

Design optimization of a scramjet under uncertainty using probabilistic learning on manifolds

R. G. Ghanem^{a,*}, C. Soize^b, C. Safta^c, X. Huan^d, G. Lacaze^e, J. C. Oefelein^f, H. N. Najm^c

^aUniversity of Southern California, 210 KAP Hall, Los Angeles, CA 90089, United States

^bUniversité Paris-Est Marne-la-Vallée, MSME, UMR 8208 CNRS, 5 Bd Descartes, 77454 Marne-La-Vallée, Cedex 2, France

^cCombustion Research Facility, Sandia National Laboratories, Livermore, CA 94551, United States

^dDepartment of Mechanical Engineering, University of Michigan, Ann Arbor, MI 48109, United States

^eSpace Exploration Technologies Corporation, Hawthorne, CA 90250, United States

^fSchool of Aerospace Engineering, Georgia Institute of Technology, 270 Ferst Drive NW., Atlanta, GA 30332, United States

Abstract

We demonstrate, on a scramjet combustion problem, a constrained probabilistic learning approach that augments physics-based datasets with realizations that adhere to underlying constraints and scatter. The constraints are captured and delineated through diffusion maps, while the scatter is captured and sampled through a projected stochastic differential equation. The objective function and constraints of the optimization problem are then efficiently framed as non-parametric conditional expectations. Different spatial resolutions of a large-eddy simulation filter are used to explore the robustness of the model to the training dataset and to gain insight into the significance of spatial resolution on optimal design.

Keywords: Scramjet simulations, Machine learning, Uncertainty quantification, Optimization under uncertainty, Sampling on manifolds, Diffusion maps

1. Introduction

Model-based design optimization is a significant capability for systems whose behavior has not yet been catalogued but where functional dependence of system behavior on component behavior can be described through conservation or other such fundamental principles. These systems typically include new designs with no corresponding legacy knowledge and complex systems for which perturbations in any component has consequences on system-level performance. Predictive models for these systems are typically discretizations of mathematical formulations of the underlying fundamental principles resulting in computational models with very significant computational requirements. It is typical for such systems, given their novelty or complexity, to invoke reductions that could be either physical, mathematical, or algorithmic in nature. The impact of these reductions on the predicted performance of any given design can put into question the optimal character of solutions obtained through traditional design optimization procedures. One rational formalism for accounting for these assumptions throughout the design process is to explore the robustness of an optimal solution to perturbations in these assumptions. Probabilistic modeling provides an effective procedure for characterizing these assumptions and has typically been implemented through either parametric procedures where model parameters are described as random variables [1, 2, 3, 4, 5] or within a non-parametric framework where the model itself is described as a random operator [6]. Irrespective of how a non-deterministic problem is implemented, it requires the numerical exploration of a statistical ensemble, thus quickly exacerbating the computational burden of an already massive exercise. Clearly, new perspectives on characterizing and evaluating performance and optimal designs of complex systems are required if physics-based modeling is indeed a necessary ingredient of the prediction and optimization process.

*Corresponding author: R. Ghanem, ghanem@usc.edu

Email addresses: ghanem@usc.edu (R. G. Ghanem), christian.soize@u-pem.fr (C. Soize), csafta@sandia.gov (C. Safta), xhuan@umich.edu (X. Huan), Guilhem.Lacaze@spacex.com (G. Lacaze), joseph.oefelein@aerospace.gatech.edu (J. C. Oefelein), hnnajm@sandia.gov (H. N. Najm)

In this paper we tackle the above challenge in the context of a hypersonic combustion process. Large-eddy simulation (LES) modeling is used to discretize the Navier-Stokes equations, with several of the filtering parameters and operational parameters described as random variables. In the process, we rely on a new methodology developed by the authors that conjoins diffusion maps and projected Itô sampling to augment a small statistical ensemble (the training set) with a large number of realizations that are consistent, in a useful sense, with both underlying mechanisms and evidence. Diffusion maps [7] are used to extract nonlinear manifolds from the training set, while the Itô samplers [8] are used to sample on these manifolds from a specified target probability measure. We construe these manifolds as encoding non-parametric characterization of constraints on the data, including physics-based constraints.

The optimization problem is typically formulated in terms of several quantities of interest (QoIs) that are themselves functionals of the field variables constituting the solution of the LES problem. The map from input parameters to QoIs is a composition of two maps. First is the map represented by the numerical solver, which encodes conservation laws as described by the LES model. This is followed by a postprocessing map that evaluates the QoIs, which enter in the evaluation of objectives and constraints. The LES map is challenged with a computational burden that limits its ability to efficiently explore parameter space, a key requirement for design optimization. The postprocessing map faces its own challenges, namely the lack of apparent physics-based constraints on the QoIs, which forces the characterization of the QoIs in terms of the input parameters and the LES map. These challenges notwithstanding, it can be expected that these two maps do indeed constrain the QoIs, albeit in an unknown and intractable manner. By discovering these constraints, the task of input space exploration would be reduced to exploring a low-dimensional manifold. For that, we propose to use the probabilistic learning on manifolds [9, 4] that relies on diffusion maps (DMAPS), with the corresponding manifold characterized through a positive definite operator and its associated orthogonal projections. Then, statistical sampling on this manifold is readily accomplished as described above, via a projected Itô equation whose invariant measure is a nonparametric statistical estimate of the measure of the training data, represented as a Gaussian mixture model. We compare optimization results obtained from three LES models representing different spatial resolutions, and explore the impact of the size of the training set on the resulting optimal designs.

The paper is organized as follows. In the next section, an overview is presented on the physics problem representing scramjet combustion. Following that, probabilistic learning on manifolds using diffusion maps and Itô sampling is described. The optimization problem under uncertainty is then illustrated and its solution analyzed. Some remarks are presented in a concluding section.

2. Overview of the hypersonic reactive flow model

Our physical application of interest stems from the HIFiRE (Hypersonic International Flight Research and Experimentation) program [10, 11], which has cultivated a mature experimental campaign with accessible data through its HIFiRE Flight 2 (HF2) project [12, 13]. The HF2 payload involves a cavity-based hydrocarbon-fueled dual-mode scramjet and was tested under flight conditions of Mach 6–8+. A ground test rig, designated the HIFiRE Direct Connect Rig (HDCR) (Figure 1(a)), was developed to duplicate the isolator/combustor layout of the flight test hardware, and to provide ground-based data for comparisons with flight measurements, verifying engine performance and operability, and designing fuel delivery schedule [14, 15]. While data from flight tests are not publicly released, HDCR ground test data are available [14, 16]. Therefore, we aim to simulate and assess behavior of reactive flows inside the HDCR, in order to facilitate future computational developments that can make use of these experimental datasets.

The computational domain for the HDCR is highlighted by red lines in Figure 1(b). The rig consists of a constant-area isolator (planar duct) attached to a combustion chamber. It includes primary injectors that are mounted upstream of flame stabilization cavities on both the top and bottom walls. Secondary injectors along both walls are positioned downstream of the cavities. Flow travels from left to right in the x_s -direction (streamwise), and the geometry is symmetric about the centerline in the y_s -direction. Numerical simulations take advantage of this symmetry by considering a domain that covers only the bottom half of this configuration. The fuel supplied through the injectors is a gaseous mixture containing 36% methane and 64% ethylene by volume, which acts as a surrogate with similar combustion properties as JP-7 [17]. A reduced, three-step mechanism [18, 19] is initially employed to characterize the combustion process. Arrhenius formulations of

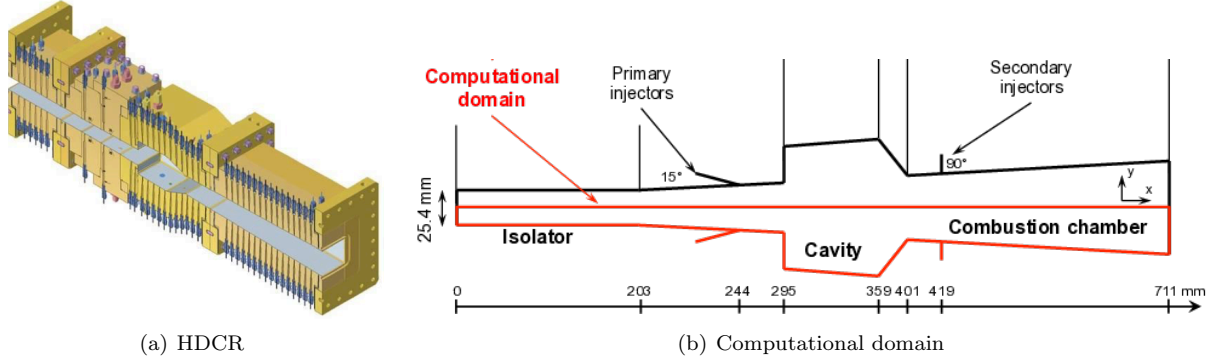


Figure 1: HDCR cut view [14] and schematic of the computational domain.

the kinetic reaction rates are adopted, and the parameters are selected to retain robust/stable combustion in the current simulations.

LES calculations are then performed using the RAPTOR code framework developed by Oefelein [20, 21, 22]. RAPTOR solves the fully coupled conservation equations of mass, momentum, total energy, and species for a chemically reacting flow. It is designed to handle high Reynolds number, high pressure, real gas and liquid conditions over a wide range of Mach numbers while accounting for detailed thermodynamics and transport processes at the molecular level. The numerical schemes in RAPTOR involve non-dissipative, discretely conservative, staggered, finite volume differencing. This scheme is well-adapted to LES simulations as it eliminates numerical contamination of the subfilter due to artificial dissipation and ensures proper discrete conservations. The subfilter closure is obtained using a mixed Smagorinsky model [23, 24].

In our numerical studies, we designate five input variables as design variables that we can directly control during the design optimization process. These design variables are the global equivalence ratio (ϕ_G), ratio of primary to secondary injector equivalence ratios (ϕ_R), location of the primary injector (x_1), location of the secondary injector (x_2), and inclination of the primary injector (θ_1). Specifically, if ϕ_1 and ϕ_2 denote the equivalence ratios from the primary and secondary injectors, respectively, then $\phi_G = \phi_1 + \phi_2$ and $\phi_R = \phi_1 / \phi_2$. The feasible design ranges are shown in Table 5. Furthermore, we allow a total of 11 model parameters to be uncertain, shown in Table 1 along with their uncertainty distributions. These distributions are assumed uniform across the ranges indicated.

The data utilized in the current analysis are from two-dimensional simulations of the scramjet computation, using grid resolutions where cell sizes are 1/8, 1/16, and 1/32 of the injector diameter $d = 3.175$ mm (respectively denoted by “d/8”, “d/16”, and “d/32” cases in this paper), corresponding to respectively around 63 thousand, 250 thousand, and 1 million grid points. The number of time steps for each run are selected to maintain an approximately equal wall-clock time, while timestep sizes are determined adaptively based on the Courant-Friedrichs-Lewy (CFL) criterion. The simulations are warm-started on solutions engineered from a quasi-steady state nominal condition simulation, and take around 1.7×10^3 , 1.1×10^4 , and 3.9×10^4 CPU hours per run for d/8, d/16, and d/32, respectively. The intense computational demand thus limits us with a total of 1053, 222, and 23 simulation runs for d/8, d/16, and d/32, respectively, where the simulation inputs are uniformly randomized jointly in the parameter and design spaces. Several QoIs are considered in this paper, either computed directly in the LES code or evaluated subsequently through postprocessing. All QoIs are time-averaged variables, where the instantaneous solutions corresponding to the second half for each run are time-averaged to generate one solution per run. The entire set of QoIs are described below, and while only a subset (in Table 4) directly enter the optimization problem definition, the rest (in Table 2) are retained to assist the discovery of the probabilistic manifold.

- **Combustion efficiency** (η_{comb}), a critical performance indicator for engines, is defined based on static enthalpy quantities [15, 25]:

$$\eta_{\text{comb}} = \frac{H(T_{\text{ref}}, Y_e) - H(T_{\text{ref}}, Y_{\text{ref}})}{H(T_{\text{ref}}, Y_{e, \text{ideal}}) - H(T_{\text{ref}}, Y_{\text{ref}})}. \quad (1)$$

Here H is the total static enthalpy, the “ref” subscript indicates a reference condition derived from the inputs, the “e” subscript is for the exit, and the “ideal” subscript is for the ideal condition where all fuel is burnt to completion. The reference condition corresponds to that of a hypothetical non-reacting mixture of all inlet air and fuel at thermal equilibrium. The numerator, $H(T_{\text{ref}}, Y_e) - H(T_{\text{ref}}, Y_{\text{ref}})$, thus reflects the global heat released during the combustion, while the denominator represents the total heat release available in the fuel-air mixture.

- **Burned equivalence ratio** (ϕ_{burn}) is defined to be equal to $\phi_{\text{burn}} \equiv \phi_G \eta_{\text{comb}}$. It represents the air excess, and high values of ϕ_{burn} results from a combination of high thermal efficiencies and stoichiometric to rich equivalence ratios, and are associated with conditions away from blowout regimes.
- **Stagnation pressure loss ratio** (P_{stagloss}) is defined as

$$P_{\text{stagloss}} = 1 - \frac{P_{s,e}}{P_{s,i}}, \quad (2)$$

where $P_{s,e}$ and $P_{s,i}$ are the wall-normal-averaged stagnation pressure quantities at the exit and inlet planes, respectively. Higher values of P_{stagloss} illustrate pressure loss across the combustor and are associated with a decrease in efficiency.

- **Maximum root-mean-square (RMS) pressure** ($\max P_{\text{rms}}$) is the maximum RMS pressure across the spatial domain:

$$\max P_{\text{rms}} = \max_{x,y} \sqrt{P(x,y)^2 - [\overline{P(x,y)}]^2}, \quad (3)$$

with \overline{P} indicating time-averaged quantity. This QoI reflects the maximum pressure oscillation amplitude, which is useful for engine structural considerations.

- **Turbulence kinetic energy (TKE)** is characterized by the RMS velocity at a given location:

$$\text{TKE} = \frac{1}{2} (u_{\text{rms}}^2 + v_{\text{rms}}^2). \quad (4)$$

In the numerical investigations of this paper, we will look at TKE from multiple streamwise locations at $x_s/d = 5, 50, 85, 110, 140, 190, 220$. TKE captures statistical signatures of the spatial heterogeneity of turbulence.

- **Initial shock location** (x_{shock}) is the most upstream shock location, which we currently compute by detecting a rapid pressure change. More upstream locations of the initial shock train are linked with higher chances of unstart and vibration, and degraded safety of the scramjet engine.

3. Probabilistic Learning on Manifolds

The initial data set, denoted by $[x_d]$, consists of N samples of \mathbb{R}^n -valued vectors, and is represented as an $n \times N$ matrix. Here, n denotes the number of observables for each sample. The data can be thought of as N points residing on a manifold in a n -dimensional ambient space. Alternatively, the data can be viewed as n points in N -dimensional space. Through a diffusion maps analysis [7, 26, 27], the manifold localization in \mathbb{R}^n is restated as subspace localization in \mathbb{R}^N , providing a more structured context for analysis. We first apply an affine transformation on the data using the eigenvectors of the empirical covariance matrix, transforming matrix $[x_d]$ into a new matrix $[\eta_d]$ whose rows are orthogonal. The purpose of this transformation is not to reduce dimensionality, but rather to improve the numerical conditioning of the data. In some instances, dimension reduction may ensue, reflecting very strong linear correlations in the data. Thus in general, matrix $[\eta_d]$ will be an $\nu \times N$ matrix, with $\nu \leq n$. In the next subsection, we provide a brief overview of the diffusion maps construction, followed by a brief overview of the projected Itô equation to sample on the associated manifolds.

Parameter	Range	Description
Inlet boundary conditions:		
p_0	$[1.406, 1.554] \times 10^6$ Pa	Stagnation pressure
T_0	$[1472.5, 1627.5]$ K	Stagnation temperature
M_0	$[2.259, 2.761]$	Mach number
I_i	$[0, 0.05]$	Turbulence intensity horizontal component
R_i	$[0.8, 1.2]$	Ratio of turbulence intensity vertical to horizontal components
L_i	$[0, 8] \times 10^{-3}$ m	Turbulence length scale
Fuel inflow boundary conditions:		
I_f	$[0, 0.05]$	Turbulence intensity magnitude
L_f	$[0, 1] \times 10^{-3}$ m	Turbulence length scale
Turbulence model parameters:		
C_R	$[0.01, 0.06]$	Modified Smagorinsky constant
Pr_t	$[0.5, 1.7]$	Turbulent Prandtl number
Sc_t	$[0.5, 1.7]$	Turbulent Schmidt number

Table 1: Uncertain input parameters. The uncertainty distributions are assumed uniform across the ranges shown.

Observable	Description
TKE_{x_s}	Turbulent kinetic energy at $x_s/d=5, 50, 85, 110, 140, 190, 220$
x_{shock}	Upstream shock location

Table 2: Additional observables from simulation output that do not directly enter the optimization formulation but are used for detecting the manifold. x_s is streamwise distance.

3.1. Diffusion Maps

The starting point for a DMAPS analysis is a symmetric, positivity preserving and positive semi-definite kernel, $k(x, y)$, used to analyze the data. We rely on a positive-definite Gaussian kernel of the form,

$$k(x, y; \epsilon) = e^{-\|x-y\|^2/\epsilon}, \quad x, y \in \mathbb{R}^\nu, \quad (5)$$

where ϵ denotes the kernel width. The $N \times N$ matrix $[K]$ is then constructed on the data $[\eta_d]$ as follows,

$$[K]_{ij}^\epsilon = k(\eta^{d,i}, \eta^{d,j}; \epsilon) \quad (6)$$

where $\eta^{d,i}$ is the i^{th} columns of $[\eta_d]$. By normalizing $[K]$ as follows,

$$[P] = [b]^{-1}[K], \quad [b]_{ij} = \sum_{\ell=1}^{\nu} [K]_{i\ell} \delta_{ij} \quad (7)$$

the resulting matrix $[P]$ can be construed as the transition matrix of a random walk on the graph associated with the data [7]. The right eigenvector g^α of $[P]$ such that $[P]g^\alpha = \lambda_\alpha g^\alpha$ can be written as $g^\alpha = [b]^{-1/2}\varphi^\alpha$ in which λ_α and φ^α are the eigenvalues and eigenvectors of the symmetric positive-definite matrix $[P_s] = [b]^{1/2}[P][b]^{-1/2} = [b]^{-1/2}[K][b]^{-1/2}$. Consequently, the family of vectors $\{g^\alpha\}_\alpha$ spans \mathbb{R}^N . More importantly, sorting the eigenvalues such that $1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_m$, it can be shown that the original data set $[\eta_d]$ is localized within the span of the first few m of these eigenvectors $\{g^\alpha, \alpha = 1 \dots, m\}$ [28, 7]. We also denote by $[g]$ the matrix whose α^{th} column is $g^\alpha \in \mathbb{R}^N$. The value of m is chosen to correspond to a noticeable reduction in the magnitude of the eigenvalues, specifically

$$m = \arg \min_{\alpha} \left\{ \frac{\lambda_\alpha}{\lambda_2} < L \right\} \quad (8)$$

with L chosen to be equal to 0.1. Typically, the ratio λ_m/λ_2 is of the order of 0.1. Clearly, these eigenvalues and hence the optimal value of m depend on ϵ , the width of the diffusion kernel. A maximum entropy argument [29] is used to simultaneously select the values of ϵ and m .

3.2. Itô Sampler on Manifolds

We construe localization to the manifold as satisfying underlying constraints, including the prevailing physics. This is justified by the fact that the information used to delineate the manifold is obtained from LES simulations. Any structure that relates the observables is thus inherited from the underlying physics. Our next task is to sample on this manifold in a manner that is consistent with the scatter originally observed in the data. In order to generate additional samples on the manifold, we first construct a probability model from the initial training data set, $[\eta_d]$. Two assumptions are required for this construction [9]. We first assume that the N columns of $[\eta_d]$ are independent realizations of an \mathbb{R}^ν -valued random vector H whose empirical covariance matrix is the identity matrix, and estimate its probability density function (PDF) as a Gaussian mixture in the form,

$$p_H(\eta) = \frac{1}{N} \sum_{j=1}^N \pi \left(\frac{\hat{s}_\nu}{s_\nu} \eta^{d,j} - \eta \right), \quad (9)$$

where π is the positive function from \mathbb{R}^ν into $]0, +\infty[$ defined, for all η in \mathbb{R}^ν , by

$$\pi(\eta) = \frac{1}{(\sqrt{2\pi} \hat{s}_\nu)^\nu} \exp \left\{ -\frac{1}{2\hat{s}_\nu^2} \|\eta\|^2 \right\}, \quad (10)$$

with $\|\eta\|$ denoting the Euclidean norm in \mathbb{R}^ν and where the positive parameters s_ν and \hat{s}_ν are defined by

$$s_\nu = \left\{ \frac{4}{N(2+\nu)} \right\}^{1/(\nu+4)}, \quad \hat{s}_\nu = \frac{s_\nu}{\sqrt{s_\nu^2 + \frac{N-1}{N}}}. \quad (11)$$

With this choice of s_ν and \hat{s}_ν the mean-squared error is minimized [30] and realizations of random vector H are normalized, a requirement consistent with their construction through an eigen-decomposition [8]. This is the PDF of random vector H characterized by equation (9), we now consider the joint occurrence of the N data points. This joint behavior is significant as it carries a signature of intrinsic structure not available in each data point separately. We are looking for structure beyond linear correlation. We thus consider matrix $[\eta_d]$ as a realization of a random matrix $[H]$, for which we next construct a probability model. We now invoke our second assumption, whereby we consider the N columns of $[H]$ as statistically independent, with the density of $[H]$ given by,

$$p_{[H]}([\eta]) = p_H(\eta^1) \times \dots \times p_H(\eta^N). \quad (12)$$

We thus obtain a nonparametric Gaussian mixture model for the PDF of random matrix $[H]$. Each realization of this random matrix will augment the initial training set $[\eta_d]$ with N new data points each of dimension ν . Alternatively, these realizations are first transformed through the eigenvectors of the empirical covariance of the original data, and are thus used to augment matrix $[x_d]$ with N new data points, each of dimension n . We next describe the procedure for generating samples of $[H]$ from the PDF specified in Equation (12).

The approach consists of constructing an Itô equation that is constrained to the manifold, through orthogonal projections, and whose invariant measure has the density specified by Equation (12). First, we introduce the orthogonal projection $[a]$ on the subspace spanned by the $[g]$,

$$[a] = [g] ([g]^T [g])^{-1}, \quad (13)$$

the $\nu \times N$ matrix $[\mathcal{N}]$ whose entries are independent standard gaussian variables, the $\nu \times N$ random matrix $[H_d]$ with known realization is $[\eta_d]$, and the $\nu \times N$ matrix $[dW(r)]$ ($r \geq 0$) whose i^{th} column is dW^i with $\{W^i, i = 1, \dots, N\}$ being independent copies of the ν -dimensional normalized Wiener process. It can then be shown that solutions $\{Z(r), r \geq 0\}$ of the following Itô stochastic differential equations [9]

$$d[Z(r)] = [Y(r)] dr, \quad (14)$$

$$d[Y(r)] = [L([Z(r)] [g]^T)] [a] dr - \frac{1}{2} f_0 [Y(r)] dr + \sqrt{f_0} [dW(r)] [a], \quad (15)$$

with the initial condition

$$[Z(0)] = [H_d][a] \quad , \quad [Y(0)] = [\mathcal{N}][a] \quad a.s. \quad , \quad (16)$$

are samples from the $\nu \times N$ random matrix $[H] = [Z][g]^T$ with PDF $p_{[H]}([u]) = c q([u])$ in which c is a constant of normalization, and where

$$[L([u])]_{k\ell} = \frac{\partial}{\partial u_k^\ell} \log\{q(u^\ell)\} \quad , \quad [u] = [u^1, \dots, u^N] \quad . \quad (17)$$

Given our choice of Gaussian mixture model for q , the expression for $[L]$ can be expanded as follows,

$$[L([u])]_{k\ell} = \frac{1}{q(u^\ell)} \frac{1}{\hat{s}_\nu^2} \frac{1}{N} \sum_{j=1}^N \left(\frac{\hat{s}_\nu}{s_\nu} \eta^{d,j} - u^\ell \right) \exp \left\{ -\frac{1}{2\hat{s}_\nu^2} \left\| \frac{\hat{s}_\nu}{s_\nu} \eta^{d,j} - u^\ell \right\|^2 \right\} \quad . \quad (18)$$

The Itô equations specified by Equations (14) and (15) are solved using a Störmer-Verlet algorithm, a symplectic scheme well-adapted to Hamiltonian non-dissipative systems [31]. In that scheme, we use a value of f_0 equal to 0.43, and an integration step of 0.02.

3.3. Conditional Expectations

All the observables are used to characterize the basis vectors for the diffusion maps, and they all contribute to its construction. The greater the number of observables, the more distinct the interdependence between them, and the more robust is the associated manifold. Only a subset of these observables (QoIs), however, is relevant to the design optimization problem. Thus, once the manifold has been characterized and the initial training set augmented as described in the previous section, new observables relevant to the objective functions and constraints are constructed from the QoIs and conjoined to the control observables to construct a conditional expectation engine to be used in solving the optimization problem.

Thus, denoting the m_q QoIs by \mathbf{Q} and the m_w control variables by \mathbf{W} , either the objective function or each of the constraints can be expressed as a conditional expectation of the generic form,

$$I = \mathbb{E}\{R | \mathbf{W} = \mathbf{w}_o\} \quad (19)$$

in which \mathbb{E} is the mathematical expectation, and where R is a scalar function of \mathbf{Q} and \mathbf{W} . We next describe efficient procedures for evaluating this expectation using the Gaussian kernel-density estimation method. We first introduce the following normalizations,

$$\hat{R} = (R - \overline{R})/\sigma_R, \quad \hat{W}_j = (W_j - \overline{W}_j)/\sigma_j, \quad j = 1, \dots, m_w \quad (20)$$

where an overline denotes the mean of a random variable, σ_R is the standard deviation of R and σ_j the standard deviation of W_j . These means and standard deviations are estimated using all the additional ν_s samples synthesized by the projected Itô sampler. These samples are also used to construct a nonparametric model of the joint density function of \hat{R} and $\hat{\mathbf{W}}$ of the following form [30]

$$p_{\hat{\mathbf{W}}, \hat{R}}(\hat{\mathbf{w}}_o, \hat{r}) \simeq \frac{1}{\nu_s} \sum_{\ell=1}^{\nu_s} \frac{1}{(\sqrt{2\pi}s)^{m_w+1}} \exp \left\{ -\frac{1}{2s^2} \left\{ \|\hat{\mathbf{w}}^\ell - \hat{\mathbf{w}}_o\|^2 + (\hat{r}^\ell - \hat{r})^2 \right\} \right\} \quad , \quad (21)$$

where the kernel width is given by [30]

$$s = \left(\frac{4}{\nu_s(2 + m_w + 1)} \right)^{1/(4+m_w+1)} \quad . \quad (22)$$

The choice of separable kernel in Equation (21) presumes local statistical independence between R and \mathbf{W} . While this assumption does not generally hold, it has diminishing influence as ν_s becomes very large, which is presently the case. A nonparametric regression of random variable R on the control variables can then be obtained in the form [4, 30],

$$\mathbb{E}\{R | \mathbf{W} = \mathbf{w}_o\} \simeq \frac{\sum_{\ell=1}^{\nu_s} \hat{r}^\ell \exp \left\{ -\frac{1}{2s^2} \|\hat{\mathbf{w}}^\ell - \hat{\mathbf{w}}_o\|^2 \right\}}{\sum_{\ell=1}^{\nu_s} \exp \left\{ -\frac{1}{2s^2} \|\hat{\mathbf{w}}^\ell - \hat{\mathbf{w}}_o\|^2 \right\}} \sigma_R + \overline{R} \quad . \quad (23)$$

<i>Resolution</i>	<i>N=Size of training set</i>							ϵ	n	ν	m
d/8	50	100	200	500	1053			441	28	24	(25-29)
d/16	50	100	150	222				441	28	25	27
d/32	10	12	14	16	18	20	23	159	9	8	9

Table 3: Spatial resolution and size of training set for all cases; ϵ =width of diffusion kernel; n = number of observables per sample; ν =dimension of de-correlated variables; m =dimension of diffusion map.

In the above, \hat{r}^ℓ and $\hat{\mathbf{w}}^\ell$ are realizations of \hat{R} and $\hat{\mathbf{W}}$, respectively. In evaluating chance constraints, it is often required to compute the probability of a random variable R exceeding some threshold value, q_f , conditional on a specified values w_o of the control variables. This probability can be expressed as the following integral,

$$h(q_f, \mathbf{w}_o) = \mathbb{P}[R > q_f \mid \mathbf{W} = \mathbf{w}_o] \simeq \int_{q_f}^{+\infty} p_{R|\mathbf{W}}(r|\mathbf{w}_o) dr, \quad (24)$$

which can be estimated as [5]

$$h(q_f, \mathbf{w}_o) \simeq \frac{\sum_{\ell=1}^{\nu_s} h^\ell(\hat{q}_f) \exp \left\{ -\frac{1}{2s^2} \|\hat{\mathbf{w}}^\ell - \hat{\mathbf{w}}_o\|^2 \right\}}{\sum_{\ell=1}^{\nu_s} \exp \left\{ -\frac{1}{2s^2} \|\hat{\mathbf{w}}^\ell - \hat{\mathbf{w}}_o\|^2 \right\}} \quad (25)$$

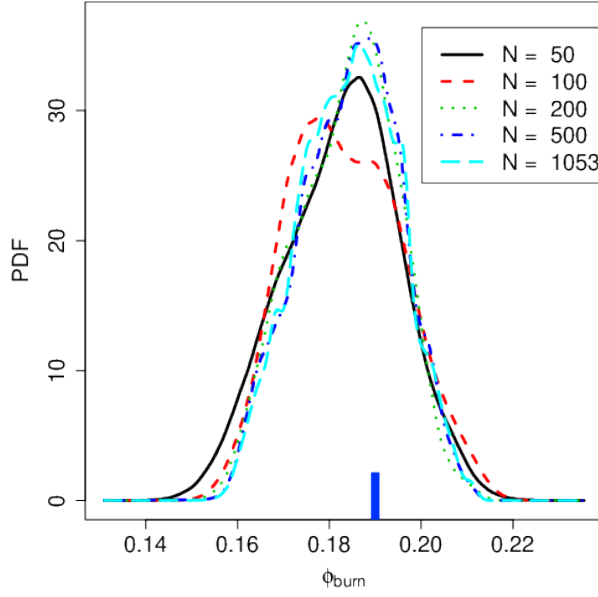
where

$$\hat{q}_f = (q_f - \bar{R})/\sigma, \quad h^\ell(\hat{q}_f) = \frac{1}{2} \left(1 - \text{erf} \left((\hat{q}_f - \hat{r}^\ell)/(s\sqrt{2}) \right) \right) \quad (26)$$

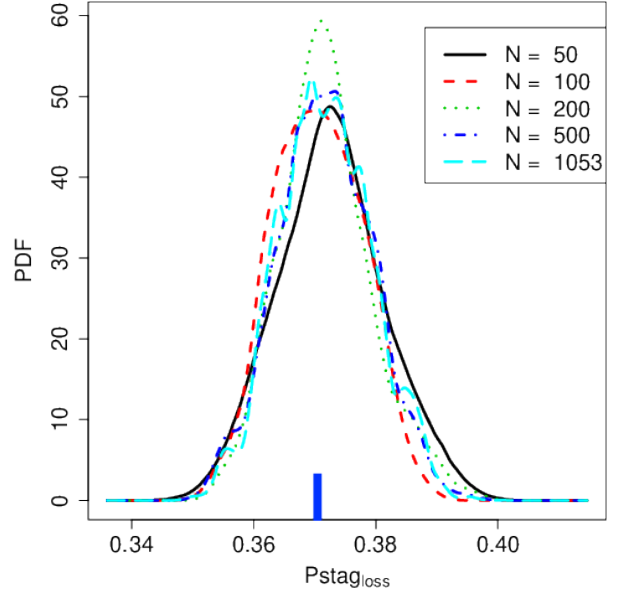
and $\text{erf}(y) = \frac{2}{\sqrt{\pi}} \int_0^y e^{-t^2} dt$ is the error function.

The augmented datasets were generated for each of the three spatial resolutions, d/8, d/16 and d/32. For each of these three cases, different values of N , the size of the training set, were used to delineate different manifolds and integrate forward the associated projected Itô sampler. The various cases are shown in Table (3). For d/8 and d/16, the 28 observables shown in Tables (1), (2), (4), and (5) are included in the manifold detection process. For the case d/32, given the small number of training samples (maximum of 23), only 9 observables were included in the analysis, consisting of the 4 QoIs and 5 controls, shown in Tables (4) and (5), respectively. The value of ν , representing the dimension of the decorrelated observations is also shown in table (3). A threshold of 99% was used for truncation in this representation, and the ensuing slight reduction in the dimension of observables is indicative of strong linear dependence of one or two variables on the remaining variables. The value of m shown in Table (3) represents the number of eigenvectors retained in the characterization of the diffusion manifold following the criterion specified by Equation (8). For each of the cases shown in this table, a total of n_{NMC} samples of $[Z(r)]$ were generated according to the manifold sampling scheme [9] summarized earlier in the paper, such that $\nu \times N \times n_{\text{NMC}} = 5 \times 10^6$. Conditional expectations and conditional probabilities expressed in Equations (23) and (25) respectively, are then estimated using these samples.

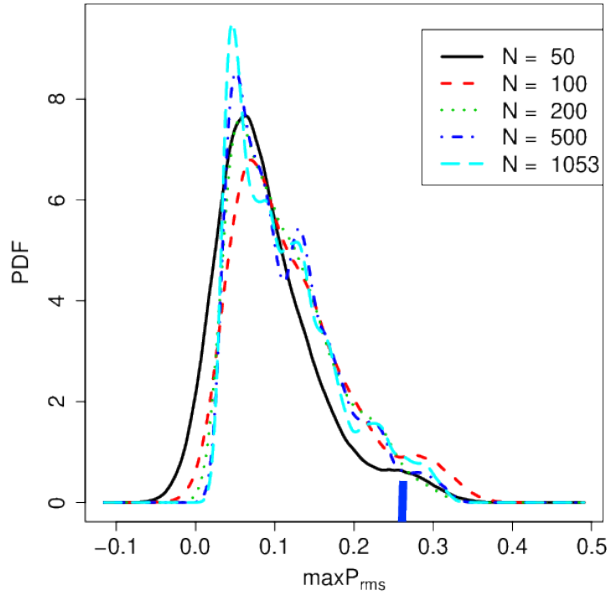
Figure (2) shows the marginal PDFs for resolution d/8 as the size N of the training set is increased monotonically over the set (50, 100, 200, 500, 1053). The vertical marker in Figure (2)-(a) indicates the lower bound on the value of ϕ_{burn} as specified subsequently for design optimization. The vertical markers in Figures (2b) and (2-c) indicate the upper bounds on P_{stagloss} and $\max P_{\text{rms}}$, respectively. The vertical bar in Figure (2-d) shows the range of the computed objective function specified in the optimization problem. This is further explained in the next section. The apparent convergence of these PDFs as N increases suggests that the probability measure is eventually supported by the manifold, with the projection inducing only a small discrepancy. Figures (3) and (4) show similar results for spatial resolutions of d/16 and d/32, respectively. It is noted that there are clear differences between the PDFs at different spatial resolutions. This indicates that further mesh refinement, beyond d/32, is necessary to arrive at a fully converged mesh solution. Some of these PDFs, for instance, exhibit bimodal behavior, symmetry, or skewness at one resolution but not at the others.



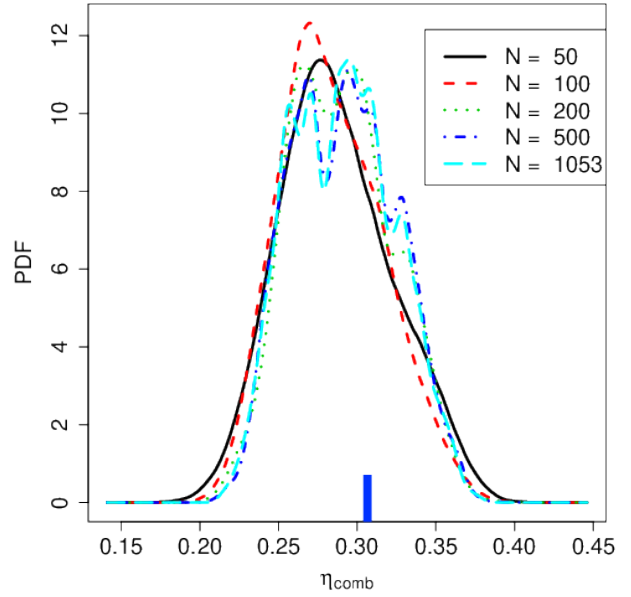
(a) Burned equivalence ratio, d/8, ϕ_{burn}



(b) Stagnation pressure loss, d/8, $P_{stagloss}$

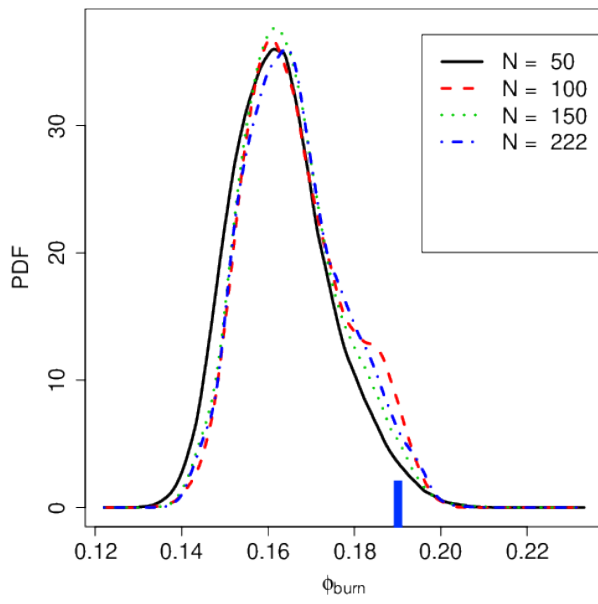


(c) Maximum RMS Pressure, d/8, $\max P_{rms}$

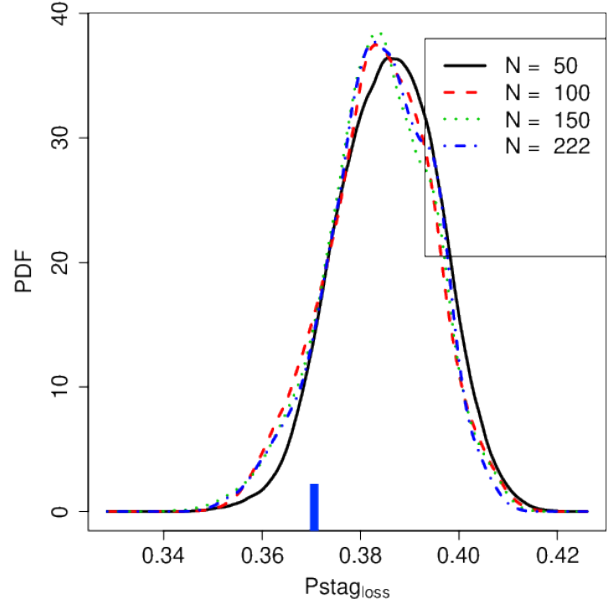


(d) Combustion efficiency, d/8, η_{comb}

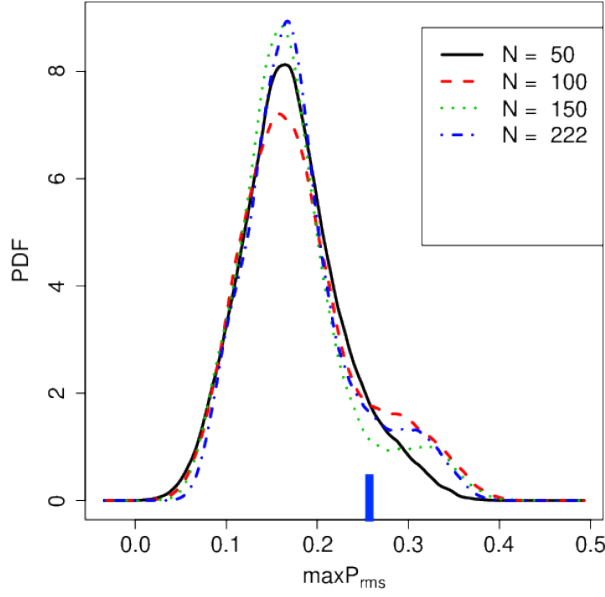
Figure 2: Probability density functions of four quantities of interest for resolution d/8. Each figure shows results for a range of values of N , the training-set size. Vertical markers indicate lower bound for optimization in subfigure a), upper bounds in subfigures b) and c) (Table (4)); vertical marker in subfigure d) shows computed range of optimal objective function as N varies.



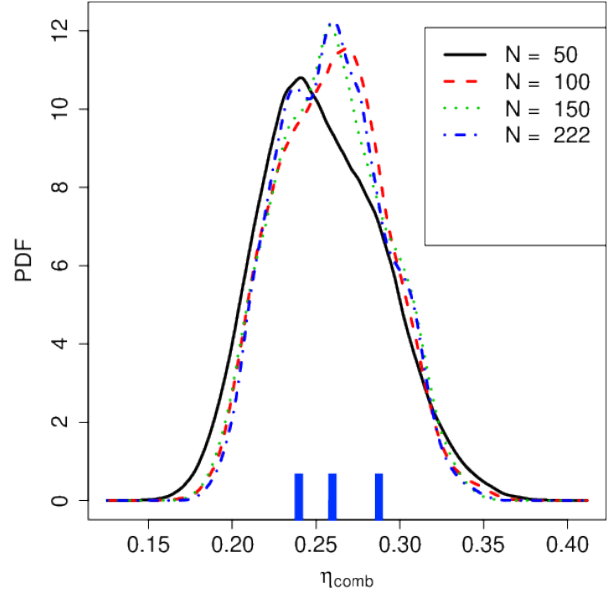
(a) Burned equivalence ratio, d/16, ϕ_{burn}



(b) Stagnation pressure Loss, d/16, $P_{stagloss}$

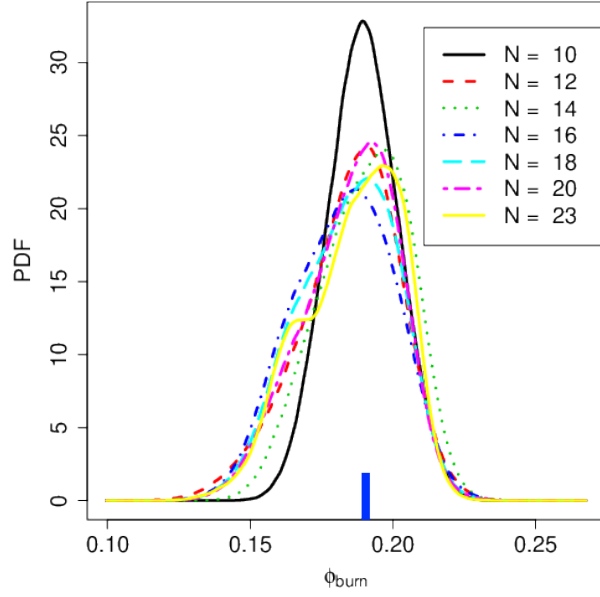


(c) Maximum RMS Pressure, d/16, $\max P_{rms}$

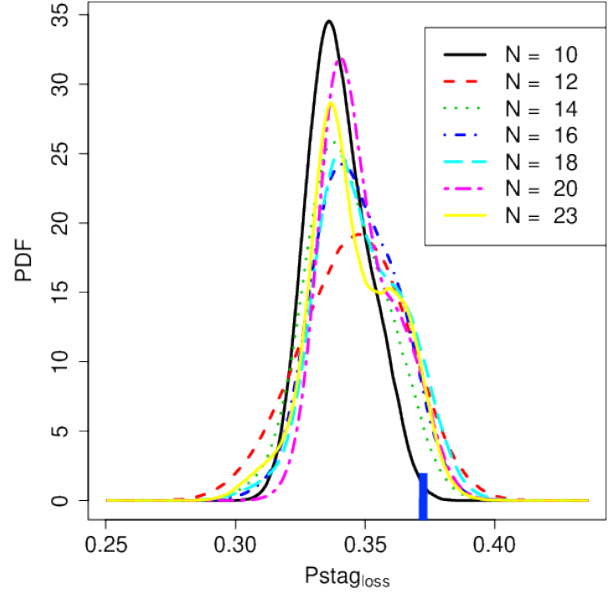


(d) Combustion efficiency, d/16, η_{comb}

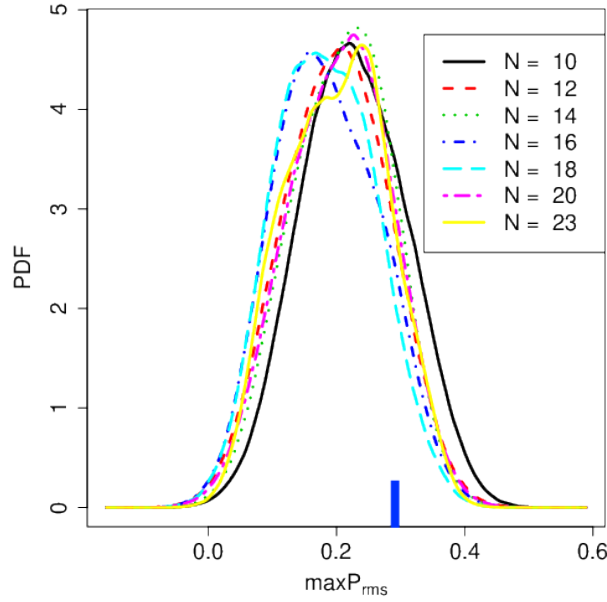
Figure 3: Probability density functions of four quantities of interest for resolution d/16. Each figure shows results for a range of values of N , the training-set size. Vertical markers indicate lower bound for optimization in subfigure a) , upper bounds in subfigures b) and c) (Table (4)); vertical marker in subfigure d) shows computed range of optimal objective function as N varies.



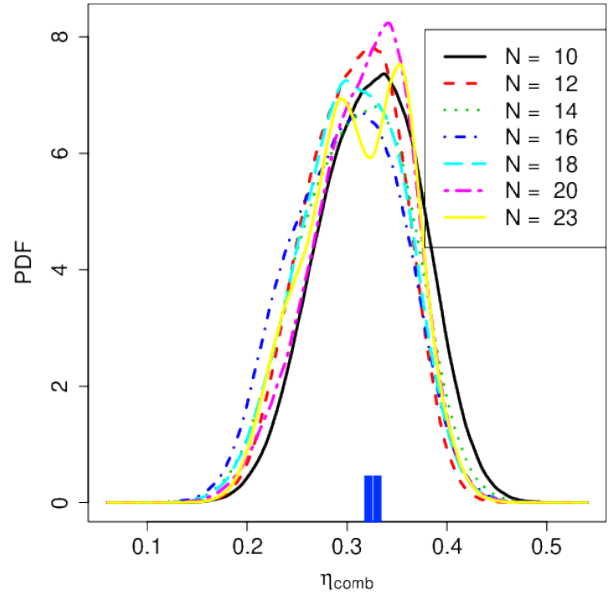
(a) Burned equivalence ratio, $d/32$, ϕ_{burn}



(b) Stagnation pressure loss, $d/32$, $P_{stagloss}$



(c) Maximum RMS Pressure, $d/32$, $\max P_{rms}$



(d) Combustion efficiency, $d/32$, η_{comb}

Figure 4: Probability density functions of four quantities of interest for resolution $d/32$. Each figure shows results for a range of values of N , the training-set size. Vertical markers indicate lower bound for optimization in subfigure a), upper bounds in subfigures b) and c) (Table (4)); vertical marker in subfigure d) shows computed range of optimal objective function as N varies.

<i>Variable name</i>	<i>Symbol</i>	<i>Lower/Upper bound</i>
Q_1 =Burned equivalence ratio	ϕ_{burn}	$\geq L_1 = 0.1900$
Q_2 =Combustion efficiency	η_{comb}	to be maximized
Q_3 =Stagnation pressure loss	P_{stagloss}	$\leq U_3 = 0.3700$
Q_4 =Maximum RMS pressure	$\max P_{\text{rms}}$	$\leq U_4 = 0.2705$

Table 4: Quantities of interest

<i>Variable name</i>	<i>Symbol</i>	<i>Range</i>
w_1 =Global equivalence ratio	ϕ_G	[0.5-1.0]
w_2 =Primary-secondary ratio	ϕ_R	[0.25-0.35]
w_3 =Primary injector location	x_1	[0.231-0.2564](m)
w_4 =Secondary injector location	x_2	[0.40755-0.43295](m)
w_5 =Primary injector angle	θ_1	[5-25] $^\circ$

Table 5: Control variables

4. Optimization Problem of Scramjet with Uncertainty

Our objective is to maximize the combustion efficiency η_{comb} , subject to constraints on the burned equivalence ratio ϕ_{burn} , the stagnation pressure loss across the combustor P_{stagloss} , and the maximum pressure RMS $\max P_{\text{rms}}$. These QoIs, which we denote by Q_i , are shown in Table (4) together with their respective ranges.

The five control variables, denoted by w_i , used to attain the target design are shown in Table (5) together with the respective feasible ranges. We also denote by $\mathbb{P}(E \mid Q_2)$ the probability of event E conditional on Q_2 associated with the current iteration in the optimization process. We reinterpret the bounds on the QoIs, shown in Table (4) to be chance constraints to be satisfied with a specified probability greater than $1 - \alpha$, and replace these bounds with the following equations,

$$\mathbb{P}\{Q_1(w) > L_1 \mid Q_2\} > 1 - \alpha \quad (27a)$$

$$\mathbb{P}\{Q_3(w) < U_3 \mid Q_2\} > 1 - \alpha \quad (27b)$$

$$\mathbb{P}\{Q_4(w) < U_4 \mid Q_2\} > 1 - \alpha \quad (27c)$$

We are thus dealing with the following optimization problem:

$$w^{\text{opt}} = \arg \min_{w \in C_w \subset \mathbb{R}^5} J(w) \quad (28a)$$

$$\text{subject to} \quad c_i < 0 \quad i = 1, 2, 3 \quad (28b)$$

where the objective function is given by the following expectation,

$$J(w) = \mathbb{E}[Q_2], \quad (29)$$

the feasible domain, C_w , is the 5-dimensional cube specified by the following ranges,

$$\begin{aligned} w_1 &\in [0.5, 1] \\ w_2 &\in [0.25, 0.35] \\ w_3 &\in [0.231, 0.2564] \\ w_4 &\in [0.40755, 0.43295] \\ w_5 &\in [5, 25] \end{aligned} \quad (30)$$

and the constraints are specified as

$$c_1(w^{\text{opt}}) = 1 - \alpha - \mathbb{P}\{Q_1(w^{\text{opt}}) > L_1 \mid Q_2\} \quad (31a)$$

$$c_2(w^{\text{opt}}) = 1 - \alpha - \mathbb{P}\{Q_3(w^{\text{opt}}) < U_3 \mid Q_2\} \quad (31b)$$

$$c_3(w^{\text{opt}}) = 1 - \alpha - \mathbb{P}\{Q_4(w^{\text{opt}}) < U_4 \mid Q_2\} \quad (31c)$$

The optimization problem is formulated and solved for each of the cases shown in Table (3). It is noted from Figures (2)-(4) that for $d/16$, the lower bound on ϕ_{burn} and the upper bound on P_{stagloss} are both in the tail area, with a small likelihood of being satisfied. For $d/8$ and $d/32$, on the other hand, the bounds are situated centrally within the support of the distribution.

To solve this non-convex constrained optimization problem we rely on GA [32], a genetic algorithm package within the R language [33]. We set the population size within each generation to 50, and the maximum number of iterations to 1000. We also set the number of unchanged iterations at convergence to 100. The probability of mutation in a parent chromosome and the number of best fitness individuals to survive at each generation (elitism) are set to 0.8 and 2.5, respectively. Finally the number of consecutive generations without improvement that characterizes the optimal solution is set to 100. In all the results shown below, an optimal solution was attained without activating a stopping criterion based on the maximum number of iterations (set to 1000). The optimization problem was solved on multicore CPUs taking advantage of the multicore feature of GA.

Following a first set of GA optimization runs using a randomly selected initial population, a second set of GA runs was performed on all cases in Table (3) using all the results from the first set as an initial population for all cases. This was done in order to test whether a unique optimal solution exists across all resolutions, which was not the case. Figure (5) shows the values of the optimal solutions as evaluated by the genetic optimization algorithm. Each of the 5 subfigures there shows the values of one of the 5 control variables. Within each subfigure, the red circles indicate solutions associated with $d/32$ resolution, the green squares show solutions for the $d/16$ resolution, and the blue losanges correspond to the $d/8$ resolution. Within each resolution, results are shown for all the values of N (size of initial training set) indicated in Table (3). We note immediately that in most cases, the optimal values of the control variables are evaluated in the neighborhood of their bounds.

The vertical markers in Figures (2(d))-4(d)) indicate the range of the objective function $\mathbb{E}(\eta_{\text{comb}})$ at the optimal solution. This range corresponds to the various values of N . It is noted that, while the optimal value of the objective function is robust to changes in N for both the $d/8$ and $d/32$ resolutions, it is very sensitive to N for the $d/16$ resolution; this in spite of the seemingly converged behavior of the marginal density functions, as function of N , in Figure (3(d)). It is clear from the foregoing observations that performance and convergence criteria based on marginal density functions may be misleading when exploring high-dimensional problems. In general, the conditional expectations used in the optimization problem probe higher order joint density functions of the QoI and the control variables. The accurate estimation of these requires significantly larger number of samples, which is facilitated by the manifold sampling approach used in the present paper.

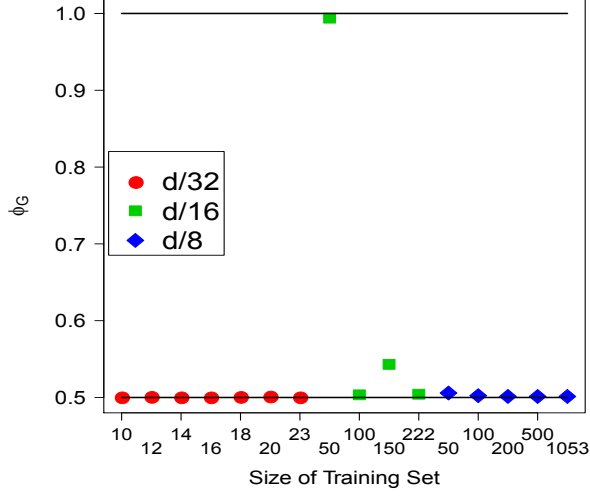
5. Conclusion

The paper presents the application of a recent machine learning algorithm to a very large scale LES-based optimization problem associated with scramjet combustion. The approach used relies on learning a joint statistical model of the observables which include QoIs, input parameters (including design parameters), as well as other observables. The scatter is construed to occur around an intrinsic structure to this dataset that relates its components. The variance of this scatter around this structure is much smaller than in ambient space, and can thus be explained with a smaller dataset. The combination of diffusion maps, Itô sampling, and non-parametric statistics come together in developing efficient algorithms that encapsulate these ideas.

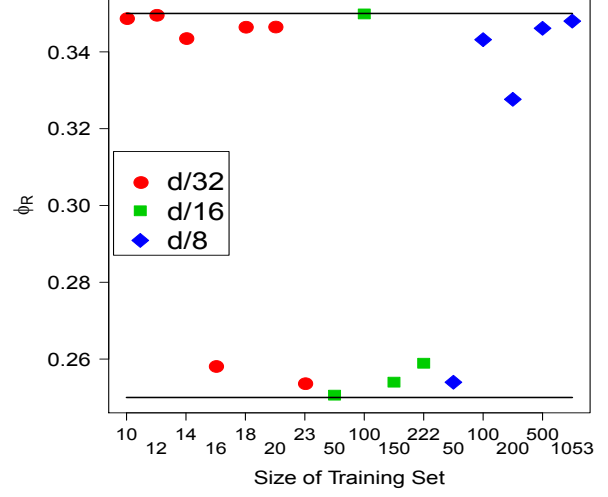
Three different LES spatial resolutions are explored, each with a different number of samples in the associated training set. Good convergence of the optimal solution is observed as function of the size of the training set, although the optimal solutions are different for each of the three resolutions. It is noted that while the optimal solution in each of the three cases tends to occur at a corner of the feasible domain (i.e. several design variables lie along their respective lower or upper bound), the specific corner depends on the specific spatial resolution. This highlights some of the challenges facing multi-fidelity/multi-resolution optimization strategies.

6. Acknowledgment

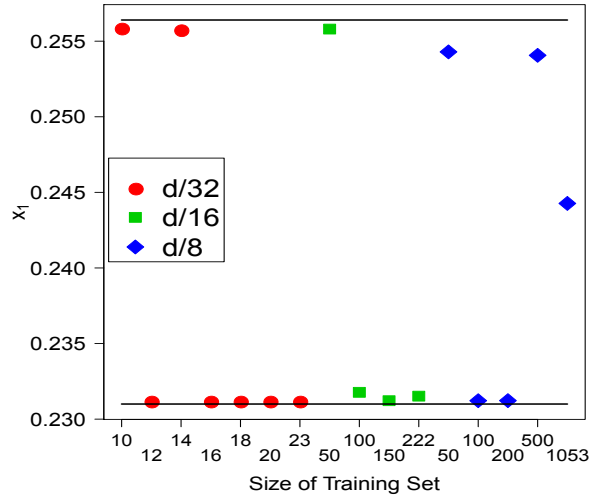
Support for this research was provided by the Defense Advanced Research Projects Agency (DARPA) program on Enabling Quantification of Uncertainty in Physical Systems (EQUiPS). This research used



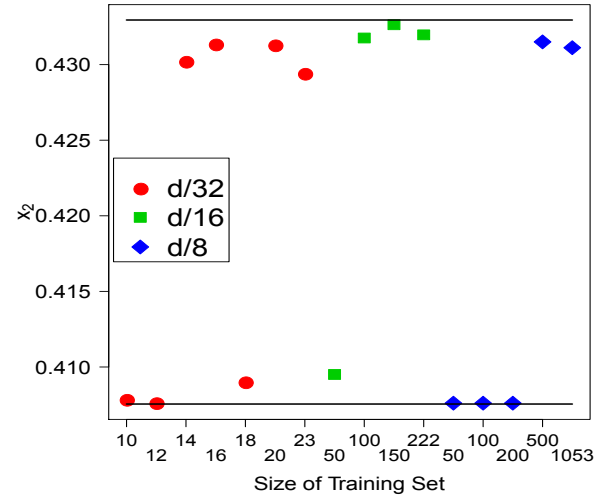
(a) Equivalence ratio, ϕ_G



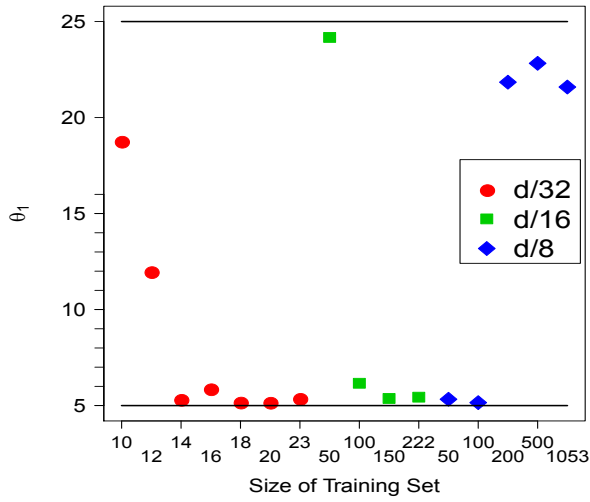
(b) Relative equivalence ratio, ϕ_R



(c) Coordinate of primary injector, x_1



(d) Coordinate of secondary injector, x_2



(e) Angle of primary injector, θ_1

Figure 5: Optimal values of the five control parameters for d/8, d/16, and d/32 resolutions and different sizes of training set.

resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operated under Contract No. DE-AC02-05CH11231. Computation for this work was also supported by the University of Southern California’s Center for High-Performance Computing (hpc.usc.edu). Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-NA-0003525. The views expressed in the article do not necessarily represent the views of the U.S. Department Of Energy or the United States Government.

7. References

- [1] R. Ghanem, P. Spanos, *Stochastic Finite Elements: A Spectral Approach*, Springer-Verlag, 1991.
- [2] C. Soize, R. Ghanem, Physical systems with random uncertainties: Chaos representations with arbitrary probability measure, *SIAM Journal of Scientific Computing* 26 (2) (2004) 395–410.
- [3] M. Eldred, Design under uncertainty employing stochastic expansion methods, *International Journal for Uncertainty Quantification* 1 (2) (2011) 119–146.
- [4] R. Ghanem, C. Soize, Probabilistic nonconvex constrained optimization with fixed number of function evaluations, *International Journal for Numerical Methods in Engineering* *Published on line 2017*. doi:10.1002/nme.5632.
- [5] R. Ghanem, C. Soize, C.-R. Thimmisetty, Optimal well-placement using a probabilistic learning, *Data-Enabled Discovery and Applications* 2 (1) (2018) 1–16. doi:10.1007/s41688-017-0014-x.
- [6] C. Soize, *Uncertainty Quantification. An Accelerated Course with Advanced Applications in Computational Engineering*, Springer, New York, 2017. doi:10.1007/978-3-319-54339-0.
- [7] R. Coifman, S. Lafon, A. Lee, M. Maggioni, B. Nadler, F. Warner, S. Zucker, Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps, *PNAS* 102 (21) (2005) 7426–7431.
- [8] C. Soize, Polynomial chaos expansion of a multimodal random vector, *SIAM/ASA Journal on Uncertainty Quantification* 3 (1) (2015) 34–60. doi:10.1137/140968495.
- [9] C. Soize, R. Ghanem, Data-driven probability concentration and sampling on manifold, *Journal of Computational Physics* 321 (2016) 242–258. doi:10.1016/j.jcp.2016.05.044.
- [10] D. J. Dolvin, Hypersonic International Flight Research and Experimentation (HIFiRE), in: 15th AIAA International Space Planes and Hypersonic Systems and Technologies Conference, no. 2008-2581, Dayton, OH, 2008. doi:10.2514/6.2008-2581.
- [11] D. J. Dolvin, Hypersonic International Flight Research and Experimentation, in: 16th AIAA/DLR/DGLR International Space Planes and Hypersonic Systems and Technologies Conference, no. 2009-7228, Bremen, Germany, 2009. doi:10.2514/6.2009-7228.
- [12] K. R. Jackson, M. R. Gruber, T. F. Barhorst, The HIFiRE Flight 2 Experiment: An Overview and Status Update, in: 45th AIAA/ASME/SAE/ASEE Joint Propulsion Conference & Exhibit, no. 2009-5029, Denver, CO, 2009. doi:10.2514/6.2009-5029.
- [13] K. R. Jackson, M. R. Gruber, S. Buccellato, HIFiRE Flight 2 Overview and Status Update 2011, in: 17th AIAA International Space Planes and Hypersonic Systems and Technologies Conference, no. 2011-2202, San Francisco, CA, 2011. doi:10.2514/6.2011-2202.
- [14] N. E. Hass, K. F. Cabell, A. M. Storch, HIFiRE Direct-Connect Rig (HDCR) Phase I Ground Test Results from the NASA Langley Arc-Heated Scramjet Test Facility, Tech. rep., NASA (2010).

- [15] A. M. Storch, M. Bynum, J. Liu, M. Gruber, Combustor Operability and Performance Verification for HIFiRE Flight 2, in: 17th AIAA International Space Planes and Hypersonic Systems and Technologies Conference, no. 2011-2249, San Francisco, CA, 2011. doi:10.2514/6.2011-2249.
- [16] K. F. Cabell, N. E. Hass, A. M. Storch, M. Gruber, HIFiRE Direct-Connect Rig (HDCR) Phase I Scramjet Test Results from the NASA Langley Arc-Heated Scramjet Test Facility, in: 17th AIAA International Space Planes and Hypersonic Systems and Technologies Conference, no. 2011-2248, San Francisco, CA, 2011. doi:10.2514/6.2011-2248.
- [17] G. L. Pellett, S. N. Vaden, L. G. Wilson, Opposed Jet Burner Extinction Limits: Simple Mixed Hydrocarbon Scramjet Fuels vs Air, in: 43rd AIAA/ASME/SAE/ASEE Joint Propulsion Conference & Exhibit, no. 2007-5664, Cincinnati, OH, 2007. doi:10.2514/6.2007-5664.
- [18] T. Lu, C. K. Law, A directed relation graph method for mechanism reduction, *Proceedings of the Combustion Institute* 30 (1) (2005) 1333–1341. doi:10.1016/j.proci.2004.08.145.
- [19] A. C. Zambon, H. K. Chelliah, Explicit reduced reaction models for ignition, flame propagation, and extinction of $C_2H_4/CH_4/H_2$ and air systems, *Combustion and Flame* 150 (1-2) (2007) 71–91. doi:10.1016/j.combustflame.2007.03.003.
- [20] J. C. Oefelein, Large eddy simulation of turbulent combustion processes in propulsion and power systems, *Progress in Aerospace Sciences* 42 (1) (2006) 2–37. doi:10.1016/j.paerosci.2006.02.001.
- [21] G. Lacaze, A. Misdariis, A. Ruiz, J. C. Oefelein, Analysis of high-pressure Diesel fuel injection processes using LES with real-fluid thermodynamics and transport, *Proceedings of the Combustion Institute* 35 (2) (2015) 1603–1611. doi:10.1016/j.proci.2014.06.072.
- [22] G. Lacaze, Z. P. Vane, J. C. Oefelein, Large Eddy Simulation of the HIFiRE Direct Connect Rig Scramjet Combustor, in: 55th AIAA Aerospace Sciences Meeting, no. 2017-0142, Grapevine, TX, 2017. doi:10.2514/6.2017-0142.
- [23] M. Germano, U. Piomelli, P. Moin, W. Cabot, A dynamic subgrid-scale eddy viscosity model, *Physics of Fluids* 3 (7) (1991) 17601765.
- [24] B. Vreman, B. Geurts, H. Kuerten, On the formulation of the dynamic mixed subgrid-scale model, *Physics of Fluids* 6 (12) (1994) 40574059.
- [25] M. R. Gruber, K. Jackson, J. Liu, Hydrocarbon-Fueled Scramjet Combustor Flowpath Development for Mach 6-8 HIFiRE Flight Experiments, Tech. rep., AFRL (2008).
- [26] R. Coifman, S. Lafon, Diffusion maps, applied and computational harmonic analysis, *Applied and Computational Harmonic Analysis* 21 (1) (2006) 5–30.
- [27] R. Coifman, I. Kevrekidis, S. Lafon, M. Maggioni, B. Nadler, Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems, *SIAM J. Multiscale Model. Simul.* 7 (2) (2008) 842–864.
- [28] M. Belkin, P. Niyogi, Laplacian eigenmaps for dimensionality reduction and data representation, *Neural Computation* 15 (2003) 1373–1396.
- [29] C. Soize, R. Ghanem, C. Safta, X. Huan, Z. Vane, J. Oefelein, G. Lacaze, H. Najm, Q. Tang, X. Chen, Entropy-based closure for probabilistic learning on manifolds, arXiv:1803.08161v1, 21 Mar2018, to appear in *Journal of Computational Physics*.
- [30] D. Scott, *Multivariate Density Estimation: Theory, Practice, and Visualization*, 2nd Edition, John Wiley and Sons, New York, 2015.
- [31] K. Burrage, I. Lenane, G. Lythe, Numerical methods for second-order stochastic differential equations, *SIAM Journal on Scientific Computing* 29 (1) (2007) 245–264.

- [32] L. Scrucca, GA: A package for genetic algorithms in R, *Journal of Statistical Software* 53 (4) (2013) 1–37.
- [33] R Core Team, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria (2018).
URL <https://www.R-project.org/>