

# Intrusive Projection Methods with Upwinding for Uncertain Nonlinear Hyperbolic Systems<sup>☆</sup>

J. Tryoen<sup>a,b,\*</sup>, O. Le Maître<sup>b,c</sup>, M. Ndjinga<sup>c</sup>, A. Ern<sup>a</sup>

<sup>a</sup>Université Paris Est, CERMICS, Ecole des Ponts, 77455 Marne la Vallée cedex 2, France

<sup>b</sup>LIMSI-CNRS, 91403 Orsay cedex, France

<sup>c</sup>CEA-Saclay, DEN/DM2S/SFME, F-91191 Gif-sur-Yvette cedex, France

---

## Abstract

This paper deals with stochastic spectral methods for uncertainty propagation and quantification in nonlinear hyperbolic systems of conservation laws. We consider problems with parametric uncertainty in initial conditions and model coefficients, whose solutions exhibit discontinuities in the spatial as well as in the stochastic variables. The stochastic spectral method relies on multi-resolution schemes where the stochastic domain is discretized using tensor-product stochastic elements supporting local polynomial bases. A Galerkin projection is used to derive a system of deterministic equations for the stochastic modes of the solution. Hyperbolicity of the resulting Galerkin system is analyzed. A finite volume scheme with a Roe-type solver is used for discretization of the spatial and time variables. An original technique is introduced for the fast evaluation of approximate upwind matrices, which is particularly well adapted to local polynomial bases. Efficiency and robustness of the overall method are assessed on the Burgers and Euler equations with shocks.

*Key words:* Uncertainty Quantification, hyperbolic systems, conservation laws, stochastic spectral methods, upwinding, Galerkin Projection

---

## 1. Introduction

In numerical simulation, accounting for uncertainties in input quantities (such as model parameters, initial and boundary conditions, and geometry) is an important issue, especially in risk analysis, safety, and design. Assuming that these input quantities can be parametrized by random variables with known distribution functions, the question is to quantify the resulting uncertainty in the numerical solution. Uncertainty Quantification (UQ) provides for instance numerical error bars that make the comparison with experimental observations easier and therefore facilitate the evaluation of physical models. Moreover, they enable to identify the uncertain parameters that should be measured or controlled with more accuracy because they have the most significant impact on the solution. Furthermore, they allow for the assessment of the reliability level that can be attached to computations.

Stochastic spectral methods provide effective tools for UQ. Such methods decompose random quantities on suitable approximation bases. Their main interest is that they provide a complete probabilistic description of the uncertain solution. A classical choice for the stochastic basis is the set of generalized Polynomial Chaos (gPC) spanned by random polynomials, continuous in the stochastic domain and truncated to some

---

<sup>☆</sup>This work is partially supported by GNR MoMaS (ANDRA, BRGM, CEA, EdF, IRSN, PACEN-CNRS). O.P. Le Maître is partially supported by the French National Research Agency (Grant ANR-08-JCJC-0022).

\*Corresponding author: Université Paris Est, CERMICS, Ecole des Ponts, 77455 Marne la Vallée cedex 2, France. Phone: (33)-1 64 15 36 65

Email addresses: tryoenj@cermics.enpc.fr (J. Tryoen), olm@limsi.fr (O. Le Maître), michael.ndjinga@cea.fr (M. Ndjinga), ern@cermics.enpc.fr (A. Ern)

degree. Polynomial Chaos (PC) methods were originally introduced by Ghanem and Spanos [11] following the Wiener Chaos theory [36] in which random processes are expanded in a Hermite polynomial basis of Gaussian random variables. The theory was then extended to the case of more general random processes that can be expanded on a basis of orthogonal polynomials associated with the chosen random variables; see among others [38]. Then, two types of resolution methods are available. The first ones are called non-intrusive and are based on the use of the numerical code solving the deterministic model (without uncertainty) as a black box to construct the spectral expansion of the solution. Two approaches can be used, either the probabilistic collocation method [25, 37, 2, 9, 30, 8], which consists in approximating the stochastic solution by a polynomial interpolation, or the non-intrusive projection method [32, 14, 22], which is based on the evaluation of the stochastic modes of the solution by numerical integration. For the two cases, the issue is to find the set of interpolation or integration points that provide the most accurate stochastic approximation. The second type of resolution methods are stochastic Galerkin methods based on a Galerkin projection of the model equations yielding a reformulated deterministic problem for the stochastic modes of the solution. Such methods are called intrusive because of the need to rewrite to some extent the simulation code. Their advantage is to rely on the weak form of the problem and thereby on a firmer theoretic background. Therefore, they are in our opinion better suited for mathematical analysis and improvements such as refinement and adaptation. In particular, stochastic Galerkin methods applied to elliptic and parabolic problems are relatively well understood. Such methods have been successfully applied in many domains (see [11] and references therein). Regarding viscous flow models, previous works have dealt with the incompressible Navier-Stokes equations [21, 22], low Mach number flows [20], and electrochemical microfluidic applications [5]. Recent reviews on uncertain fluid flows can be found in [15, 27, 28].

The application of stochastic spectral methods to hyperbolic systems of conservation laws (in particular inviscid flows) poses additional challenges. The main difficulty is that solutions can exhibit discontinuities (in the spatial domain) in finite time due to the development of shock waves and contact discontinuities (in the spatial variables). Although these discontinuities concern the spatial variables, their propagation speed can be affected by uncertainty, thereby leading to discontinuities in the stochastic variables as well. As a result, bases of continuous polynomials in the stochastic domain become inappropriate, because of aliasing errors [3] and Gibbs-type phenomena [17]. To overcome this issue, Multi-Resolution Analysis (MRA) methods using stochastic finite elements [4], multi-element gPC (ME-gPC) [35], and multi-wavelet expansions [17, 18, 19] can be used to make the spectral representation more local by decomposing the stochastic domain into different regions or different scales. Another difficulty originates from the nonlinearities in the physical fluxes of the stochastic hyperbolic system raising the issue of computing such fluxes in the context of Galerkin projections. Indeed, all mathematical operations must be applied to the stochastic expansions that represent the variables. One attractive approach is to use pseudo-spectral techniques [6].

Polynomial collocation methods have already been applied by Mathelin et al. [26] to the Euler equations but in the continuous case. Other non-intrusive approaches include that of Abgrall [1] based on ENO-like reconstructions for the convection, Burgers, and Euler equations, and that of Lin et al. [24] based on multi-element probabilistic collocation methods for supersonic flows past a wedge with random roughness. Concerning intrusive methods, most of the approaches found in the literature are in fact pseudo-intrusive because the fluxes in the Galerkin system are computed in a non-intrusive way by quadrature methods, as for instance in Ge et al. [10] for the shallow-water equations and in Poette et al. [31] for the Burgers and Euler equations. One nice feature of this latter approach is that the polynomial expansion is carried on suitable entropic variables and not on the original conservative variables, so that it can be proven that the Galerkin projection leads to a hyperbolic system; however the numerical algorithm requires a minimization procedure to recover the solution expansion that can be time consuming.

To our knowledge, very few intrusive stochastic spectral methods have been investigated for uncertain hyperbolic problems. The scalar wave equation has been treated with gPC methods by Gottlieb et al. [13]. The case of nonlinear hyperbolic systems is obviously more difficult. Supersonic flows past a wedge with random inflow fluctuations or random wedge oscillations around its apex have been studied using ME-gPC methods by Lin et al. [23]. In the context of intrusive methods, a crucial question is the design of a suitable scheme to approximate in the spatial and time domains the evolution problem associated with the Galerkin projection. Typically, one would like to use a Finite Volume (FV) scheme with appropriate upwinding.

For instance, Lin et al. [23] considered upwinding using the mean values (in the stochastic domain) of the eigenvectors of the Galerkin Jacobian matrix. As mentioned in [23], this approach is only justified in the case of relatively small fluctuations of the random quantities. The present paper improves on this point both theoretically and numerically, by using full spectral information on the eigenvectors of the Galerkin Jacobian matrix and by proposing a cost-effective method to approximate the absolute value of this matrix.

The purpose of the present work is to investigate intrusive methods for nonlinear stochastic hyperbolic systems. To this end we discretize as in [35] the stochastic domain using tensor-product stochastic elements supporting local polynomial bases. A Stochastic Galerkin projection is then used to derive the Galerkin system, that is, the set of deterministic equations coupling the stochastic modes of the solution on the selected basis. The nonlinear fluxes in the Galerkin system are computed in a pseudo-spectral way with the tools described in [6]. A FV method with a Roe-type solver is used to approximate the Galerkin system in the spatial and time domains. At the theoretical level, our main result is that the Galerkin system is proven to be hyperbolic in two specific cases, namely when the original stochastic problem has a symmetric Jacobian (the main application being scalar conservation laws) and when its eigenvectors are independent of the uncertainty (the main application being linear hyperbolic systems with uncertainty only on initial or boundary conditions). Moreover, in the general case, we identify an approximate Galerkin Jacobian matrix which is shown to be  $\mathbb{R}$ -diagonalizable and whose eigenvalues can be easily determined from those of the original stochastic problem. These eigenvalues are solely used as data in the determination of a (low-degree) fitting polynomial that can be applied to the actual Galerkin Jacobian matrix to compute an upwinding matrix for the Roe solver. This new methodology for computing approximate upwind matrices is particularly well adapted to the stochastic discretization since the computation of the fitting polynomial can be localized to each stochastic element, thereby making the procedure more robust and efficient.

The paper is organized as follows. In Section 2, the stochastic hyperbolic framework is presented, including the stochastic approximation spaces and the stochastic Galerkin projection. The hyperbolicity of the Galerkin system is investigated in Section 3. Numerical methods are described in Section 4. Finally, simulation results are presented in Section 5.

We adopt the following notation: lower case symbols represent deterministic quantities, whereas upper case symbols represent stochastic quantities.

## 2. Galerkin projection of stochastic hyperbolic systems

### 2.1. Probabilistic framework and parametric uncertainty

We are interested in uncertainty propagation and quantification in nonlinear hyperbolic problems. The uncertainty is treated in a probabilistic framework. We rely on an abstract probability space  $\mathcal{P} = (\Theta, \Sigma, d\mu)$ , where  $\Theta$  is the set of random events,  $\Sigma$  the associated  $\sigma$ -algebra, and  $d\mu$  the probability measure. For any random variable  $H(\theta)$  defined on  $\mathcal{P}$ , the expectation of  $H$  is

$$E[H] = \int_{\Theta} H(\theta) d\mu(\theta). \quad (1)$$

We denote by  $L^2(\Theta, d\mu)$  the space of second-order random variables on  $\mathcal{P}$ . We assume hereafter that all random quantities are second-order.

In view of stochastic discretization, we introduce a finite set of  $N$  random variables  $\xi(\theta) := \{\xi_1(\theta), \dots, \xi_N(\theta)\}$  defined on  $\mathcal{P}$  with known distributions. These random variables will be used to parametrize the uncertain coefficients or initial conditions of the hyperbolic problem. For simplicity, we consider  $\xi_i(\theta)$  as real-valued independent identically distributed random variables, such that the joined density function of  $\xi(\theta)$  factorizes, namely

$$p_{\xi}(y) = \prod_{i=1}^N p(y_i), \quad (2)$$

where  $p(y_i)$  is the probability density function of  $\xi_i(\theta)$ . We further denote by  $\Xi$  the range of  $\xi$  and by  $\mathcal{P}_{\xi}$  the image probability space,  $\mathcal{P}_{\xi} := (\Xi, \mathcal{B}_{\Xi}, p_{\xi})$ , where  $\mathcal{B}_{\Xi}$  is the Borel set of  $\Xi$ . Similarly,  $L^2(\Xi, p_{\xi})$  is the

space of second-order random variables defined on the image space. The expectation operator in the image space is denoted using brackets and is related to the expectation on  $\mathcal{P}$  through the identity

$$E[H] = \int_{\Theta} H(\xi(\theta))d\mu(\theta) = \int_{\Xi} H(y)p_{\xi}(y)dy =: \langle H \rangle. \quad (3)$$

### 2.2. Stochastic hyperbolic systems

We consider conservative systems of nonlinear hyperbolic PDE's. The uncertainty can result from a variability of the initial condition and/or of some coefficients in the model. For simplicity, we focus on one-dimensional spatial domains. The extension to higher spatial dimension is straightforward at least concerning the stochastic aspects. We seek for  $U(x, t, \xi)$  solving almost surely the following conservative system

$$\begin{cases} \frac{\partial}{\partial t}U(x, t, \xi) + \frac{\partial}{\partial x}F(U(x, t, \xi); \xi) = 0, \\ U(x, t = 0, \xi) = U^0(x, \xi). \end{cases} \quad (4)$$

Let  $\Omega \subset \mathbb{R}$  be the bounded spatial domain over which the problem is posed and let  $\mathcal{A}_U \subset \mathbb{R}^m$ ,  $m \geq 1$ , be the set of admissible values for the solutions which we assume independent of the random event. For instance, for the Burgers equation, we can take  $\mathcal{A}_U = \mathbb{R}$ , whereas for the Euler equations,  $\mathcal{A}_U$  is the set of states with positive density and pressure. Then,  $U : (x, t, \xi) \in \Omega \times [0, T] \times \Xi \mapsto U(x, t, \xi) \in \mathcal{A}_U \otimes L^2(\Xi, p_{\xi})$  denotes the uncertain state vector of conservative variables parametrized by  $\xi$ ,  $U^0(x, \xi)$  is a parametrization by  $\xi$  of the uncertain initial condition, and  $F : (U; \xi) \in \mathcal{A}_U \otimes L^2(\Xi, p_{\xi}) \times \Xi \mapsto F(U; \xi) \in \mathbb{R}^m \otimes L^2(\Xi, p_{\xi})$  is the uncertain flux function, involving some random coefficients parametrized again by  $\xi$ . Moreover, since the domain  $\Omega$  is bounded, appropriate boundary conditions have to be enforced at the boundary  $\partial\Omega$ ; they will be specified in Section 5 when presenting the test cases.

For smooth  $U$ , the system (4) can also be written in the non-conservative form

$$\begin{cases} \frac{\partial}{\partial t}U(x, t, \xi) + \nabla_U F(U(x, t, \xi); \xi) \frac{\partial}{\partial x}U(x, t, \xi) = 0, \\ U(x, t = 0, \xi) = U^0(x, \xi). \end{cases} \quad (5)$$

This stochastic system is assumed to be hyperbolic in the sense that the stochastic Jacobian matrix  $\nabla_U F \in \mathbb{R}^{m,m} \otimes L^2(\Xi, p_{\xi})$  is  $\mathbb{R}$ -diagonalizable almost surely, that is, for almost every  $\xi \in \Xi$ , there exist  $m$  eigenvalues  $\Lambda^1(\cdot; \xi), \dots, \Lambda^m(\cdot; \xi)$  and  $m$  associated eigenvectors  $W^1(\cdot; \xi), \dots, W^m(\cdot; \xi)$  forming a complete basis of  $\mathbb{R}^m$ , such that

$$\nabla_U F(\cdot; \xi) = P^{-1}(\cdot; \xi)D(\cdot; \xi)P(\cdot; \xi), \quad (6)$$

with

$$D(\cdot; \xi) = \text{diag}(\Lambda^k(\cdot; \xi))_{k=1, \dots, m} \quad \text{and} \quad P(\cdot; \xi) = ( W^1(\cdot; \xi) \quad \dots \quad W^m(\cdot; \xi) ). \quad (7)$$

The matrices  $D(\cdot; \xi)$  and  $P(\cdot; \xi)$  are in  $\mathbb{R}^{m,m} \otimes L^2(\Xi, p_{\xi})$ . To alleviate the notation, the dependence of the eigenvalues and eigenvectors on  $U$  is omitted in the sequel.

### 2.3. Stochastic discretization

To approximate the solution in  $L^2(\Xi, p_{\xi})$ , we need a stochastic discretization of the problem. This is obtained by considering an appropriate Hilbertian basis of random functionals in  $\xi$  spanning  $L^2(\Xi, p_{\xi})$ ,

$$L^2(\Xi, p_{\xi}) = \overline{\text{span}\{\Psi_1(\xi), \Psi_2(\xi), \dots\}}, \quad \langle \Psi_{\alpha} \Psi_{\beta} \rangle = \delta_{\alpha\beta}, \quad (8)$$

where  $\delta_{\alpha\beta}$  denotes the Kronecker symbol. The discrete solution is sought in a finite dimensional subspace  $\mathcal{S}^P$  constructed by truncating the Hilbertian basis:

$$\mathcal{S}^P = \text{span}\{\Psi_1(\xi), \Psi_2(\xi), \dots, \Psi_P(\xi)\} \subset L^2(\Xi, p_{\xi}), \quad \dim(\mathcal{S}^P) =: P. \quad (9)$$



We assume for simplicity that  $\xi$  is a uniform random vector in  $[0, 1]^N$  (an isoprobabilistic transformation can be used to map the original independent random variables to this random vector [18, 19]). The image probability space is then  $\mathcal{P}_\xi := ([0, 1]^N, \mathcal{B}_{[0,1]^N}, 1)$ , where  $\mathcal{B}_{[0,1]^N}$  is the Borel set of  $[0, 1]^N$ .

We decompose the stochastic domain  $[0, 1]^N$  dyadically and approximate the stochastic solution by piecewise polynomial functions. In addition to the number  $N$  of random variables  $\xi_i$  in the parametrization, this approximation depends on the resolution level  $Nr \geq 0$  (controlling the minimal size of the stochastic elements, that is, the discretization cells in the stochastic domain) and on the expansion order  $No \geq 0$  (controlling the degree of the piecewise polynomial approximation). Let  $\mathbf{i} = (i_1, \dots, i_N) \in \{1, \dots, 2^{Nr}\}^N$  be a multi-index and let  $K_{\mathbf{i}} = \{\xi \in [0, 1]^N, \forall 1 \leq j \leq N, \xi_j \in [2^{-Nr}(i_j - 1), 2^{-Nr}i_j]\}$  be the associated stochastic element. Thus, we define  $\mathcal{S}^{No, Nr}$  as the stochastic approximation space of piecewise polynomial functions

$$\mathcal{S}^{No, Nr} := \{f : [0, 1]^N \rightarrow \mathbb{R}, \forall \mathbf{i} \in \{1, \dots, 2^{Nr}\}^N, f|_{K_{\mathbf{i}}} \in \mathbb{Q}_{No}^N[\xi]\}, \quad (10)$$

where  $\mathbb{Q}_{No}^N[\xi]$  denotes the vector space of real polynomials in  $\mathbb{R}^N$  with degree  $\leq No$  in each variable  $\xi_i$ . The space  $\mathcal{S}^{No, Nr}$  has dimension

$$\dim \mathcal{S}^{No, Nr} = (No + 1)^N 2^{NNr} =: P_\pi P_\sigma =: P, \quad (11)$$

where  $P_\pi := (No + 1)^N$  is the dimension of the local polynomial basis on each stochastic element, and  $P_\sigma := 2^{NNr}$  is the number of stochastic elements. The spaces  $\mathcal{S}^{No, Nr}$  form a hierarchical family of stochastic spaces since  $\mathcal{S}^{No, Nr} \subset \mathcal{S}^{No', Nr}$  for  $No \leq No'$  and  $\mathcal{S}^{No, Nr} \subset \mathcal{S}^{No, Nr'}$  for  $Nr \leq Nr'$ . It is also possible to work with smaller stochastic approximation spaces, for instance spanned by polynomials of total degree  $\leq No$ , that is, using sparse polynomial tensorization instead of full polynomial tensorization. The resulting changes in the numerical method will be indicated whenever relevant.

Two kinds of basis can be considered. Firstly,  $\mathcal{S}^{No, Nr}$  can be spanned by the hierarchical Multi-Wavelet (MW) system of order  $No$  and resolution level  $Nr$  introduced in [18]. Alternatively,  $\mathcal{S}^{No, Nr}$  can be spanned by local Legendre polynomial bases, where each function of  $\mathcal{S}^{No, Nr}$  is expanded in each Stochastic Element (SE) of size  $2^{-Nr}$  on a local fully tensorized set with dimension  $(No + 1)^N$  of Legendre polynomials. For convenience, Legendre polynomials are henceforth defined with respect to the reference interval  $[0, 1]$ . The case  $Nr = 0$  corresponds to the classical continuous approximation (Wiener-Legendre expansion), while the choice  $Nr > 0$  and  $No = 0$  leads to the Wiener-Haar expansion (piecewise constant approximation). In view of adaptive algorithms, the MW basis provides a natural framework. The SE basis is more convenient for theoretical analysis and implementation. Therefore, unless stated explicitly, we focus in this work on the SE basis, which we denote by  $\{\Psi_\alpha(\xi)\}_{\alpha=1, \dots, P}$ . In practice,  $\alpha$  is a double index,  $\alpha = \{\alpha_\sigma, \alpha_\pi\}$ , the first index ( $\alpha_\sigma$ ) referring to the stochastic element and the second ( $\alpha_\pi$ ) referring to the polynomial function within the stochastic element.

The approximate solution in  $\mathcal{S}^P := \mathcal{S}^{No, Nr}$  is expanded as a series in the form

$$U(x, t, \xi) \approx U^P(x, t, \xi) = \sum_{\alpha=1}^P u_\alpha(x, t) \Psi_\alpha(\xi). \quad (12)$$

The deterministic  $\mathbb{R}^m$ -valued fields  $u_\alpha(x, t)$  are called the stochastic modes of the solution (in  $\mathcal{S}^P$ ). If  $U^P(x, t, \xi)$  is known, then  $u_\alpha = \langle \Psi_\alpha U^P \rangle$ . The knowledge of the stochastic modes allows one to compute interesting statistic quantities, such as expectation, variance, higher moments, density functions, cross correlations, etc. . . , relying either on analytic expressions or on a sampling of  $\Xi$ .

#### 2.4. The Galerkin system

The computation of the stochastic modes  $u_\alpha(x, t)$  is based on a weak interpretation, or Galerkin projection, of (4). Projecting (4) on the basis of  $\mathcal{S}^P$  and accounting for orthonormality, we obtain

$$\begin{cases} \frac{\partial}{\partial t} u_\alpha(x, t) + \frac{\partial}{\partial x} \langle \Psi_\alpha F(U^P; \cdot) \rangle = 0, & \forall \alpha = 1, \dots, P, \\ u_\alpha(x, t = 0) = \langle \Psi_\alpha U^0 \rangle, & \forall \alpha = 1, \dots, P. \end{cases} \quad (13)$$

Equation (13) shows that the  $\alpha$ -th stochastic mode of the approximate solution is governed by an equation that generally couples all the stochastic modes in the term  $\langle \Psi_\alpha F(U^P; \cdot) \rangle$ . It is convenient to define the vectors of stochastic modes and flux

$$u(x, t) = \begin{pmatrix} u_1(x, t) \\ \vdots \\ u_P(x, t) \end{pmatrix}, \quad f(u(x, t)) = \begin{pmatrix} f_1(u) \\ \vdots \\ f_P(u) \end{pmatrix}, \quad (14)$$

with

$$f_\alpha(u) := \langle \Psi_\alpha F(U^P; \cdot) \rangle, \quad \alpha = 1, \dots, P, \quad \text{and} \quad U^P = \sum_{\beta=1}^P u_\beta \Psi_\beta(\xi). \quad (15)$$

The component vector  $u$  must belong to the admissible set  $\mathcal{A}_u \subset \mathbb{R}^{m(P+1)}$  such that  $u \in \mathcal{A}_u \Leftrightarrow U^P(\xi) = \sum_{\alpha=1}^P u_\alpha \Psi_\alpha(\xi) \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi)$ . With obvious notation for  $u^0$ , the deterministic Galerkin system takes the simple form

$$\begin{cases} \frac{\partial}{\partial t} u(x, t) + \frac{\partial}{\partial x} f(u(x, t)) = 0, \\ u(x, t = 0) = u^0(x). \end{cases} \quad (16)$$

Thus, the problem on  $u$  has the same form as the original stochastic problem (4), except that the state vector is now of size  $mP$ .

### 3. Hyperbolicity of the Galerkin system

Before detailing the construction of a numerical method to approximate the Galerkin system (16), we investigate the hyperbolicity of this system. Introducing the Galerkin Jacobian matrix of order  $mP$  such that

$$(\nabla_u f(u))_{\alpha, \beta=1, \dots, P} = \langle \nabla_U F(U^P; \cdot) \Psi_\alpha \Psi_\beta \rangle_{\alpha, \beta=1, \dots, P}, \quad (17)$$

we aim at understanding whether this matrix is  $\mathbb{R}$ -diagonalizable.

The advantage of using SE bases rather than MW bases for investigating the  $\mathbb{R}$ -diagonalization of the Galerkin Jacobian matrix  $\nabla_u f$  is that owing to the adopted index convention, this matrix has a diagonal block structure. Indeed,  $(\nabla_u f)_{\alpha\beta} = 0$  whenever  $Supp(\Psi_\alpha) \cap Supp(\Psi_\beta)$  has zero measure. Consequently,  $\nabla_u f$  is diagonalizable if and only if each block in the diagonal is diagonalizable. Such blocks are of size  $m(\text{No} + 1)^N \times m(\text{No} + 1)^N$  and correspond to a given stochastic element. The issue of the hyperbolicity of the Galerkin system can then be studied for the case  $\text{Nr} = 0$ .

An interesting point is that the two different representations of  $\nabla_u f$  using the MW basis or the SE basis for the stochastic discretization are equivalent in view of  $\mathbb{R}$ -diagonalization. Indeed, let  $\{\psi_\alpha^{MW}(\xi)\}_{\alpha=1, \dots, P}$  denote the MW basis and let  $B \in \mathbb{R}^{P, P}$  denote the transition matrix between the two bases, such that  $\Psi_\alpha^{MW}(\xi) = \sum_{\gamma=1}^P B_{\alpha\gamma} \Psi_\gamma(\xi)$ , for all  $\alpha = 1, \dots, P$ , that is,  $(B)_{1 \leq \alpha, \gamma \leq P} = \langle \Psi_\alpha^{MW} \Psi_\gamma \rangle$ . Let  $\nabla_u f^{MW} \in \mathbb{R}^{mP, mP}$  be the representation of the Galerkin Jacobian matrix using the MW basis. Then, for all  $\alpha, \beta = 1, \dots, P$ ,

$$\begin{aligned} (\nabla_u f^{MW})_{\alpha\beta} &= \langle \nabla_U F(U^P; \cdot) \Psi_\alpha^{MW} \Psi_\beta^{MW} \rangle = \sum_{\gamma, \delta} \langle \nabla_U F(U^P; \cdot) B_{\alpha\gamma} \Psi_\gamma B_{\beta\delta} \Psi_\delta \rangle \\ &= \sum_{\gamma, \delta} B_{\alpha\gamma} (\nabla_u f)_{\gamma\delta} B_{\beta\delta} = (B \nabla_u f B^T)_{\alpha\beta}. \end{aligned} \quad (18)$$

Moreover,  $B$  is orthogonal owing to the orthonormality of the two bases, which implies that  $\nabla_u f^{MW}$  and  $\nabla_u f$  are similar and therefore proves the equivalence of the two representations with respect to  $\mathbb{R}$ -diagonalization.

### 3.1. Stochastic symmetric hyperbolic systems

**Theorem 1.** Consider either sparse or full polynomial tensorization for the stochastic space  $\mathcal{S}^{\text{No}, \text{Nr}}$ . If the stochastic Jacobian matrix  $\nabla_U F(\cdot; \xi)$  is symmetric, then the Galerkin Jacobian matrix  $\nabla_u f$  is  $\mathbb{R}$ -diagonalizable. In particular, the Galerkin projection of a scalar conservation law always leads to a hyperbolic system.

*Proof.* If  $\nabla_U F(\cdot; \xi)$  is symmetric, then the Galerkin matrix  $\nabla_u f$  defined by (17) is also symmetric and therefore  $\mathbb{R}$ -diagonalizable.  $\square$

### 3.2. Stochastic eigenvectors independent of the uncertainty

**Theorem 2.** Consider either sparse or full polynomial tensorization for the stochastic space  $\mathcal{S}^{\text{No}, \text{Nr}}$ . If the eigenvectors of the stochastic Jacobian matrix  $\nabla_U F(\cdot; \xi)$  are independent of the uncertainty, then the Galerkin Jacobian matrix  $\nabla_u f$  is  $\mathbb{R}$ -diagonalizable.

*Proof.* If the eigenvectors of  $\nabla_U F(\cdot; \xi)$  are independent of  $\xi$ , then the spectral decomposition (6) becomes

$$\nabla_U F(\cdot; \xi) = p_0^{-1} D(\xi) p_0, \quad \text{with } p_0 = (w_0^1 \quad \dots \quad w_0^m), \quad (19)$$

where  $w_0^1, \dots, w_0^m$  are independent of  $\xi$ . A generic element in  $\nabla_u f$  can be identified with the multi-index  $(\alpha i, \beta j)$  with  $i, j = 1, \dots, m$  and  $\alpha, \beta = 1, \dots, P$ , in such a way that

$$\begin{aligned} (\nabla_u f(u))_{\alpha i, \beta j} &= \left\langle (\nabla_U F(U^P; \cdot))_{ij} \Psi_\alpha \Psi_\beta \right\rangle = \sum_k \left\langle (p_0^{-1})_{ik} \Lambda^k (p_0)_{kj} \Psi_\alpha \Psi_\beta \right\rangle \\ &= \sum_k (p_0^{-1})_{ik} \langle \Lambda^k \Psi_\alpha \Psi_\beta \rangle (p_0)_{kj} \\ &= \sum_{k, k'} \sum_{\gamma, \gamma'} \{ \delta_{\alpha\gamma} (p_0^{-1})_{ik} \} \{ \delta_{kk'} \langle \Lambda^k \Psi_\gamma \Psi_{\gamma'} \rangle \} \{ \delta_{\gamma'\beta} (p_0)_{k'j} \} \\ &= \sum_{k, k'} \sum_{\gamma, \gamma'} (q)_{\alpha i, \gamma k} (d)_{\gamma k, \gamma' k'} (r)_{\gamma' k', \beta j} = (q d r)_{\alpha i, \beta j} \end{aligned} \quad (20)$$

where  $d$  is the block-diagonal matrix of size  $mP \times mP$  such that

$$(d)_{\gamma k, \gamma' k'} = \delta_{kk'} \langle \Lambda^k \Psi_\gamma \Psi_{\gamma'} \rangle, \quad (21)$$

and  $q$  and  $r$  are  $mP \times mP$  matrices such that

$$(q)_{\alpha i, \gamma k} = \delta_{\alpha\gamma} (p_0^{-1})_{ik}, \quad (r)_{\gamma k, \beta j} = \delta_{\gamma\beta} (p_0)_{kj}. \quad (22)$$

Each block of the diagonal of  $d$  is symmetric, and therefore  $\mathbb{R}$ -diagonalizable so that  $d$  is  $\mathbb{R}$ -diagonalizable. Besides,

$$(q^r)_{\alpha i, \beta j} = \sum_{\alpha, k} (q)_{\alpha i, \gamma k} (r)_{\gamma k, \beta j} = \sum_{\alpha, k} (p_0^{-1})_{ik} \delta_{\alpha\gamma} \delta_{\gamma\beta} (p_0)_{kj} = \delta_{\alpha\beta} \delta_{ij}, \quad (23)$$

which means that  $q = r^{-1}$ . This concludes the proof.  $\square$

**Remark.** Theorem 2 provides another proof of the fact that the Galerkin system derived from a stochastic scalar conservation law is hyperbolic. Indeed,  $W(\xi) = 1$  is the eigenvector of  $\nabla_U F \in \mathbb{R} \otimes L^2(\Xi, p_\xi)$ . A relevant application of Theorem 2 is the scalar wave equation with uncertain sound velocity. Theorem 2 can also be applied to linear hyperbolic systems with uncertainty only on initial or boundary conditions.

### 3.3. An approximate Galerkin Jacobian matrix

In the most general case, the Galerkin Jacobian matrix  $\nabla_u f$  is not guaranteed to be  $\mathbb{R}$ -diagonalizable. However, we can identify an  $\mathbb{R}$ -diagonalizable approximation of  $\nabla_u f$  obtained by quadrature, denoted by  $\overline{\nabla_u f}$ , for which explicit expressions of the eigenvalues can be derived. The main application of this result (see section 4) is to use the spectrum of  $\overline{\nabla_u f}$  to compute a fitting polynomial that can then be applied to the Galerkin Jacobian matrix  $\nabla_u f$  to approximate its absolute value in the context of upwind matrices for Roe-type solvers.

To study the  $\mathbb{R}$ -diagonalization of  $\overline{\nabla_u f}$ , we can assume that  $N_r = 0$  since the extension to  $N_r \geq 1$  is straightforward owing to the block diagonal structure of the Galerkin Jacobian matrix. Moreover, it is sufficient to consider the one-dimensional stochastic case ( $N = 1$ ). In this case, we denote by  $\{\xi_\gamma\}_{\gamma=0,\dots,\text{No}}$  the set of  $P = \text{No} + 1$  Gauss points in  $[0, 1]$ , i.e., the  $(\text{No} + 1)$  zeroes of the Legendre polynomial of degree  $(\text{No} + 1)$ , and by  $\{\omega_\gamma\}_{\gamma=0,\dots,\text{No}}$  the associated quadrature weights. The extension to  $N > 1$  is straightforward, owing to the tensorized structure of the polynomial basis so that the multidimensional Gauss points are simply obtained by tensorization of the one-dimensional Gauss points.

**Theorem 3.** *Assume that the stochastic Jacobian matrix  $\nabla_U F(\cdot; \xi)$  is defined at the  $(\text{No} + 1)$  Gauss points in  $[0, 1]$ . Consider the matrix  $\overline{\nabla_u f}$  obtained by approximating the coefficients of the Galerkin Jacobian matrix  $\nabla_u f$  by the above Gauss quadrature, namely*

$$(\overline{\nabla_u f}(u))_{\alpha,\beta=0,\dots,\text{No}} = \left( \sum_{\gamma=0}^{\text{No}} \omega_\gamma \nabla_U F(U^P(\xi_\gamma); \xi_\gamma) \Psi_\alpha(\xi_\gamma) \Psi_\beta(\xi_\gamma) \right)_{\alpha,\beta=0,\dots,\text{No}}. \quad (24)$$

Then,  $\overline{\nabla_u f}$  is  $\mathbb{R}$ -diagonalizable with eigenvalues  $\{\Lambda^k(\xi_\eta)\}_{k=1,\dots,m,\eta=0,\dots,\text{No}}$ , and eigenvectors  $\{v_\eta^k\}_{k=1,\dots,m,\eta=0,\dots,\text{No}}$  defined by

$$(v_\eta^k)_{\beta=0,\dots,\text{No}} = \langle V_\eta^k \Psi_\beta \rangle_{\beta=0,\dots,\text{No}}, \quad (25)$$

where  $V_\eta^k(\xi) \in \mathbb{R}^m \otimes \mathcal{S}^P$  is the polynomial of degree  $\leq \text{No} + 1$  in  $\xi$  such that

$$V_\eta^k(\xi_{\eta'}) = \delta_{\eta\eta'} W^k(\xi_\eta), \quad \eta' = 0, \dots, \text{No}. \quad (26)$$

Here,  $\{\Lambda^k(\xi)\}_{k=1,\dots,m}$  and  $\{W^k(\xi)\}_{k=1,\dots,m}$  are the eigenvalues and eigenvectors of the stochastic Jacobian matrix  $\nabla_U F(\cdot; \xi)$  defined by (7).

*Proof.* We observe that  $\sum_{\eta=0}^{\text{No}} V_\eta^k(\xi)$  is the interpolation polynomial of  $W^k(\xi)$  at the  $(\text{No} + 1)$  Gauss points. Since the order of the quadrature is  $(2\text{No} + 1)$ , for all  $V(\xi) \in \mathcal{S}^P$ ,  $\langle \Psi_\beta V \rangle$  is exact for all  $\beta = 0, \dots, \text{No}$  if evaluated using the quadrature. Hence, for all  $\beta = 0, \dots, \text{No}$ ,

$$(v_\eta^k)_\beta = \langle \Psi_\beta V_\eta^k \rangle = \sum_{\gamma=0}^{\text{No}} \omega_\gamma V_\eta^k(\xi_\gamma) \Psi_\beta(\xi_\gamma) = \omega_\eta W^k(\xi_\eta) \Psi_\beta(\xi_\eta). \quad (27)$$

Furthermore, observe that for all  $\xi$ ,

$$\sum_{\beta=0}^{\text{No}} (v_\eta^k)_\beta \Psi_\beta(\xi) = V_\eta^k(\xi), \quad (28)$$

since the basis is orthonormal. As a result,

$$\begin{aligned}
(\overline{\nabla_u f}(u)v_\eta^k)_\alpha &= \sum_{\beta=0}^{\text{No}} \left( \sum_{\gamma=0}^{\text{No}} \omega_\gamma \nabla_U F(U^P(\xi_\gamma); \xi_\gamma) \Psi_\alpha(\xi_\gamma) \Psi_\beta(\xi_\gamma) \right) (v_\eta^k)_\beta \\
&= \sum_{\gamma=0}^{\text{No}} \omega_\gamma \nabla_U F(U^P(\xi_\gamma); \xi_\gamma) \Psi_\alpha(\xi_\gamma) \left( \sum_{\beta=0}^{\text{No}} (v_\eta^k)_\beta \Psi_\beta(\xi_\gamma) \right) \\
&= \sum_{\gamma=0}^{\text{No}} \omega_\gamma \nabla_U F(U^P(\xi_\gamma); \xi_\gamma) \Psi_\alpha(\xi_\gamma) V_\eta^k(\xi_\gamma) \\
&= \omega_\eta \Psi_\alpha(\xi_\eta) \nabla_U F(U^P(\xi_\eta); \xi_\eta) W^k(\xi_\eta) = \omega_\eta \Psi_\alpha(\xi_\eta) \Lambda^k(\xi_\eta) W^k(\xi_\eta) = \Lambda^k(\xi_\eta) (v_\eta^k)_\alpha. \quad (29)
\end{aligned}$$

Moreover, the eigenvectors defined by (25) form a complete basis of  $\mathbb{R}^{m(\text{No}+1)}$ . Indeed, let  $m(\text{No}+1)$  reals  $(a_{k\eta})_{k=1, \dots, m, \eta=0, \dots, \text{No}}$  be such that  $\sum_{k, \eta} a_{k\eta} v_\eta^k = 0$ . This yields the stochastic vector  $\sum_\alpha \sum_{k, \eta} a_{k\eta} (v_\eta^k)_\alpha \Psi_\alpha(\xi)$ . Evaluating it at  $\xi_{\eta'}$  for any  $\eta' = 0, \dots, \text{No}$ , we obtain

$$\sum_\alpha \sum_{k, \eta} a_{k\eta} (v_\eta^k)_\alpha \Psi_\alpha(\xi_{\eta'}) = \sum_{k, \eta} a_{k\eta} \left( \sum_\alpha (v_\eta^k)_\alpha \Psi_\alpha(\xi_{\eta'}) \right) = \sum_{k, \eta} a_{k\eta} V_\eta^k(\xi_{\eta'}) = 0,$$

that is,  $\sum_k a_{k\eta'} W^k(\xi_{\eta'}) = 0$ . Since for each  $\eta'$ , the stochastic eigenvectors  $W^k(\xi_{\eta'})$  form a complete basis of  $\mathbb{R}^m$ , we infer  $a_{k\eta'} = 0$  for all  $k = 1, \dots, m$ . Since  $\eta'$  is arbitrary, the proof is complete.  $\square$

**Remark.** Theorem 3 can be exploited when working with partial polynomial tensorization since the approximate eigenvalues only serve as data to compute a fitting polynomial. This point will be further discussed in section 4.3 and illustrated numerically in section 5.2.3.

## 4. Numerical method

The Galerkin system (16) is discretized using a FV method [12, 34]. Consider for simplicity a uniform spatial step  $\Delta x$  and discrete times  $t^n$  with time step  $\Delta t^n = t^{n+1} - t^n$  verifying a CFL condition specified below. The FV scheme takes the form

$$u_i^{n+1} = u_i^n - \frac{\Delta t^n}{\Delta x} (\varphi(u_i^n, u_{i+1}^n) - \varphi(u_{i-1}^n, u_i^n)), \quad (30)$$

where  $u_i^n$  is an approximation to the mean value in space of the solution  $u$  in the cell of center  $i\Delta x$  with width  $\Delta x$  at the time  $t^n$  and  $\varphi(\cdot, \cdot)$  is the numerical flux. On a given interface  $LR$  separating left and right states indexed by  $L$  and  $R$  respectively, the numerical flux is chosen in the form

$$\varphi(u_L, u_R) = \frac{f(u_L) + f(u_R)}{2} - a \frac{u_R - u_L}{2}, \quad (31)$$

where  $\frac{1}{2}(f(u_L) + f(u_R))$  is the centered part of the flux and  $a \in \mathbb{R}^{m(P+1), m(P+1)}$  is a (nonnegative) upwind matrix whose construction will be discussed in section 4.3.

### 4.1. Roe matrix and Roe state

We assume that the original stochastic problem (4) possesses a Roe matrix  $A^{\text{Roe}}(U_L, U_R; \xi) \in \mathbb{R}^{m, m} \otimes L^2(\Xi, p_\xi)$  almost surely. Recall that  $A^{\text{Roe}}(U_L, U_R; \xi)$  verifies the following properties:

- $A^{\text{Roe}}(U_L, U_R; \xi)$  is  $\mathbb{R}$ -diagonalizable,  $\forall U_L, U_R \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi)$ .
- Consistency with the stochastic Jacobian matrix  $\nabla_U F$ ,

$$A^{\text{Roe}}(U, U; \xi) = \nabla_U F(U; \xi), \quad \forall U \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi).$$

- Conservativity through shocks,

$$F(U_R; \xi) - F(U_L; \xi) = A^{\text{Roe}}(U_L, U_R; \xi)(U_R - U_L), \quad \forall U_L, U_R \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi).$$

**Theorem 4.** Consider either sparse or full polynomial tensorization for the stochastic space  $\mathcal{S}^{\text{No}, \text{Nr}}$ . Under the above hypotheses,  $\forall u_L, u_R \in \mathcal{A}_u$ , the matrix  $a(u_L, u_R) \in \mathbb{R}^{m^P, m^P}$  defined by

$$a(u_L, u_R) := \langle A^{\text{Roe}}(U_L^P, U_R^P; \cdot) \Psi_\alpha \Psi_\beta \rangle_{\alpha, \beta=1, \dots, P} \quad (32)$$

with  $U_L^P(\xi) = \sum_{\alpha=1}^P (u_L)_\alpha \Psi_\alpha(\xi)$  and  $U_R^P(\xi) = \sum_{\alpha=1}^P (u_R)_\alpha \Psi_\alpha(\xi)$ , verifies the following properties:

- Consistency with the Galerkin Jacobian matrix  $\nabla_u f$ ,

$$a(u, u) = \nabla_u f(u), \quad \forall u \in \mathcal{A}_u.$$

- Conservativity through shocks,

$$f(u_R) - f(u_L) = a(u_L, u_R)(u_R - u_L), \quad \forall u_L, u_R \in \mathcal{A}_u.$$

*Proof.* To prove the consistency with the Galerkin Jacobian matrix, observe that  $\forall u \in \mathcal{A}_u$ , letting  $U^P = \sum_{\alpha=1}^P u_\alpha \Psi_\alpha(\xi)$ ,

$$a(u, u) = \langle \langle A^{\text{Roe}}(U^P, U^P; \cdot) \Psi_\alpha \Psi_\beta \rangle \rangle_{\alpha, \beta=1, \dots, P} = \langle \langle \nabla_U F(U^P; \cdot) \Psi_\alpha \Psi_\beta \rangle \rangle_{\alpha, \beta=1, \dots, P} = \nabla_u f(u).$$

To prove the conservativity through shocks, observe that  $\forall u_L, u_R \in \mathcal{A}_u$  and  $\forall \alpha = 1, \dots, P$ , letting  $U_L^P = \sum_{\alpha=1}^P (u_L)_\alpha \Psi_\alpha(\xi)$  and  $U_R^P = \sum_{\alpha=1}^P (u_R)_\alpha \Psi_\alpha(\xi)$ ,

$$\begin{aligned} (f(u_R) - f(u_L))_\alpha &= \langle (F(U_R^P; \cdot) - F(U_L^P; \cdot)) \Psi_\alpha \rangle = \langle A^{\text{Roe}}(U_L^P, U_R^P; \cdot) (U_R^P - U_L^P) \Psi_\alpha \rangle \\ &= \left\langle A^{\text{Roe}}(U_L^P, U_R^P; \cdot) \sum_{\beta=1}^P (\langle \Psi_\beta U_R^P \rangle - \langle \Psi_\beta U_L^P \rangle) \Psi_\alpha \Psi_\beta \right\rangle \\ &= \sum_{\beta=1}^P \langle A^{\text{Roe}}(U_L^P, U_R^P; \cdot) \Psi_\alpha \Psi_\beta \rangle (\langle \Psi_\beta U_R^P \rangle - \langle \Psi_\beta U_L^P \rangle) \\ &= \sum_{\beta=1}^P (a)_{\alpha, \beta} \left( (u_R)_\beta - (u_L)_\beta \right). \end{aligned}$$

This completes the proof.  $\square$

Assume furthermore that for all  $U_L, U_R \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi)$ , there exists a Roe state  $U_{LR}^{\text{Roe}} \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi)$  almost surely such that

$$A^{\text{Roe}}(U_L, U_R; \xi) = \nabla_U F(U_{LR}^{\text{Roe}}; \xi). \quad (33)$$

Then, for all  $U_L^P, U_R^P \in \mathcal{A}_U \otimes \mathcal{S}^P$ , introducing  $U_{LR}^{\text{Roe}} \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi)$  such that  $A^{\text{Roe}}(U_L^P, U_R^P; \xi) = \nabla_U F(U_{LR}^{\text{Roe}}; \xi)$ , we set

$$a_{LR}^{\text{Roe}} := a(u_L, u_R) := \langle \nabla_U F(U_{LR}^{\text{Roe}}; \cdot) \Psi_\alpha \Psi_\beta \rangle_{\alpha, \beta=1, \dots, P}. \quad (34)$$

Moreover, if  $a_{LR}^{\text{Roe}}$  is  $\mathbb{R}$ -diagonalizable, this matrix is a Roe linearized matrix.

#### 4.2. An efficient method for approximating the absolute value of a matrix

Let  $A$  be a deterministic  $\mathbb{R}$ -diagonalizable matrix of size  $N_A$ . The method presented here holds for a general matrix  $A$ ; its application to stochastic hyperbolic systems is detailed in section 4.3 below. By definition,  $|A|$  is the co-diagonalizable matrix with  $A$  whose eigenvalues are the absolute values of the eigenvalues of  $A$ ,

$$|A| = \sum_{\gamma=1}^{N_A} |\lambda_\gamma| l_\gamma \otimes r_\gamma, \quad (35)$$

where  $\{\lambda_\gamma\}_{\gamma=1,\dots,N_A}$  are the real eigenvalues of  $A$ ,  $\{l_\gamma\}_{\gamma=1,\dots,N_A}$  the left eigenvectors, and  $\{r_\gamma\}_{\gamma=1,\dots,N_A}$  the right eigenvectors. It is possible to diagonalize  $A$  and to compute  $|A|$  using (35), but in practice this method is extremely costly. A more interesting method has been proposed in [29], which consists in computing a sequence of polynomial iterations based on the exact knowledge of the eigenvalues (or at least an explicit bound), and converging to the matrix sign if all the eigenvalues are real. However, this method also becomes costly when  $N_A$  grows. Another method has been proposed in [7], relying on the computation of a polynomial which interpolates some absolute values of the eigenvalues of  $A$ . We derive here a new method based on a single computation of a low-degree polynomial. Our method is clearly less costly, and it is also better adapted to the situations where only approximations of the eigenvalues are known. Denote by  $\{\lambda'_\gamma\}_{\gamma=1,\dots,N_A}$  the approximate eigenvalues of  $A$ . The method consists in finding a polynomial  $q_{d,\{\lambda'\}}$  with degree  $d$  ( $d$  is fixed a priori) which minimizes the least-squares error between  $|\lambda'_\gamma|$  and  $q_{d,\{\lambda'\}}(\lambda'_\gamma)$ , and then applying this polynomial to the matrix  $A$  in order to approximate  $|A|$ .

Let  $q(X) = \sum_{j=0}^d c_j X^j$  be a polynomial. We seek  $q_{d,\{\lambda'\}}$  which minimizes the error  $\sum_{\gamma=1}^{N_A} (|\lambda'_\gamma| - q_{d,\{\lambda'\}}(\lambda'_\gamma))^2$ . It is well-known that this minimization problem is equivalent to solving a linear system with the polynomial coefficients  $(c_j)_{j=0,\dots,d}$  as unknowns. This system of size  $(d+1) \times (d+1)$  can be written as

$$\begin{pmatrix} \sum_{\gamma=1}^{N_A} \lambda'^0_\gamma \lambda'^0_\gamma & \dots & \sum_{\gamma=1}^{N_A} \lambda'^0_\gamma \lambda'^d_\gamma \\ \vdots & \ddots & \vdots \\ \sum_{\gamma=1}^{N_A} \lambda'^d_\gamma \lambda'^0_\gamma & \dots & \sum_{\gamma=1}^{N_A} \lambda'^d_\gamma \lambda'^d_\gamma \end{pmatrix} \begin{pmatrix} c_0 \\ \vdots \\ c_d \end{pmatrix} = \begin{pmatrix} \sum_{\gamma=1}^{N_A} |\lambda'_\gamma| \lambda'^0_\gamma \\ \vdots \\ \sum_{\gamma=1}^{N_A} |\lambda'_\gamma| \lambda'^d_\gamma \end{pmatrix}. \quad (36)$$

Solving this linear system yields the coefficients  $(c_j)_{j=0,\dots,d}$  that define the polynomial  $q_{d,\{\lambda'\}}$ . We then apply this polynomial to  $A$  and obtain an approximation to  $|A|$ . For efficiency, Hörner's method [33, p. 44] can be used:  $q_{d,\{\lambda'\}}(A)$  can be rewritten as

$$q_{d,\{\lambda'\}}(A) = c_0 I + (c_1 I + (c_2 I + \dots + (c_{d-1} I + c_d A) \dots A) A). \quad (37)$$

The number of matrix-matrix products is thus reduced to  $d$  instead of  $d(d-1)/2$  if all the powers of the matrix were computed independently. Therefore, the computational cost is proportional to  $d$  instead of being of order  $d^2$ . We can further reduce the computational cost in the present case since we only evaluate the product of  $|A|$  times a given vector  $x$ . By computing directly  $q_{d,\{\lambda'\}}(A)x$ , the cost is reduced to  $d$  matrix-vector products instead of  $d$  matrix-matrix products.

#### 4.3. The upwind scheme

We apply the method presented in the previous section to approximate the absolute value of  $a_{LR}^{\text{Roe}} \in \mathbb{R}^{m_P, m_P}$  at each interface  $LR$  in the spatial domain. To this purpose, we assume to have at our disposal explicit expressions of the eigenvalues  $\Lambda^1(\cdot; \xi), \dots, \Lambda^m(\cdot; \xi)$  of the stochastic Jacobian matrix  $\nabla_U F$  and we evaluate these eigenvalues at  $U_{LR}^{\text{Roe}}(\xi)$  and at the Gauss points of each stochastic element. This yields the approximate eigenvalues  $\{\lambda'_\gamma\}_{\gamma=1,\dots,m_P}$ . In other words, we use the eigenvalues of the matrix  $\overline{\nabla_u f}$  identified in Theorem 3. In the case of full polynomial tensorization, the number of approximate and exact eigenvalues is the same. In the case of sparse polynomial tensorization, there are more approximate eigenvalues than exact eigenvalues. The resulting fitting polynomial is still expected to catch relatively well the exact eigenvalues. Indeed, because of localization on each stochastic element, the eigenvalues are expected to be clustered around the eigenvalues of the stochastic Jacobian matrix  $\nabla_U F$  in each stochastic

element. We refer to the end of section 5.2.3 for an example. Choosing a degree  $d$  then yields a polynomial  $q_{d,\{\lambda'\}}$ . The linear system (36) can be singular if the number of distinct approximate eigenvalues is less than  $d$ . In particular, this occurs in the deterministic case. To properly handle this issue, we use a Singular Value Decomposition method. Moreover, an important point is that we exploit the diagonal block structure of the Galerkin Jacobian matrix  $\nabla_u f$  by evaluating a fitting polynomial on each stochastic element. The key advantage is that the polynomial has to fit less points, so that computations are at the same time more efficient and more accurate. As a result, (38) is applied separately on each stochastic element using a specific polynomial.

The numerical flux in the Finite Volume scheme (30) is chosen in the form

$$\varphi(u_L, u_R) = \frac{f(u_L) + f(u_R)}{2} - q_{d,\{\lambda'\}}(a_{LR}^{\text{Roe}}) \frac{u_R - u_L}{2}. \quad (38)$$

We emphasize that this flux is a numerical (upwind) flux associated with the Galerkin system (16), and not the projection of a numerical flux associated with the original stochastic problem (4) as some methods discussed in the introduction propose. Specifically, the constructed flux is not equivalent in general to the flux that would result from a non-intrusive projection using deterministic Roe fluxes at some collocation or quadrature points. In the present method, some collocative information is used to calculate the polynomial  $q_{d,\{\lambda'\}}$ , but this polynomial is applied to the Galerkin Jacobian matrix, so that we refer to our method as intrusive.

Finally, the time-step  $\Delta t^n$  is selected from a CFL-condition based on the highest characteristic velocity over the spatial and stochastic discretization cells. In practice,  $\Delta t^n$  is computed such that

$$\frac{\Delta t^n}{\Delta x} = \frac{C}{\max_{LR \in \mathcal{I}, \gamma=1, \dots, m_P} |\lambda'_\gamma|}, \quad (39)$$

where  $\mathcal{I}$  denotes the set of interfaces  $LR$  and  $\{\lambda'_\gamma\}_{\gamma=1, \dots, m_P}$  are the (deterministic) approximate eigenvalues identified above. In other words, the maximum of the eigenvalues over the stochastic domain is evaluated by considering the eigenvalues at the Gauss points of all the stochastic elements. In the sequel, we set the CFL constant  $C$  to 0.95.

We observe that the matrix  $q_{d,\{\lambda'\}}(a_{LR}^{\text{Roe}})$  is not guaranteed to control the eigenvalues of  $a_{LR}^{\text{Roe}}$  (this matrix is not even guaranteed to be nonnegative). Indeed, approximate eigenvalues have been used to build  $q_{d,\{\lambda'\}}$ , and, in addition, this polynomial only provides a least-square fit to the eigenvalues. This issue can possibly be handled by tightening the stochastic resolution or increasing the polynomial degree  $d$ ; these aspects will be further explored numerically in the next section. Furthermore, it is also possible to lower the CFL constant  $C$  as a safeguard in the case where the eigenvalues are underestimated.

## 5. Results

The methodology presented in the previous sections is assessed on three test cases. The first two deal with the Burgers equation and the third one with the Euler equations. Unless specified, full polynomial tensorization is used to span the stochastic approximation space. Furthermore, it is well-known that Roe solvers need to be supplemented with an entropy corrector to prevent non-entropic shocks across sonic points. The present test cases are designed so as to avoid this situation. The extension of the present Roe solvers to include entropy correctors is possible. Details are reported elsewhere [? ].

### 5.1. Test case 1: Burgers equation with positive wave speeds

The goal of this first test case is to assess the proposed methodology for a stochastic scalar conservation law (the Burgers equation) so that the Galerkin system is guaranteed to be hyperbolic from theorems 1 or 2, and involving only a positive wave speed so that the computation of  $|a_{LR}^{\text{Roe}}|$  is trivial.



### 5.1.1. Problem definition

We consider a one-dimensional spatial domain  $\Omega = [0, 1]$  with periodic boundary conditions. The governing equation, in conservative form, is

$$\frac{\partial U}{\partial t} + \frac{\partial F(U)}{\partial x} = 0, \quad F(U) = \frac{U^2}{2}, \quad (40)$$

and we consider an uncertain initial condition  $U^0(x, \xi)$  consisting of three piecewise constant deterministic states in  $x$ . Specifically, the three states are  $\bar{u}^1 = 1$ ,  $\bar{u}^2 = 1/2$ , and  $\bar{u}^3 = 1/6$ , and the position of some jumps is uncertain: the jump from states  $\bar{u}^1$  to  $\bar{u}^2$  occurs at a random location  $X_{1,2}$  having a uniform distribution in  $[0.1, 0.2]$ , while the jump from states  $\bar{u}^2$  to  $\bar{u}^3$  occurs at a random location  $X_{2,3}$  having a uniform distribution in  $[0.3, 0.4]$ . Finally, the jump from states  $\bar{u}^3$  to  $\bar{u}^1$  is at  $x_{31} = 0.6$ . The random locations  $X_{1,2}$  and  $X_{2,3}$  are independent and parameterized using two independent random variables  $\xi_1$  and  $\xi_2$  respectively, both with uniform distribution in  $[0, 1]$ :

$$X_{1,2} = 0.1 + 0.1\xi_1, \quad X_{2,3} = 0.3 + 0.1\xi_2, \quad \xi_1, \xi_2 \sim \mathcal{U}[0, 1]. \quad (41)$$

Therefore, the problem has two stochastic dimensions ( $N = 2$ ), and the dimension of the approximation space for expansion order  $N_o$  and resolution level  $N_r$  is  $\dim \mathcal{S}^{N_o, N_r} = (N_o + 1)^2 2^{2N_r}$ .

The initial condition is discretized on the spatial mesh by taking cell averaged random states as initial values. At the stochastic level, the discretization uses piecewise continuous bilinear approximations over the  $2^{2N_r}$  stochastic elements for  $N_o \geq 1$ , or the stochastic element averaged state for  $N_o = 0$ . The bilinear approximation uses nodal interpolation at the vertices at the stochastic elements, so that initial discrete states are continuous in the stochastic domain. This procedure prevents the presence of overshoots in the initial data. However, no particular treatment is applied to enforce the stochastic continuity during time integration. Figure 1 provides an illustration of the random initial condition for a spatial discretization with  $N_c = 200$  uniform cells in the spatial domain. The plot shows a sample set of 20 realizations of the random initial condition  $U^0(x, \xi)$ , with its expectation and standard deviation. It can be observed that the realizations present slightly inclined shocks, an effect caused by the cell average procedure and which can be reduced by taking a finer spatial mesh.

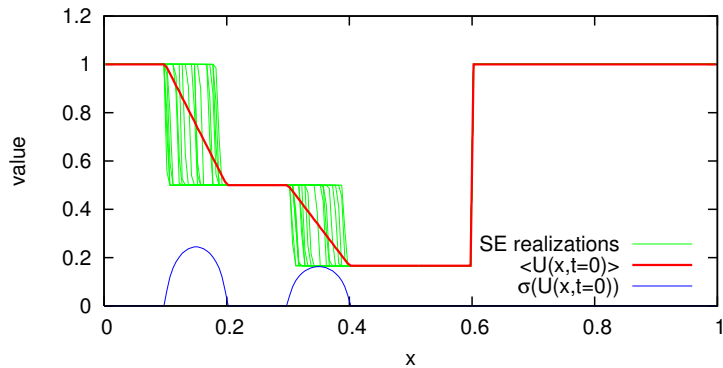


Figure 1: Random initial condition for test case 1: sample set of 20 random realizations, mean, and standard deviation.

### 5.1.2. Time integration

The stochastic Burgers equation is time-integrated using the Roe solver described above. We recall that since  $U$  is a scalar, the Galerkin problem is hyperbolic. The evaluation of the stochastic expansion of the nonlinear flux  $F(U)$  relies on an exact Galerkin projection using the third-order multiplication tensor  $\mathcal{M}_{\alpha\beta\delta}$  defined below by equation (55).

It is well-known that for the deterministic Burgers equation, the eigenvalue of the stochastic Jacobian matrix  $\nabla_U F$  is  $U$ . Because in the present setting the initial condition is almost surely positive for any  $x$ , we

expect  $U > 0$  with probability one, for all  $(x, t)$ . Therefore, the spectrum of the Galerkin Jacobian matrix is expected to be strictly positive, so that the upwinding matrix of the Galerkin problem reduces to the Galerkin Jacobian matrix (the polynomial transformation is in fact the identity).

Figure 2 shows the stochastic solution at times  $t = 0.2, 0.4, 0.6,$  and  $0.8$ . The computation uses  $N_r = 3$  and  $N_o = 3$ , so that the dimension of the stochastic space is  $16 \times 64 = 1024$ . The solution expectation and standard deviation, together with a random sample set of realizations, are also plotted. The realizations are reconstructed from the stochastic expansions of the solutions, using a unique set of randomly generated realizations of  $\xi \in [0, 1]^2$ .

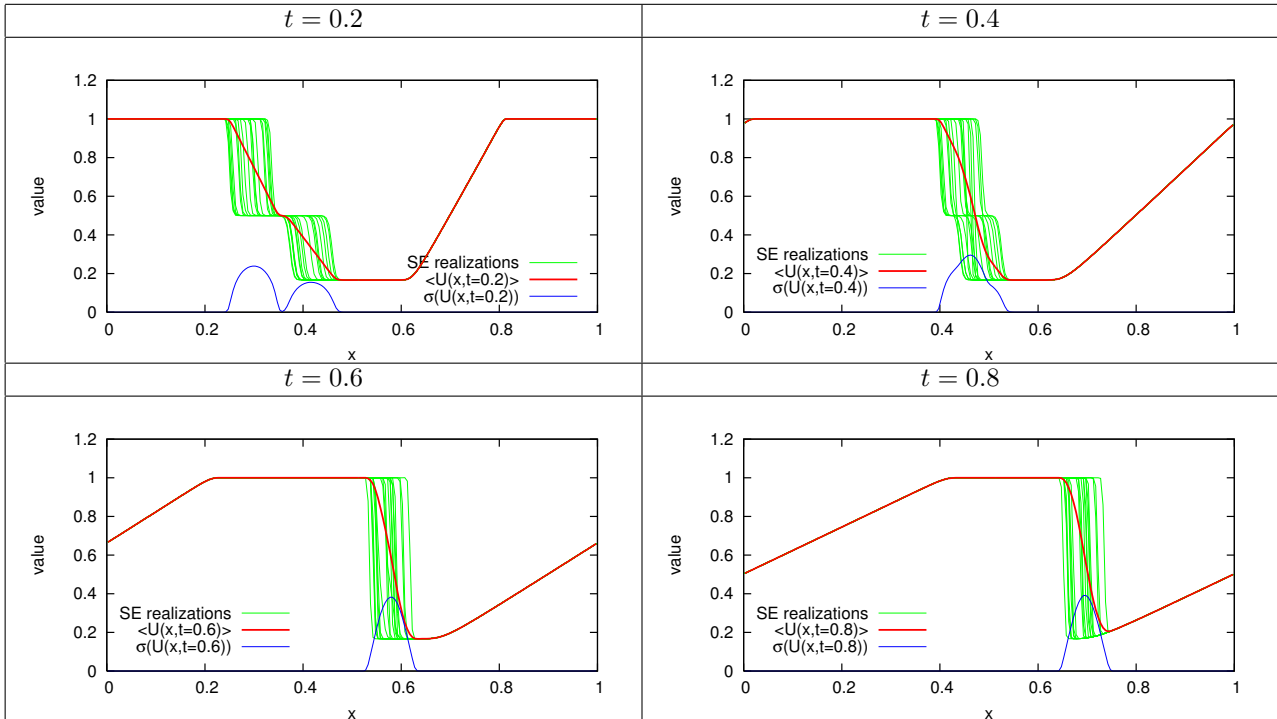


Figure 2: Solution of the stochastic Burgers equation at different times. The solution mean (red) and standard deviation (blue) are plotted as a function of  $x$ , together with a reconstruction of 20 randomly generated realizations (green). Computations with  $N_r = 3$  and  $N_o = 3$ .

Focusing first on the stochastic solution, we observe that the proposed method correctly captures the dynamics of the Burgers equation. The shocks are transported with the correct velocity and the discontinuities remain sharp as time evolves. For  $t = 0.2$ , the first shock whose velocity is  $0.75$  has not yet reached the second shock whose velocity is  $1/3$ . At  $t = 0.4$ , a fraction of the realizations corresponds to a situation where the first and second shocks have merged. At  $t = 0.6$ , the shocks have merged for nearly all realizations, a situation which is achieved at  $t = 0.8$ . It can be observed that the realizations, although corresponding to the same sample set of  $\xi$  in all plots, present a different distribution before and after the shocks have merged. Indeed, since the merging happens at different times depending on the initial locations of the two shocks and the shock velocities are different before and after merging, the location of the shock at later times is not expected to be uniform.

The uncertain shock dynamics can also be analyzed from the standard deviations of the stochastic solution: not only the maximum standard deviation is larger at  $t = 0.8$ , denoting the higher amplitude of the discontinuity, but the profiles are different. The expectation plots confirm the previous observations. While the uncertainty in shock location induces an affine evolution of  $\langle U \rangle$  when the two shocks are distinct, a variable slope of  $\langle U \rangle$  with  $x$  is observed after the shocks have merged: this indicates a non-uniform distribution of the shock location after merging. Similarly, the dynamics of the (deterministic) rarefaction

wave is well captured.

In addition to the analysis of the uncertain shock dynamics, Figure 2 also demonstrates that the Roe solver for the Galerkin system does not create spurious uncertainty in the solution, through numerical diffusion for instance. This can be better seen from Figure 3 where the space-time diagrams of the solution expectation and standard deviation are plotted over the larger period of time  $t \in [0, 2]$ . For time  $t > 0.7$ , in a moving frame attached to the remaining shock, the standard deviation reaches a maximum at  $t \approx 1$  where it peaks at  $\sigma(U) \approx 0.42$ , and then slowly decays. This decay is not a numerical artifact, but is induced by the rarefaction wave which has grown up to occupy the whole domain, as seen from the expectation plot where the plateau  $U = 1$  has disappeared for  $t > 1$ .

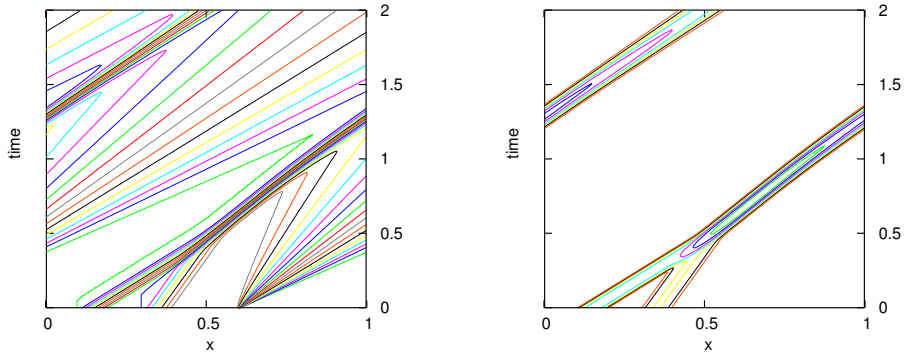


Figure 3: Space-time diagrams of the expectation (left) and standard deviation (right) of the stochastic Burgers solution. Contours are in the range  $[0, 1]$  with a constant spacing 0.05. Computations with  $Nr = 3$  and  $No = 3$ .

For analysis purpose, we define a moving observation point  $x_o(t) = 0.25 + 0.5t$ . The observation point is initially located between the two stochastic shocks. Since the velocity of  $x_o$  is lower than 0.75,  $x_o$  will be caught-up by the first random shock. Moreover, since  $x_o$  moves faster than the second shock, there is a time interval for which the stochastic solution at  $x_o$  corresponds to a set of events  $\xi$  with different configurations of the shocks. This is seen from Figure 4 where the stochastic solution  $U(x_o(t), t, \xi)$  is plotted as a function of  $\xi = (\xi_1, \xi_2)$  for various times  $t \in [0.2, 0.7]$ . For  $t = 0.2$ , the observation point starts to be caught-up by some events corresponding to the largest realizations of  $X_{1,2}$ : the solution is a function of  $\xi_1$  only. At  $t = 0.3$ , a larger fraction (roughly 1/4) of the first shock has overrun the observation point, and the stochastic solution exhibits two plateaus. At  $t = 0.4$ , the observation point starts to reach the second shock, introducing some dependence on  $\xi_2$ , while a fraction of events corresponds to shocks having merged. This creates a stochastic solution with three distinct plateaus with respective values 1, 1/2, and 1/6, whose configuration evolves in time. At  $t = 0.7$ , the solution at the observation point is essentially constant and equal to 1, with only a small fraction of events for which  $U = 1/6$ .

These results demonstrate the ability of the proposed method to account for nonlinear dynamics and complex interaction between random shocks. However, plots in Figure 4 deserve more comments. Firstly, although the numerical scheme allows for discontinuities across the stochastic discretization cells, the solutions reported here appear essentially continuous. While the initialization procedure ensures stochastic continuity of the initial condition, the numerical method maintains satisfactorily this property as time advances, as expected from the properties of the Burgers equation, provided that the resolution is fine enough. Secondly, the transitions between the states are smooth. This is due to the numerical diffusion of the Roe method which is known to spread the shocks over a few spatial cells. The smoothness of the stochastic solution reflects this spatial numerical diffusion. This point will be further evidenced below, where we show that the smooth transitions in the stochastic domain have a characteristic thickness independent of the stochastic resolution. In addition, we can observe that the smooth transitions are thicker along the second ( $\xi_2$ ) stochastic direction than along the first ( $\xi_1$ ). This is due to the different shock velocities (effects of different local CFLs).

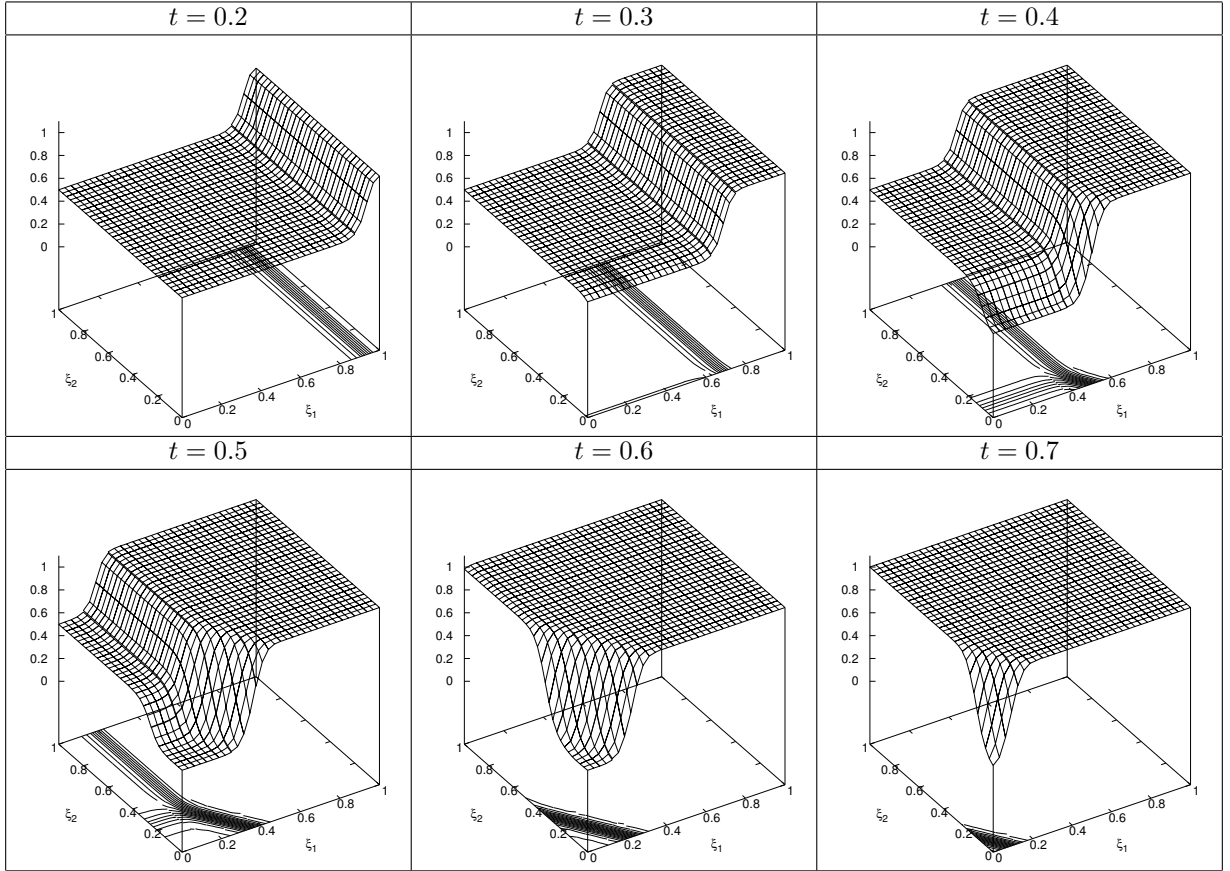


Figure 4: Stochastic solution of the Burgers equation at observation point  $x_o(t)$  as a function of  $(\xi_1, \xi_2)$  and for different times as indicated. Computations with  $N_o = 3$  and  $N_r = 3$ .

### 5.1.3. Convergence analysis

We present in Figure 5 the stochastic solutions at the observation point  $x = 0.5$  and time  $t = 0.5$  for different stochastic discretizations. The plots of the first line illustrate the convergence of the approximation with the expansion order  $N_o$ , while those of the second line highlight the convergence with the resolution level  $N_r$ . It is seen that when the stochastic discretization is too coarse, the solution exhibits significant discontinuities between stochastic discretization cells. Moreover, as claimed above, the transition thicknesses in the stochastic domain become independent of  $N_o$  and  $N_r$  as they increase.

## 5.2. Test case 2: Burgers equation with positive and negative wave speeds

The purpose of this test case is to assess the method still for the Burgers equation (so that the Galerkin system is guaranteed to be hyperbolic), but in a situation involving positive and negative wave speeds thereby requiring the calculation of  $|a_{LR}^{\text{Roe}}|$  as outlined in sections 4.2 and 4.3.

### 5.2.1. Problem definition

We still consider the Burgers equation, but with stochastic initial condition  $U^0(x, \xi)$  defined using two uncertain states,  $U^+(\xi_1)$  and  $U^-(\xi_2)$ , the first one almost surely positive and the second one almost surely

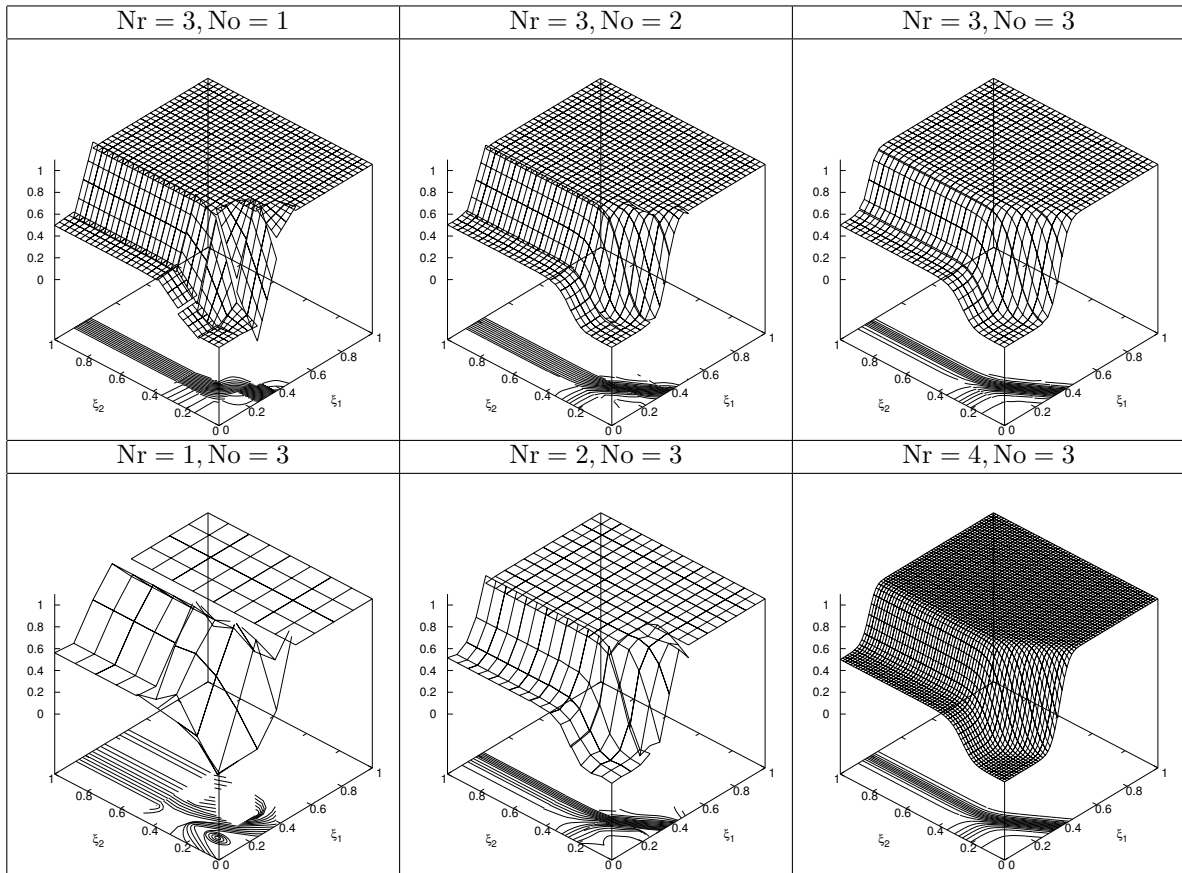


Figure 5: Stochastic solutions of the Burgers equation as a function of  $(\xi_1, \xi_2)$  at  $x = 0.5$  and time  $t = 0.5$  for different stochastic discretization parameters  $N_r$  and  $N_o$  as indicated.

negative. We take for  $x \in [0, 1]$ ,

$$U^0(x, \xi) = \begin{cases} U^+(\xi_1) & x < 1/3, \\ U^-(\xi_2) & x > 2/3, \\ U^+(\xi_1)(2 - 3x) + U^-(\xi_2)(3x - 1) & 1/3 \leq x \leq 2/3, \end{cases} \quad (42)$$

such that  $U^0(x, \xi)$  is continuous for any  $\xi \in [0, 1]^2$ . We define the stochastic states as

$$\begin{aligned} U^+(\xi_1) &= 1 + 0.1(2\xi_1 - 1), & \xi_1 &\sim \mathcal{U}[0, 1] \rightarrow U^+ \sim \mathcal{U}[0.9, 1.1], \\ U^-(\xi_2) &= -1 + 0.05(2\xi_2 - 1), & \xi_2 &\sim \mathcal{U}[0, 1] \rightarrow U^- \sim \mathcal{U}[-1.05, -0.95], \end{aligned} \quad (43)$$

and we solve the stochastic Burgers equation with Dirichlet boundary conditions,  $U = U^+$  at  $x = 0$  and  $U = U^-$  at  $x = 1$ . The initial condition is illustrated in Figure 6.  $N_c = 200$  cells are used for the spatial discretization.

### 5.2.2. Time integration

Although initially continuous, the stochastic solution will develop in finite time a discontinuity with a stochastic jump  $|U^+ - U^-|$  and a stochastic propagation velocity  $(U^+ + U^-)/2$ . The stochastic character of the shock magnitude and velocity has to be contrasted with the situation of the previous test case, where the jumps and shock velocity were certain. This yields a more complex situation as illustrated in Figure 7 where

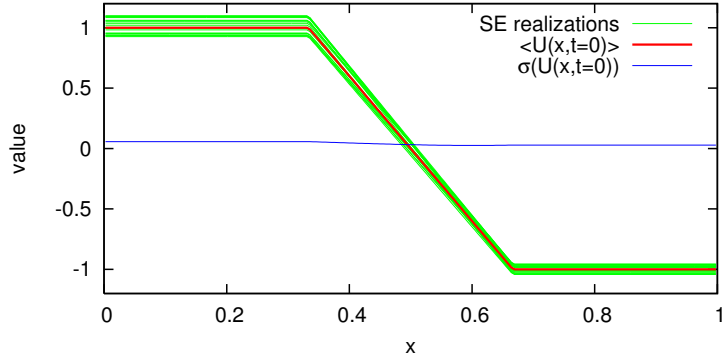


Figure 6: Random initial condition for test case 2: sample set of 20 random realizations, mean, and standard deviation.

the solution is plotted at different times for the stochastic discretization parameters  $N_o = 3$  and  $N_r = 3$  so that  $\dim \mathcal{S}^{N_o, N_r} = 1024$ .

From the realizations in Figure 7, we observe the appearance of overshoots which are related to the Gibbs phenomenon. We emphasize that no instability occurs as time increases. In order to verify this assertion, we compare in Figure 8 randomly generated realizations of the solution obtained by our method with realizations obtained by a non-intrusive projection method. Specifically, we rely on tensorized quadratures with  $N_Q^2$  points on each stochastic element to compute the solution modes at the selected analysis time. At each quadrature point, the stochastic solution is computed using a deterministic code also based on a Roe solver. Since  $U$  is not a polynomial in  $\xi$ , the number of quadrature points cannot be selected *a priori*; for all the results shown, we used  $N_Q = 16$ , a value for which the quadratures are sufficiently accurate for all expansion orders  $N_o$  and resolution levels  $N_r$  investigated. Note that the non-intrusive projection requires the resolution of a fairly large number ( $16^2 \times 2^{2N_r} = 1024 \times 2^{N_r}$ ) of independent deterministic Burgers equations. We see that overshoots are also present for the realizations of the stochastic expansion obtained by the non-intrusive method. Moreover, for both approaches, overshoots can be reduced by increasing the stochastic resolution. In addition, we compare in Figure 9 the maximum value of the stochastic solution of the Burgers equation as a function of time for  $N_r = 4$  and  $N_o = 3$  using the present method and a non-intrusive projection method. We observe that the maximum value increases with time in both cases until reaching a plateau at time  $t \approx 2.35$ . This again corroborates the stability of our method. Further test cases (omitted for brevity) show that the size of the overshoots can be reduced by increasing  $N_r$  at  $N_o$  fixed or vice versa.

To get further insight, we present in Figure 10 the evolution of the solution at a fixed point  $x_o = 0.5$  and different times. The plots show the evolution from the initially smooth solution to a shocked solution with states  $U^+$  or  $U^-$  according to the sign of  $2(\xi_1 - 1/2) - (\xi_2 - 1/2)$ . In addition, it is seen that overshoots occur only in a neighborhood of the discontinuity, namely in stochastic elements containing the developing discontinuity. Using a finer stochastic discretization (increasing  $N_r$ ) delays the emergence of the overshoots and reduces the portion of the stochastic domain affected by them.

### 5.2.3. Validation of the method used to evaluate the upwinding matrix

Another interesting property of the present test case is that contrary to the previous one, there exist spatial cells where the solution  $U$  can take positive and negative values. As a result, the eigenvalues of  $a_{LR}^{Roe}$  are no longer always positive, and the polynomial transformation  $q$  to approach the absolute value of  $a_{LR}^{Roe}$  is no longer trivial as in the previous example. We then investigate the impact of the selected polynomial degree  $d$  of  $q$  on the computed solution. In the example presented previously, we used polynomials with degree  $d = 3$ . In Figure 11 we report the stochastic solution at  $x = 0.5$  and  $t = 0.5$  computed using increasing polynomial degree  $d$ . It is seen that for  $d = 1$ , the solution exhibits spurious discontinuities and overshoots across the stochastic discretization cells containing the developing shock (where the solution changes sign) meaning that the eigenvalues of the upwinding matrix are not approximated with enough

accuracy. When  $d = 2$ , the overshoots and discontinuities are greatly reduced compared to the case  $d = 1$ . Increasing further  $d$  does not bring significant improvement in the solution. In fact, at that stage the error in the solution is essentially dominated by the stochastic and spatial discretization error, whereby the error in the approximation of  $|a_{LR}^{\text{Roe}}|$  for  $d > 3$  is negligible.

To measure more precisely the error on the approximation of  $|a_{LR}^{\text{Roe}}|$ , we compute the set of exact eigenvalues  $\{\lambda_\alpha\}_{\alpha=1,\dots,P}$  of  $a_{LR}^{\text{Roe}}$ . We then compare the quantities  $|\lambda_\alpha|$  with their respective polynomial approximation  $q(\lambda_\alpha)$ . The error is quantified using the following quantities

$$\epsilon_2^2 = \frac{1}{P} \sum_{\alpha=1}^P (|\lambda_\alpha| - q(\lambda_\alpha))^2 \quad \text{and} \quad \epsilon_\infty = \max_{1 \leq \alpha \leq P} ||\lambda_\alpha| - q(\lambda_\alpha)|. \quad (44)$$

We recall that the fitting polynomial  $q$  is actually different on each stochastic element. In Figure 12 we present the error measures at  $t = 0.4$  as a function of  $x$ . We first remark that the error is limited to the portion of the spatial domain where the stochastic shock can be present, and diminishes as  $d$  increases. Both error measures appear to stagnate when  $d$  increases beyond 5 as can be expected since the estimated eigenvalues (at the tensorized Gauss points) used for the determination of  $q$  are not the actual eigenvalues of  $a_{LR}^{\text{Roe}}$ .

Finally, we have verified that the present procedure to compute approximate upwind matrices can be applied when working with sparse polynomial tensorization. To this purpose, we have proceeded as described in section 4.3. The profiles of the stochastic solution as a function of  $(\xi_1, \xi_2)$  are similar to those reported in Figure 10, indicating that the approximate upwind matrix is sufficient to yield stable computations.

#### 5.2.4. Convergence of the stochastic error

We take advantage of this simple problem setting to investigate the convergence of the stochastic solution. Indeed, for this Riemann problem, we can easily derive the exact solution  $U(x, t, \xi)$  for any given  $\xi$ , hereafter denoted  $U^{ex}$ , as long as the shock has not reached one of the domain boundaries [12]. We rely on a Monte-Carlo sampling strategy to estimate the two first moments of  $U^{ex}$ . We proceed as follows. Firstly, a random sample set of  $M$  realizations of  $\xi$  is generated by sampling uniformly  $[0, 1]^2$ . Secondly, for each element  $\xi^{(i)}$  of the sample set, we define  $u^{(i)}(x, t) := U^{ex}(x, t, \xi^{(i)})$  for  $i = 1, \dots, M$ . The sample set estimate of the mean is

$$\langle U^{ex} \rangle(x, t) \approx \frac{1}{M} \sum_{i=1}^M u^{(i)}(x, t) = E_s(U^{ex})(x, t), \quad (45)$$

while the sample set estimate of the standard deviation is

$$\sigma^2(U^{ex})(x, t) \approx \frac{1}{M} \sum_{i=1}^M \left( u^{(i)} - E_s(U^{ex}) \right)^2(x, t) = \sigma_s^2(U^{ex})(x, t). \quad (46)$$

To minimize the random sampling error in the empirical estimate, we use  $M = 100000$ .

In Figure 13, we compare the mean and standard deviation of the exact and computed solution for  $\text{No} = 2$  and  $\text{Nr} = 4$  at  $t = 0.6$  on a mesh with  $\text{Nc} = 201$  cells. It is seen that the means of the computed and exact solutions are in excellent agreement. For the standard deviations, computed and exact solutions are in good agreement, although the computed solution slightly under-estimates the standard deviation with less than 5% of relative error. The top panel of Figure 14 presents the spatial distribution of the error of the standard deviation for various resolution levels of the spatial grid and for fixed stochastic discretization parameters  $\text{No} = 2$  and  $\text{Nr} = 4$ . We observe that refining the spatial grid improves the accuracy. The bottom panel of Figure 14 displays for  $\text{No} = 2$  and various  $\text{Nr}$  the spatially integrated error defined as  $S_h^2 = \Delta x \sum_{i=1}^{\text{Nc}} (\sigma(U_i^{ex}) - \sigma(U_i^P))^2$ , where the subscript  $i$  refers to the spatial discretization cell. The results show that except for the lowest stochastic resolution level and the finest spatial grid, the error is dominated by the spatial discretization error.

To further analyze the stochastic convergence of the method, we monitor the convergence on a fixed spatial mesh with  $N_c = 201$ . We consider the error measure

$$\epsilon_h^2(t) := \frac{1}{M} \sum_{i=1}^M \int_{\Omega} \left( U_h^{\text{No}, \text{Nr}}(x, t, \xi^{(i)}) - U_h^{MC}(x, t, \xi^{(i)}) \right)^2 dx, \quad (47)$$

where  $U_h^{\text{No}, \text{Nr}}(x, t, \xi^{(i)})$  and  $U_h^{MC}(x, t, \xi^{(i)})$  are evaluated for each element  $\xi^{(i)}$  in a sample set from the stochastic expansion of the computed solution and by solving the corresponding deterministic (discrete) Burgers problem respectively. We use a sample set dimension  $M = 10000$ . Figure 15 reports the stochastic error  $\epsilon_h^2$  at  $t = 0.6$  and  $t = 1.8$ , as a function of the resolution level Nr and for expansion orders No = 1, 2, and 3. In these simulations, the approximation of the upwind matrix uses a polynomial degree defined as  $d = \min(8, (\text{No} + 1)^2)$ . For both times, we observe a similar decay rate of the stochastic error as a function of resolution level. The errors are larger for longer times, since the shocks have expanded on a larger portion of the spatial domain.



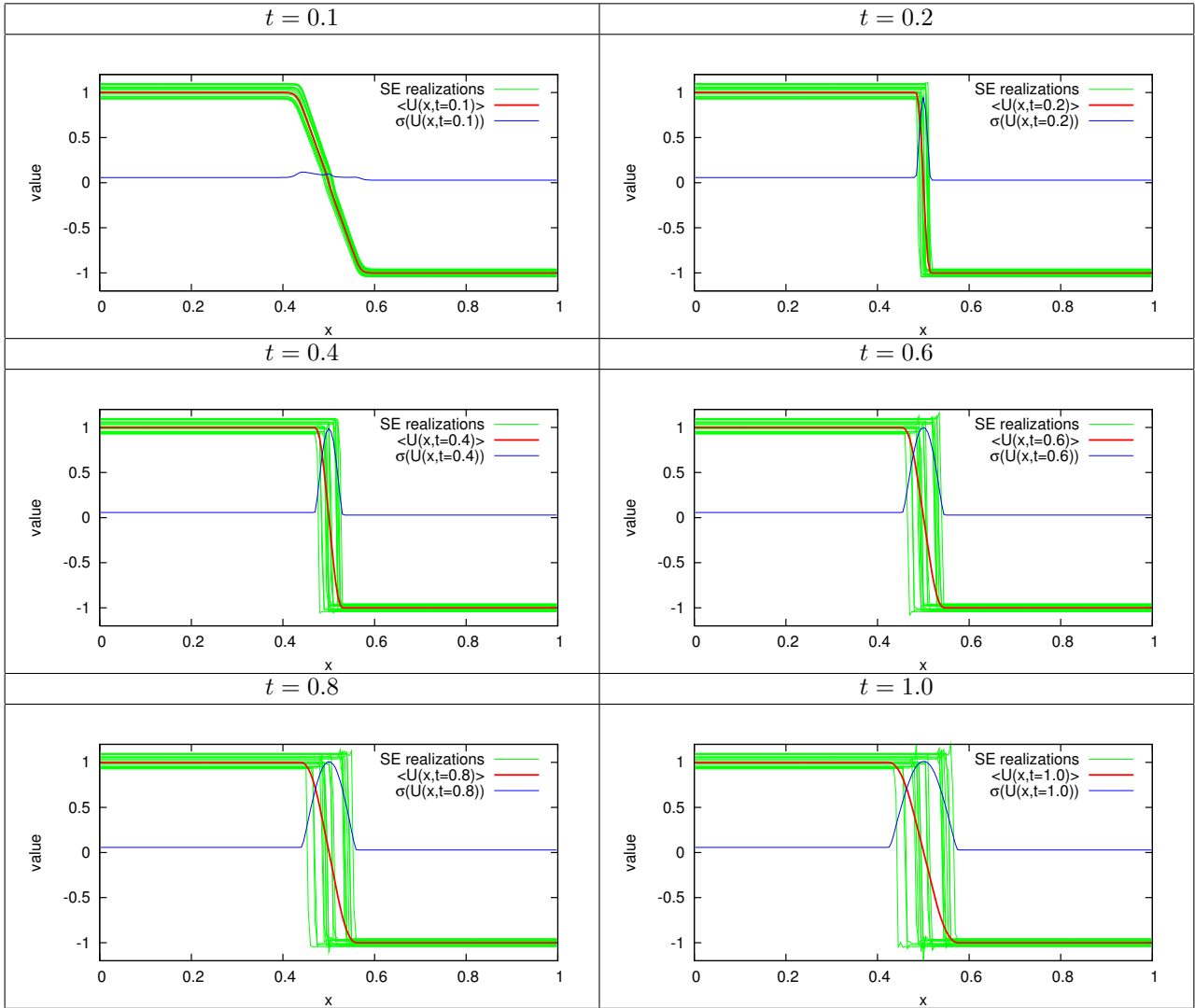


Figure 7: Stochastic solution of the Burgers equation at different times. The solution mean (red) and standard deviation (blue) are plotted as a function of  $x$ , together with a reconstruction of 20 randomly generated realizations of the solution (green). Computations with  $N_r = 3$  and  $N_o = 3$ .

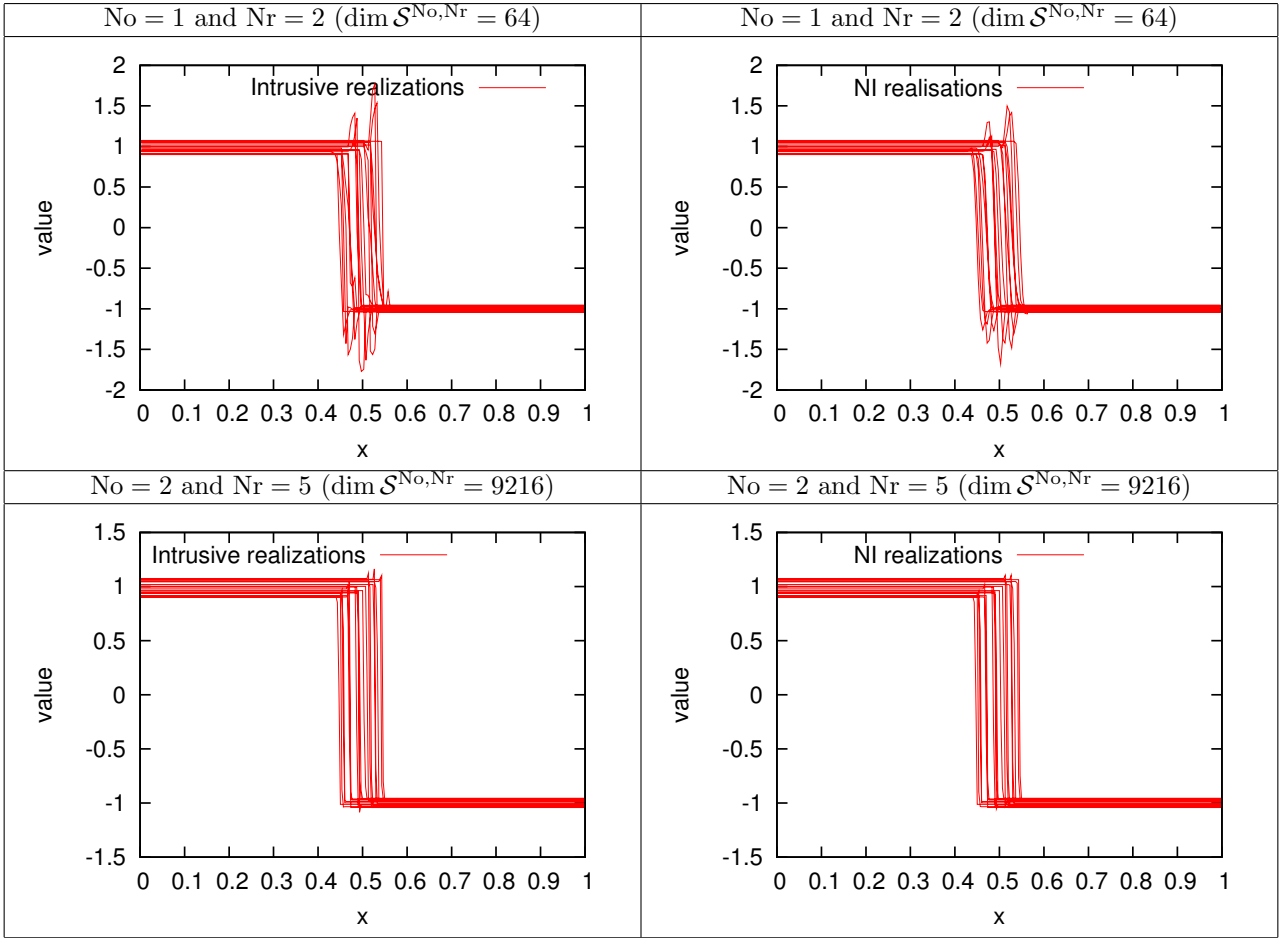


Figure 8: Reconstruction of 20 randomly generated realizations of the stochastic solution of the Burgers equation at  $t = 1.0$  s. Computations with different Nr and No using the present method (left) and a non-intrusive projection method (right).

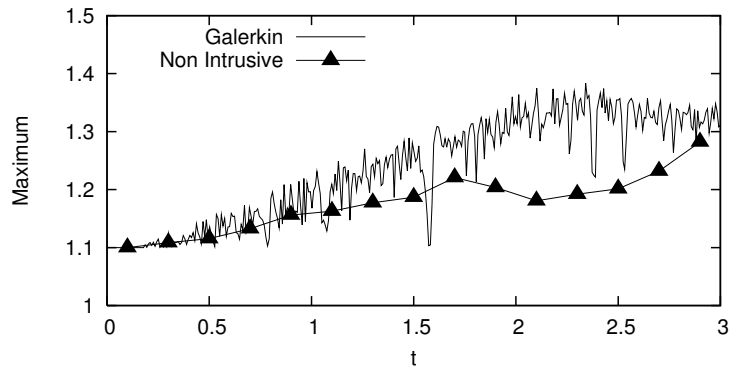


Figure 9: Maximum value of the stochastic solution of the Burgers equation as a function of time for Nr = 4 and Nr = 3 using the present method and a non-intrusive projection method.

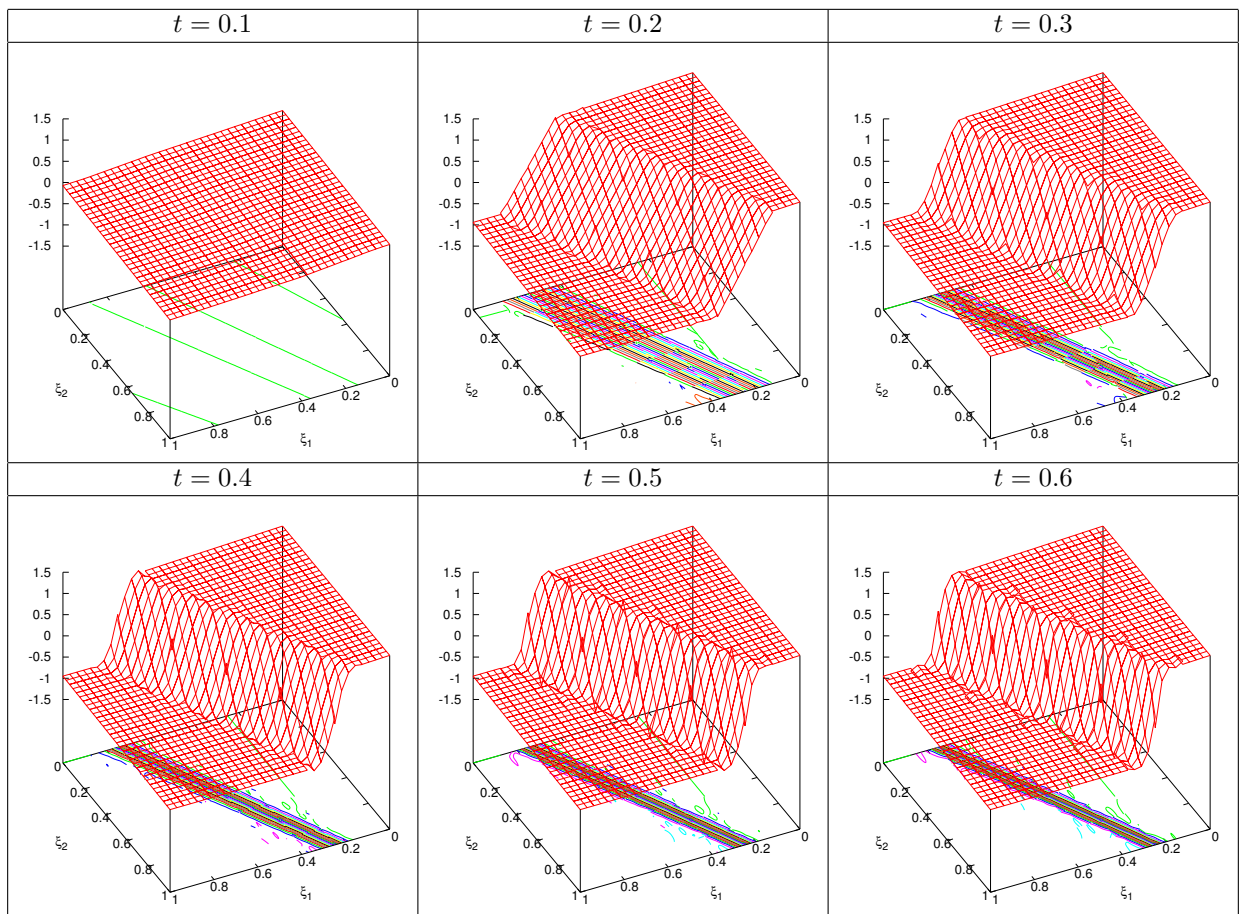


Figure 10: Stochastic solution of the Burgers equation as a function of  $(\xi_1, \xi_2)$  at  $x = 0.5$  and different times as indicated. Computations with  $N_r = 3$  and  $N_o = 3$ .

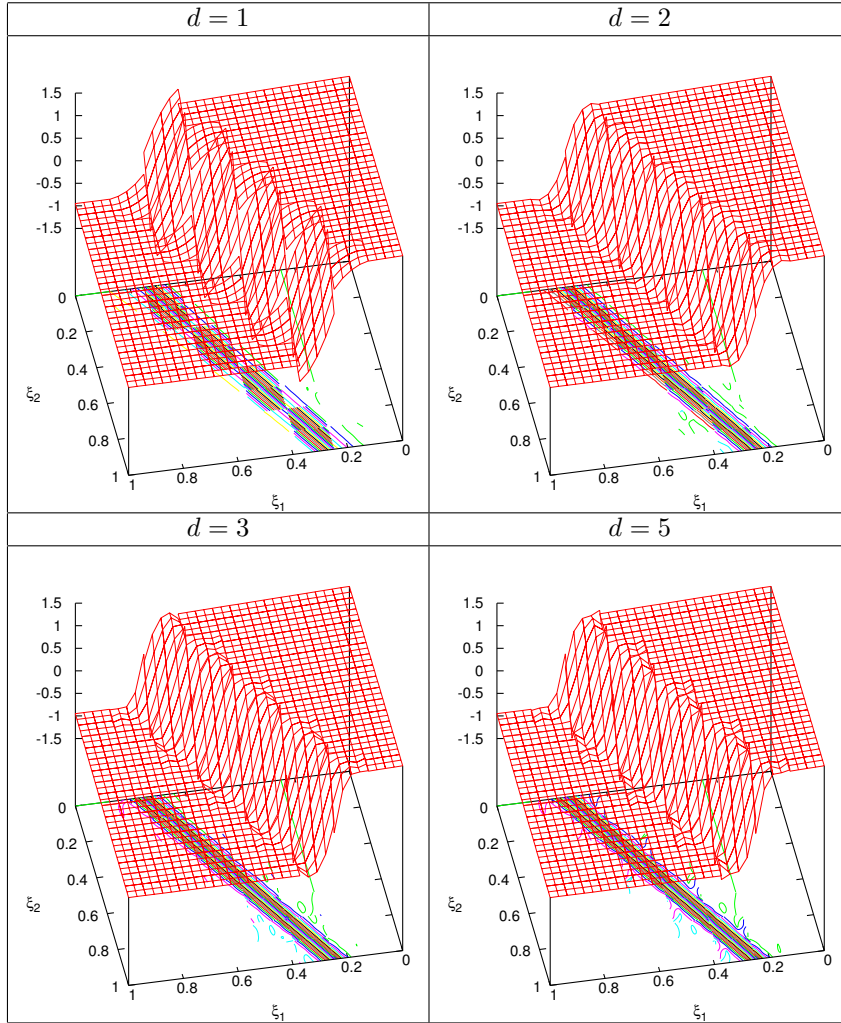


Figure 11: Stochastic solution of the Burgers equation as a function of  $(\xi_1, \xi_2)$  at  $x = 0.5$  and  $t = 0.5$  and for different degrees  $d$  of the polynomial  $q$  to approximate the absolute value of  $a_{LR}^{\text{Roe}}$ . Computations with  $N_r = 3$  and  $N_o = 3$ .

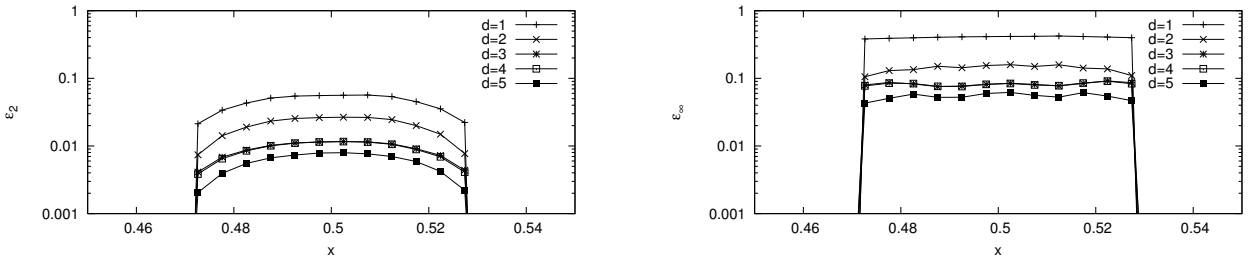


Figure 12: Measures  $\epsilon_2$  (left) and  $\epsilon_\infty$  (right) of the errors on the eigenvalues of the absolute value of  $a_{LR}^{\text{Roe}}$  at time  $t = 0.4$  and different degrees  $d$  of the polynomial  $q$ . Computations with  $N_r = 3$  and  $N_o = 3$ .

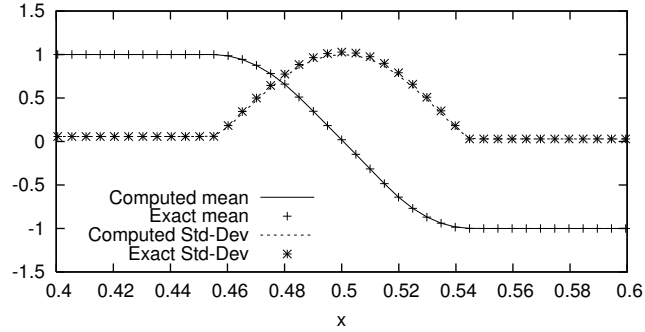


Figure 13: Comparison of the mean and standard deviation of the numerical solution at  $t = 0.6$ , computed with  $N_o = 2$ ,  $N_r = 4$  and  $N_c = 201$ , with the corresponding MC estimates of the mean and standard deviation of the exact solution of the stochastic Burgers equation. Only a portion of the computational domain is shown for clarity.

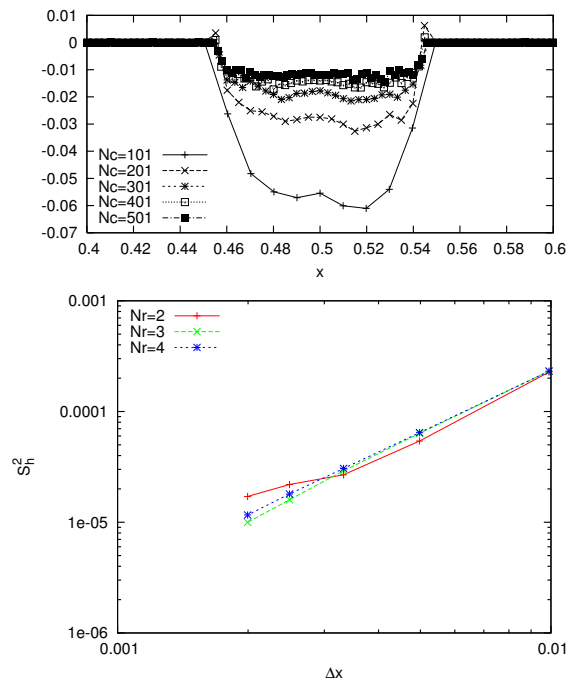


Figure 14: Errors in the standard deviation at  $t = 0.6$  for different spatial meshes. Top: Error distribution in spatial domain for  $N_o = 2$  and  $N_r = 4$ . Bottom: Spatially integrated error for  $N_o = 2$  and various  $N_r$ .

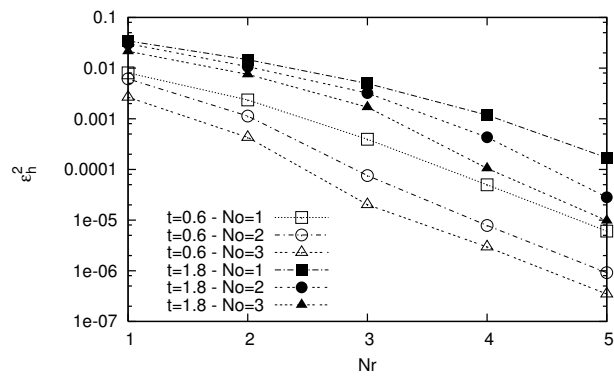


Figure 15: Stochastic error  $\epsilon_h^2$  at  $t = 0.6$  and  $t = 1.8$  as a function of the resolution level  $Nr$  and for different expansion orders  $No$  as indicated. Computations with  $Nc = 201$ .

### 5.3. Test case 3: Euler equations

In this section, the method is tested on the stochastic Euler equations with one random parameter. The goal of this test case is to assess the method on a nonlinear hyperbolic system of conservation laws, so that the obtained Galerkin system is not guaranteed to be hyperbolic. We consider the one-dimensional Sod shock tube problem, where the flow of an ideal gas is governed by the Euler equations. Conventional thermodynamic notation is used instead of the lower/upper case convention adopted previously. The conserved quantities are the fluid density  $\rho$ , the impulse  $q = \rho v$  (with  $v$  the velocity), and the total energy  $E = 1/2\rho v^2 + \rho e$ , where the first term is the kinetic energy and the second one the internal energy (per unit volume). The tube extends over one unit of length and is closed by two rigid walls at  $x = 0$  and  $x = 1$ . Boundary conditions are  $q = 0$  and  $\frac{\partial \rho}{\partial x} = \frac{\partial E}{\partial x} = 0$  at the solid walls. The discretization uses  $N_c = 250$  cells in the spatial domain.

#### 5.3.1. Problem definition

We consider an uncertainty on the adiabatic coefficient  $\gamma$  which is parametrized using a unique random variable  $\xi$  having a uniform distribution in  $[0, 1]$ . We consider a uniform probability distribution of  $\gamma$  in the range  $[1.4, 1.6]$ , so that the parametrization is

$$\gamma(\xi) = 1.4 + 0.2 \xi, \quad \xi \sim U[0, 1]. \quad (48)$$

Consistently with the notation introduced above, we set

$$U(x, t, \xi) = (\rho(x, t, \xi), q(x, t, \xi), E(x, t, \xi)) \in \mathcal{A}_U \otimes L^2(\Xi, p_\xi), \quad (49)$$

where  $\mathcal{A}_U \subset \mathbb{R}^3$  is the set of admissible states such that the density and the pressure are positive, and

$$F(U; \xi) = (F_\rho(U; \xi), F_q(U; \xi), F_E(U; \xi)) = (q(\xi), (q^2/\rho + p)(\xi), (v(E + p))(\xi)) \in \mathbb{R}^3 \otimes L^2(\Xi, p_\xi), \quad (50)$$

with the pressure  $p$  given by the ideal gas law

$$p(\rho, q, E) = (\gamma - 1) \left( E - \frac{1}{2} \rho v^2 \right). \quad (51)$$

The initial conditions are

$$\rho^0(x) = \begin{cases} 1 & x \in [0, 1/2], \\ 0.125 & x \in ]1/2, 1], \end{cases} \quad v^0(x) = 0, \quad \text{and} \quad p^0(x) = \begin{cases} 1 & x \in [0, 1/2], \\ 0.125 & x \in ]1/2, 1]. \end{cases} \quad (52)$$

#### 5.3.2. Numerical solver

*Computation of the Galerkin flux*  $f(u) \in \mathbb{R}^{3P}$ . Let  $u = (\rho_\alpha, q_\alpha, E_\alpha)_\alpha \in \mathbb{R}^{3P}$  yielding the expansions  $\rho^P(\xi) = \sum_{\alpha=1}^P \rho_\alpha \Psi_\alpha(\xi)$ ,  $q^P(\xi) = \sum_{\alpha=1}^P q_\alpha \Psi_\alpha(\xi)$ , and  $E^P(\xi) = \sum_{\alpha=1}^P E_\alpha \Psi_\alpha(\xi)$ . Then, the Galerkin flux has components  $(f_{\rho\alpha}, f_{q\alpha}, f_{E\alpha})_\alpha$  such that  $(f_{\rho\alpha})_\alpha$ ,  $(f_{q\alpha})_\alpha$ , and  $(f_{E\alpha})_\alpha$  are the stochastic modes of the components  $F_\rho(U^P; \xi)$ ,  $F_q(U^P; \xi)$ , and  $F_E(U^P; \xi)$  of the flux  $F(U^P; \xi)$  with  $U^P = (\rho^P, q^P, E^P)$ .

Contrary to approaches that compute the Galerkin flux in a non-intrusive way by quadrature formulae, we consider an approximation of the Galerkin projection of the flux  $F(U^P; \xi)$  on  $\mathcal{S}^P$ . Exact Galerkin projections of the stochastic Euler fluxes can hardly be envisioned since they would result in unnecessary complex nonlinear operations. As motivated in [6], pseudo-spectral computations allow for significant computational savings, while introducing negligible numerical errors as long as the stochastic resolution is fine enough. Furthermore, proceeding step by step in the approximation of the nonlinear fluxes yields intermediate quantities, such as kinetic energy and sound velocity, that can be re-used at different steps of the numerical scheme, in particular when computing the Roe state.

Tools for accurate evaluations of polynomial and non-polynomial functions of variables represented by stochastic expansions are described in [6]. Letting  $a(\xi) = \sum_{\alpha} a_\alpha \Psi_\alpha(\xi) \in \mathcal{S}^P$  and  $b(\xi) = \sum_{\beta} b_\beta \Psi_\beta(\xi) \in \mathcal{S}^P$ , the product  $ab$  can be expanded as

$$(ab)(\xi) = \left( \sum_{\alpha=1}^P a_\alpha \Psi_\alpha \right) \left( \sum_{\beta=1}^P b_\beta \Psi_\beta \right) = \sum_{\alpha, \beta=1}^P a_\alpha b_\beta \Psi_\alpha \Psi_\beta. \quad (53)$$

Generally,  $(ab) \notin \mathcal{S}^P$  since this function possesses terms with degree  $> P$ . Its Galerkin projection on  $\mathcal{S}^P$  is given by

$$(a * b) := \sum_{\alpha=1}^P (a * b)_\alpha \Psi_\alpha, \quad (a * b)_\alpha = \sum_{\beta, \delta=1}^P a_\beta b_\delta \mathcal{M}_{\alpha\beta\delta}, \quad (54)$$

where we have introduced the third-order multiplication tensor

$$\mathcal{M}_{\alpha\beta\delta} := \langle \Psi_\alpha \Psi_\beta \Psi_\delta \rangle. \quad (55)$$

This third-order tensor depends only on the stochastic basis, can be computed once and for all at the beginning of the simulation, and its sparse character in the stochastic space is exploited for its storage. The so-called Galerkin product  $(a * b)$  is the building block for evaluating the projection on  $\mathcal{S}^P$  of more general nonlinearities. For instance, the Galerkin projection on  $\mathcal{S}^P$  of  $1/a$ , which we denote by  $a^{-*}$ , and the Galerkin projection on  $\mathcal{S}^P$  of  $\sqrt{a}$ , which we denote by  $a^{*/2}$ , are obtained from the resolution of the linear system  $a * a^{-*} = 1$  and from the resolution with Newton's method of the nonlinear system  $a^{*/2} * a^{*/2} = a$ , respectively.

For the Euler equations, we define the components  $F_\rho^*(U^P; \xi)$ ,  $F_q^*(U^P; \xi)$ , and  $F_E^*(U^P; \xi)$  of the stochastic flux  $F^*(U^P; \xi)$  as follows:

$$F_\rho^*(U^P; \cdot) = q^P, \quad F_q^*(U^P; \cdot) = (q^P * q^P) * (\rho^P)^{-*} + p^*, \quad F_E^*(U^P; \cdot) = v^* * (E^P + p^*), \quad (56)$$

where  $v^* := q^P * (\rho^P)^{-*}$  and  $p^* := (\gamma - 1) * (E^P - (q^P * q^P) * (\rho^P)^{-*}/2)$ . We observe that  $F^*(U^P; \xi)$  is only an approximation of the Galerkin projection of  $F(U^P; \xi)$  since the composition of the elementary Galerkin operations (product, inversion) introduces a so-called pseudo-spectral approximation. Consistently, the pseudo-spectral Galerkin flux is denoted by

$$f^*(u) = (f_\alpha^*(u))_\alpha = (f_{\rho\alpha}^*, f_{q\alpha}^*, f_{E\alpha}^*)_\alpha, \quad (57)$$

where  $(f_{\rho\alpha}^*)_\alpha$ ,  $(f_{q\alpha}^*)_\alpha$ , and  $(f_{E\alpha}^*)_\alpha$  are the stochastic modes of  $F^*(U^P; \xi)$ . All in all, the computation of the pseudo-spectral Galerkin flux amounts to four Galerkin products and a Galerkin inversion.

*Computation of the Galerkin Jacobian matrix.* The pseudo-spectral Galerkin Jacobian matrix  $\nabla_u f^*(u) \in \mathbb{R}^{3P, 3P}$  is given by

$$\nabla_u f^*(u) = \left( \sum_{\delta=1}^P (\nabla_U F^*(U^P; \cdot))_\delta \mathcal{M}_{\alpha\beta\delta} \right)_{\alpha\beta}. \quad (58)$$

Letting  $H^* := (E^P + p^*) * \rho^{-*}$ ,  $\nabla_U F^*(U^P; \cdot)$  is defined as

$$\nabla_U F^*(U^P; \cdot) = \begin{pmatrix} 0 & 1 & 0 \\ 1/2(\gamma - 3) * (v^* * v^*) & -(\gamma - 3) * v^* & \gamma - 1 \\ 1/2(\gamma - 1) * (v^* * (v^* * v^*)) - v^* * H^* & H^* - (\gamma - 1) * (v^* * v^*) & \gamma * v^* \end{pmatrix}. \quad (59)$$

In terms of computational costs, this approach requires nine Galerkin products and a Galerkin inversion owing, in particular, to the re-use of the Galerkin product  $(v^* * v^*)$ .

*Computation of  $a_{LR}^{\text{Roe}}$  and its absolute value.* The Roe density, velocity, enthalpy, and corresponding sound velocity are approximated on  $\mathcal{S}^P$  for a given interface  $LR$  as

$$\rho_{LR}^{\text{Roe},*} := (\rho_L^P)^{*/2} * (\rho_R^P)^{*/2}, \quad v_{LR}^{\text{Roe},*} := ((\rho_L^P)^{*/2} * v_L^* + (\rho_R^P)^{*/2} * v_R^*) * ((\rho_L^P)^{*/2} + (\rho_R^P)^{*/2})^{-*}, \quad (60)$$

$$H_{LR}^{\text{Roe},*} := ((\rho_L^P)^{*/2} * H_L^* + (\rho_R^P)^{*/2} * H_R^*) * ((\rho_L^P)^{*/2} + (\rho_R^P)^{*/2})^{-*}, \quad (61)$$

$$(c_{LR}^{\text{Roe},*})^2 := (\gamma - 1) * (H_{LR}^{\text{Roe},*} - (v_{LR}^{\text{Roe},*} * v_{LR}^{\text{Roe},*})/2). \quad (62)$$



This yields  $U_{LR}^{\text{Roe},*}$  and the pseudo-spectral Roe matrix  $a_{LR}^{\text{Roe},*}$  is given by

$$a_{LR}^{\text{Roe},*} := \left( \sum_{\delta=1}^P (\nabla_U F^*(U_{LR}^{\text{Roe},*}; \cdot))_{\delta} \mathcal{M}_{\alpha\beta\delta} \right)_{\alpha\beta}. \quad (63)$$

Finally, the absolute value of  $a_{LR}^{\text{Roe},*}$  is computed as described in sections 4.2 and 4.3, using the approximation in  $\mathcal{S}^P$  of the stochastic eigenvalues of  $\nabla_U F^*(U_{LR}^{\text{Roe},*}(\xi); \xi)$  at Gauss points in each stochastic element, that is,  $(v_{LR}^{\text{Roe},*} \pm c_{LR}^{\text{Roe},*})(\xi_{\eta})_{\eta=0, \dots, \text{No}}$  and  $v_{LR}^{\text{Roe},*}(\xi_{\eta})_{\eta=0, \dots, \text{No}}$ .

### 5.3.3. Results

In this section we present and analyze the results for the shock tube problem with uncertainty in the adiabatic coefficient. We begin with a general analysis of the results, taking  $\text{No} = 2$  and  $\text{Nr} = 3$  as stochastic discretization parameters, so that the dimension of the stochastic space is 24.

In the deterministic case and for the initial condition (52) for a certain realization of  $\gamma(\xi)$ , a shock wave generated at the discontinuity travels to the right with velocity  $v + c$ , while a slower rarefaction fan travels to the left with velocity  $v - c$ , and a contact discontinuity wave travels to the right with velocity  $v$ . When the waves reach the solid wall, they are reflected inside the spatial domain and so propagate toward each other, merge, and interact. When the waves have crossed, they continue to propagate up to the point where they again reflect on a wall, and so on.

Here, the uncertain sound velocity will affect the propagation velocity of the shock, contact discontinuity, and rarefaction fan. Solutions for different realizations of  $\gamma(\xi)$  exhibit similar patterns as in the deterministic case, but with different slopes for the shock, contact discontinuity, and rarefaction fan in the space-time diagram. This is verified in Figure 16 where the density in the deterministic case (with adiabatic coefficient set to  $\langle \gamma \rangle$ ) and its expectation in the stochastic case are plotted. The spreading of the location of both the shock and the contact discontinuity when time increases is clearly visible, while for the rarefaction fan, which is already smooth in the deterministic case, the impact of the uncertain sound velocity is less pronounced.

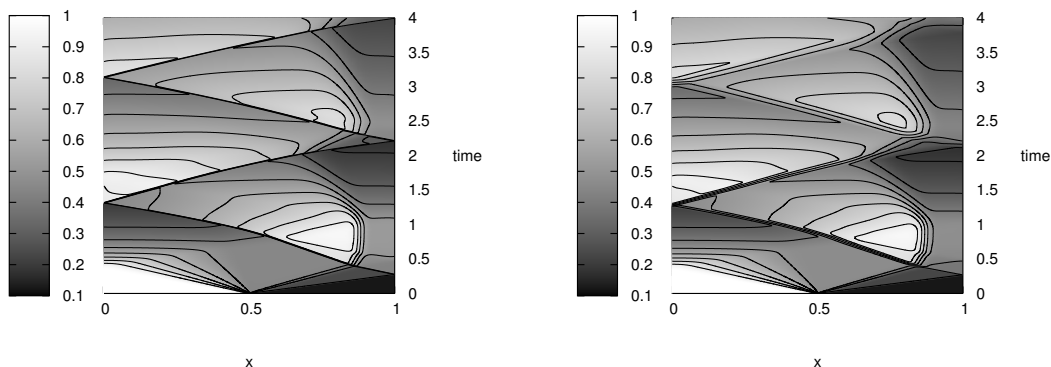


Figure 16: Space-time diagram of the deterministic density (left) and the expected density computed with parameters  $\text{No} = 2$  and  $\text{Nr} = 3$  (right).

The impact of the uncertainty can also be appreciated from the standard deviations of the density, reported in Figure 17 for early and longer times. The highest values of the standard deviations are observed along the path of the shock wave, the maximum values corresponding to times at which the shock wave reflects on the tube walls. For early times ( $t \leq 0.25$ ), uncertainty is present only in areas where the shocks can depend on the sound velocity in the prescribed uncertainty range. The first process leading to larger uncertainty levels is the shock-wall interaction since the arrival of the shock at a wall causes an abrupt increase of the density over a short time interval. The uncertainty in the arrival time of the shock therefore induces a large variability in the solution. The second process leading to larger uncertainty levels is the interaction between the uncertain shock, contact discontinuity, and rarefaction fan.

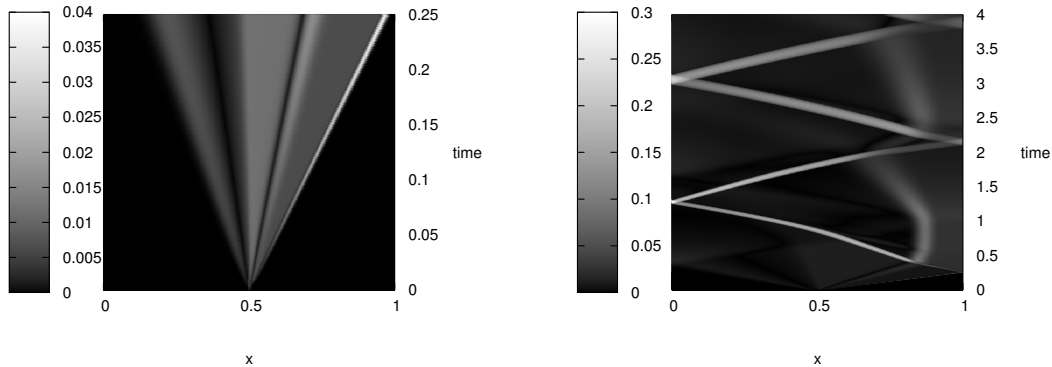


Figure 17: Space-time diagram of the standard deviations in the density for early (left) and longer times (right) computed with parameters  $No = 2$  and  $Nr = 3$ . Different color scales are used.

To assess the validity of the stochastic expansion, we show in Figure 18 a reconstruction of the stochastic density  $\rho(x, t, \xi)$  at selected times. The discontinuity in  $\rho(x, \cdot, \xi)$  is initially in the  $x$ -direction. As time increases, the density becomes discontinuous in both  $x$ - and  $\xi$ -directions since the stochastic shock wave propagates with an uncertain velocity. In the  $(x, \xi)$ -plane, the discontinuity becomes more and more oblique reflecting a monotone dependence of the shock velocity on  $\xi$ . For points  $(x, \xi)$  not too close to the discontinuity, the solution is smooth and appears to be accurately approximated by the stochastic expansion. In the neighborhood of the discontinuity, the solution exhibits small unphysical oscillations that are triggered by the well-known Gibbs phenomenon so that the density takes values slightly outside its expected range. Such oscillations appear to be caused by an insufficient stochastic resolution and can be reduced by increasing the resolution level and/or the polynomial order of the stochastic approximation. This is illustrated in Figures 19 and 20, which show the convergence of the density field as the value of  $No$  or  $Nr$  is increased. Oscillations become smaller as the level of stochastic resolution increases.

We define the error measure on the density as

$$\epsilon_h(x, t) := \left( \frac{1}{M} \sum_{i=1}^M \left( \rho_h^{No, Nr}(x, t, \xi^{(i)}) - \rho_h^{MC}(x, t, \xi^{(i)}) \right)^2 \right)^{1/2}, \quad (64)$$

where  $\rho_h^{No, Nr}(x, t, \xi^{(i)})$  and  $\rho_h^{MC}(x, t, \xi^{(i)})$  are evaluated for each element  $\xi^{(i)}$  in a sample set from the stochastic expansion of the computed solution and by solving the corresponding deterministic (discrete) Euler problem respectively. We use a sample set dimension  $M = 10000$ . Figure 21 reports the error  $\epsilon_h(x, t)$  for early and longer times using parameters  $Nc = 250$ ,  $Nr = 3$ , and  $No = 2$ . For early time, uncertainty has not propagated all over the spatial domain. We can distinguish three zones of error corresponding to the neighborhoods of the three waves. We notice as expected that the error hits its maximum in the neighborhood of the shock. As time increases, uncertainties propagate all over the domain. The zone of error corresponding to the shock spreads and the discontinuity in the  $(x, \xi)$ -plane becomes more oblique. Furthermore, after several reflexions of the waves on the walls ( $t = 6.5$ ), the error remains small indicating that no instability occurs at longer times. In Figure 22, we examine the convergence of the error  $\epsilon_h(x, t = 6.5)$  as the value of  $No$  or  $Nr$  is increased, confirming that both parameters can be used to improve stochastic resolution.

To complete the discussion, we provide a brief estimate of the computational efficiency of the numerical method for the Euler problem with random  $\gamma$ . We show in Table 1 the evolution of the computational times for different stochastic discretizations. CPU times ( $T_{CPU}$ ) are reported for an integration of the Euler equations up to  $t = 3$  on a fixed grid with  $Nc = 250$  cells, and are normalized by the computational time using  $No = Nr = 0$ , *i.e.* for the deterministic problem. Since the time step for the integration is based on a fixed CFL, it also depends on the stochastic discretization: the measured times correspond to different numbers of iterations performed. However, we observed roughly 0.5% variability in the number of time

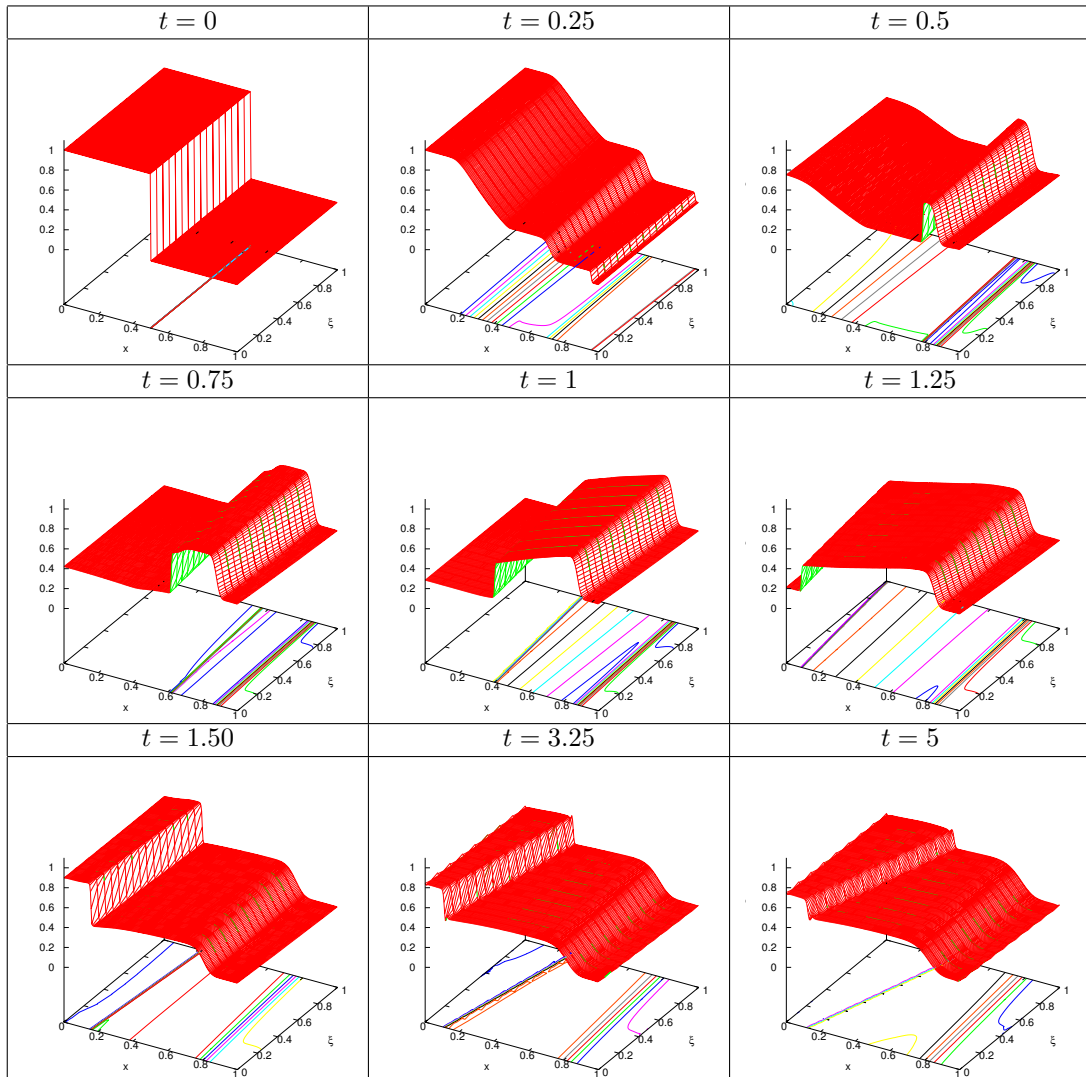


Figure 18: Reconstruction of the stochastic density  $\rho(x, t, \xi)$  at selected times.

iterations between the most and least refined simulations, hence, CPU times can also be interpreted as times to perform a fixed number of iterations.

Inspection of Table 1 shows a linear scaling with the number  $2^{N_r}$  of stochastic elements for fixed polynomial order  $N_o$ . This scaling was expected and achieved owing to the decoupling of the Galerkin problems over the stochastic elements for the projection on the SE basis. This linear scaling with respect to the number of stochastic elements is expected to hold also for problems with higher stochastic dimension ( $N > 1$ ); however, the overall CPU time increases significantly with  $N$  owing to the exponentially growing dimension of the local stochastic basis (except for  $N_o = 0$ ). For fixed resolution level  $N_r$ , the scaling of  $T_{CPU}$  with the dimension of  $\mathcal{S}^{N_o, N_r}$  is roughly linear at least for  $N_o \leq 4$ . There are two effects. Firstly, the spectral evaluation of the nonlinearities in the flux and Roe's states has a complexity (number of operations) that essentially scales with the number of nonzero terms in the third-order multiplication tensor  $\mathcal{M}_{\alpha\beta\delta}$ , which itself increases exponentially with  $N_o$ . Secondly, as  $N_o$  increases, a higher degree  $d$  has to be used for the polynomial approximation of the upwind matrixes, resulting in higher computational costs. The second effect can be tempered, based on the numerical experiments on the Burgers equation, by limiting  $d$  to a

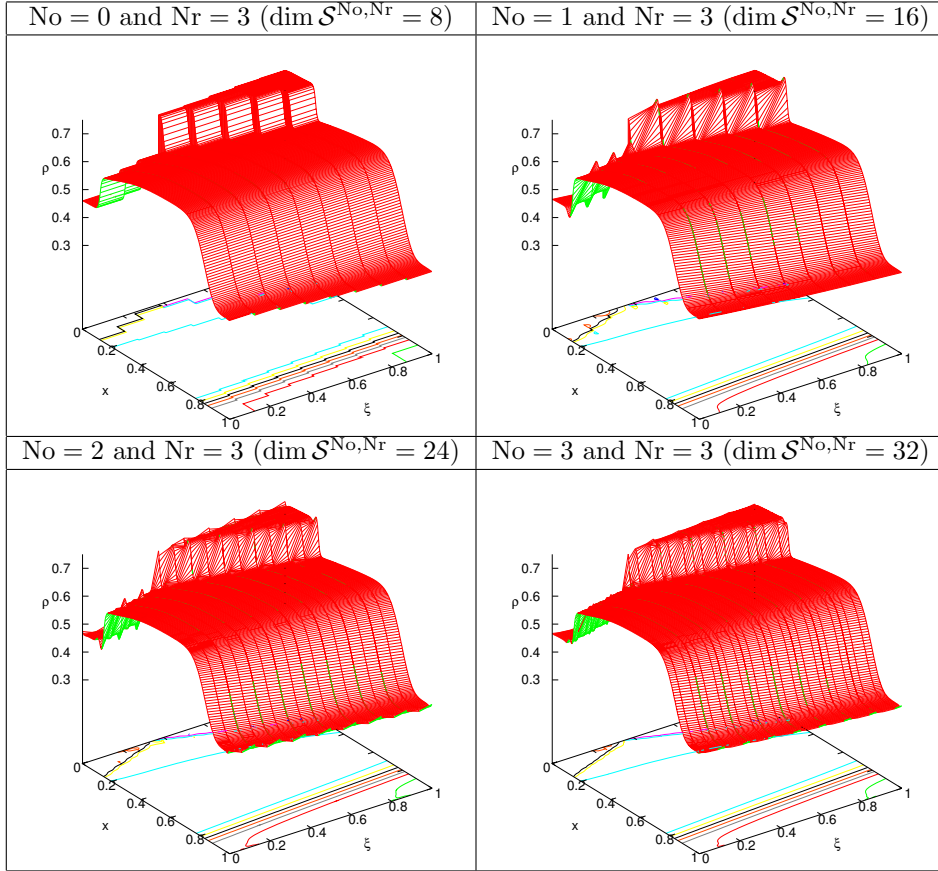


Figure 19: Convergence of the stochastic density  $\rho(x, t, \xi)$  with  $N_r = 3$ ,  $t = 6.5$ .

low value; the present simulations actually used a definition  $d + 1 = \min(9, 3(N_o + 1))$ , since allowing for higher degree  $d$  for  $N_o \geq 3$  was found to have no significant effect on the solution. All in all, the complexity of the nonlinearity resolution appears to be the most limiting factor of the present method, and this effect is expected to be worse for problems with higher stochastic dimension  $N$  (see for instance [16]): this trend pleads for using stochastic approximation spaces with low-degree polynomials for non-smooth stochastic problems.

	$N_r = 2$		$N_r = 3$		$N_r = 4$	
	$T_{CPU}$	$\dim \mathcal{S}^{N_r, N_o}$	$T_{CPU}$	$\dim \mathcal{S}^{N_r, N_o}$	$T_{CPU}$	$\dim \mathcal{S}^{N_r, N_o}$
$N_o = 0$	4.0	(4)	8.1	(8)	16.1	(16)
$N_o = 1$	6.9	(8)	13.9	(16)	27.8	(32)
$N_o = 2$	11.8	(12)	23.2	(24)	46.5	(48)
$N_o = 3$	17.1	(16)	34.1	(32)	68.1	(64)
$N_o = 4$	24.8	(20)	49.3	(40)	98.0	(80)

Table 1: Normalized computational times  $T_{CPU}$  for different stochastic discretization parameters  $N_r$  and  $N_o$ .

## 6. Conclusion

In this paper we have investigated theoretically and numerically an intrusive projection method for stochastic hyperbolic systems of conservation laws, exhibiting discontinuities in both spatial and stochastic

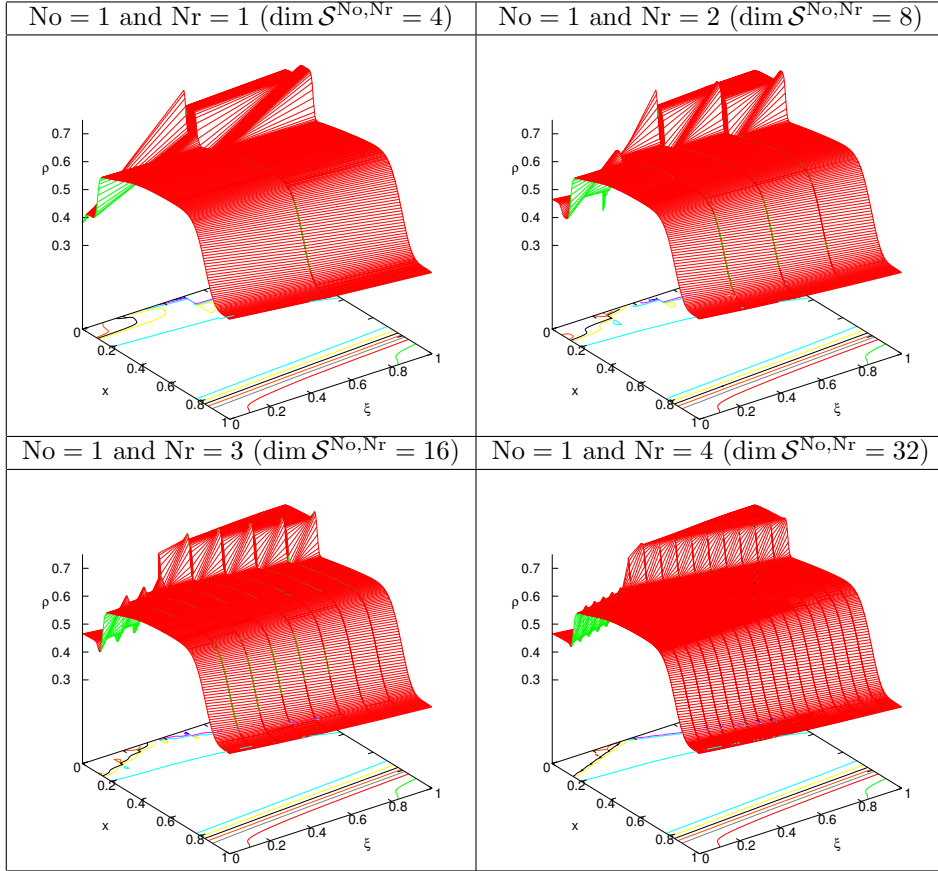


Figure 20: Convergence of the stochastic density  $\rho(x, t, \xi)$  with Nr. No = 1,  $t = 6.5$ .

domains. The method is based on the Galerkin projection of the original stochastic problem on a space of piecewise polynomials and uses a Roe-type solver with upwind matrices that are efficiently computed by an original and fast method. Numerical tests on the stochastic Burgers and Euler equations in one spatial dimension and, respectively, in two and one stochastic dimensions indicate that the method is accurate and robust while maintaining moderate computational costs. Despite these improvements, Galerkin projection methods remain expensive, especially to explore problems with higher stochastic dimensions. To reach their full potential, Galerkin projection methods need further developments, in particular stochastic adaptivity. This is the focus of ongoing efforts.

## References

- [1] R. Abgrall. A simple, flexible and generic deterministic approach to uncertainty quantifications in non linear problems: application to fluid flow problems. *Rapport de recherche INRIA*, 00325315, 2008.
- [2] Ivo Babuška, Fabio Nobile, and Raúl Tempone. A stochastic collocation method for elliptic partial differential equations with random input data. *SIAM J. Numer. Anal.*, 45(3):1005–1034, 2007.
- [3] Alexandre Joel Chorin. Gaussian fields and random flow. *J. Fluid Mech.*, 63:21–32, 1974.
- [4] Manas K. Deb, Ivo M. Babuška, and J. Tinsley Oden. Solution of stochastic partial differential equations using Galerkin finite element techniques. *Comput. Methods Appl. Mech. Engrg.*, 190(48):6359–6372, 2001.
- [5] B. Debusschere, H.N. Najm, A. Matta, O.M. Knio, R.G. Ghanem, and O.P. Le Maître. Protein Labeling Reactions in Electrochemical Microchannel Flow: Numerical Prediction and Uncertainty Propagation. *Physics of Fluids*, 15(8):2238–2250, 2003.
- [6] Bert J. Debusschere, Habib N. Najm, Philippe P. Pébay, Omar M. Knio, Roger G. Ghanem, and Olivier P. Le Maître. Numerical challenges in the use of polynomial chaos representations for stochastic processes. *SIAM J. Sci. Comput.*, 26(2):698–719, 2004.

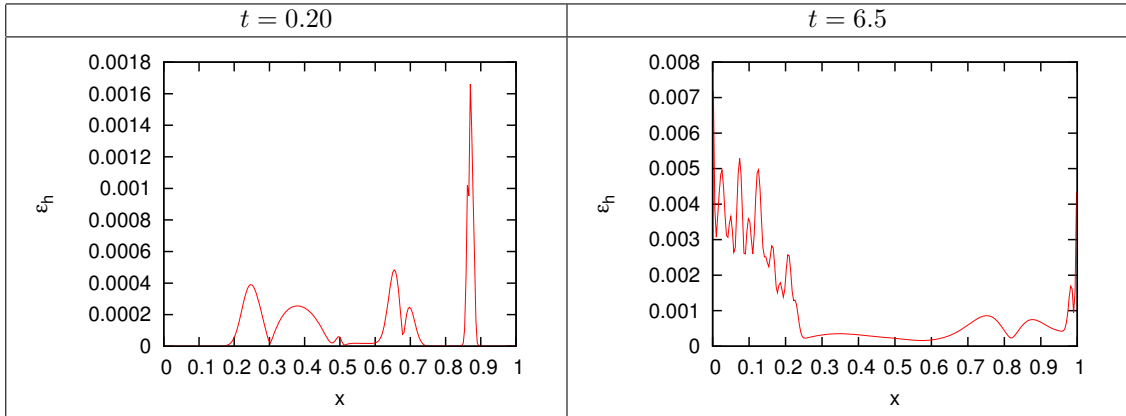


Figure 21: Stochastic error  $\epsilon_h(x, t)$  for early (left) and longer (right) times. Computations with  $N_c = 250$ ,  $N_r = 3$ , and  $N_o = 2$ .

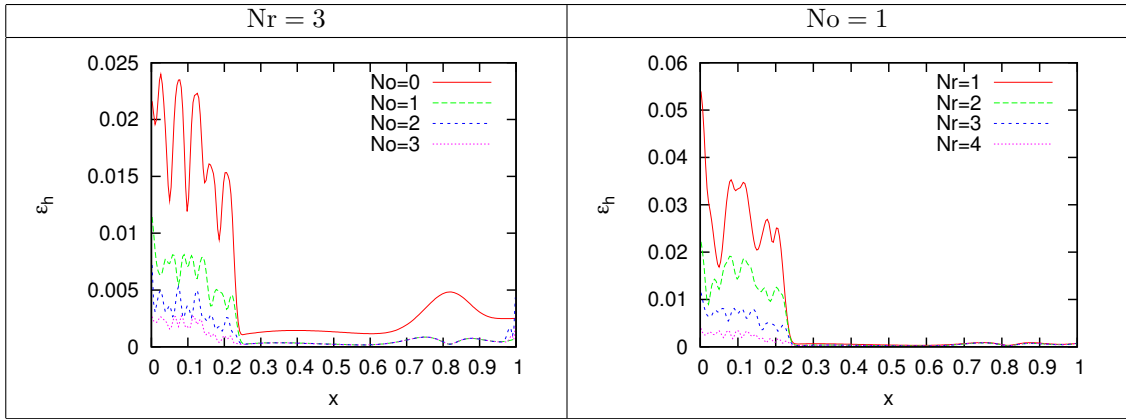


Figure 22: Stochastic error  $\epsilon_h(x, t = 6.5)$  for various  $N_o$  and  $N_r$ . Computations with  $N_c = 250$ .

- [7] Pierre Degond, Pierre-François Peyrard, Giovanni Russo, and Philippe Villedieu. Polynomial upwind schemes for hyperbolic systems. *C. R. Acad. Sci. Paris Sér. I Math.*, 328(6):479–483, 1999.
- [8] Jasmine Foo, Xiaoliang Wan, and George Em Karniadakis. The multi-element probabilistic collocation method (ME-PCM): error analysis and applications. *J. Comput. Phys.*, 227(22):9572–9595, 2008.
- [9] Baskar Ganapathysubramanian and Nicholas Zabaras. Sparse grid collocation schemes for stochastic natural convection problems. *J. Comput. Phys.*, 225(1):652–685, 2007.
- [10] L. Ge, K.F. Cheung, and M.H. Kobayashi. Stochastic Solution for Uncertainty Propagation in Nonlinear Shallow-Water Equations. *Journal of Hydraulic Engineering*, Decembre 2008:1732–1743, 2008.
- [11] Roger G. Ghanem and Pol D. Spanos. *Stochastic finite elements: a spectral approach*. Dover, 2003.
- [12] Edwige Godlewski and Pierre-Arnaud Raviart. *Numerical approximation of hyperbolic systems of conservation laws*, volume 118 of *Applied Mathematical Sciences*. Springer-Verlag, New York, 1996.
- [13] David Gottlieb and Dongbin Xiu. Galerkin method for wave equations with uncertain coefficients. *Commun. Comput. Phys.*, 3(2):505–518, 2008.
- [14] A. Keese and H.G. Matthies. Numerical methods and Smolyak quadrature for nonlinear stochastic partial differential equations. Technical report, Institute of Scientific Computing TU Braunschweig Brunswick, 2003.
- [15] O. M. Knio and O. P. Le Maître. Uncertainty propagation in CFD using polynomial chaos decomposition. *Fluid Dynam. Res.*, 38(9):616–640, 2006.
- [16] O. Le Maître. A Newton method for the resolution of steady stochastic Navier-Stokes equations. *Comput Fluids*, 2009. in press, doi:10.1016/j.compfluid.2009.01.001.
- [17] O. P. Le Maître, O. M. Knio, H. N. Najm, and R. G. Ghanem. Uncertainty propagation using Wiener-Haar expansions. *J. Comput. Phys.*, 197(1):28–57, 2004.
- [18] O. P. Le Maître, H. N. Najm, R. G. Ghanem, and O. M. Knio. Multi-resolution analysis of Wiener-type uncertainty propagation schemes. *J. Comput. Phys.*, 197(2):502–531, 2004.
- [19] O. P. Le Maître, H. N. Najm, P. P. Pébay, R. G. Ghanem, and O. M. Knio. Multi-resolution-analysis scheme for uncertainty quantification in chemical systems. *SIAM J. Sci. Comput.*, 29(2):864–889, 2007.

- [20] Olivier Le Maître, M. T. Reagan, B. Debusschere, H. N. Najm, R. G. Ghanem, and O. M. Knio. Natural convection in a closed cavity under stochastic non-Boussinesq conditions. *SIAM J. Sci. Comput.*, 26(2):375–394, 2004.
- [21] Olivier P. Le Maître, Omar M. Knio, Habib N. Najm, and Roger G. Ghanem. A stochastic projection method for fluid flow. I. Basic formulation. *J. Comput. Phys.*, 173(2):481–511, 2001.
- [22] Olivier P. Le Maître, Matthew T. Reagan, Habib N. Najm, Roger G. Ghanem, and Omar M. Knio. A stochastic projection method for fluid flow. II. Random process. *J. Comput. Phys.*, 181(1):9–44, 2002.
- [23] G. Lin, C.-H. Su, and G. E. Karniadakis. Predicting shock dynamics in the presence of uncertainties. *J. Comput. Phys.*, 217(1):260–276, 2006.
- [24] G. Lin, C.-H. Su, and G. E. Karniadakis. Stochastic modeling of random roughness in shock scattering problems: theory and simulations. *Comput. Methods Appl. Mech. Engrg.*, 197(43-44):3420–3434, 2008.
- [25] L. Mathelin and M. Hussaini. A stochastic collocation algorithm for uncertainty analysis. Technical Report NASA/CR-2003-212153, NASA Langley Research Center, 2003.
- [26] Lionel Mathelin, M. Yousuff Hussaini, and Thomas A. Zang. Stochastic approaches to uncertainty quantification in CFD simulations. *Numer. Algorithms*, 38(1-3):209–236, 2005.
- [27] H.N. Najm. Uncertainty quantification and polynomial chaos techniques in computational fluid dynamics. *Ann. Rev. Fluid Mech.*, 41:35–52, 2009.
- [28] H.N. Najm, B.J. Debusschere, Y.M. Marzouk, S. Widmer, and O.P. Le Maître. Uncertainty quantification in chemical systems. *Int. J. Numer. Meth. Engrg.*, 2009. in press, doi: 10.1002/nme.2551.
- [29] Michaël Ndjinga. Computing the matrix sign and absolute value functions. *C. R. Math. Acad. Sci. Paris*, 346(1-2):119–124, 2008.
- [30] F. Nobile, R. Tempone, and C. G. Webster. A sparse grid stochastic collocation method for partial differential equations with random input data. *SIAM J. Numer. Anal.*, 46(5):2309–2345, 2008.
- [31] G. Poette, B. Després, and D. Lucor. Uncertainty quantification for systems of conservation laws. *J. Comput. Phys.*, 228(7):2443–2467, 2009.
- [32] M.T. Reagan, H.N. Najm, R.G. Ghanem, and O.M. Knio. Uncertainty quantification in reacting flow simulations through non-intrusive spectral projection. *Combustion and Flame*, 132:545–555, 2003.
- [33] J. Stoer and R. Bulirsch. *Introduction to numerical analysis*, volume 12 of *Texts in Applied Mathematics*. Springer-Verlag, New York, third edition, 2002. Translated from the German by R. Bartels, W. Gautschi and C. Witzgall.
- [34] Eleuterio F. Toro. *Riemann solvers and numerical methods for fluid dynamics: A practical introduction*. Springer-Verlag, Berlin, second edition, 1999.
- [35] Xiaoliang Wan and George Em Karniadakis. Multi-element generalized polynomial chaos for arbitrary probability measures. *SIAM J. Sci. Comput.*, 28(3):901–928 (electronic), 2006.
- [36] Norbert Wiener. The Homogeneous Chaos. *Amer. J. Math.*, 60(4):897–936, 1938.
- [37] Dongbin Xiu and Jan S. Hesthaven. High-order collocation methods for differential equations with random inputs. *SIAM J. Sci. Comput.*, 27(3):1118–1139, 2005.
- [38] Dongbin Xiu and George Em Karniadakis. The Wiener-Askey polynomial chaos for stochastic differential equations. *SIAM J. Sci. Comput.*, 24(2):619–644 (electronic), 2002.