, where $\Phi(\chi) := (\phi_{\beta^{(1)}}, \dots, \phi_{\beta^{(N)}})^T$ denotes the
45 corresponding polynomial basis vector.

46
47
48
49
50 *Proof.* Let $\hat{\Psi}(\xi)$ be the monomial basis vector. Note that $\{\psi_{\alpha}(\xi)\}_{|\alpha|=0}^p$ and $\{\phi_{\beta}(\xi)\}_{|\beta|=0}^p$ are two sets of
51 basis in Π_p^d . There exists transfer matrices M_{ψ} and $M_{\phi} \in \mathbb{R}^{N \times N}$ such that

52
53
54
55
56
57
58
59
60
61
62
63
64
65
$$\Psi(\xi) = M_{\psi} \hat{\Psi}(\xi), \quad \Phi(\chi) = M_{\phi} \hat{\Psi}(\xi).$$

With $\chi = Q\xi$, $\Phi(Q\xi)$ is also a basis in Π_p^d . Then there exists an invertible matrix $T \in \mathbb{R}^{N \times N}$ such that

66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1300
1301
1302
1303
1304
130

which gives $\Phi(\chi) = \mathcal{U}\Psi(\xi)$, where $\mathcal{U} = \mathcal{T}\mathcal{M}_\psi^{-1}$. Recall $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ and $\{\phi_\beta(\chi)\}_{|\beta|=0}^p$ are orthonormal basis with respect to $\rho(\xi)$ and $\rho'(\chi)$, we have

$$\mathbf{I} = \int \Phi(\chi)\Phi(\chi)^T d\rho'(\chi) = \int \mathcal{U}\Psi(\xi)\Psi(\xi)^T \mathcal{U}^T d\rho(\xi) = \mathcal{U}\mathcal{U}^T \quad (3.8)$$

□

We do not need further assumptions on $\rho(\xi)$ because Theorem 3.1 holds both when $\rho(\xi)$ is a measure on the continuous random vector ξ (with probability density function $\omega(\xi)$) or a discrete measure $\nu_S(\xi)$ on a sample set S . Furthermore, it is straightforward to show the following Corollary.

Corollary 1. Let $S_1 := \{\xi^{(k)}\}_{k=1}^M$ and $S_2 := \{\chi^{(k)}\}_{k=1}^M$ be two sets of random sampling points where $\chi^{(k)} = \mathbf{Q}\xi^{(k)}$ with invertible \mathbf{Q} . Let \mathbf{G}_ξ and \mathbf{G}_χ be the Gram matrix constructed by $\Psi(\xi)$ and $\Phi(\chi)$ defined in Theorem 3.1, i.e., $\mathbf{G}_\xi := \sum_{k=1}^M \Psi(\xi^{(k)})\Psi(\xi^{(k)})^T/M$ and $\mathbf{G}_\chi := \sum_{k=1}^M \Phi(\chi^{(k)})\Phi(\chi^{(k)})^T/M$. Then \mathbf{G}_\cdot has invariant l_2 norm, that is, $\|\mathbf{G}_\xi\|_2 = \|\mathbf{G}_\chi\|_2$. Moreover, $\|\mathbf{G}_\xi - \mathbf{I}\|_2 = \|\mathbf{G}_\chi - \mathbf{I}\|_2$.

In general, the l_2 norm of $\|\mathbf{G}_\xi - \mathbf{I}\|_2$ is independent of specific monomial order of α and invariant under linear transformations of the random vector. The basis functions $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ constructed by Equations (3.4) and (3.5) provide an appropriate candidate for representing the surrogate model $f(\xi)$ via CS.

3.1.2. Near-orthonormal basis

When $\rho(\xi)$ is implicitly represented by a sample set S , we employ the discrete measure ν_S to construct $\{\psi_\alpha(\xi)\}_{|\beta|=0}^p$. However, we note that the training set that queries $f(\cdot)$, denoted by S_f , may not be a subset of S . In practice, the sample set S and the training set S_f are usually collected in a sequential manner or directly from different experiments, although individual sampling points of both S and S_f follow the same distribution. Since S only contains a finite number of samples of ξ , basis $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ constructed by (3.4) and (3.5) is not the “exact orthonormal” basis with respect to $\rho(\xi)$. Especially, let $S' = \{\xi'_k\}_{k=1}^{N_s}$ be another sample set following the same distribution $\rho(\cdot)$ and $\nu_{S'}(\cdot)$ be the discrete measure defined on S' . For the orthonormal **amdP** basis functions $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ constructed on S , we have $\mathbb{E}[\Psi(\xi)\Psi(\xi)^T] \neq \mathbf{I}$ under the discrete measure $\nu_{S'}(\xi)$ and vice versa.

The above observation forces us to re-examine the orthonormal condition imposed by (3.5). Since the pre-constructed basis $\psi_\alpha(\xi)$ does not retain the exact orthonormal condition when later being applied to approximate $f(\xi)$, we may relax the condition when determining the coefficients f_β^α in (3.4). In the present study, we propose the following heuristic criterion

$$\arg \min_{\hat{\mathbf{f}}^\alpha} \|\hat{\mathbf{f}}^\alpha\|_2 \text{ subject to } \left| \int \psi_\alpha(\xi)\psi_\beta(\xi) d\nu_S(\xi) - \delta_{\alpha,\beta} \right| < \zeta_{\alpha,\beta}, \quad \beta \leq \alpha, \quad (3.9)$$

4 where $\hat{\mathbf{f}}^\alpha$ is the coefficient vector of ψ_α when represented using monomial basis functions, i.e., $\psi_\alpha(\xi) =$
 5 $\sum_{\beta \leq \alpha} \hat{f}_\beta^\alpha \hat{\psi}_\beta(\xi)$. $\hat{\mathbf{f}}^\alpha$ is related to \mathbf{f}^α through the linear transformation
 6
 7

$$8 \quad 9 \quad 10 \quad 11 \quad 12 \quad 13 \quad 14 \quad 15 \quad 16 \quad 17 \quad 18 \quad 19 \quad 20 \quad 21 \quad 22 \quad 23 \quad 24 \quad 25 \quad 26 \quad 27 \quad 28 \quad 29 \quad 30 \quad 31 \quad 32 \quad 33 \quad 34 \quad 35 \quad 36 \quad 37 \quad 38 \quad 39 \quad 40 \quad 41 \quad 42 \quad 43 \quad 44 \quad 45 \quad 46 \quad 47 \quad 48 \quad 49 \quad 50 \quad 51 \quad 52 \quad 53 \quad 54 \quad 55 \quad 56 \quad 57 \quad 58 \quad 59 \quad 60 \quad 61 \quad 62 \quad 63 \quad 64 \quad 65$$

$$\hat{\mathbf{f}}^\alpha = \begin{pmatrix} \mathbf{F} & 0 \\ 0 & 1 \end{pmatrix} \mathbf{f}^\alpha, \quad (3.10)$$

where \mathbf{F} is an upper triangle matrix determined by pre-computed $\hat{\mathbf{f}}^\beta, \beta \prec \alpha$, i.e.,

$$\mathbf{F}_{I_\beta' I_\beta} = \begin{cases} \hat{f}_{\beta'}^\beta & \beta' \preceq \beta \\ 0 & \beta' \succ \beta, \end{cases} \quad (3.11)$$

240 where I_β represents the mapping from multi-index to single index.

The parameter $\zeta_{\alpha, \beta}$ quantifies the relaxation of the orthonormal condition. We split the sample set S equally into two parts $S := S_1 \cup S_2$. Denote $\{\psi_\alpha^{(1)}(\xi)\}_{|\alpha|=0}^p$ and $\{\psi_\alpha^{(2)}(\xi)\}_{|\alpha|=0}^p$ the orthonormal bases constructed by Equations (3.4) and (3.5) on the discrete measures $\nu_{S_1}(\xi)$ and $\nu_{S_2}(\xi)$, respectively. Inspired by cross-validation, we have chosen $\zeta_{\alpha, \beta} = \frac{|\zeta_1| + |\zeta_2|}{2\sqrt{2}}$

$$\zeta_1 = \int \psi_\alpha^{(1)}(\xi) \psi_\beta^{(1)}(\xi) d\nu_{S_2}(\xi), \quad \zeta_2 = \int \psi_\alpha^{(2)}(\xi) \psi_\beta^{(2)}(\xi) d\nu_{S_1}(\xi). \quad (3.12)$$

33 ALGORITHM 3: Construct the near-orthonormal **amdP** basis $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ on discrete sample set S .

35 1: Collect samples of ξ from sample set $S = \{\xi^{(k)}\}_{k=1}^{N_s}$, split S equally into two disjoint subsets, i.e.,
 36 $S = S_1 \cup S_2$, $S_1 \cap S_2 = \emptyset$.
 37 2: Given fixed monomial index order $\{\alpha^{(l)}\}_{l=1}^N$, construct the orthonormal **amdP** basis $\{\psi_\alpha^{(1)}(\xi)\}_{|\alpha|=0}^p$ and
 38 $\{\psi_\alpha^{(2)}(\xi)\}_{|\alpha|=0}^p$ on set S_1 and S_2 by **Algorithm 1**.
 39 3: **for** $l = 1$ to N **do**
 40 4: Let $\alpha = \alpha^{(l)}$, construct $\psi_\alpha(\xi) = f_\alpha^\alpha \hat{\psi}_\alpha(\xi) - \sum_{\beta \prec \alpha} f_\beta^\alpha \hat{\psi}_\beta(\xi)$ on by Equations (3.9), (3.10), and (3.12).
 41 5: **end for**

47 Algorithm 3 describes construction for a set of near-orthonormal **amdP** basis functions on the sample
 48 set S . When applied to the sample set S' to approximate $f(\xi)$, the basis shows comparable orthonormal
 49 conditions with the basis constructed by (3.5). Such results can be partially understood by the theoretical
 50 bound from Theorem 2.5 on the number of samples M for exact recovery in orthonormal polynomial systems,
 51 $M \geq C_1 K^2 s \log^3(s) \log(N)$, where $s = \|c\|_0$ and $K = \sup_\alpha \|\psi_\alpha\|_\infty$. Theoretical analysis of the recovery error
 52 under different basis functions is out of the scope of the present work and is left for future investigation.
 53 However, we note that the accuracy of the surrogate model $f(\xi)$ can be further improved by enhancing
 54
 55
 56
 57
 58

the sparsity of \mathbf{c} . This can be achieved through the ideas presented in our previous work [1] which will be extended to general distributions below.

Remark 3.2. We emphasize that (3.9) provides a heuristic approach to construct the near-orthonormal **amdP** basis functions $\psi_{\alpha}(\xi)$ with a smaller basis bound. In practice, (3.9) can be further relaxed to

$$\begin{aligned} \arg \min_{\hat{\mathbf{f}}^{\alpha}} \|\hat{\mathbf{f}}^{\alpha}\|_2 \text{ subject to } & \sum_{|\beta|=r, \beta<\alpha} \left| \int \psi_{\alpha}(\xi) \psi_{\beta}(\xi) d\nu_S(\xi) \right|^2 < \sum_{|\beta|=r, \beta<\alpha} \zeta_{\alpha, \beta}^2, \\ & \left| \int \psi_{\alpha}(\xi) \psi_{\alpha}(\xi) d\nu_S(\xi) - 1 \right| < \zeta_{\alpha, \alpha}, \quad r = 0, \dots, |\alpha|, \end{aligned} \quad (3.13)$$

which shows similar numerical performance. There is no theoretical guarantee yet that Equations (3.9) and (3.13) yield a smaller basis bound than (3.5) on S_f , S or the entire domain of ξ . We numerically compare some properties of different bases in Section 4.1, which illustrate the performance of the near-orthonormal **amdP** basis constructed above. There may exist other numerical approaches to optimize $\psi_{\alpha}(\xi)$ that can lead to an even smaller basis bound. We also note that the threshold $\zeta_{\alpha, \beta}$ is determined by directly splitting S into two disjoint sets. In practice, it is possible to design more sophisticated strategies to optimize the choice of $\zeta_{\alpha, \beta}$ and the basis construction procedure. We leave such studies for future work.

3.2. Sparsity enhancement

For the linear system in (2.6), the numerical accuracy of the recovered $\tilde{\mathbf{c}}$ via l_1 -minimization depends on the sparsity of \mathbf{c} . This dependence motivates us to develop a numerical approach to further enhance the sparsity of \mathbf{c} through the variability analysis of $f(\xi)$ [1]. If we know $f(\xi)$ explicitly, the (sorted) directions of variance in $f(\xi)$ under the distribution of ξ can be found based on the active subspace method [85, 66]. In particular, we define the gradient matrix \mathbf{G} by

$$\mathbf{G} = \mathbb{E} \left[\nabla f(\xi) \nabla f(\xi)^T \right] \quad (3.14)$$

where $\nabla f(\xi)$ is the gradient vector defined by $\nabla f(\xi) = \left(\frac{\partial f}{\partial \xi_1}, \frac{\partial f}{\partial \xi_2}, \dots, \frac{\partial f}{\partial \xi_d} \right)^T$. Eigendecomposition of \mathbf{G} ,

$$\mathbf{G} = \mathbf{Q} \mathbf{K} \mathbf{Q}^T, \quad \mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_d], \quad (3.15a)$$

$$\mathbf{K} = \text{diag}(k_1, \dots, k_d), \quad k_1 \geq \dots \geq k_d \geq 0, \quad (3.15b)$$

yields the sorted variability directions $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_d$. Accordingly, we may define a new random vector χ following the sorted variability directions via linear transformation

$$\chi = \mathbf{Q}^T \xi. \quad (3.16)$$

$f(\xi) = f((\mathbf{Q}^T)^{-1} \chi) = f(\mathbf{Q} \chi)$ can be approximated by expansion in an orthonormal polynomial basis χ with a coefficient vector \mathbf{c} which is sparser than the $f(\xi)$ being expanded by orthonormal basis of ξ . For the

1
2
3
4 remainder of this paper, we use \mathbf{Q} to denote the rotation matrix to transform ξ to χ and $g(\chi)$ to represent
5 $f(\mathbf{Q}\chi)$.
6

7 In practice, $f(\xi)$ is usually not explicitly known. We may numerically approximate \mathbf{G} by
8

9
$$\mathbf{G} \approx \mathbb{E} \left[\nabla \tilde{f}(\xi) \nabla \tilde{f}(\xi)^T \right], \quad (3.17)$$

10
11

12 where $\tilde{f}(\xi)$ represents the approximation of $f(\xi)$ by the orthonormal polynomial basis functions $\psi_\alpha(\xi)$
13 as proposed in [1] or obtained via solving (2.6) with the data-driven basis approach (i.e., basis functions
14 constructed with respect to an arbitrary measure) described in Section 3.1. In particular, if ξ is a random
15 vector with i.i.d. Gaussian components, χ is also a random vector with i.i.d. Gaussian components. Thus,
16 $\tilde{f}(\xi)$ and $\tilde{g}(\chi) := \tilde{f}(\mathbf{Q}\chi)$ can be represented by the orthonormal basis functions of the same form, e.g.,
17 tensor products of univariate Hermite polynomials. Without loss of generality, from now on, we use $\tilde{g}(\chi)$
18 to represent $\tilde{f}(\mathbf{Q}\chi)$.
19

20 However, if $\rho(\xi)$ is not i.i.d. Gaussian, χ and ξ do not generally have the same distribution. Therefore,
21

22 an orthonormal polynomial basis $\psi(\cdot)$ with respect to ξ cannot be directly applied to χ . The general approach
23 presented in Section 3.1 enables us to construct the **amdP** basis with respect to the probability measure of
24 the rotated vector χ . The two orthonormal bases associated with ξ and χ respectively are related to each
25 other via a unitary transformation as shown in Theorem 3.1. In particular, if $\rho(\xi)$ is implicitly described by
26 a sample set $S = \{\xi^{(k)}\}_{k=1}^{N_s}$, \mathbf{G} can be easily evaluated by representing $\psi_\alpha(\xi)$ via the monomial basis, i.e.,
27 $\psi_\alpha(\xi) = \sum_{\beta \leq \alpha} \hat{f}_\beta^\alpha \hat{\psi}_\beta(\xi)$ via Equation (3.10) and then integrating with discrete measure ν_S . By transforming
28 S and S_f into $\{\chi^{(k)}\}_{k=1}^{N_s}$ and $\{\chi'^{(k)}\}_{k=1}^M$, the orthonormal and near-orthonormal **amdP** basis functions with
29 respect to χ can be constructed by Eqs. (3.5) (3.13). The surrogate model $\tilde{g}(\chi)$ can then be constructed by
30 solving (2.6).
31

32 The entire DSRAR procedure is presented in **Algorithm 4**. Compared with $\tilde{f}(\xi)$, $\tilde{g}(\chi)$ shows smaller
33 numerical error in general. The additional cost of sparsity enhancement procedure in Step 4 - 6 is less than
34 0.6 CPU (3.7 GHz Quad-Core Intel Xeon E5) hour for the numerical examples considered in this study. For
35 realistic applications, the overhead of Step 4 - 6 could be relatively small if sampling of QoI is expensive or
36 the available training set is limited.
37

38 The DSRAR framework described above is also applicable to systems with standard density distributions,
39

40 where $\rho(\xi)$ is known explicitly. Without loss of generality, we assume that an orthonormal polynomial basis
41 $\{\psi_\alpha(\xi)\}_{|\alpha|=0}^p$ is known. Evaluation of \mathbf{G} by (3.17) on $\rho(\xi)$ is straightforward. The surrogate model of f can
42 be constructed via l_1 minimization with enhanced sparsity through Algorithm 5.
43

44 The procedures for random vector rotation and surrogate construction presented in Algorithms 4 and
45 5 can be conducted in an iterative manner. We have investigated this issue [86] by applying a previously
46 developed rotation procedure [1] successively to systems with underlying Gaussian distributions. For the
47

ALGORITHM 4: DSRAR: Surrogate model construction with discrete sample set S and training set S_f .

- 1: Collect the sample set within the random space $S = \{\xi^{(k)}\}_{k=1}^{N_s}$.
 - 2: Generate evaluations of f on training set $S_f = \{\xi'^{(k)}\}_{k=1}^M$ with M outputs f_1, f_2, \dots, f_M .
 - 3: Construct the data-driven **amdP** basis $\{\psi_i(\xi)\}_{i=1}^N$ on discrete measure $\nu_S(\xi)$ as the exact orthonormal basis by Algorithm 1 or the near orthonormal basis by Algorithm 3.
 - 4: Evaluate the measurement matrix $A_{ij} = \psi_j(\xi'^{(i)})$, $1 \leq i \leq M$, $1 \leq j \leq N$; construct surrogate model $\tilde{f}(\xi) = \sum_{|\alpha|=0}^p c_\alpha \psi_\alpha(\xi)$ by solving the l_1 -minimization problem.
 - 5: Evaluate the gradient matrix $\mathbf{G} \approx \mathbb{E} [\nabla \tilde{f}(\xi) \nabla \tilde{f}(\xi)^T]$ on $\nu_S(\xi)$. Find the eigendecomposition $\mathbf{G} = \mathbf{Q} \mathbf{K} \mathbf{Q}^T$, define sample set $\{\chi^{(k)}\}_{k=1}^{N_s}$ and training set $\{\chi'^{(k)}\}_{k=1}^M$ by $\chi^{(k)} = \mathbf{Q}^T \xi^{(k)}$, $\chi'^{(k)} = \mathbf{Q}^T \xi'^{(k)}$.
 - 6: Reconstruct the data-driven **amdP** basis $\{\phi_\alpha(\chi)\}_{|\alpha|=0}^p$ by Algorithm 3 and surrogate model $\tilde{g}(\chi)$ with enhanced sparsity following Step 3 and Step 4.
-

ALGORITHM 5: DSRAR: Surrogate model construction with training set S_f and probability measure $\rho(\xi)$.

- 1: Evaluate f on training set $S_f = \{\xi'^{(k)}\}_{k=1}^M$ with M outputs f_1, f_2, \dots, f_M .
 - 2: Evaluate the measurement matrix $A_{ij} = \psi_j(\xi'^{(i)})$, $1 \leq i \leq M$, $1 \leq j \leq N$; construct surrogate model $\tilde{f}(\xi) = \sum_{|\alpha|=0}^p c_\alpha \psi_\alpha(\xi)$ by solving l_1 minimization problem.
 - 3: Evaluate the gradient matrix $\mathbf{G} = \mathbb{E} [\nabla \tilde{f}(\xi) \nabla \tilde{f}(\xi)^T]$ on $\rho(\xi)$. Conduct eigendecomposition $\mathbf{G} = \mathbf{Q} \mathbf{K} \mathbf{Q}^T$ and define training set $\{\chi'^{(k)}\}_{k=1}^M$, $\chi'^{(k)} = \mathbf{Q}^T \xi'^{(k)}$.
 - 4: Re-construct the orthonormal **amdP** basis $\{\phi_\alpha(\chi)\}_{|\alpha|=0}^p$ with respect to $\rho'(\chi)$ by **Algorithm 2**. Construct the surrogate model $\tilde{g}(\chi)$ with enhanced sparsity following Step 3.
-

systems studied in the present work, the improvement of the numerical accuracy is marginal after the first rotation procedure. Therefore, the numerical results with only one rotation procedure will be presented in this manuscript.

4. Results

This section presents the numerical results of the present DSRAR framework for surrogate model construction with arbitrary underlying distributions. For numerical examples where the probability measure $\rho(\xi)$ (with density function $\omega(\xi)$) is not known explicitly and is represented by a discrete data set $S = \{\xi^{(k)}\}_{k=1}^{N_s}$, we split S equally into two subsets $S = S_1 \cup S_2$. We use S_1 to construct the data-driven amdP basis and split S_2 into two disjoint subset $S_2 = S_{2,1} \cup S_{2,2}$, where $S_{2,1}$ is the training set for surrogate model construction and $S_{2,2}$ is the test set to evaluate the accuracy of the constructed surrogate model. The size of the training set is $O(10^2) - O(10^3)$ and size of the test set is $O(10^5)$.

4.1. Accurate recovery of linear systems with data-driven bases

In this test, we collected a sample set $S = \{\xi^{(k)}\}_{k=1}^{N_s}$ with $N_s = 2 \times 10^5$. The random vector ξ followed the Gaussian mixture distribution

$$\omega(\xi) = \sum_{i=1}^{N_m} a_i \mathcal{N}(\mu_i, \Sigma_i) \quad (4.1)$$

where N_m is the number of Gaussian modes. We set $N_m = 3$, $a_i > 0$ for $i = 1, 2, 3$ and $\sum_{i=1}^3 a_i = 1$. For each Gaussian mode, μ_i is a 25-dimensional i.i.d. random vector with uniform distribution $\mathcal{U}[-2.5, 2.5]$ on each dimension and then shifted such that $\sum_{i=1}^3 a_i \mu_i = 0$. The matrices Σ_i were chosen such that

$$\Sigma_i = (\Upsilon_i \Upsilon_1^T + \mathbf{I})/4, \quad (4.2)$$

where Υ_i is a random matrix with i.i.d. entries from $\mathcal{U}[0, 1]$ for $i = 1, 2, 3$.

We considered a linear system

$$A\mathbf{c} = \mathbf{b} + \boldsymbol{\varepsilon}$$

and recovered \mathbf{c} using M training points by solving the l_1 minimization problem defined by (2.6) where

$$[A]_{i,j} = \psi_j(\xi^{(i)}), \quad b_i = \sum_{k=1}^N c_k \psi_k(\xi^{(i)}), \quad (4.3)$$

with $1 \leq i \leq M$, $1 \leq j \leq N$, and $\boldsymbol{\varepsilon}$ is noise with $\|\boldsymbol{\varepsilon}\|_2 \leq 10^{-7}$. We set $d = 25$, $p = 2$ and $N = \begin{pmatrix} d+p \\ p \end{pmatrix} = 351$. The basis functions $\psi_\alpha(\xi)$ were constructed on the set S_1 by the following approaches:

1. the orthonormal amdP basis subject to Equations (3.4) and (3.5);
2. the near-orthonormal amdP basis subject to Equation (3.9);

4 3. tensor product of univariate normalized Legendre polynomials (both sampling points and training
 5 points are scaled to lie in $[-1, 1]$ on each dimension accordingly).

8 310 Training points from set S_2 were used to examine the recovery accuracy of \mathbf{c} .

10 4.1.1. Sparse linear systems

12 First, we considered the scenario where \mathbf{c} is a s -sparse vector and employed the following theoretical
 13 bound to examine the recovery accuracy via l_1 -minimization.

16 **Theorem 4.1.** *Given a matrix $\Psi \in \mathbb{R}^{M \times N}$ and set T_α with $s = |T_\alpha|$, a s -sparse vector \mathbf{c} with non-zero
 17 entries on T_α can be exactly recovered via l_1 -minimization if $\frac{\theta_s}{1-\delta_s} < 0.5$, where δ_s and θ_s are defined by*

$$20 \quad \begin{aligned} \delta_s &:= \inf [\delta : (1-\delta)\|\mathbf{y}\|_2^2 \leq \|\Psi_t \mathbf{y}\|_2^2 \leq (1+\delta)\|\mathbf{y}\|_2^2], \forall t \subseteq \mathbf{T}, \forall \mathbf{y} \in \mathbb{R}^{|t|} \\ 21 \quad \theta_s &:= \inf [\theta : |\langle \Psi_{t'}, \Psi_t \mathbf{y}' \rangle| \leq \theta \|\mathbf{y}'\|_2 \|\mathbf{y}\|_2], \forall t \subseteq \mathbf{T}, t' \not\subseteq \mathbf{T}, |t'| \leq s, \forall \mathbf{y} \in \mathbb{R}^{|t|}, \mathbf{y}' \in \mathbb{R}^{|t'|} \end{aligned} \quad (4.4)$$

24 where Ψ_t and $\Psi_{t'}$ denote the sub-matrices of Ψ with column indices in t and t' respectively.

26 315 Theorem 4.1 (see Appendix A for proof) provides a sufficient condition to exactly recover \mathbf{c} with non-
 27 zero entries on index set T_α . For numerical study, we randomly chose an index set T_α from Λ_p^d with
 28 $|T_\alpha| = 3$, where Λ_p^d is defined by (2.4). For each training set, we constructed the measurement matrix \mathbf{A}
 29 with different bases and computed $\theta_s / (1 - \delta_s)$ by (4.4). Figure 1(a) shows the mean value $\mathbb{E}[\theta_s / (1 - \delta_s)]$
 30 on 200 independent training sets chosen from S_2 for each M . The exact and near-orthonormal bases yield
 31 similar results: $\mathbb{E}[\theta_s / (1 - \delta_s)]$ becomes smaller than 0.5 as M approaches 210, which is also shown in the
 32 inset plot of Figure 1(a). In contrast, $\mathbb{E}[\theta_s / (1 - \delta_s)]$ obtained from Legendre polynomial basis shows worse
 33 performance due to the loss of orthonormality.

39 In our numerical experiments, we were able to recover \mathbf{c} using fewer samples than the number M —as
 40 suggested by the sufficient condition (Theorem 2.5) originally given by Rauhut [52]—since this number is
 41 based on the worst case scenario and is not, in general, a sharp bound. Figure 1(b) shows the numerical
 42 results of a test case with $c_{T_\alpha} = 1$, $c_{T_\alpha^c} = 0$, $|T_\alpha| = 5$. For each M , 200 CS implementations were conducted
 43 to compute the average of the relative error $\|\mathbf{c} - \tilde{\mathbf{c}}\|_1 / \|\mathbf{c}\|_1$. The exact and near-orthonormal **amdP** bases
 44 show similar performance, where \mathbf{c} can be accurately recovered (up to $\|\varepsilon\|_2$) using $M = 45$ training points.
 45 In contrast, the Legendre basis yields larger relative error in ℓ_1 -norm. The relative error of the recovered
 46 coefficients from one CS implementation with Legendre basis is shown in the inset plot of Figure 1(b).

52 4.1.2. Non-sparse linear systems

54 We also tested the recovery performance when the exact representation is not sparse. The vector \mathbf{c} is
 55 chosen with a random non-zero index set T_α with $|T_\alpha| = 120$. Individual components of \mathbf{c}_{T_α} are i.i.d. log-
 56 normal, such that $\log c_{T_\alpha} \sim \mathcal{N}(0, 2)$. For each size (M) of the training set, 200 CS implementations were

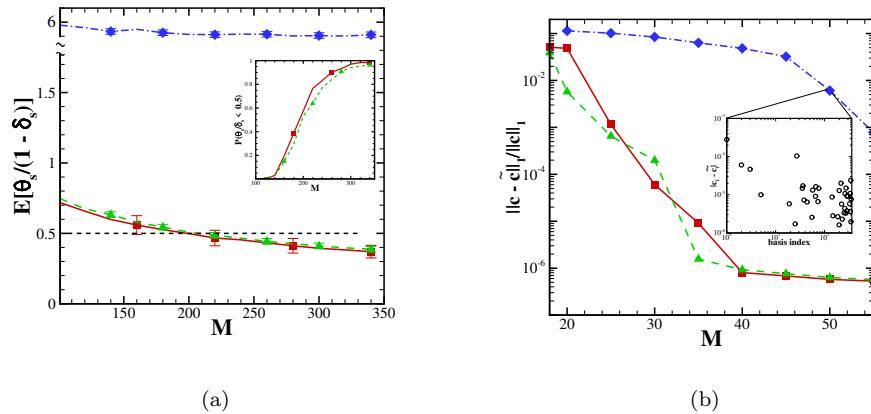


Figure 1: The measurement matrices constructed by the exact and near-orthonormal bases exhibit similar performance in the theoretical (sufficient) bound and numerical results for recovery of sparse vector. Both bases outperform the Legendre basis. “—■—”: the exact orthonormal `amdP` basis; “—▲—”: the near-orthonormal `amdP` basis; “—◆—”: Legendre basis. (a) Mean value of the theoretical bound $\mathbb{E}[\theta_s / (1 - \delta_s)]$ of exact recovery for measurement matrices \mathbf{A} constructed by various bases for the chosen non-zero index T_α with $s = 3$. The error bar represents the standard deviation. The inset plot shows the theoretical prediction of the exact recovery probability. (b) Relative l_1 error of the recovered sparse vector ($s = 5$) using different training set size M . The inset plot shows the recovery error $\|\mathbf{c} - \tilde{\mathbf{c}}\|_1$ of one training set for the Legendre basis system.

conducted to compute the average of the numerical error $\|\mathbf{c} - \tilde{\mathbf{c}}\|_2$, as shown in Figure 2(a). Similar to the previous example, the Legendre basis exhibits the largest approximation error. The near-orthonormal basis shows smaller error than the exact orthonormal basis.

We also computed the density distribution of individual component $|\mathbf{c}_{i'} - \tilde{\mathbf{c}}_{i'}|$, where i' refers to single index sorted by the magnitude in descending order. Figure 2(b-d) shows that, compared with the exact orthonormal basis and the Legendre basis, the distribution of $\log |\mathbf{c}_{i'} - \tilde{\mathbf{c}}_{i'}|$ obtained from near-orthonormal basis is biased toward the smallest magnitudes for error of individual i' . This result can be interpreted as that the average of $\|\mathbf{c} - \tilde{\mathbf{c}}\|_2$ of the near orthogonal basis is smaller than that of the exact orthogonal basis and also outperforms the Legendre basis.

4.2. Systems with explicit knowledge of density function

In this subsection, we demonstrate the proposed method in systems with common non-Gaussian randomness with analytical density function $\omega(\xi)$. We show that the present method based on orthonormal basis construction and rotation of the random variables exploits the sparser representation of QoI while retaining proper orthogonality with respect to rotated variables. Therefore, it yields more accurate surrogate models than other approaches based on the direct recovery of c without the sparsity enhancement rotation procedure and/or directly applying the rotation procedure without reconstruction of the orthonormal **amdP** basis.

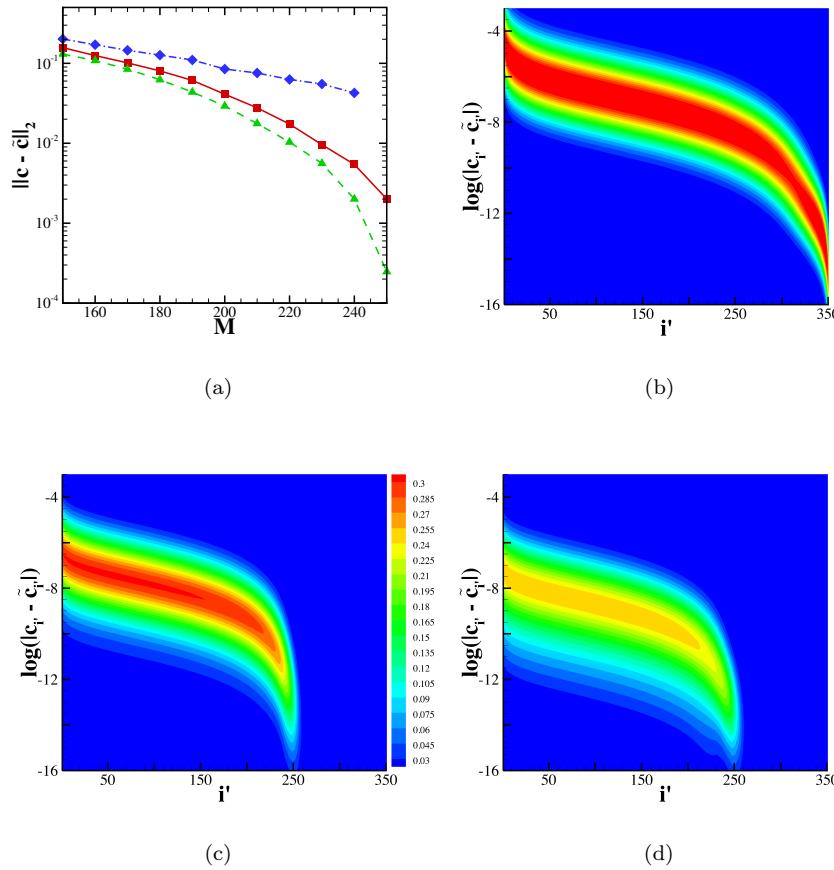


Figure 2: The measurement matrices constructed by different bases show different numerical performance for the recovery of non-sparse vector. The near-orthonormal basis shows the most accurate result. (a) l_1 error of the recovered vector \mathbf{c} with different bases. “ $\text{---} \blacksquare \text{---}$ ”: the exact orthonormal `amdP` basis; “ $\text{---} \blacktriangle \text{---}$ ”: the near-orthonormal `amdP` basis; “ $\text{---} \blacklozenge \text{---}$ ”: Legendre basis. (b-d) Contours of $|\mathbf{c}_{\alpha'} - \tilde{\mathbf{c}}_{\alpha'}|$ (sorted by magnitude) from training sets of size $M = 230$ with Legendre (top right), orthonormal (bottom left) and near-orthonormal bases (bottom right).

1
2
3
4 4.2.1. High-dimensional polynomial
5
6

For the first numerical example, we consider a high-dimensional polynomial function

$$f(\xi) = \sum_{|\alpha| \leq 3} \hat{c}_\alpha \hat{\psi}_\alpha(\xi) = \sum_{i=1}^N \frac{\eta_i}{|i|^{1.5}} \hat{\psi}_i(\xi), \quad (4.5)$$

where $\hat{\psi}_\alpha$ and $\hat{\psi}_i$ represent monomial basis functions, η_i represents uniform random variables $\mathcal{U}[0, 1]$. We employed this polynomial function with sparse coefficients as a benchmark problem to examine the recovery accuracy of the present method. ξ is a random vector consisting of 20 i.i.d. random variables. The density function of the i -th variable ξ_i is given by

$$\omega(\xi_i) = e^{-\xi_i}, \quad (4.6)$$

where the corresponding orthonormal basis are given by the Laguerre polynomials. Accordingly, we construct a 3rd-order polynomial expansion $\tilde{f}(\xi)$ with $N = 1771$ multivariate basis functions, which are the tensor product of the univariate Laguerre polynomials. Figure 3 shows the relative l_2 error of \tilde{f} computed by level 4 sparse grid integration. Similar to the previous example, the probability density function (PDF) of χ does not retain the form $\omega'(\chi) = \prod_{i=1}^d \exp(-\chi_i)$ after the rotation. Iteratively employing the multivariate Laguerre polynomials to represent $\tilde{g}(\chi)$ may result in erroneous prediction (the red dash-dotted curve). Alternatively, such a problem can be addressed by using the reconstructed orthonormal **amdP** basis with respect to χ , which yields a smaller error than $\tilde{f}(\xi)$ (the blue dashed curve).

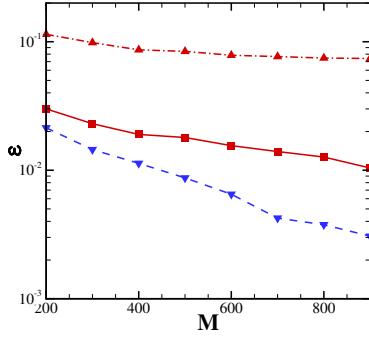


Figure 3: Sparsity-enhancing rotation with the reconstructed orthonormal **amdP** basis yields the most accurate recovery of a high-dimensional polynomial function of random vectors following density function given by Equation (4.6). Directly applying the rotation procedure without reconstructing the orthonormal basis yields erroneous prediction. “—■—”: Laguerre polynomial basis with respect to ξ ; “-·-▲-·-”: Laguerre polynomial basis with respect to rotated vector χ ; “-·-▼-·-”: the reconstructed **amdP** orthonormal basis with respect to rotated vector χ .

1
2
3
4 360 4.2.2. One-dimensional elliptic PDEs with high-dimensional random inputs
5
6

We applied the proposed method to model the solution to a one-dimensional (1D) elliptic PDE with high dimensional random input

$$-\frac{d}{dx} \left(D(x; \boldsymbol{\xi}) \frac{du(x; \boldsymbol{\xi})}{dx} \right) = 1, \quad x \in (0, 1) \quad (4.7)$$

$$u(0) = u(1) = 0,$$

where $a(x; \boldsymbol{\xi}) := \log D(x; \boldsymbol{\xi})$ is the stochastic input and $a(x; \boldsymbol{\xi})$ was a stationary process with correlation function

$$K(x, x') = \exp \left(\frac{|x - x'|}{l_c} \right), \quad (4.8)$$

where l_c is the correlation length. We constructed $a(x; \boldsymbol{\xi})$ by the Karhunen-Loève (KL) expansion:

$$a(x; \boldsymbol{\xi}) = a_0(x) + \sigma \sum_{i=1}^d \sqrt{\lambda_i} \phi_i(x) \xi_i, \quad (4.9)$$

where $\{\lambda_i\}_{i=1}^d$, and $\{\phi_i(x)\}_{i=1}^d$ are the d largest eigenvalues and the corresponding eigenfunctions of $K(x, x')$. The values of λ_i and the analytical expressions for ϕ_i were available from the literature [87]. The ξ_i are i.i.d. random variables on $[-1, 1]$. The density function of ξ_i is given by

$$\omega(\xi_i) = \frac{1}{\pi \sqrt{1 - \xi_i^2}}, \quad (4.10)$$

where the corresponding orthonormal basis consists of Chebyshev polynomials of the first kind. For this example, we set $a_0(x) \equiv 1$, $\sigma = 0.8$, $l_c = 0.14$ and $d = 16$. We chose the quantity of interest as $u(x; \boldsymbol{\xi})$ at $x = 0.45$ and constructed a 3rd-order polynomial expansion with $N = 969$ basis functions. Figure 4 shows the relative l_2 error of the constructed $\tilde{f}(\boldsymbol{\xi})$ and $\tilde{g}(\boldsymbol{\chi})$. For the density function $\omega(\xi_i)$ given by (4.10), $\tilde{f}(\boldsymbol{\xi})$ can be represented by a multivariate basis constructed by the tensor products of univariate Chebyshev polynomials. However, in general, the PDF of $\boldsymbol{\chi}$ does not retain the form $\omega'(\boldsymbol{\chi}) = \prod_{i=1}^d \frac{1}{\pi \sqrt{1 - \chi_i^2}}$. As shown in Figure 4, iteratively employing the multivariate Chebyshev polynomials to represent $\tilde{g}(\boldsymbol{\chi})$ (the red dash-dotted curve)—as done in previous studies [88]—resulted in a larger error than $\tilde{f}(\boldsymbol{\xi})$. Representing $\tilde{g}(\boldsymbol{\chi})$ by the reconstructed orthonormal **amdP** basis (the blue dashed curve) further decreases the numerical error compared to $\tilde{f}(\boldsymbol{\xi})$ (the solid red curve).

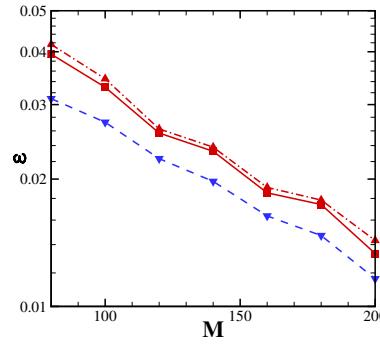
47 4.3. Systems with implicit knowledge of density function
48

In this suite of benchmark examples, we investigated the applicability and efficiency of the developed DSRAR framework based on data-driven orthonormal bases construction and sparsity enhanced rotation.

53 4.3.1. High-dimensional polynomials
54

We studied the ability of the data-driven method to recover a high-dimensional polynomial function

$$f(\boldsymbol{\xi}) = \sum_{\alpha \in T_\alpha} \hat{\psi}_\alpha(\boldsymbol{\xi}), \quad (4.11)$$



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure 4: Sparsity-enhancing rotation with reconstructed orthonormal basis yield the most accurate surrogate models for a 1D elliptical PDE with random permeability coefficient modeled by Equations (4.9) and (4.10). Directly applying the rotation procedure without reconstructing the orthonormal basis yields increased numerical error. “—■—”: Chebyshev polynomial basis with respect to ξ ; “—▲—”: Chebyshev polynomial basis with respect to rotated vector χ ; “—▼—”: the reconstructed orthonormal `amdP` basis with respect to rotated vector χ .

where $\hat{\psi}_\alpha$ represents the monomial basis function, T_α represents a set containing 50 indices randomly chosen from Λ_p^d with $d = 25$ and $p = 3$. The sample set S of random vector ξ for basis construction was generated from the Gaussian mixture model specified in (4.1) with $|S| = 2 \times 10^5$.

We approximated $f(\xi)$ by a 3rd-order polynomial expansion $\tilde{f}(\xi) = \sum_{i=1}^N \tilde{c}_i \psi_i(\xi)$ with $N = 3276$. Figure 5(a) shows the relative l_2 error of the constructed surrogate model \tilde{f} defined by

$$\epsilon = \left(\int (f(\xi) - \tilde{f}(\xi))^2 d\nu_{S_2}(\xi) / \int f(\xi)^2 d\nu_{S_2}(\xi) \right)^{\frac{1}{2}}, \quad (4.12)$$

where 20 implementations were utilized for each training sample size number M . As shown in Figure 5(a), $\tilde{f}(\xi)$ constructed by the near-orthonormal `amdP` basis yielded the smallest error while the tensor product of Legendre basis functions yielded the largest error. Accordingly, the magnitudes of the recovered coefficients $|\tilde{c}_i|$ by the exact and near-orthonormal bases decayed more quickly than those recovered using the Legendre basis functions, as shown in Figure 5(b). Furthermore, $\tilde{f}(\xi)$ allowed us to define a new random vector χ , which further enhanced the sparsity of c , as shown in Figures 5(c) and (d). Following Step 5 in Algorithm 4, we defined a new random χ through rotation. The associated representation coefficient vector c has enhanced sparsity.

However, for the exact and near-orthonormal basis, the $\tilde{g}(\chi)$ gave smaller errors (the dashed curve) than $\tilde{f}(\xi)$ (the solid curve), as shown in Figure 5(a). Thus, enhancing the sparsity of c alone does not guarantee enhanced accuracy of \tilde{f} . In particular, $\tilde{g}(\chi)$ constructed by the Legendre basis yielded larger error than $\tilde{f}(\xi)$ as demonstrated in Figure 5(a); although, the sparsity of c was greater, as seen in Figure 5(d). This behavior indicates that retaining the orthonormal condition can be crucial for the accurate construction of \tilde{f} . The basis bound (see Table B.2 in Appendix B) provides a metric to understand why the near-orthonormal basis

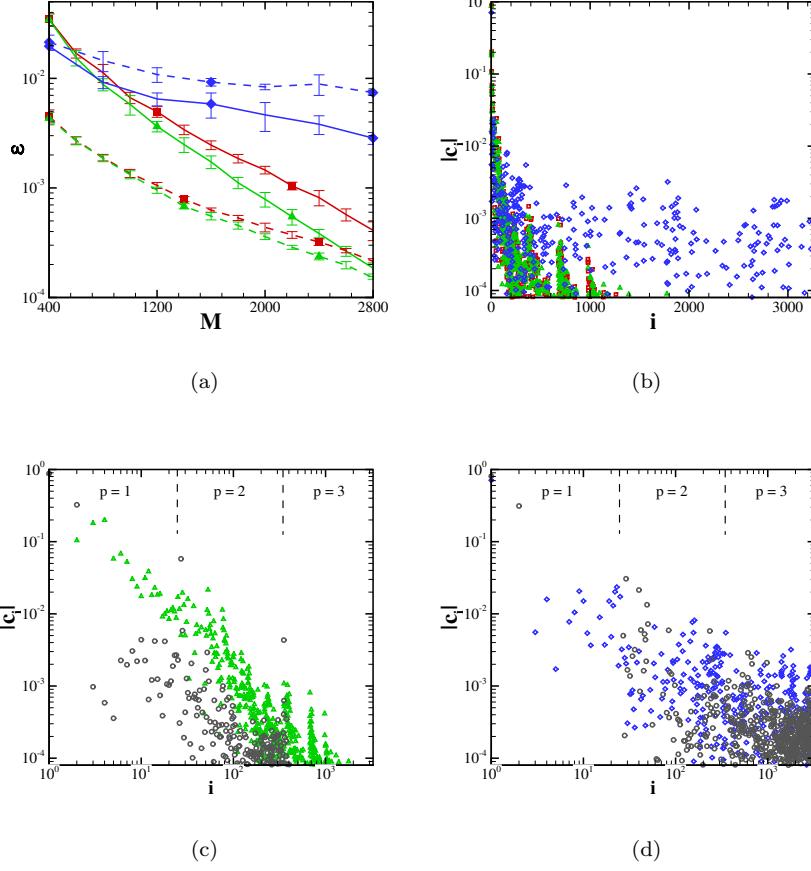
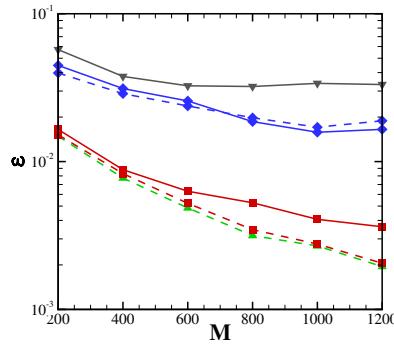


Figure 5: Numerical results for recovery of a high-dimensional polynomial function. The combination of near-orthonormal basis construction with the sparsity enhancement rotation procedure yields the most accurate results. Directly applying the rotation procedure to the Legendre basis may lead to increased error despite increased sparsity in \mathbf{c} . (a) Relative l_2 error of the recovered polynomial function with different bases: the exact orthonormal amdP basis with respect to ξ (“ \blacksquare ”) and χ (“ $\blacksquare\blacksquare$ ”); the near-orthonormal amdP basis with respect to ξ (“ \blacktriangle ”) and χ (“ $\blacktriangle\blacktriangle$ ”); Legendre basis with respect to ξ (“ \blacklozenge ”) and χ (“ $\blacklozenge\blacklozenge$ ”). (b) Coefficients magnitude $|c_i|$ recovered using different bases. “ \blacksquare ”: the exact orthonormal amdP basis with respect to ξ ; “ \blacktriangle ”: the near-orthonormal amdP basis with respect to ξ ; “ \blacklozenge ”: Legendre basis with respect to ξ . (c) Recovered coefficient magnitude $|c_i|$ using the near orthogonal basis with respect to ξ (“ \blacktriangle ”) and χ (“ \circ ”). The dashed vertical lines indicate the separation between different polynomial orders p . (d) Recovered coefficient magnitude $|c_i|$ using the Legendre basis with respect to ξ (“ \blacklozenge ”) and χ (“ \circ ”).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure 6: The combination of near-orthonormal basis construction and sparsity enhancement rotation yields the most accurate results, as shown through the relative l_2 error of the constructed surrogate model for the 1D elliptic PDE with random permeability coefficient: the exact orthonormal `amdP` basis with respect to ξ (“ $\text{---} \blacksquare \text{---}$ ”) and χ (“ $\text{---} \blacksquare \text{---}$ ”); Legendre basis with respect to ξ (“ $\text{---} \blacklozenge \text{---}$ ”) and χ (“ $\text{---} \blacklozenge \text{---}$ ”); Hermite basis with respect to ξ (“ $\text{---} \blacktriangledown \text{---}$ ”); the near-orthonormal `amdP` basis with respect to χ (“ $\text{---} \blacktriangle \text{---}$ ”).



performs better than the exact orthonormal basis.

4.3.2. 1D elliptic PDEs with high-dimensional random inputs

In this example, we revisited the 1D elliptic PDE (4.7) with random coefficient given by Equation (4.9).

Here we set $a_0(x) \equiv 1$, $\sigma = 1$, $l_c = 0.12$ and $d = 20$ such that $\sum_{i=1}^d \lambda_i > 0.91 \sum_{i=1}^\infty \lambda_i$.

Similar to the work by Zabaras et al. [70], a non-Gaussian multivariate distribution was used for $\xi = (\xi_1, \xi_2, \dots, \xi_d)$. We generated a sample set $\{\tilde{\xi}^{(k)}\}_{k=1}^{N_s}$, where $N_s = 2 \times 10^5$ and $\tilde{\xi}$ came from the Gaussian mixture distribution specified in (4.1). We used PCA to transform $\tilde{\xi}$ to ξ such that $\mathbb{E}[\xi_i] = 0$ and $\mathbb{E}[\xi_i \xi_j] = \delta_{ij}$. For each input sample $\xi^{(k)}$, a and u only depended on x and the solution of the deterministic elliptic equation is given by [54]

$$\begin{aligned} u(x) &= u(0) + \int_0^x \frac{a(0)u(0)' - y}{a(y)} dy \\ a(0)u(0)' &= \left(\int_0^1 \frac{y}{a(y)} dy \right) / \left(\int_0^1 \frac{1}{a(y)} dy \right). \end{aligned} \quad (4.13)$$

We chose the QoI to be $u(x; \xi)$ at $x = 0.35$ and constructed a 3rd-order polynomial expansion with $N = 1771$ basis functions. Figure 6 shows the relative l_2 error of $\tilde{f}(\xi)$ (solid curve) and $\tilde{g}(\chi)$ (dashed curve) constructed by different bases. The data-driven bases (both exact orthonormal basis and near-orthonormal basis) showed more accurate results than the Legendre basis and the Hermite basis. In particular, the near-orthonormal basis with respect to the rotated variable χ yielded the most accurate result (the green dashed curve). In contrast, directly employing the Legendre basis to the rotated variable χ without reconstructing the basis function led to increased l_2 error, although c shows more sparsity in terms of χ (the gray dashed curve) than ξ (the gray solid curve).

1
2
3
4 4.4. UQ study of a molecule system under Non-Gaussian conformational distributions
5

6 405 We demonstrated the proposed method on a physical system exploring conformational uncertainty in
7 a small molecule system. Molecular properties, such as solvation energies or solvent-accessible surface ar-
8 eas (SASAs), are often calculated using single molecular conformations. However, due to thermal energy, a
9 molecule undergoes conformational fluctuations which can induce significant uncertainty in properties calcu-
10 lated from single structures. Our previous work [1] was focused on quantifying this uncertainty using a simple
11 multivariate Gaussian model for conformational fluctuations: the elastic network model [89]. However, it is
12 well known that the conformational fluctuations are often non-Gaussian due to the complicated structure
13 of the underlying energy landscape. Therefore, in the current study, we construct the data-driven basis
14 *directly* from the samples of molecular trajectories collected from molecular dynamics (MD) simulations,
15 thus eliminating the *over-simplified* Gaussian assumption.
16
17

18 We simulated the dynamics of the small molecule benzyl bromide under equilibrium (see Appendix E for
19 details) and collected a sample set of the instantaneous molecular structure $\{\mathbf{r}^{(k)}\}_{k=1}^{N_s}$ from MD simulation
20 trajectories over $20\mu\text{s}$. In what follows, $N_s = 2 \times 10^5$ and \mathbf{r} represent the positions of individual atoms. As
21 a pre-processing step, we transformed $\{\mathbf{r}^{(k)}\}_{k=1}^{N_s}$ into a set of uncorrelated random vectors $S = \{\boldsymbol{\xi}^{(k)}\}_{k=1}^{N_s}$
22 via PCA:
23
24

$$\begin{aligned}\Sigma &= \mathbb{E} \left[(\mathbf{r} - \bar{\mathbf{r}})(\mathbf{r} - \bar{\mathbf{r}})^T \right] \\ \Sigma &= \mathbf{Q}\boldsymbol{\Gamma}\mathbf{Q}^T \quad \boldsymbol{\xi} = \boldsymbol{\Gamma}^{-1/2}\mathbf{Q}^T\mathbf{r},\end{aligned}\tag{4.14}$$

25 415 where the average $\mathbb{E}[\cdot]$ is taken over the entire sample set and $\boldsymbol{\xi} \in \mathbb{R}^{12}$ is the normalized random vector
26 that represents 99.99% of the observed variance. Figures 7(a) and (b) show the joint distributions of (ξ_1, ξ_2)
27 and (ξ_1, ξ_3) . Although the individual components of $\boldsymbol{\xi}$ are uncorrelated, the joint density distributions are
28 mutually dependent and deviate from the standard Gaussian distributions.
29
30

31 41 We chose the polar solvation energy and SASA as the target QoIs for this system. The polar solvation
32 energy was modeled by the Poisson-Boltzmann equation [90, 91]
33
34

$$-\nabla \cdot (\epsilon_f(\mathbf{x}; \boldsymbol{\xi}) \nabla \varphi(\mathbf{x}; \boldsymbol{\xi})) = \rho_f(\mathbf{x}; \boldsymbol{\xi})\tag{4.15}$$

35 44 which relates the electrostatic potential φ to a dielectric coefficient ϵ_f and a fixed charge distribution ρ_f .
36 Equation (4.15) is typically solved with Dirichlet boundary conditions set to an analytical asymptotic solution
37 of the equation for an infinite domain. The dielectric coefficient ϵ_f implicitly represents the boundary
38 between the atoms of the molecule and the surrounding solvent: the coefficient changes rapidly across this
39 boundary from a low dielectric value in the molecular interior to a high dielectric value in the solvent. The
40 charge distribution ρ_f is generally modeled as a collection of δ -like functions centered on the atoms of the
41 molecule with magnitudes proportional to the atomic partial charges. Both ϵ_f and ρ_f are dependent on the
42
43

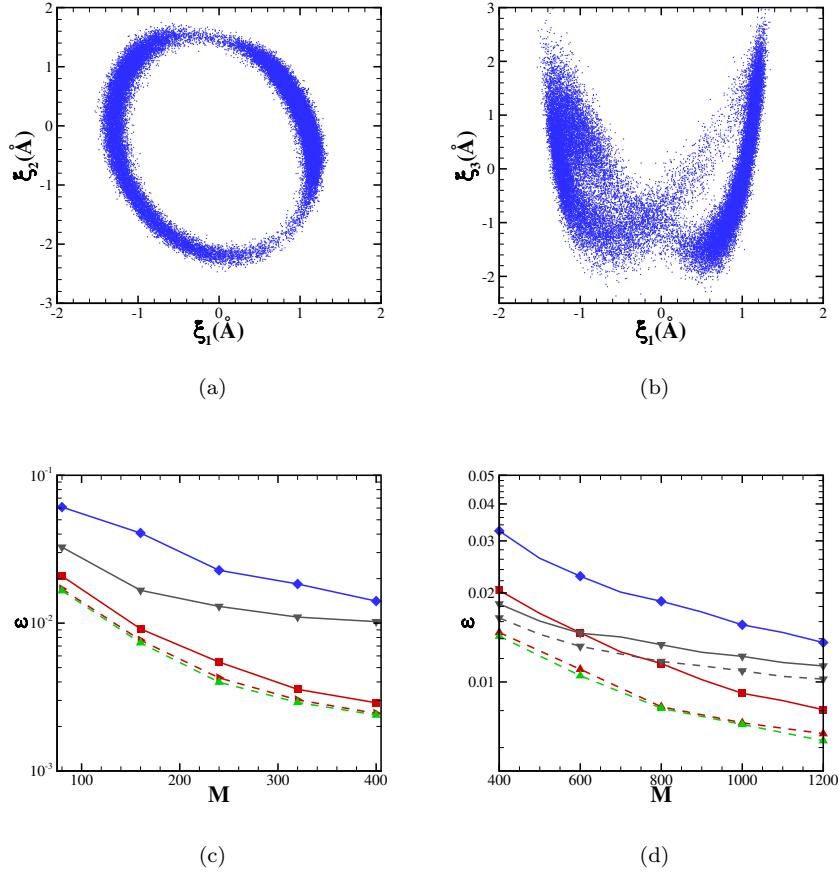


Figure 7: The present method based on data-driven basis construction and sparsity enhancement rotation yields the most accurate surrogate model for molecular systems with mutually dependent non-Gaussian density distributions. (a-b) Sampling points representing the joint distributions (ξ_1, ξ_2) (left) and (ξ_1, ξ_3) (right). (c-d) Relative l_2 error of the polar solvation energy (left) and the local SASA (right) of an individual atom (the H9 atom attached to the ortho-carbon atom) obtained with different numbers of training data M : the exact amdP orthonormal basis with respect to ξ (“—■—”) and χ (“—■—”); Hermite basis with respect to ξ (“—▼—”) and χ (“—▽—”); Legendre basis with respect to ξ (“—◆—”); the near-orthonormal amdP basis with respect to χ (“—▲—”).

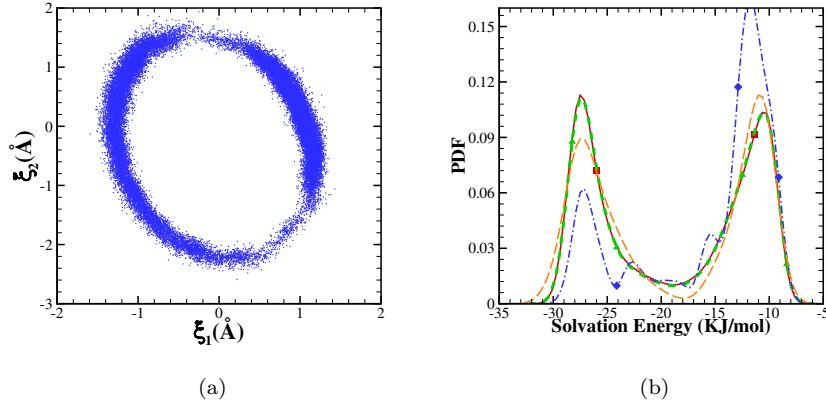


Figure 8: The present method yields the most accurate prediction on the PDF of the QoI for the molecular systems. Direct fitting of the underlying density $\omega(\xi)$ using Gaussian Mixture model may induce biased error to the PDF prediction. (a) Fitted random variables (ξ_1, ξ_2) with Gaussian mixture models. (b) PDF of the solvation energy obtained with the Gaussian Mixture model and the present data-driven approach. “■”: reference solution obtained from 2×10^5 MC samples; “◆”: direct MC sampling using the same set of 200 samples; “▲”: present method using the same set of 200 samples; “—”: fitting Gaussian Mixture model using 800 samples.

instantaneous molecular structure (i.e., ξ). The polar solvation energy was calculated from

$$G_p(\xi) = \int \rho_f(\mathbf{x}; \xi) (\varphi(\mathbf{x}; \xi) - \varphi_h(\mathbf{x}; \xi)) d\mathbf{x} \quad (4.16)$$

where φ_h is a reference potential obtained from solution of

$$-\epsilon_h \nabla^2 \varphi_h(\mathbf{x}; \xi) = \rho_f(\mathbf{x}; \xi) \quad (4.17)$$

where ϵ_h is a constant reference dielectric value. We used the Adaptive Poisson-Boltzmann Solver (APBS) software to solve the equations above [92]. Besides the solvation energy of the whole molecule, we also studied a local property like the SASA of an individual atom (the H9 atom attached to the ortho-carbon atom of the benzyl bromide molecule, see Figure E.11) by the Shake-Rupley algorithm [93] using APBS. Details of the APBS calculations are presented in Appendix E.

Figures 7(c) and (d) show the relative l_2 error of the constructed surrogate model $\tilde{f}(\xi)$ for the solvation energy and SASA using a 4th-order gPC expansion with $N = 1820$ basis functions. For both QoIs, the near-orthonormal and orthonormal bases with respect to the rotated variable χ (dashed curves) yield similar error which is much smaller than the error of Legendre and Hermite bases. A possible explanation for the similar performance of the near-orthonormal and orthonormal bases is the closeness of the basis bound estimates for these two bases (see Table B.3 in Appendix B).

Instead of the direct construction of $\tilde{f}(\xi)$ using data-driven basis functions, another possible approach to characterize the uncertainty of the molecular system is to fit the distribution density $\omega(\xi)$ with a distribution

model such as a Gaussian Mixture model. Figure 8 (a) shows a scatter plot of the joint distribution (ξ_1, ξ_2) extracted from the fitted Gaussian mixture distribution $\tilde{\omega}(\xi)$ using 7 Gaussian modes. Accordingly, we can construct the surrogate model for each Gaussian mode using standard Hermite basis function. However, it is well-known that accurate construction of $\omega(\xi)$ is a numerically challenging problem for $d > 4$. As shown in Figure 8(b), direct fitting $\omega(\xi)$ by $\tilde{\omega}(\xi)$ induces non-negligible error and leads to biased prediction of the PDF of the solvation energy. Furthermore, we lose the one-to-one mapping between the individual conformation state ξ and the QoIs through the constructed surrogate model $\tilde{f}(\xi)$.

5. Summary

In this study, we have developed a DSRAR framework for constructing surrogate models irrespective of the mutual dependence between the components of random inputs using limited training points. To the best of our knowledge, this problem has not been addressed by previous UQ studies based on polynomial chaos expansions. The DSRAR framework does not assume mutual independence between the components of random inputs and therefore can be applied to UQ in complex systems where information about the underlying random distribution can be implicit. To construct the surrogate model, this framework uses data-driven `amdP` basis construction and a sparsity-enhancing rotation procedure which leads to more accurate recovery of the sparse representation of the target function. The method benefits from both the orthonormal basis expansion and the enhanced sparsity of the expansion coefficients. With the assumption that there exists a sparse representation of the surrogate model, the DSRAR approach can be applied to challenging UQ problems under two widely encountered situations: (I) probability measure implicitly represented by a large collection of samples and (II) non-Gaussian probability measures with explicit (analytical) forms. For systems with explicit knowledge of the probability measure, our method exploits sparser representations of QoIs while retaining proper orthogonality with respect to rotated variables. For systems with randomness implicitly represented by a large collection of random samples, we also proposed a heuristic method to construct a *near-orthonormal* basis in addition to the exact orthonormal basis with respect to the discrete measure. The near-orthonormal basis shows a smaller basis bound and empirically yields more accurate representations. The numerical examples show the effectiveness of our method for realistic problems on quantifying uncertainty propagation in molecular system under conformational fluctuations as well as PDEs with arbitrary underlying probability measures.

For future study, we note that several issues not considered in the present work could further improve the performance of the present DSRAR framework. The heuristic approach to constructing near-orthonormal basis introduced in this study yields smaller basis bounds and more accurate representations than existing methods. However, we do not have the theoretical analysis to formally show that the near-orthonormal basis is optimal and to establish the conditions under which it outperforms the exact orthonormal basis. It

would be interesting to investigate different approaches of data-driven basis construction to further improve the properties of measurement matrix for CS purposes. For instance, if new data becomes available after the surrogate construction, it is worth exploring how to use the new information to design more sophisticated (cross-validation) strategies to optimize the orthonormal threshold values and the basis construction procedure. Furthermore, our study used a standard ℓ_1 minimization approach for relaxing the CS problem and recovering a sparse solution of the under-determined system. However, other optimization approaches can be employed when the measurement matrix is highly coherent when ℓ_1 minimization is not necessarily optimal. Finally, it would be interesting to employ the developed DSRAR approach for UQ study in other complex biological systems [94, 95]. Such results will be presented in a future publication.

Appendix A. Proof of Theorem 4.1

Proof. Let $\mathbf{v} \in \text{Ker } \mathbf{A}$ and $\mathbf{x} \neq \mathbf{c}$ another solution of $\mathbf{Ax} = \mathbf{b}$. To show that \mathbf{c} is the unique ℓ_1 minimizer of $\mathbf{Ac} = \mathbf{b}$, it is sufficient if

$$\|\mathbf{v}_{T_\alpha}\|_1 < \|\mathbf{v}_{T_\alpha^c}\|_1, \quad (\text{A.1})$$

which gives

$$\begin{aligned} \|\mathbf{c}\|_1 &\leq \|\mathbf{c} - \mathbf{x}_{T_\alpha}\|_1 + \|\mathbf{x}_{T_\alpha}\|_1 = \|\mathbf{c}_{T_\alpha} - \mathbf{x}_{T_\alpha}\|_1 + \|\mathbf{x}_{T_\alpha}\|_1 = \|\mathbf{v}_{T_\alpha}\|_1 + \|\mathbf{x}_{T_\alpha}\|_1 \\ &< \|\mathbf{v}_{T_\alpha^c}\|_1 + \|\mathbf{x}_{T_\alpha}\|_1 = \|\mathbf{x}\|_1. \end{aligned} \quad (\text{A.2})$$

To satisfy (A.1), we partition T_α^c into $T_\alpha^c = T_{\alpha,1}^c \cup T_{\alpha,2}^c \cup \dots$, where $T_{\alpha,1}^c$ is the index set of s largest absolute entries of \mathbf{v} in T_α^c , $T_{\alpha,2}^c$ is the index set of s largest absolute entries of \mathbf{v} in $T_\alpha^c \setminus T_{\alpha,1}^c$. Accordingly,

$$\|\mathbf{v}_{T_\alpha}\|_2^2 \leq \frac{1}{1-\delta_s} \|\mathbf{Av}_{T_\alpha}\|_2^2 = \frac{1}{1-\delta_s} \sum_{k=1} \langle \mathbf{Av}_{T_\alpha}, \mathbf{A}(-\mathbf{v}_{T_{\alpha,k}^c}) \rangle \leq \frac{\theta_s}{1-\delta_s} \sum_{k=1} \|\mathbf{v}_{T_\alpha}\|_2 \|\mathbf{v}_{T_{\alpha,k}^c}\|_2, \quad (\text{A.3})$$

which gives $\|\mathbf{v}_{T_\alpha}\|_2 \leq \frac{\theta_s}{1-\delta_s} \sum_{k=1} \|\mathbf{v}_{T_{\alpha,k}^c}\|_2$. The remaining of the proof is straightforward and follows Theorem 2.6 of Rauhut [96]. By the Cauchy-Schwarz inequality, we obtain

$$\|\mathbf{v}_{T_\alpha}\|_1 \leq \frac{\theta_s}{1-\delta_s} (\|\mathbf{v}_{T_\alpha}\|_1 + \|\mathbf{v}_{T_\alpha^c}\|_1). \quad (\text{A.4})$$

Equation (A.1) follows if $\frac{\theta_s}{1-\delta_s} < 0.5$. \square

Remark Appendix A.1. We emphasize that Theorem 4.1 holds only for the given index set T_α ; it provides a metric to examine the recovery accuracy with respect to measurement matrix \mathbf{A} and should not be viewed as the sufficient condition for exact recovery of *arbitrary s-sparse vector* via ℓ_1 -minimization (see canonical references [82, 83, 96] for details). Theorem 4.1 also indicates that, for the given index set T_α , small $\|\mathbf{A}_{T_\alpha}^* \mathbf{A}_{T_\alpha} - I\|_2$ will promote the recover of \mathbf{v}_{T_α} .

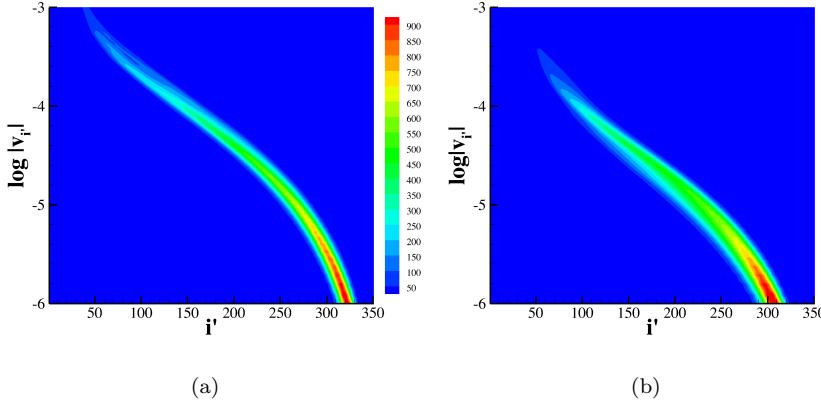


Figure B.9: The null spaces of measurement matrices constructed by the exact and near-orthonormal bases are different under $\|\mathbf{v}_{T_\alpha}\|_1 > \|\mathbf{v}_{T_\alpha^c}\|_1$ —a necessary condition for \mathbf{c} not being recoverable exactly. Density contour of the normalized null space vector component $\log |\mathbf{v}_{i'}|$ (sorted by magnitude) of the measurement matrix \mathbf{A} constructed by orthogonal (a) and near-orthogonal basis functions (b) that satisfy $\|\mathbf{v}_{T_\alpha}\|_1 > \|\mathbf{v}_{T_\alpha^c}\|_1$ and $\|\mathbf{v}\|_2 = 1$.

Appendix B. Measurement matrix and basis bounds

Appendix B.1. Null space of measurement matrix from Section 4.1

Let $\tilde{\mathbf{c}} = \mathbf{c} + \mathbf{v}$, $\mathbf{v} \in \text{Ker } \mathbf{A}$ where \mathbf{A} is the measurement matrix defined in (4.3). From the null space property [96], $\tilde{\mathbf{c}}$ does not fully recover \mathbf{c} by ℓ_1 minimization (i.e., equation (2.6)) only if $\|\tilde{\mathbf{c}}\|_1 < \|\mathbf{c}\|_1$. As a *necessary condition* for the failure of recovery, it requires

$$\|\mathbf{v}_{T_\alpha}\|_1 > \|\mathbf{v}_{T_\alpha^c}\|_1, \quad (\text{B.1})$$

where T_α^c refers to the complement of T_α . Accordingly, different null space of measurement matrix \mathbf{A} generally leads to different recovery error.

We examined the above necessary condition (B.1) for different measurement matrices by randomly choosing a non-zero index set T_α with $|T_\alpha| = 50$ and $M = 180$. For \mathbf{A} constructed by both basis sets, we collected 1000 normalized $\mathbf{v} \in \text{Ker } \mathbf{A}$ that satisfy $\|\mathbf{v}_{T_\alpha}\|_1 > \|\mathbf{v}_{T_\alpha^c}\|_1$. Figure B.9 shows the density contour of individual component $|\mathbf{v}_{i'}|$ in log-scale, where i' refers to the index sorted by magnitude in descending order. The two basis sets demonstrate different distributions of $\log |\mathbf{v}_{i'}|$, which likely contribute to the different recovery errors shown in Figure 2.

Appendix B.2. Basis bounds

The lower bound of the required number of samples M given in Theorem 2.5 suggests that bases with smaller basis bounds K are preferred. We expect that smaller basis bounds will correlate with higher

accuracy representations. For the constructed basis set $\psi_i(\xi), i = 1, \dots, N$, we define the basis bound \tilde{K} on the given data set S by

$$\tilde{K} := \frac{1}{|S_{M_\sigma}|} \sum_{\xi \in S_{M_\sigma}} |k(\xi)|, \quad (\text{B.2})$$

where the set S_{M_σ} is defined by $S_{M_\sigma} = \{\xi \mid |k(\xi) - \mathbb{E}[k]| > M_\sigma \sigma[k], \xi \in S\}$. Here $k(\xi) := \max_i |\psi_i(\xi)|$ denotes the maximum magnitude for an individual sampling point ξ , $\mathbb{E}[k]$ and $\sigma[k]$ represent the mean and the standard deviation of $k(\xi)$ on S with respect to the discrete measure ν_S . In this study, we present \tilde{K} as an indication of the difference between the exact and near-orthonormal basis function. In compressive sensing, the measurement matrix only consists of limited number of samples. Therefore, we employ the mean of the tails in the basis bounds as an indicator of the upper bound of the largest entry values from the measurement matrix. M_σ defines the range of this tail set. We choose $M_\sigma = 5$ if not specified otherwise.

Table B.1: \tilde{K} of constructed basis set for Gaussian mixture system $d = 25$, $p = 2$ and $N_s = 1 \times 10^5$.

M_σ	3	4	5	6	$\max_{\xi \in S} k(\xi)$
\tilde{K}_{orth}	10.359	12.048	13.895	15.513	22.208
$\tilde{K}_{\text{near-orth}}$	9.622	11.196	12.867	14.448	18.790

Following the definition by Equation (B.2), we examine the basis bound \tilde{K} of the numerical examples

presented in this study. Table B.1 shows the results of Gaussian mixture system $\{\xi^{(i)}\}, i = 1, \dots, N_s$ with $N_s = 1 \times 10^5$, $d = 25$ and $p = 2$ which is defined in Section 4.1. For different values of M_σ , \tilde{K} of the near orthogonal basis shows consistently smaller values than the values of the exact orthogonal basis set.

Table B.2 shows the basis bound \tilde{K} of the Gaussian mixture system which is studied in Section 4.3.1 with $N_s = 2 \times 10^5$, $d = 25$ and $p = 3$. The values of \tilde{K} for the near orthogonal basis are consistently smaller than the value for the exact orthogonal basis set no matter on the original random sample set or the rotated sample set. Furthermore, we present the basis bounds on the rotated sampling set $\{\chi_M^{(i)}\}_{i=1}^{N_s}$, where the subscript “ M ” refers to the different number of training points utilized to construct the surrogate model $X(\xi)$. The near-orthogonal basis yields smaller \tilde{K} than the exact orthogonal basis in each case.

Similarly, Table B.3 shows \tilde{K} of the constructed basis for uncertainty quantification of the molecular solvation energy ($d = 12$, $p = 4$ and $N_s = 2 \times 10^5$), which is studied in Section 4.4. The near-orthogonal basis yields smaller values consistently for different number (χ_M) of training points.

Appendix C. Other metrics for the surrogate model

Besides the relative l_2 error, we have also computed the predictivity coefficients Q_2 for the test cases of Gaussian mixture (with $d = 25$ and $p = 3$) and molecular systems. Similar to Marrel et al. [97], Q_2 is defined

Table B.2: \tilde{K} of constructed basis set for Gaussian mixture system $d = 25$, $p = 3$ and $N_s = 2 \times 10^5$.

	ξ	$\chi_{M=400}$	$\chi_{M=1200}$	$\chi_{M=1600}$	$\chi_{M=2400}$
\tilde{K}_{orth}	32.497	32.522	32.079	33.142	32.308
$\tilde{K}_{\text{near-orth}}$	28.320	29.811	29.407	29.512	29.192

Table B.3: \tilde{K} of constructed basis set for molecular system $d = 12$, $p = 4$ and $N_s = 2 \times 10^5$.

	$\chi_{M=80}$	$\chi_{M=160}$	$\chi_{M=240}$	$\chi_{M=320}$	$\chi_{M=400}$
\tilde{K}_{orth}	40.596	39.914	39.789	39.218	39.142
$\tilde{K}_{\text{near-orth}}$	39.970	39.278	39.290	38.528	38.631

by

$$Q_2 = 1 - \int (f(\xi) - \bar{f}(\xi))^2 d\nu_{S_2}(\xi) / \int (f(\xi) - \bar{f})^2 d\nu_{S_2}(\xi), \quad (\text{C.1})$$

where \bar{f} represents the mean of QoI on S_2 . The results are shown in Tab. C.4, where the surrogate models are constructed by the our data-driven basis approach.

Table C.4: The predictivity coefficient Q_2 for polynomial function with Gaussian Mixture measure ($d = 25$ and $p = 3$) and the molecular system for solvation energy and SASA of atom H9.

molecule solvation	M	80	160	240	320	400
	Q_2	0.995715	0.999132	0.999731	0.999864	0.999911
molecule SASA	M	200	300	400	500	600
	Q_2	0.988675	0.996069	0.998272	0.998709	0.999027
Gaussian Mixture	M	200	300	400	500	600
	Q_2	0.998372	0.999347	0.999844	0.999892	0.999941

With the constructed surrogate model, we can compute the Sobol sensitivity indices for QoI with dependent random variables. In brief, $f(\xi)$ is expanded by

$$f(\xi) = \eta_0(\xi) + \sum_{\beta \in \Theta^d} \eta_\beta(\xi), \quad (\text{C.2})$$

where Θ^d represents the collection of all subsets of $[1 : d]$ and $\eta_\beta(\xi)$ satisfies $\mathbb{E}[\eta_\alpha, \eta_\beta] = 0$, if $\alpha \subset \beta$. The sensitivity index S_β is given by

$$S_\beta = \frac{\mathbb{V}(\eta_\beta) + \sum_{\alpha \cap \beta \neq \alpha, \beta} \text{Cov}(\eta_\alpha, \eta_\beta)}{\mathbb{V}(f)} \quad (\text{C.3})$$

where $\mathbb{V}(\cdot)$ refers to the variance on ν_S . We refer to Chastaing et al. [98] for the details. Fig. C.4 shows the first-order sensitivity indices for the test cases of Gaussian mixture systems ($d = 25$, $p = 3$) and the

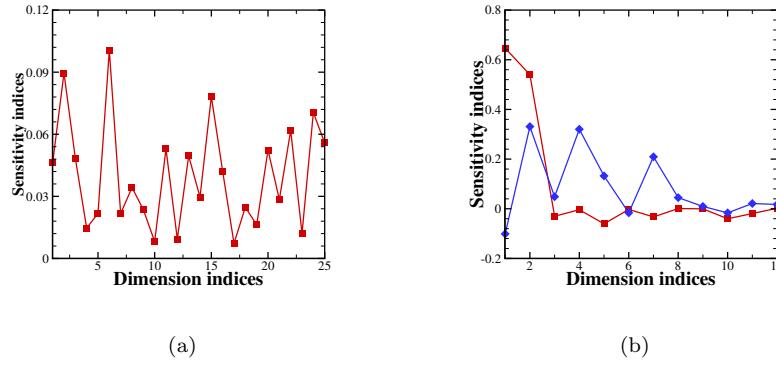


Figure C.10: The first-order sobol sensitivity indices for (a) polynomial function with Gaussian mixture measure ($d = 25, p = 3$) (b) molecular system for solvation energy (“ \blacksquare ”) and SASA of atom H9(“ \blacktriangleleft ”).

molecular systems, where the surrogate models are constructed by the data-driven basis approach using $M = 800$, $M = 240$ and $M = 600$ training points, respectively. The dominant components are on the dimensions $(1, 2, 3, 6, 11, 13, 14, 15, 16, 20, 22, 24, 25)$, $(1, 2, 5)$ and $(1, 2, 4, 5, 7)$ (90% of total variance).

520 Appendix D. Generation of the Gaussian mixture data set

We used Matlab to generate the Gaussian mixture data set in Sec. 4.1 by calling the function `gmdistribution(μ , $\{\Sigma_i\}_{i=1}^3, \mathbf{a}$)` with $\mathbf{a} = (0.5358, 0.1281, 0.3361)$. μ is a 25×3 random matrix with i.i.d. entries on $U[-2.5, 2.5]$. $\{\Sigma_i\}_{i=1}^3$ is a $25 \times 25 \times 3$ array where Σ_i is defined by

$$\Sigma_i = (\Upsilon_i \Upsilon_i^T + \mathbf{I})/4, \quad (D.1)$$

where Υ_i is a random matrix with i.i.d. entries from $\mathcal{U}[0, 1]$ for $i = 1, 2, 3$. μ and Υ_i are generated by calling the Matlab function `rand()` with random number seed 200.

44 Appendix E. Molecular Dynamics simulation and calculation details

We performed all-atom MD simulation of benzyl bromide in water using GROMACS 5.1.2 [99]. The simulation system included a benzyl bromide molecule (see Figure E.11 for the molecular structure) and 1011 water molecules. The General AMBER Force Field (GAFF) [100] was used for the benzyl bromide parameters. The partial charges of benzyl bromide molecule were calculated by RESP method [101]. Bond lengths of benzyl bromide were constrained using the LINCS algorithm [102]. The water molecule was modeled with the rigid TIP3P water model [103]. The bond lengths and angles were held constant through the SETTLE algorithm [104]. The system was equilibrated in the isothermal-isobaric ensemble for 10 ns at 300K and 1 bar after energy minimization. The van der Waals cut-off radii was 1.0 nm. Long-range

1
2
3
4 electrostatics were calculated using a Particle Mesh Ewald (PME) summation with grid spacing of 0.12 nm.
5 The time step was 2 fs. Isobaric-isothermal simulations were equilibrated using a V-rescale thermostat and
6 Berendsen barostat. Following equilibration, the simulation was run for a production period of 20 μ s in a
7 NVT ensemble with a Nosé-Hoover thermostat. The trajectory was stored every 10000 time steps.
8
9

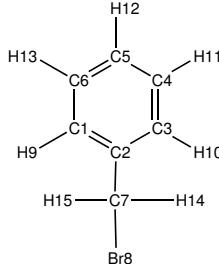


Figure E.11: Sketch of the molecule benzyl bromide with labeled atoms.

APBS calculations [105, 92] were performed with 129^3 grid points over a $40 \times 40 \times 40 \text{ \AA}^3$ coarse grid domain with focusing to a $14 \times 14 \times 14 \text{ \AA}^3$ fine grid domain with the grid origin located at the geometric center of the molecule. The Poisson equation was solved with Dirichlet boundary conditions based on the asymptotic behavior of multiple point charges in a homogeneous dielectric medium. The dielectric coefficient inside the domain used a van der Waals molecular volume definition with a dielectric value of 2.0 inside the molecule and 78.0 outside the molecule. Charges were modeled by Dirac delta functions but discretized to the finite difference grid points using a cubic spline approximation.

Acknowledgements

We thank Ling Guo (Shanghai Normal University), Lei Wu (Princeton University), Wen Zhou (Colorado State University), and David Sept (University of Michigan, ORCID:0000-0003-3719-2483) for helpful discussions. This work was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research as part of the Collaboratory on Mathematics for Mesoscopic Modeling of Materials (CM4) and by the National Institutes of Health grant R01 GM069702. The research was performed using resources available through Research Computing at Pacific Northwest National Laboratory HL acknowledges grant support from AMS Simons Post-doctoral Travel Grant and PNNL Laboratory Directed Research & Development (LDRD) under project “Development of physics-compatible stochastic models for multiphysics systems”.

References

- [1] H. Lei, X. Yang, B. Zheng, G. Lin, N. A. Baker, Constructing surrogate models of complex systems with enhanced sparsity: Quantifying the influence of conformational uncertainty in biomolecular solvation,

- SIAM Multiscale Model. Simul. 13 (4) (2015) 1327–1353.
- [2] S. Oladyshkin, W. Nowak, Data-driven uncertainty quantification using the arbitrary polynomial chaos expansion, Reliability Engineering & System Safety 106 (2012) 179 – 190.
- [3] A. Saltelli, [Global sensitivity analysis: the primer](#), John Wiley, 2008.
- [4] C. M. Bishop, [Pattern Recognition and Machine Learning \(Information Science and Statistics\)](#), Springer-Verlag, Berlin, Heidelberg, 2006.
- [5] A. Laio, M. Parrinello, Escaping free-energy minima, Proceedings of the National Academy of Sciences 99 (20) (2002) 12562–12566.
- [6] G. Fishman, Monte Carlo: Concepts, Algorithms, and Applications, Springer-Verlag New York, Inc., 1996.
- [7] S. Kucherenko, D. Albrecht, A. Saltelli, [Exploring multi-dimensional spaces: a Comparison of Latin Hypercube and Quasi Monte Carlo Sampling Techniques](#) (2015). [arXiv:arXiv:1505.02350](#).
- [8] M. B. Giles, [Multilevel Monte Carlo methods](#), Acta Numerica 24 (2015) 259 – 328.
- [9] S. Heinrich, [Multilevel Monte Carlo Methods](#), in: S. Margenov, J. Waśniewski, P. Yalamov (Eds.), Large-Scale Scientific Computing, Springer Berlin Heidelberg, Berlin, Heidelberg, 2001, pp. 58–67.
- [10] M. Pisaroni, F. Nobile, P. Leyland, [A Continuation Multi Level Monte Carlo \(C-MLMC\) method for uncertainty quantification in compressible inviscid aerodynamics](#), Computer Methods in Applied Mechanics and Engineering 326 (2017) 20 – 50.
- [11] P. Koutsourelakis, [Accurate Uncertainty Quantification Using Inaccurate Computational Models](#), SIAM Journal on Scientific Computing 31 (5) (2009) 3274–3300.
- [12] B. Peherstorfer, K. Willcox, M. Gunzburger, [Optimal Model Management for Multifidelity Monte Carlo Estimation](#), SIAM Journal on Scientific Computing 38 (5) (2016) A3163–A3194.
- [13] B. Fox, Strategies for Quasi-Monte Carlo, Kluwer Academic Pub., 1999.
- [14] H. Niederreiter, Random number generation and Quasi-Monte Carlo methods, SIAM, 1992.
- [15] H. Niederreiter, P. Hellekalek, G. Larcher, P. Zinterhof, Monte Carlo and Quasi-Monte Carlo methods 1996, Springer-Verlag, 1998.
- [16] M. D. McKay, R. J. Beckman, W. J. Conover, [Comparison of Three Methods for Selecting Values of Input Variables in the Analysis of Output from a Computer Code](#), Technometrics 21 (2) (1979) 239–245.

- [17] M. Stein, Large sample properties of simulations using Latin Hypercube Sampling, *Technometrics* 29 (2) (1987) 143–151.
- [18] W. Loh, On Latin hypercube sampling, *Ann. Stat.* 24 (5) (1996) 2058–2080.
- [19] J. Sacks, W. J. Welch, T. J. Mitchell, H. P. Wynn, [Design and Analysis of Computer Experiments](#), *Statist. Sci.* 4 (4) (1989) 409–423.
- [20] M. C. Kennedy, A. O'Hagan, [Bayesian calibration of computer models](#), *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63 (3) (2001) 425–464.
- [21] C. Rasmussen, C. Williams, [Gaussian Processes for Machine Learning](#), Adaptive Computation and Machine Learning, MIT Press, Cambridge, MA, USA, 2006.
- [22] N. Wiener, The homogeneous chaos, *Amer. J. Math.* 60 (1938) 897–936.
- [23] R. Ghanem, P. Spanos, [Stochastic Finite Elements: A Spectral Approach](#), Springer-Verlag, 1991.
- [24] D. Xiu, G. E. Karniadakis, The wiener-askey polynomial chaos for stochastic differential equations, *SIAM J. Sci. Comput.* 24 (2002) 619–644.
- [25] P. Z. G. Qian, C. F. J. Wu, [Bayesian Hierarchical Modeling for Integrating Low-Accuracy and High-Accuracy Experiments](#), *Technometrics* 50 (2) (2008) 192–204.
- [26] B. Williams, D. Higdon, J. Gattiker, L. Moore, M. McKay, S. Keller-McNulty, [Combining experimental data and computer simulations, with an application to flyer plate experiments](#), *Bayesian Anal.* 1 (4) (2006) 765–792. doi:10.1214/06-BA125.
- [27] J. Oakley, A. O'Hagan, [Bayesian inference for the uncertainty distribution of computer model outputs](#), *Biometrika* 89 (4) (2002) 769–784.
- [28] B. A. Lockwood, M. Anitescu, [Gradient-Enhanced Universal Kriging for Uncertainty Propagation](#), *Nuclear Science and Engineering* 170 (2) (2012) 168–195. doi:10.13182/NSE10-86.
- [29] D. Xiu, G. Karniadakis, Modeling uncertainty in flow simulations via generalized polynomial chaos, *J. Comput. Phys.* 187 (2003) 137–167.
- [30] R. Ghanem, S. Masri, M. Pellissetti, R. Wolfe, Identification and prediction of stochastic dynamical systems in a polynomial chaos basis, *Comput. Meth. Appl. Math. Engrg.* 194 (2005) 1641–1654.
- [31] O. Knio, O. Le Maître, Uncertainty propagation in CFD using polynomial chaos decomposition, *Fluid Dyn. Res.* 38 (9) (2006) 616–640.

- [32] B. Sudret, Global sensitivity analysis using polynomial chaos expansions, *Reliability Engineering & System Safety* 93 (7) (2008) 964 – 979.
- [33] J. Li, D. Xiu, A generalized polynomial chaos based ensemble Kalman filter with high accuracy, *J. Comput. Phys.* 228 (2009) 5454–5469.
- [34] Y. Marzouk, D. Xiu, A stochastic collocation approach to bayesian inference in inverse problems, *Communications in Computational Physics* 6 (4) (2009) 826–847. doi:10.4208/cicp.2009.v6.p826.
- [35] J. Li, D. Xiu, Evaluation of failure probability via surrogate models, *J. Comput. Phys.* 229 (2010) 8966–8980.
- [36] J. Li, P. Stinis, Mori-zwanzig reduced models for uncertainty quantification, arXiv:1803.02826.
- [37] R. Schobi, B. Sudret, J. Wiart, [Polynomial-chaos-based Kriging](#), *International Journal for Uncertainty Quantification* 5 (2) (2015) 171–193.
- [38] L. L. Gratiet, S. Marelli, B. Sudret, [Metamodel-Based Sensitivity Analysis: Polynomial Chaos Expansions and Gaussian Processes](#), Springer International Publishing, Cham, 2016, pp. 1–37.
- [39] N. Owen, P. Challenor, P. Menon, S. Bennani, [Comparison of Surrogate-Based Uncertainty Quantification Methods for Computationally Expensive Simulators](#), *SIAM/ASA Journal on Uncertainty Quantification* 5 (1) (2017) 403–435. doi:10.1137/15M1046812.
- [40] P. T. Roy, N. El Moçayd, S. Ricci, J.-C. Jouhaud, N. Goutal, M. De Lozzo, M. C. Rochoux, [Comparison of polynomial chaos and Gaussian process surrogates for uncertainty quantification and correlation estimation of spatially distributed open-channel steady flows](#), *Stochastic Environmental Research and Risk Assessment* 32 (6) (2018) 1723–1741.
- [41] L. Mathelin, M. Hussaini, A stochastic collocation algorithm for uncertainty analysis, *Tech. Rep. NASA/CR-2003-212153*, NASA Langley Research Center (2003).
- [42] D. Xiu, J. Hesthaven, High-order collocation methods for differential equations with random inputs, *SIAM J. Sci. Comput.* 27 (3) (2005) 1118–1139.
- [43] I. Babuška, F. Nobile, R. Tempone, A stochastic collocation method for elliptic partial differential equations with random input data, *SIAM J. Numer. Anal.* 45 (3) (2007) 1005–1034.
- [44] F. Nobile, R. Tempone, C. G. Webster, A sparse grid stochastic collocation method for partial differential equations with random input data, *SIAM J. Numer. Anal.* 46 (5) (2008) 2309–2345.

- [45] X. Ma, N. Zabaras, An adaptive hierarchical sparse grid collocation algorithm for the solution of stochastic differential equations, *J. Comput. Phys.* 228 (8) (2009) 3084–3113.
- [46] J. Foo, G. E. Karniadakis, Multi-element probabilistic collocation method in high dimensions, *J. Comput. Phys.* 229 (5) (2010) 1536 – 1557.
- [47] P. G. Constantine, M. S. Eldred, E. T. Phipps, Sparse pseudospectral approximation method, *Computer Methods in Applied Mechanics and Engineering* 229-232 (2012) 1 – 12. doi:<https://doi.org/10.1016/j.cma.2012.03.019>.
- [48] J. D. Jakeman, S. G. Roberts, Local and dimension adaptive stochastic collocation for uncertainty quantification, in: *Sparse grids and applications*, Springer, 2013, pp. 181–203.
- [49] J. Li, P. Stinis, A unified framework for mesh refinement in random and physical space, *Journal of Computational Physics* 323 (2016) 243 – 264. doi:<https://doi.org/10.1016/j.jcp.2016.07.027>.
- [50] A. Doostan, H. Owhadi, A non-adapted sparse approximation of pdes with stochastic inputs, *J. Comput. Phys* 230 (2011) 3015–3034.
- [51] L. Yan, L. Guo, D. Xiu, Stochastic collocation algorithms using ℓ^1 minimization, *Inter. J. Uncertain Quantification* 2 (2012) 279–293.
- [52] H. Rauhut, R. Ward, Sparse legendre expansions via ℓ_1 -minimization, *J. Approx. Theory* 164 (2012) 517–533.
- [53] L. Mathelin, K. A. Gallivan, A compressed sensing approach for partial differential equations with random input data, *Communications in Computational Physics* 12 (4) (2012) 919?954. doi:[10.4208/cicp.151110.090911a](https://doi.org/10.4208/cicp.151110.090911a).
- [54] X. Yang, G. E. Karniadakis, Reweighted ℓ_1 minimization method for stochastic elliptic differential equations, *J. Comput. Phys.* 248 (2013) 87–108.
- [55] J. Hampton, A. Doostan, Compressive sampling of polynomial chaos expansions: Convergence analysis and sampling strategies., *J. Comput. Phys* 280 (2015) 363–386.
- [56] J. Peng, J. Hampton, A. Doostan, On polynomial chaos expansion via gradient-enhanced ℓ_1 -minimization, *J. Comput. Phys* 310 (2016) 440–458.
- [57] L. Yan, Y. Shin, D. Xiu, Sparse approximation using $\ell_1 - \ell_2$ minimization and its application to stochastic collocation, *SIAM Journal on Scientific Computing* 39 (1) (2017) A229–A254. doi:[10.1137/15M103947X](https://doi.org/10.1137/15M103947X).

- [58] Y. L. Liu, L. Guo, Stochastic collocation via l1-minimisation on low discrepancy point sets with application to uncertainty quantification, *EAJAM* 6 (2016) 171–191.
- [59] H. Lei, X. Yang, Z. Li, G. E. Karniadakis, Systematic parameter inference in stochastic mesoscopic modeling, *J. Comput. Phys.* 330 (4) (2017) 571–593.
- [60] N. Alemazkoor, H. Meidani, Divide and conquer: An incremental sparsity promoting compressive sampling approach for polynomial chaos expansions, *Computer Methods in Applied Mechanics and Engineering* 318 (2017) 937 – 956. doi:<https://doi.org/10.1016/j.cma.2017.01.039>.
- [61] P. Diaz, A. Doostan, J. Hampton, Sparse polynomial chaos expansions via compressed sensing and d-optimal design, *Computer Methods in Applied Mechanics and Engineering* 336 (2018) 640 – 666. doi:<https://doi.org/10.1016/j.cma.2018.03.020>.
- [62] P. Rai, K. Sargsyan, H. Najm, Compressed sparse tensor based quadrature for vibrational quantum mechanics integrals, *Computer Methods in Applied Mechanics and Engineering* 336 (2018) 471 – 484. doi:<https://doi.org/10.1016/j.cma.2018.02.026>.
- [63] K. Huang, [Lectures on Statistical Physics and Protein Folding](#), WORLD SCIENTIFIC, 2005. doi: [10.1142/5741](https://doi.org/10.1142/5741).
- [64] S. Kumar, J. M. Rosenberg, D. Bouzida, R. H. Swendsen, P. A. Kollman, The weighted histogram analysis method for free-energy calculations on biomolecules. i. the method, *Journal of Computational Chemistry* 13 (8) (1992) 1011–1021.
- [65] L. Maragliano, E. Vanden-Eijnden, Single-sweep methods for free energy calculations, *The Journal of Chemical Physics* 128 (18) (2008) 184110. doi:[10.1063/1.2907241](https://doi.org/10.1063/1.2907241).
- [66] P. G. Constantine, E. Dow, Q. Wang, Active subspace methods in theory and practice: Applications to kriging surfaces, *SIAM J. Sci. Comput.* 36 (4) (2014) A1500–A1524.
- [67] W. Li, G. Lin, An adaptive importance sampling algorithm for bayesian inversion with multimodal distributions, *Journal of Computational Physics* 294 (2015) 173 – 190.
- [68] V. Vittaldev, R. P. Russell, R. Linares, Spacecraft uncertainty propagation using gaussian mixture models and polynomial chaos expansions, *Journal of Guidance, Control, and Dynamics* 39 (12) (2016) 2615–2626. doi:[10.2514/1.g001571](https://doi.org/10.2514/1.g001571).
- [69] J. Feinberg, H. P. Langtangen, Chaospy: An open source tool for designing methods of uncertainty quantification, *Journal of Computational Science* 11 (2015) 46 – 57. doi:<https://doi.org/10.1016/j.jocs.2015.08.008>.

- [70] J. Wan, N. Zabaras, A probabilistic graphical model based stochastic input model construction, *Journal of Computational Physics* 272 (2014) 664 – 685. doi:<https://doi.org/10.1016/j.jcp.2014.05.002>.
- [71] X. Wan, G. Karniadakis, Multi-element generalized polynomial chaos for arbitrary probability measures, *SIAM J. Sci. Comput.* 28 (2006) 901–928.
- [72] J. A. S. Witteveen, H. Bijl, Modeling arbitrary uncertainties using gram-schmidt polynomial chaos, in: 44th AIAA Aerospace Sciences Meeting and Exhibit, American Institute of Aeronautics and Astronautics, 2006, pp. 1706–1713. doi:[doi:10.2514/6.2006-896](https://doi.org/10.2514/6.2006-896).
- [73] M. Zheng, X. Wan, G. E. Karniadakis, Adaptive multi-element polynomial chaos with discrete measure: Algorithms and application to spdes, *Applied Numerical Mathematics* 90 (2015) 91–110.
- [74] S. Yin, D. Yu, Z. Luo, B. Xia, An arbitrary polynomial chaos expansion approach for response analysis of acoustic systems with epistemic uncertainty, *Computer Methods in Applied Mechanics and Engineering* 332 (2018) 280 – 302. doi:<https://doi.org/10.1016/j.cma.2017.12.025>.
- [75] R. Ahlfeld, B. Belkouchi, F. Montomoli, Samba: Sparse approximation of moment-based arbitrary polynomial chaos, *Journal of Computational Physics* 320 (2016) 1 – 16. doi:<https://doi.org/10.1016/j.jcp.2016.05.014>.
- [76] C. F. Dunkl, Y. Xu, *Orthogonal Polynomials of Several Variables*, 2nd Edition, Encyclopedia of Mathematics and its Applications, Cambridge University Press, 2014. doi:[10.1017/CBO9781107786134](https://doi.org/10.1017/CBO9781107786134).
- [77] E. Candès, M. Rudelson, T. Tao, R. Vershynin, Error correction via linear programming, in: 46th Annual IEEE Symposium on Foundations of Computer Science (FOCS’05), 2005, pp. 668–681.
- [78] E. J. Candès, The restricted isometry property and its implications for compressed sensing, *C. R. Acad. Sci. Paris Sér. I Math.* 346 (2008) 589–592.
- [79] M. E. Davies, R. Gribonval, Restricted isometry constants where ℓ_p sparse recovery can fail for $0 < p \leq 1$, *IEEE Trans. Inf. Theory* 55 (2010) 2203–2214.
- [80] D. Donoho, M. Elad, V. Temlyakov, Stable recovery of sparse overcomplete representations in the presence of noise, *Information Theory, IEEE Transactions on* 52 (1) (2006) 6–18.
- [81] E. V. D. Berg, M. Friedlander, Spgl1: A solver for large-scale sparse reconstruction, <http://www.cs.ubc.ca/labs/scl/spgl>.
- [82] E. J. Candès, T. Tao, Decoding by linear programming, *IEEE Trans. Inform. Theory* 51 (2005) 4203–4215.

- [83] E. J. Candès, J. Romberg, T. Tao, Stable signal recovery from incomplete and inaccurate measurements, *Comm. Pure Appl. Math.* 56 (2006) 1207–1223.
- [84] B. K. Natarajan, Sparse approximate solutions to linear systems, *SIAM J. Sci. Comput.* 2 (1995) 227–234.
- [85] T. M. Russi, [Uncertainty Quantification with Experimental Data and Complex System Models](#), Ph.D. thesis, University of California, Berkeley (2001).
- [86] X. Yang, H. Lei, N. A. Baker, G. Lin, Enhancing sparsity of hermite polynomial expansions by iterative rotations, *J. Comput. Phys.* 307 (2016) 94 – 109.
- [87] M. Jardak, C.-H. Su, G. E. Karniadakis, Spectral polynomial chaos solutions of the stochastic advection equation, *J. Sci. Comput.* 17 (1-4) (2002) 319–338.
- [88] X. Yang, X. Wan, L. Lin, A general framework of enhancing sparsity of generalized polynomial chaos expansions (2017). [arXiv:arXiv:1707.02688](#).
- [89] A. R. Atilgan, S. R. Durell, R. L. Jernigan, M. C. Demirel, O. Keskin, I. Bahar, Anisotropy of fluctuation dynamics of proteins with an elastic network model, *Biophysical Journal* 80 (1) (2001) 505–515.
- [90] P. Ren, J. Chun, D. G. Thomas, M. J. Schnieders, M. Marucho, J. Zhang, N. A. Baker, Biomolecular electrostatics and solvation: a computational perspective, *Quarterly reviews of biophysics* 45 (4) (2012) 427–491. doi:[10.1017/S003358351200011X](https://doi.org/10.1017/S003358351200011X).
- [91] N. A. Baker, Biomolecular Applications of Poisson?Boltzmann Methods, in: K. B. Lipkowitz, R. Larter, T. R. Cundari (Eds.), *Reviews in Computational Chemistry*, John Wiley & Sons, Inc., 2005, pp. 349–379. doi:[10.1002/0471720895.ch5](https://doi.org/10.1002/0471720895.ch5).
- [92] E. Jurrus, D. Engel, K. Star, K. Monson, J. Brandi, L. E. Felberg, D. H. Brookes, L. Wilson, J. Chen, K. Liles, M. Chun, P. Li, D. W. Gohara, T. Dolinsky, R. Konecny, D. R. Koes, J. E. Nielsen, T. Head-Gordon, W. Geng, R. Krasny, G.-W. Wei, M. J. Holst, J. A. McCammon, N. A. Baker, Improvements to the APBS biomolecular solvation software suite, *Protein Science* 27 (1) (2018) 112–128. doi:[10.1002/pro.3280](https://doi.org/10.1002/pro.3280).
- [93] A. Shrake, J. Rupley, Environment and exposure to solvent of protein atoms. lysozyme and insulin, *Journal of Molecular Biology* 79 (2) (1973) 351 – 371.
- [94] M. Rasheed, N. Clement, A. Bhowmick, C. Bajaj, [Statistical Framework for Uncertainty Quantification in Computational Molecular Modeling](#), in: Proceedings of the 7th ACM International Conference on

- 1
2
3
4 Bioinformatics, Computational Biology, and Health Informatics, ACM, New York, NY, USA, 2016,
5 pp. 146–155.
6
7 [95] N. Clement, M. Rasheed, C. L. Bajaj, [Viral Capsid Assembly: A Quantified Uncertainty Approach](#),
8 Journal of Computational Biology 25 (1) (2018) 51–71. doi:[10.1089/cmb.2017.0218](https://doi.org/10.1089/cmb.2017.0218).
9
10 [96] H. Rauhut, Compressive sensing and structured random matrices, Radon Series Comp. Appl. Math. 9
11 (2010) 1–92.
12
13 [97] A. Marrel, B. Iooss, B. Laurent, O. Roustant, [Calculations of Sobol indices for the Gaussian process](#)
14 [metamodel](#), Reliability Engineering & System Safety 94 (3) (2009) 742 – 751.
15
16 [98] G. Chastaing, F. Gamboa, C. Prieur, [Generalized Sobol sensitivity indices for dependent variables: numerical methods](#), Journal of Statistical Computation and Simulation 85 (7) (2015) 1306–1333. doi:
17 [10.1080/00949655.2014.960415](https://doi.org/10.1080/00949655.2014.960415).
18
19 [99] H. J. C. Berendsen, D. van der Spoel, R. van Drunen, GROMACS: A message-passing parallel molecular
20 dynamics implementation, Computer Physics Communications 91 (1) (1995) 43–56. doi:[10.1016/0010-4655\(95\)00042-E](https://doi.org/10.1016/0010-4655(95)00042-E).
21
22
23 [100] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, D. A. Case, Development and testing of a general
24 amber force field, Journal of Computational Chemistry 25 (9) (2004) 1157–1174.
25
26
27 [101] C. I. Bayly, P. Cieplak, W. Cornell, P. A. Kollman, A well-behaved electrostatic potential based method
28 using charge restraints for deriving atomic charges: the resp model, The Journal of Physical Chemistry
29 97 (40) (1993) 10269–10280.
30
31
32 [102] B. Hess, H. Bekker, J. C. Berendsen Herman, G. E. M. Fraaije Johannes, Lincs: A linear constraint
33 solver for molecular simulations, Journal of Computational Chemistry 18 (12) (1998) 1463–1472.
34
35
36 [103] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, M. L. Klein, Comparison of simple
37 potential functions for simulating liquid water, The Journal of Chemical Physics 79 (2) (1983) 926–935.
38
39
40 [104] S. Miyamoto, A. Kollman Peter, Settle: An analytical version of the shake and rattle algorithm for
41 rigid water models, Journal of Computational Chemistry 13 (8) (2004) 952–962.
42
43
44 [105] N. A. Baker, D. Sept, S. Joseph, M. J. Holst, J. A. McCammon, Electrostatics of nanosystems: ap-
45 plication to microtubules and the ribosome, Proceedings of the National Academy of Sciences 98 (18)
46 (2001) 10037–10041.
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

LaTeX Source Files

[Click here to download LaTeX Source Files: revision.zip](#)