

# Hyperbolic Stochastic Galerkin Formulation for the $p$ -System

Stephan Gerster, Michael Herty and Aleksey Sikstel

Institut für Geometrie und Praktische Mathematik  
Templergraben 55, 52062 Aachen, Germany

---

IGPM, RWTH Aachen University, Templergraben 55, D-52062 Aachen

This work is supported by DFG HE5386/13–15, HE6954/4, BMBF 05M18PAA.

# Hyperbolic Stochastic Galerkin Formulation for the $p$ -System

Stephan Gerster\*, Michael Herty, Aleksey Sikstel

*RWTH Aachen University, IGPM, Templergraben 55, 52062 Aachen, Germany*

---

## Abstract

We analyze properties of stochastic hyperbolic systems using a Galerkin formulation, which reformulates the stochastic system as a deterministic one that describes the evolution of polynomial chaos modes. We investigate conditions such that the resulting systems are hyperbolic. We state the eigendecompositions in closed form. A Roe flux is presented and theoretical results are illustrated numerically.

*Keywords:* Hyperbolic Partial Differential Equations, Uncertainty Quantification, Stochastic Galerkin Method, Euler Equations, Roe Variable Transform

---

## 1. Introduction

Recently, the representation of stochastic processes by orthogonal polynomials has gained interest in the mathematical and engineering community [1]. This idea has been employed in uncertainty quantification and inverse problems. The first work in this direction by Wiener [2] used Hermite polynomials to represent Brownian motion. This approach has been extended by Cameron, Martin [3], Ghanem, Spanos [4] and Xiu, Karniadakis [5] for non-Gaussian processes using polynomials from the Askey scheme. It is known today under the name generalized polynomial chaos (gPC). If the solution depends sufficiently regularly on the stochastic input, spectral convergence is observed [5].

Non-intrusive methods compute statistics directly using numerical quadrature or Monte-Carlo methods. However, in the case of an intrusive approach gPC expansions of the stochastic input are substituted into the governing equations. Then, they are projected by a Galerkin method to obtain deterministic evolution equations for the gPC coefficients. Applications of this procedure have been proven successful for diffusion [6, 7] and kinetic equations [8, 9, 10].

So far, results for general hyperbolic systems are not available [11, Sec. 10.2]. A problem is posed by the fact that the deterministic Jacobian of the projected system differs from the random Jacobian of the original system and hence hyperbolicity is not guaranteed. The loss of hyperbolicity prevents existence, uniqueness results and consequently the use of robust numerical schemes. Applications to hyperbolic conservation laws are in general limited to linear and scalar hyperbolic equations. Then, the resulting gPC systems remain hyperbolic, i.e. the Jacobian in the quasilinear form is diagonalizable with real eigenvalues [12, 13, 14, 15]. Hyperbolicity of nonlinear systems is shown

---

\*Corresponding author.

E-mail addresses: [gerster@igpm.rwth-aachen.de](mailto:gerster@igpm.rwth-aachen.de), [stephan.gerster@gmail.com](mailto:stephan.gerster@gmail.com)

in the case of a symmetric Jacobian and when eigenvectors are deterministic [16]. For the classical fluid-dynamic equations, however, eigenvectors are uncertain as well.

Also for quasilinear hyperbolic systems the resulting gPC system may not be hyperbolic [17, 18]. However, in [18] an approach to regain hyperbolicity has been proposed. This approach is limited to quasilinear forms and solvers which are developed for conservative formulations are not directly applicable. Like in the deterministic case, a non-conservative formulation may lead to wrong shock speeds.

Therefore, operator-splitting methods have been developed for Euler [19] and shallow water equations [20]. The idea is to split the underlying system into subsystems such that for each of them the gPC method does not lead to complex eigenvalues. The splitted system may differ from the original one. Numerical experiments suggest, however, that at least in the deterministic case both systems are similar. Also for this approach, standard solvers are not directly applicable. Still these methods, as well as semi-intrusive methods, which approximate numerically the Galerkin projection, see e.g. [21, 22], may yield complex eigenvalues unless positive densities are guaranteed [22].

For some fluid-dynamic equations, like Euler equations, we may first transform the partial differential equation (PDE) into non-conserved variables and then apply the intrusive method. Resulting systems for entropy variables are hyperbolic [23, 17]. To obtain a conservative formulation, an optimization problem needs be solved at each spatial point and in each time step, which makes the method numerically expensive. A similar method, but with a different variable transform is proposed in [24]. There, Roe variables are used to guarantee real wave speeds. Computational experiments suggest appropriate computational cost, see [24, Table 1]. The Roe formulation in [24] is so far restricted to particular expansions, including the Wiener-Haar basis and linear multiwavelets.

Besides hyperbolicity, another open problem [11, Sec. 10.2] is the representation of strictly positive quantities, e.g. density of a gas or height of the water. This problem is frequently related to truncation errors of the gPC expansion, when the underlying quantity is assumed to be almost surely (a.s.) positive as long as the order of truncation is large enough [24, 25].

The main result of this paper is Theorem 3.2. Therein, we present a gPC formulation for the  $p$ -system, which is a  $2 \times 2$  hyperbolic system that includes shallow water and isothermal Euler equations. The introduced systems are based on the Roe variable transform in [26, 24]. The assumption of strictly positive quantities in the truncated gPC expansion is weakened to a positive definiteness condition (A2). Then, hyperbolicity for some expansions is shown in Corollary 1. In particular for isothermal Euler equations, Corollary 2 guarantees hyperbolicity even for arbitrary gPC expansions.

Furthermore, we establish novel results on the mathematical structure and properties of the gPC system, see [11, Sec. 10.2]. We state eigendecompositions of the gPC systems in closed form. In particular for isothermal Euler equations, we deduce a Roe matrix for arbitrary gPC bases. This shows the relation of the Roe variable transform [24] to the original deterministic Roe matrix for Euler equations [26]. This also yields a numerical flux. Finally, we illustrate the hyperbolic character of the introduced systems numerically. We discuss an efficient implementation of the numerical flux function and the eigenvalue decomposition. Furthermore, an eigenvalue estimate is presented.

The reader finds a summary of the discussed hyperbolic systems with the corresponding assumptions and a list of all symbols in the appendix.

## 2. Hyperbolic Conservation Laws

In this section we introduce different formulations of hyperbolic conservation laws, which result from the Roe variable transform in [26, 24].

### 2.1. Conservative Formulation

We consider a system of hyperbolic differential equations

$$\frac{\partial}{\partial t} y(t, x) + \frac{\partial}{\partial x} f(y(t, x)) = 0 \quad (1)$$

with time and space variables  $(t, x) \in [0, T] \times [0, x_{\text{end}}]$  for a fixed  $T < \infty$ . For isentropic Euler equations the conserved quantities  $y := (\rho, q)^T$  are density and momentum, for shallow water equations water depth and momentum. In both cases the conservation law (1) is a  $2 \times 2$  system with strictly hyperbolic flux function  $f \in C^2(\mathbb{R}^2; \mathbb{R}^2)$ . We consider the  $p$ -system

$$\frac{\partial}{\partial t} \begin{pmatrix} \rho(t, x) \\ q(t, x) \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} q(t, x) \\ \frac{q^2(t, x)}{\rho(t, x)} + p(\rho(t, x)) \end{pmatrix} = 0, \quad (2)$$

where  $p(\rho)$  describes a pressure law, which satisfies  $p(\rho) > 0$ ,  $p'(\rho) > 0$  and  $p''(\rho) \geq 0$ . The pressure law reads in the isothermal case  $p(\rho) = a^2 \rho$  with sound speed  $a > 0$ , in the isentropic case  $p(\rho) \sim \rho^\kappa$  with  $\kappa > 0$  and for shallow water equations  $p(\rho) = \frac{g}{2} \rho^2$  with gravitational constant  $g > 0$ . To simplify notation we omit time and space arguments. For smooth solutions the  $p$ -system (2) is equivalent to solving the quasilinear form

$$y_t + D_y f(y) y_x = 0 \quad \text{with Jacobian} \quad D_y f(y) = \begin{pmatrix} 0 & 1 \\ p'(\rho) - u^2(y) & 2u(y) \end{pmatrix}, \quad u(y) := \frac{q}{\rho}.$$

The eigenvalues read  $\lambda^\pm(y) = u(y) \pm \sqrt{p'(\rho)}$  and satisfy  $\lambda^-(y) < 0 < \lambda^+(y)$  for subsonic flows with velocity  $u(y) < |\sqrt{p'(\rho)}|$ . The eigendecomposition of the Jacobian is

$$D_y f(y) = T(y) \Lambda(y) T^{-1}(y), \quad T(y) := \begin{pmatrix} 1 & 1 \\ \lambda^+(y) & \lambda^-(y) \end{pmatrix}, \quad \Lambda(y) := \text{diag}\{\lambda^+(y), \lambda^-(y)\}. \quad (3)$$

Note that eigenvectors are expressed in terms of eigenvalues. For isothermal Euler equations with pressure law  $p(\rho) := a^2 \rho$  both eigenvalues depend only on velocity  $u(y)$ .

### 2.2. Roe Formulation

Following [26, 27, 24] we introduce **Roe variables**.

**Definition 2.1** (Roe Variables). Let the Roe variables  $\omega(t, x)$  and the mapping into conserved variables  $y(t, x)$  be defined by

$$\omega := \begin{pmatrix} \alpha \\ \beta \end{pmatrix} := \begin{pmatrix} \sqrt{\rho} \\ \frac{q}{\sqrt{\rho}} \end{pmatrix} \quad \text{and} \quad \mathcal{Y} : \mathbb{R}^+ \times \mathbb{R} \rightarrow \mathbb{R}^+ \times \mathbb{R}, \quad \omega \mapsto \begin{pmatrix} \alpha^2 \\ \alpha\beta \end{pmatrix} = y. \quad (4)$$

The pressure law in Roe variables reads  $\pi(\alpha) := p(\alpha^2)$  and the velocity is  $\nu(\omega) := \beta/\alpha$ . The flux function for conserved variables depending on Roe variables is denoted by

$$F : \mathbb{R}^+ \times \mathbb{R} \rightarrow \mathbb{R}^2, \omega \mapsto \begin{pmatrix} \alpha\beta \\ \beta^2 + \pi(\alpha) \end{pmatrix} = f(\mathcal{Y}(\omega)). \quad (5)$$

The following Lemma states the relationship between conserved and Roe variables.

**Lemma 2.2.** *For smooth solutions with  $\rho > 0$  and  $\alpha > 0$  the following systems are equivalent:*

$$y_t + f(y)_x = 0, \quad (\mathcal{C})$$

$$\mathcal{Y}(\omega)_t + F(\omega)_x = 0, \quad (\mathcal{R})$$

$$y_t + D_y f(y) y_x = 0 \quad \text{with} \quad D_y f(y) := D_\omega F(\omega) [D_\omega \mathcal{Y}]^{-1}(\omega), \quad (6)$$

$$\omega_t + D_\omega F_{\text{Roe}}(\omega) \omega_x = 0 \quad \text{with} \quad D_\omega F_{\text{Roe}}(\omega) := [D_\omega \mathcal{Y}]^{-1}(\omega) D_\omega F(\omega) \quad (7)$$

*Proof.* For  $\rho > 0$  the mapping  $\mathcal{Y}(\cdot)$  is bijective on  $\mathbb{R}^+ \times \mathbb{R}$ , so equations  $(\mathcal{C})$  and  $(\mathcal{R})$  are equal by definition. Since the Jacobian  $D_\omega \mathcal{Y}(\omega)$  is then invertible, we obtain

$$D_y f(y) = D_y F(\mathcal{Y}^{-1}(y)) = D_\omega F(\omega) D_y \mathcal{Y}^{-1}(y) = D_\omega F(\omega) [D_\omega \mathcal{Y}]^{-1}(\omega).$$

We obtain for the quasilinear form (7)

$$\begin{aligned} 0 &= \mathcal{Y}(\omega)_t + F(\omega)_x = D_\omega \mathcal{Y}(\omega) \omega_t + D_\omega F(\omega) \omega_x \\ \Leftrightarrow \quad 0 &= \omega_t + D_\omega F_{\text{Roe}}(\omega) \omega_x. \end{aligned}$$

□

Although these systems are equivalent in the deterministic case, we will see that they result in different stochastic Galerkin formulations. We obtain the Jacobian  $D_\omega F_{\text{Roe}}(\omega)$  by calculating

$$D_\omega \mathcal{Y}(\omega) = \begin{pmatrix} 2\alpha & 0 \\ \beta & \alpha \end{pmatrix}, \quad D_\omega F(\omega) = \begin{pmatrix} \beta & \alpha \\ \pi'(\alpha) & 2\beta \end{pmatrix} \Rightarrow D_\omega F_{\text{Roe}}(\omega) = \begin{pmatrix} \frac{\nu(\omega)}{2} & \frac{1}{2} \\ -\frac{\nu^2(\omega)}{2} + \frac{\pi'(\alpha)}{\alpha} & \frac{3}{2}\nu(\omega) \end{pmatrix}.$$

However, the flux function  $F_{\text{Roe}}(\omega)$  itself is not specifiable in closed form, which would be necessary for the analysis of shocks. In any case, the Roe formulation does not preserve the physical correct shock speed in terms of conserved variables. So we conclude that the Roe formulation (7) is not relevant for further analysis.

### 3. Random Hyperbolic Conservation Laws

The hyperbolic problem (1) is extended in the way that initial and boundary conditions, but *not the pressure law* are allowed to depend on a random parameter  $\xi$ , i.e. a measurable mapping  $\xi : \Omega \rightarrow \mathbb{R}$  on a probability space  $(\Omega, \mathcal{F}(\Omega), \mathbb{P})$ . Similarly to [11], we call  $\xi$  **germ**. For our purposes it does not matter if the germ is one- or multidimensional. Therefore, we consider for simplicity a one-dimensional germ. An example, where hyperbolicity depends on the dimension of stochastic input, is found in [14, Sec. 3.3].

We briefly recall basic results from [11, 28, 29, 30, 31, 32]. Interpreting integrals and inner products component wise, the space of second-order random variables defined on the probability space  $(\Omega, \mathcal{F}(\Omega), \mathbb{P})$  equipped with an inner product is defined by

$$\mathbb{L}^2(\Omega, \mathbb{P}) := \left\{ \xi \mid \xi : \Omega \rightarrow \mathbb{R} \text{ measurable, } \|\xi\|_{\mathbb{P}} < \infty \right\} \quad \text{with} \quad \langle \xi_1, \xi_2 \rangle_{\mathbb{P}} := \int \xi_1 \xi_2 \, d\mathbb{P}.$$

The expected value is  $\mathbb{E}[\xi_1 \xi_2] := \langle \xi_1, \xi_2 \rangle_{\mathbb{P}}$  and the covariance is  $\text{Cov}[\xi_1, \xi_2] := \mathbb{E}[\xi_1 \xi_2] - \mathbb{E}[\xi_1] \mathbb{E}[\xi_2]$ . A **generalized polynomial chaos (gPC)** is a set of orthogonal subspaces  $\hat{\mathcal{S}}_k \subset \mathbb{L}^2(\Omega, \mathbb{P})$  with

$$\mathcal{S}_K := \bigoplus_{k=0}^K \hat{\mathcal{S}}_k \rightarrow \mathbb{L}^2(\Omega, \mathbb{P}) \quad \text{for } K \rightarrow \infty.$$

We refer to an orthogonal basis of  $\mathcal{S}_K$  as a **gPC basis**  $\{\phi_k(\xi)\}_{k=0}^K$  with germ  $\xi$ . Its distribution is given by the probability measure  $\mathbb{P}$ , we write  $\xi \sim \mathbb{P}$  for brevity. Common choices are Legendre and Hermite polynomials with uniformly and normally distributed germs. A stochastic process  $y(t, x; \xi)$ , which is for each fixed  $(t, x)$  square-integrable, admits the truncated series expansion

$$\mathcal{G}_K[y](t, x; \xi) := \sum_{k=0}^K \hat{y}_k(t, x) \phi_k(\xi), \quad \hat{y}_k(t, x) := \frac{\langle y(t, x; \cdot), \phi_k(\cdot) \rangle_{\mathbb{P}}}{\|\phi_k\|_{\mathbb{P}}^2}, \quad (8)$$

where  $\mathcal{G}_K$  denotes the projection operator of the stochastic process  $y(t, x; \xi)$  onto the gPC basis of degree  $K \in \mathbb{N}_0$ . It converges in the sense  $\|\mathcal{G}_K[y](t, x; \cdot) - y(t, x; \cdot)\|_{\mathbb{P}} \rightarrow 0$  for  $K \rightarrow \infty$  [3, 5]. We use normalized basis functions such that  $\|\phi_k\|_{\mathbb{P}} = 1$ . In the case  $\|\tilde{\phi}_k\|_{\mathbb{P}} \neq 1$ , we rescale

$$\phi_k(\xi) := \frac{\tilde{\phi}_k(\xi)}{\|\tilde{\phi}_k\|_{\mathbb{P}}} \quad \text{and} \quad \mathcal{G}_K[y](t, x; \xi) := \sum_{k=0}^K \hat{y}_k(t, x) \phi_k(\xi) \quad \text{with} \quad \hat{y}_k(t, x) := \frac{\langle y(t, x; \cdot), \tilde{\phi}_k(\cdot) \rangle_{\mathbb{P}}}{\|\tilde{\phi}_k\|_{\mathbb{P}}}.$$

The reason for this normalization is that we have observed numerical instabilities for Hermite polynomials with  $\|\phi_k\|_{\mathbb{P}}^2 = k!$ . Then, the gPC modes  $\hat{y}_k$  give the expected value and the variance

$$\mathbb{E}[\mathcal{G}_K[y]](t, x) = \hat{y}_0(t, x) \quad \text{and} \quad \text{Var}[\mathcal{G}_K[y]](t, x) = \sum_{k=1}^K \hat{y}_k^2(t, x).$$

A straightforward analogue for the product of  $\mathcal{G}_K[y](t, x; \xi)$  and  $\mathcal{G}_K[z](t, x; \xi)$  would be

$$\left( \mathcal{G}_K[y] \mathcal{G}_K[z] \right)(t, x; \xi) = \sum_{i,j=0}^K \hat{y}_i(t, x) \hat{z}_j(t, x) \phi_i(\xi) \phi_j(\xi). \quad (9)$$

However, this leads to basis functions up to order  $2K$ , so equation (9) requires an additional projection to recover the degree  $K$ . Therefore, we define the **pseudo-spectral product**

$$\hat{\mathcal{G}}_K[y, z](t, x; \xi) := \sum_{k=0}^K (\hat{y} * \hat{z})_k(t, x) \phi_k(\xi), \quad (\hat{y} * \hat{z})_k(t, x) := \sum_{i,j=0}^K \hat{y}_i(t, x) \hat{z}_j(t, x) \langle \phi_i \phi_j, \phi_k \rangle_{\mathbb{P}}. \quad (10)$$

Similar to [24, 29], we introduce the symmetric matrix

$$\mathcal{P}(\hat{\alpha}) := \sum_{\ell=0}^K \hat{\alpha}_\ell \mathcal{M}_\ell \quad \text{with} \quad \mathcal{M}_\ell := \left( \langle \phi_\ell, \phi_i \phi_j \rangle_{\mathbb{P}} \right)_{i,j=0,\dots,K} \quad (11)$$

such that  $\mathcal{P}(\hat{\alpha})\hat{\beta} = \hat{\alpha} * \hat{\beta}$ ,  $D_{\hat{\alpha}}(\hat{\alpha} * \hat{\beta}) = \mathcal{P}(\hat{\beta})$ .

The pseudo-spectral product is exact for  $(yz) \in \mathcal{S}_K$ , but in general we have  $\hat{\mathcal{G}}_K[y, z] \neq \mathcal{G}_K[yz]$ . It introduces a truncation error by disregarding the components of  $(yz)$  which are orthogonal to  $\mathcal{S}_K$  and it is not associative [28, 31]. The operator  $(*)$  is called **Galerkin product**.

### 3.1. Non-Intrusive Formulation

The strong formulations for the random systems in conserved ( $\mathcal{C}$ ) and Roe variables ( $\mathcal{R}$ ) are

$$\frac{\partial}{\partial t} \begin{pmatrix} \rho(t, x; \xi) \\ q(t, x; \xi) \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} q(t, x; \xi) \\ \frac{q^2(t, x; \xi)}{\rho(t, x; \xi)} + p(\rho(t, x; \xi)) \end{pmatrix} = 0 \quad \mathbb{P}\text{-a.s.} \quad (\mathcal{C}(\xi))$$

$$\frac{\partial}{\partial t} \begin{pmatrix} \alpha^2(t, x; \xi) \\ (\alpha\beta)(t, x; \xi) \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} (\alpha\beta)(t, x; \xi) \\ \beta^2(t, x; \xi) + \pi(\alpha(t, x; \xi)) \end{pmatrix} = 0 \quad \mathbb{P}\text{-a.s.} \quad (\mathcal{R}(\xi))$$

To ensure hyperbolicity, the Jacobian  $D_y f(y(t, x; \xi))$  must be diagonalizable with real eigenvalues at least  $\mathbb{P}$ -a.s. Therefore, we require the assumption

$$\rho(t, x; \xi) > 0 \quad \mathbb{P}\text{-a.s.} \quad \text{and} \quad \alpha(t, x; \xi) > 0 \quad \mathbb{P}\text{-a.s.} \quad (\text{A1})$$

### 3.2. Intrusive Formulations

We substitute the truncated expansions into the systems ( $\mathcal{C}(\xi)$ ) and ( $\mathcal{R}(\xi)$ ) to obtain

$$\frac{\partial}{\partial t} \begin{pmatrix} \mathcal{G}_K[\rho](t, x; \xi) \\ \mathcal{G}_K[q](t, x; \xi) \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} \mathcal{G}_K[q](t, x; \xi) \\ \frac{\mathcal{G}_K^2[q](t, x; \xi)}{\mathcal{G}_K[\rho](t, x; \xi)} + p(\mathcal{G}_K[\rho](t, x; \xi)) \end{pmatrix} = 0, \quad (\mathcal{C}_K(\xi))$$

$$\frac{\partial}{\partial t} \begin{pmatrix} \hat{\mathcal{G}}_K[\alpha, \alpha](t, x; \xi) \\ \hat{\mathcal{G}}_K[\alpha, \beta](t, x; \xi) \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} \hat{\mathcal{G}}_K[\alpha, \beta](t, x; \xi) \\ \hat{\mathcal{G}}_K[\beta, \beta](t, x; \xi) + \pi(\mathcal{G}_K[\alpha](t, x; \xi)) \end{pmatrix} = 0. \quad (\mathcal{R}_K(\xi))$$

The truncated systems ( $\mathcal{C}_K(\xi)$ ) and ( $\mathcal{R}_K(\xi)$ ) should be solved for the gPC modes in  $\mathbb{L}^2(\Omega, \mathbb{P})$ -sense. The solution, however, does in general not exist due to the possible loss of hyperbolicity. Another open problem is the convergence of the solution for  $K \rightarrow \infty$  in the general case [17, 10]. A straight forward approach is to describe the evolution of the gPC modes for the random system ( $\mathcal{C}_K(\xi)$ ). A conservative formulation for the isothermal case with pressure law  $\hat{p}(\hat{\rho}) = a^2 \hat{\rho}$ , presented in [33], is

$$\frac{\partial}{\partial t} \begin{pmatrix} \hat{\rho} \\ \hat{q} \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} \hat{q} \\ \mathcal{P}(\hat{q})\mathcal{P}^{-1}(\hat{\rho})\hat{q} + a^2 \hat{\rho} \end{pmatrix} = 0 \quad \text{with Jacobian} \quad (\hat{\mathcal{C}}_K)$$

$$D_{\hat{y}} \hat{f}(\hat{y}) = \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ -\mathcal{P}(\hat{q})\mathcal{P}^{-1}(\hat{\rho})\mathcal{P}(\mathcal{P}^{-1}(\hat{\rho})\hat{q}) + a^2 \mathbb{1} & \mathcal{P}(\hat{q})\mathcal{P}^{-1}(\hat{\rho}) + \mathcal{P}(\mathcal{P}^{-1}(\hat{\rho})\hat{q}) \end{pmatrix}$$

for  $\mathbb{O} := \text{diag}\{0, \dots, 0\} \in \mathbb{R}^{(K+1) \times (K+1)}$  and  $\mathbb{1} := \text{diag}\{1, \dots, 1\} \in \mathbb{R}^{(K+1) \times (K+1)}$ . However, it is shown in [33] that the formulation ( $\hat{\mathcal{C}}_K$ ) is in general *not hyperbolic*. We use this formulation as

a comparison only. Instead, we consider the random formulation  $(\mathcal{R}_K(\xi))$ . The idea to project a system of the form  $(\mathcal{R}_K(\xi))$  instead of  $(\mathcal{C}_K(\xi))$  is borrowed from [24]. Our main motivation is to analyze hyperbolicity and in Theorem 3.2 we derive a conservative formulation for the gPC modes in system  $(\mathcal{R}_K(\xi))$ . In Corollary 1 and Corollary 2 we present conditions that guarantee hyperbolicity of the derived systems for fixed time and space variables  $(t, x)$ . Therefore, we drop the dependency on the variables  $(t, x)$  in this section and we extend Definition 2.1 for both conserved and Roe variables.

**Definition 3.1** (Intrusive Variables). The conserved variables on the admissible set  $\mathbb{A}_K := \mathbb{R}^+ \times \mathbb{R}^K$  are denoted by  $\hat{y} := (\hat{\rho}, \hat{q})^T \in \mathbb{A}_K \times \mathbb{R}^{K+1}$  and Roe variables are  $\hat{\omega} := (\hat{\alpha}, \hat{\beta})^T \in \mathbb{A}_K \times \mathbb{R}^{K+1}$ . The mapping into conserved variables is defined as

$$\hat{\mathcal{Y}} : \mathbb{A}_K \times \mathbb{R}^{K+1} \rightarrow \mathbb{A}_K \times \mathbb{R}^{K+1}, \quad \hat{\omega} \mapsto \begin{pmatrix} \hat{\alpha} * \hat{\alpha} \\ \hat{\alpha} * \hat{\beta} \end{pmatrix} = \hat{y}. \quad (12)$$

The projected pressure law in terms of Roe variables is denoted by  $\hat{\pi}(\hat{\alpha})$ . The flux function for conserved variables depending on Roe variables is defined as

$$\hat{F} : \mathbb{A}_K \times \mathbb{R}^{K+1} \rightarrow \mathbb{R}^{2(K+1)}, \quad \hat{\omega} \mapsto \begin{pmatrix} \hat{\alpha} * \hat{\beta} \\ \hat{\beta} * \hat{\beta} + \hat{\pi}(\hat{\alpha}) \end{pmatrix}.$$

In contrast to the deterministic case (4) we cannot claim that the Galerkin root  $\hat{\alpha} * \hat{\alpha} = \hat{\rho}$  in (12) is uniquely solvable for any fixed  $\hat{\rho} \in \mathbb{A}_K$ . A system of  $K + 1$  nonlinear equations has to be solved, but there is no guarantee that positive quantities are well represented [28, Sec. 2.2]. In fact, we will only make use of a local invertibility that follows from the implicit function theorem.

**Remark 1.** Non-polynomial pressure laws can be pseudo-projected using a Taylor series at some point  $\check{\alpha}_0 \in \mathbb{R}^+$  with gPC modes  $\check{\alpha} := (\check{\alpha}_0, 0, \dots, 0)^T$ . Thus, we assume that the Taylor expansion up to order  $m \in \mathbb{N}_0$  in terms of Roe variables exists. For example, an expansion at the mean  $\check{\alpha}_0 := \hat{\alpha}_0 > 0$  is proposed in [28], where  $(\mathcal{G}_K[\alpha] - \check{\alpha}_0)$  describes the stochastic deviation. The truncated Taylor series reads

$$\pi_m(\mathcal{G}_K[\alpha](t, x; \xi)) := \sum_{\ell=0}^m \frac{\pi^{(\ell)}(\check{\alpha}_0)}{\ell!} \left( \mathcal{G}_K[\alpha](t, x; \xi) - \check{\alpha}_0 \right)^\ell.$$

The  $\ell$ -th moment  $y^\ell(t, x; \xi)$  is approximated componentwise by the recursion

$$\hat{\mathcal{G}}_k[y, \dots, y](t, x; \xi) := \sum_{k=0}^K (\hat{y}^{\ell*})_k(t, x) \phi_k(\xi), \quad \hat{y}^{\ell*} := ((\hat{y} * \hat{y}) * \dots * \hat{y}) * \hat{y}, \quad \hat{y}^{0*} := (1, 0, \dots, 0)^T. \quad (13)$$

As emphasised in [28], projections in the repeated Galerkin multiplications (13) are essentially truncations, which introduce additional approximation errors, since each one reduces a projection of order  $2K$  to  $K$ . Therefore, the order  $K$  must be sufficiently large compared to the power  $m$ . Then, the gPC modes of the Taylor expansion are

$$\widehat{\pi}_m(\hat{\alpha}) := \sum_{\ell=0}^m \frac{\pi^{(\ell)}(\check{\alpha}_0)}{\ell!} (\hat{\alpha} - \check{\alpha})^{\ell*}. \quad (14)$$



### 3.3. Hyperbolic Stochastic Galerkin Formulation

We introduce the stochastic Galerkin formulation as a conservation law describing gPC modes.

**Theorem 3.2** (Stochastic Galerkin Formulation). *Let  $\bar{y} := (\bar{\rho}, \bar{q})^T \in \mathbb{A}_K \times \mathbb{R}^{K+1}$  be given such that there is  $\bar{\alpha} \in \mathbb{A}_K$  satisfying  $\bar{\alpha} * \bar{\alpha} = \bar{\rho}$  and*

$$\mathcal{P}(\bar{\alpha}) \text{ is strictly positive definite.} \quad (\text{A2})$$

*Then, there is an open set  $\bar{\mathbb{A}}_K \subset \mathbb{A}_K$ , containing  $\bar{\rho}$ , such that for all  $\hat{y} = (\hat{\rho}, \hat{q})^T \in \bar{\mathbb{A}}_K \times \mathbb{R}^{K+1}$  the transform  $\hat{\mathcal{Y}}(\hat{\omega}) = \hat{y}$ ,  $\hat{\omega} = (\hat{\alpha}, \hat{\beta})^T$  is invertible with  $\hat{\omega} = \hat{\mathcal{Y}}^{-1}(\hat{y})$  and the gPC modes of the stochastic system  $(\mathcal{R}_K(\xi))$  are given by the conservation law*

$$\frac{\partial}{\partial t} \begin{pmatrix} \hat{\rho} \\ \hat{q} \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} \hat{q} \\ \hat{\beta} * \hat{\beta} + \hat{\pi}(\hat{\alpha}) \end{pmatrix} = 0. \quad (\hat{\mathcal{R}}_K)$$

We denote by  $\hat{f}(\hat{y}) := \hat{F}(\hat{\mathcal{Y}}^{-1}(\hat{y}))$  the flux function of the system  $(\hat{\mathcal{R}}_K)$ . Its Jacobian reads

$$\begin{aligned} D_{\hat{y}} \hat{f}(\hat{y}) &= \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ \frac{1}{2} D_{\hat{\alpha}} \hat{\pi}(\hat{\alpha}) \mathcal{P}^{-1}(\hat{\alpha}) - \mathcal{P}_{\hat{\beta}}^2(\hat{\omega}) & 2\mathcal{P}_{\hat{\beta}}(\hat{\omega}) \end{pmatrix} \quad \text{for } \mathcal{P}_{\hat{\beta}}(\hat{\omega}) := \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1}(\hat{\alpha}) \\ \text{and } \mathbb{O} &= \text{diag}\{0, \dots, 0\} \in \mathbb{R}^{(K+1) \times (K+1)}, \quad \mathbb{1} = \text{diag}\{1, \dots, 1\} \in \mathbb{R}^{(K+1) \times (K+1)}. \end{aligned} \quad (D\hat{\mathcal{R}}_K)$$

In particular, the projected pressure laws read

$$\begin{aligned} \hat{\pi}(\hat{\alpha}) &= a^2 (\hat{\alpha} * \hat{\alpha}) \quad \text{for isothermal Euler equations with pressure law } p(\rho) = a^2 \rho, \\ \hat{\pi}(\hat{\alpha}) &= \frac{g}{2} ((\hat{\alpha} * \hat{\alpha}) * (\hat{\alpha} * \hat{\alpha})) \quad \text{for shallow water equations with pressure law } p(\rho) = \frac{g}{2} \rho^2. \end{aligned} \quad (15)$$

*Proof.* The function  $\tilde{F}(\hat{\alpha}) := \mathcal{P}(\hat{\alpha})\hat{\alpha} - \hat{\rho}$  is continuously differentiable and under assumption (A2) its Jacobian is invertible at  $\hat{\alpha} = \bar{\alpha}$ . Therefore, we conclude with equation (11) and the implicit function theorem that the mapping

$$\hat{\mathcal{Y}}(\hat{\omega}) = \begin{pmatrix} \mathcal{P}(\hat{\alpha})\hat{\alpha} & \mathcal{P}(\hat{\alpha})\hat{\beta} \end{pmatrix}^T = \hat{y} \quad \text{and the Jacobian} \quad D_{\hat{\omega}} \hat{\mathcal{Y}}(\hat{\omega}) = \begin{pmatrix} 2\mathcal{P}(\hat{\alpha}) & \mathbb{O} \\ \mathcal{P}(\hat{\beta}) & \mathcal{P}(\hat{\alpha}) \end{pmatrix} \quad (16)$$

are also invertible on the open set  $\bar{\mathbb{A}}_K$ . Projections onto the  $k$ -th basis function read

$$\begin{aligned} \langle \hat{\mathcal{G}}_K[\alpha, \alpha], \phi_k \rangle_{\mathbb{P}} &= (\hat{\alpha} * \hat{\alpha})_k, & \langle \hat{\mathcal{G}}_K[\alpha, \beta], \phi_k \rangle_{\mathbb{P}} &= (\hat{\alpha} * \hat{\beta})_k, \\ \langle \hat{\mathcal{G}}_K[\beta, \beta], \phi_k \rangle_{\mathbb{P}} &= (\hat{\beta} * \hat{\beta})_k, & \langle \pi(\mathcal{G}_K[\alpha]), \phi_k \rangle_{\mathbb{P}} &= \hat{\pi}(\hat{\alpha})_k. \end{aligned}$$

Projection of the equations  $(\mathcal{R}_K(\xi))$  yields for the  $k$ -th component

$$\begin{aligned} &\left\langle \frac{\partial}{\partial t} \begin{pmatrix} \hat{\mathcal{G}}_K[\alpha, \alpha] \\ \hat{\mathcal{G}}_K[\alpha, \beta] \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} \hat{\mathcal{G}}_K[\alpha, \beta] \\ \hat{\mathcal{G}}_K[\beta, \beta] + \pi(\mathcal{G}_K[\alpha]) \end{pmatrix}, \phi_k \right\rangle_{\mathbb{P}} = 0 \\ \Leftrightarrow &\frac{\partial}{\partial t} \begin{pmatrix} (\hat{\alpha} * \hat{\alpha})_k \\ (\hat{\alpha} * \hat{\beta})_k \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} (\hat{\alpha} * \hat{\beta})_k \\ (\hat{\beta} * \hat{\beta})_k + \hat{\pi}(\hat{\alpha})_k \end{pmatrix} = 0. \end{aligned}$$

Thus, equations  $(\hat{\mathcal{R}}_K)$  hold. For smooth solutions we obtain

$$\begin{aligned}\hat{f}(\hat{y})_x &= \hat{F}(\hat{\omega})_x = D_{\hat{\omega}} \hat{F}(\hat{\omega}) \hat{\omega}_x = D_{\hat{\omega}} \hat{F}(\hat{\omega}) [\hat{\mathcal{Y}}^{-1}(\hat{y})]_x = D_{\hat{\omega}} \hat{F}(\hat{\omega}) [D_{\hat{\omega}} \hat{\mathcal{Y}}]^{-1}(\hat{\omega}) \hat{y}_x \\ &= \begin{pmatrix} \mathcal{P}(\hat{\beta}) & \mathcal{P}(\hat{\alpha}) \\ D_{\hat{\alpha}} \hat{\pi}(\hat{\alpha}) & 2\mathcal{P}(\hat{\beta}) \end{pmatrix} \begin{pmatrix} \frac{1}{2} \mathcal{P}^{-1}(\hat{\alpha}) & \mathbb{O} \\ -\frac{1}{2} \mathcal{P}^{-1}(\hat{\alpha}) \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1}(\hat{\alpha}) & \mathcal{P}^{-1}(\hat{\alpha}) \end{pmatrix} \hat{y}_x \\ &= \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ \frac{1}{2} D_{\hat{\alpha}} \hat{\pi}(\hat{\alpha}) \mathcal{P}^{-1}(\hat{\alpha}) - \mathcal{P}_{\hat{\nu}}^2(\hat{\omega}) & 2\mathcal{P}_{\hat{\nu}}(\hat{\omega}) \end{pmatrix} \hat{y}_x.\end{aligned}$$

Using the pseudo-spectral products, we obtain as projected pressure laws

$$\begin{aligned}\langle \pi(\mathcal{G}_K[\alpha]), \phi_k \rangle_{\mathbb{P}} &= \langle a^2 \hat{\mathcal{G}}_K[\alpha, \alpha], \phi_k \rangle_{\mathbb{P}} = a^2 (\hat{\alpha} * \hat{\alpha})_k && \text{for isothermal flows,} \\ \langle \pi(\mathcal{G}_K[\alpha]), \phi_k \rangle_{\mathbb{P}} &= \langle \frac{g}{2} \hat{\mathcal{G}}_K[\alpha, \dots, \alpha], \phi_k \rangle_{\mathbb{P}} = \frac{g}{2} ((\hat{\alpha} * \hat{\alpha}) * (\hat{\alpha} * \hat{\alpha}))_k && \text{for shallow water equations.}\end{aligned}$$

□

Similarly to [24], the Jacobian  $(D\hat{\mathcal{R}}_K)$  allows to prove hyperbolicity of the conservative formulation  $(\hat{\mathcal{R}}_K)$ . In particular for isothermal flows, this formulation allows to use any gPC expansion. Furthermore, we recover for  $K = 0$  the deterministic system  $(\mathcal{C})$ . To state conditions that ensure real eigenvalues of the Jacobian  $(D\hat{\mathcal{R}}_K)$ , we prove Lemma 3.3 first.

**Lemma 3.3.** *We assume the positive definiteness assumption (A2) and we define the matrices*

$$\mathcal{P}_{\hat{\nu}}(\hat{\omega}) := \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1}(\hat{\alpha}) \quad \text{and} \quad \mathcal{P}_2(\hat{\omega}) := \mathcal{P}^{-1/2}(\hat{\alpha}) \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1/2}(\hat{\alpha}).$$

We denote eigenvalue decompositions as

$$\mathcal{P}(\hat{\alpha}) = V(\hat{\alpha}) D_{\mathcal{P}}(\hat{\alpha}) V^T(\hat{\alpha}) \quad \text{and} \quad \mathcal{P}_{\hat{\nu}}(\hat{\omega}) = Q(\hat{\omega}) D_{\hat{\nu}}(\hat{\omega}) Q^{-1}(\hat{\omega}).$$

Then, it holds:

- (i) *There exists an orthogonal eigenvector matrix  $V(\hat{\alpha})$ . The symmetric square root, defined as  $\mathcal{P}^{1/2}(\hat{\alpha}) := V(\hat{\alpha}) D_{\mathcal{P}}^{1/2}(\hat{\alpha}) V^T(\hat{\alpha})$ , exists and it is unique.*
- (ii) *The matrix  $\mathcal{P}_{\hat{\nu}}(\hat{\omega})$  is diagonalizable with real eigenvalues for all  $\hat{\beta} \in \mathbb{R}^{K+1}$ . These eigenvalues coincide with those of  $\mathcal{P}_2(\hat{\omega})$ .*
- (iii) *The matrices  $a\mathcal{P}(\hat{\alpha}) \pm \mathcal{P}(\hat{\beta})$  are strictly positive definite if and only if  $a\mathbb{1} \pm D_{\hat{\nu}}(\hat{\omega}) > 0$  holds.*
- (iv) *The matrix  $\mathcal{P}_2(\hat{\omega})$  is strictly positive definite and symmetric if and only if  $\mathcal{P}(\hat{\beta})$  is strictly positive definite.*

*Proof.* Statement (i) holds according to [34]. Sylvester's law of inertia states that two congruent symmetric matrices have the same number of strictly positive eigenvalues.

- (ii) Under assumption (A2) there exists an invertible and symmetric square root  $\mathcal{P}^{1/2}(\hat{\alpha})$ . The matrix  $\mathcal{P}_2(\hat{\omega})$  is symmetric and hence diagonalizable with real eigenvalues. Due to

$$\mathcal{P}^{-1/2}(\hat{\alpha}) \mathcal{P}_{\hat{\nu}}(\hat{\omega}) \mathcal{P}^{1/2}(\hat{\alpha}) = \mathcal{P}_2(\hat{\omega})$$

eigenvalues of the nonsymmetric, but similar matrix  $\mathcal{P}_\nu(\hat{\omega})$  coincide with those of  $\mathcal{P}_2(\hat{\omega})$ .

- (iii) The matrices  $a\mathcal{P}(\hat{\alpha}) \pm \mathcal{P}(\hat{\beta})$  and  $a\mathbb{1} \pm \mathcal{P}_2(\hat{\omega}) = a\mathbb{1} \pm \mathcal{P}^{-1/2}(\hat{\alpha})\mathcal{P}(\hat{\beta})\mathcal{P}^{-1/2}(\hat{\alpha})$  are congruent and thus have the same number of strictly positive eigenvalues. Due to (ii) the matrix  $a\mathbb{1} \pm \mathcal{P}_2(\hat{\omega})$  is similar to  $a\mathbb{1} \pm \mathcal{P}(\hat{\beta})\mathcal{P}^{-1}(\hat{\alpha})$  which is strictly positive definite if and only if  $a\mathbb{1} \pm D_\nu(\hat{\omega}) > 0$ .
- (iv) Symmetry follows from the symmetry of the square root. Sylvester's law of inertia states that all eigenvalues of  $\mathcal{P}(\hat{\beta})$  are strictly positive if and only if all eigenvalues of  $\mathcal{P}_2(\hat{\omega})$  are strictly positive.

□

The Jacobian ( $D\hat{\mathcal{R}}_K$ ) is in general not diagonalizable with real eigenvalues. In the deterministic case and in the absence of vacuum states the system is always *strictly* hyperbolic. We will establish in the following corollaries two sets of real eigenvalues  $\{\hat{\Lambda}_j^\pm(\hat{\omega}) \in \mathbb{R} \mid j = 0, \dots, K\}$ , which will simplify in the deterministic case to  $\lambda^\pm(y)$ . The notation is inspired by viewing a **stochastic subsonic flow** as states with eigenvalues satisfying

$$\hat{\Lambda}_i^-(\hat{\omega}) < 0 < \hat{\Lambda}_j^+(\hat{\omega}) \quad \text{for all } i, j = 0, \dots, K. \quad (\text{S})$$

Still eigenvalues of the sets  $\hat{\Lambda}^\pm(\hat{\omega})$  may *coincide* and condition (S) is as assumption for the following results not needed. With additional assumptions on the gPC basis we can show hyperbolicity.

**Corollary 1** (Constant Eigenvectors). *Assume there exists an eigenvalue decomposition with constant eigenvectors, i.e.  $\mathcal{P}(\hat{\alpha}) = VD_{\mathcal{P}}(\hat{\alpha})V^T$ . Under assumption (A2) and for an eigenvalue decomposition of the pressure law with positive eigenvalues, denoted as  $D_{\hat{\alpha}}\hat{\pi}(\hat{\alpha}) = VD_{\hat{\pi}}(\hat{\alpha})V^T$ , the Jacobian ( $D\hat{\mathcal{R}}_K$ ) has the eigenvalue decomposition  $D_{\hat{y}}\hat{f}(\hat{y}) = [\mathcal{V}\hat{T}(\hat{\omega})]\hat{\Lambda}(\hat{\omega})[\mathcal{V}\hat{T}(\hat{\omega})]^{-1}$  with*

$$\begin{aligned} \hat{\Lambda}^\pm(\hat{\omega}) &:= D_\nu(\hat{\omega}) \pm \sqrt{\frac{1}{2}D_{\hat{\pi}}(\hat{\alpha})D_{\mathcal{P}}^{-1}(\hat{\alpha})}, & \hat{\Lambda}(\hat{\omega}) &:= \text{diag}\{\hat{\Lambda}^+(\hat{\omega}), \hat{\Lambda}^-(\hat{\omega})\}, \\ \hat{T}(\hat{\omega}) &:= \begin{pmatrix} \mathbb{1} & \mathbb{1} \\ \hat{\Lambda}^+(\hat{\omega}) & \hat{\Lambda}^-(\hat{\omega}) \end{pmatrix}, & \mathcal{V} &:= \text{diag}\{V, V\}. \end{aligned}$$

In particular, we have  $D_\nu(\hat{\omega}) = D_{\mathcal{P}}(\hat{\beta})D_{\mathcal{P}}^{-1}(\hat{\alpha})$  and we obtain

$$\hat{\Lambda}^\pm(\hat{\omega}) = D_\nu(\hat{\omega}) \pm a\mathbb{1} \quad \text{for isothermal Euler equations with pressure law } p(\rho) = a^2\rho.$$

*Proof.* The projected pressure law for isothermal Euler equations satisfies

$$D_{\hat{\alpha}}\hat{\pi}(\hat{\alpha}) = 2a^2\mathcal{P}(\hat{\alpha}) = V[2a^2D_{\mathcal{P}}(\hat{\alpha})]V^T, \quad \text{i.e. } D_{\hat{\pi}}(\hat{\alpha}) = 2a^2D_{\mathcal{P}}(\hat{\alpha}).$$

The eigenvalue decompositions of the blockmatrices of the Jacobian ( $D\hat{\mathcal{R}}_K$ ) read

$$\begin{aligned} \mathcal{P}_\nu(\hat{\omega}) &= [VD_{\mathcal{P}}(\hat{\beta})V^T][VD_{\mathcal{P}}(\hat{\alpha})V^T]^{-1} = V[D_{\mathcal{P}}(\hat{\beta})D_{\mathcal{P}}^{-1}(\hat{\alpha})]V^T = VD_\nu(\hat{\omega})V^T, \\ D_{\hat{\alpha}}\hat{\pi}(\hat{\alpha})\mathcal{P}^{-1}(\hat{\alpha}) &= [VD_{\hat{\pi}}(\hat{\alpha})V^T][VD_{\mathcal{P}}(\hat{\alpha})V^T]^{-1} = V[D_{\hat{\pi}}(\hat{\alpha})D_{\mathcal{P}}^{-1}(\hat{\alpha})]V^T. \end{aligned}$$

Assumption (A2) and  $D_{\hat{\pi}}(\hat{\alpha}) > 0$  guarantee  $\hat{\Lambda}^+(\hat{\omega}) - \hat{\Lambda}^-(\hat{\omega}) = \sqrt{2D_{\hat{\pi}}(\hat{\alpha})D_{\hat{\mathcal{P}}}^{-1}(\hat{\alpha})} > 0$ . Then, the claim follows from Theorem 3.2, the orthogonal matrix  $\mathcal{V}^T = \mathcal{V}^{-1}$  and

$$\begin{aligned} D_{\hat{y}}\hat{f}(\hat{y}) &= \mathcal{V} \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ \frac{1}{2}D_{\hat{\pi}}(\hat{\alpha})D_{\hat{\mathcal{P}}}^{-1}(\hat{\alpha}) - D_{\hat{\nu}}^2(\hat{\omega}) & 2D_{\hat{\nu}}(\hat{\omega}) \end{pmatrix} \mathcal{V}^T = \mathcal{V} \left[ \hat{T}(\hat{\omega})\hat{\Lambda}(\hat{\omega})\hat{T}^{-1}(\hat{\omega}) \right] \mathcal{V}^T \\ \text{with } \hat{T}^{-1}(\hat{\omega}) &= \begin{pmatrix} -(\hat{\Lambda}^+(\hat{\omega}) - \hat{\Lambda}^-(\hat{\omega}))^{-1}\hat{\Lambda}^-(\hat{\omega}) & (\hat{\Lambda}^+(\hat{\omega}) - \hat{\Lambda}^-(\hat{\omega}))^{-1} \\ (\hat{\Lambda}^+(\hat{\omega}) - \hat{\Lambda}^-(\hat{\omega}))^{-1}\hat{\Lambda}^+(\hat{\omega}) & -(\hat{\Lambda}^+(\hat{\omega}) - \hat{\Lambda}^-(\hat{\omega}))^{-1} \end{pmatrix}. \end{aligned}$$

□

The assumption of constant eigenvectors is taken from [24, Lemma 1]. There, a Wiener-Haar expansion and linear multiwavelets, which satisfy this assumption, are used. Moderate stochastic variations are assumed as well. The interested reader finds a proof in [24, Appendix B, C] and an introduction into these expansions in [11, 32]. However, this assumption is not true for general basis functions, e.g. Legendre polynomials. For isothermal Euler equations with projected flux function

$$\hat{F}(\hat{\omega}) = \left( \mathcal{P}(\hat{\alpha})\hat{\beta}, \mathcal{P}(\hat{\beta})\hat{\beta} + a^2\mathcal{P}(\hat{\alpha})\hat{\alpha} \right)^T, \quad (17)$$

we can extend Corollary 1 to arbitrary gPC expansions.

**Corollary 2** (Isothermal Euler Equations). *Under assumption (A2) the Jacobian  $(D\hat{\mathcal{R}}_K)$  with the pressure law (15) for isothermal Euler equations has real eigenvalues and it is given by*

$$\begin{aligned} D_{\hat{y}}\hat{f}(\hat{y}) &= \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ a^2\mathbb{1} - \mathcal{P}_{\hat{\nu}}^2(\hat{\omega}) & 2\mathcal{P}_{\hat{\nu}}(\hat{\omega}) \end{pmatrix} = [\mathcal{Q}(\hat{\omega})\hat{T}(\hat{\omega})]\hat{\Lambda}(\hat{\omega})[\mathcal{Q}(\hat{\omega})\hat{T}(\hat{\omega})]^{-1} \quad \text{with} \\ \hat{\Lambda}^{\pm}(\hat{\omega}) &:= D_{\hat{\nu}}(\hat{\omega}) \pm a\mathbb{1}, & \hat{\Lambda}(\hat{\omega}) &:= \text{diag}\{\hat{\Lambda}^+(\hat{\omega}), \hat{\Lambda}^-(\hat{\omega})\}, \\ \hat{T}(\hat{\omega}) &:= \begin{pmatrix} \mathbb{1} & \mathbb{1} \\ \hat{\Lambda}^+(\hat{\omega}) & \hat{\Lambda}^-(\hat{\omega}) \end{pmatrix}, & \mathcal{Q}(\hat{\omega}) &:= \text{diag}\{Q(\hat{\omega}), Q(\hat{\omega})\}, \end{aligned}$$

where  $\mathcal{P}_{\hat{\nu}}(\hat{\omega}) = \mathcal{P}(\hat{\beta})\mathcal{P}^{-1}(\hat{\alpha})$  has the eigenvalue decomposition  $\mathcal{P}_{\hat{\nu}}(\hat{\omega}) = Q(\hat{\omega})D_{\hat{\nu}}(\hat{\omega})Q^{-1}(\hat{\omega})$ .

*Proof.* Provided that the assumption (A2) holds, Lemma 3.3 ensures an eigenvalue decomposition  $\mathcal{P}_{\hat{\nu}}(\hat{\omega}) = Q(\hat{\omega})D_{\hat{\nu}}(\hat{\omega})Q^{-1}(\hat{\omega})$  with real eigenvalues. So  $\hat{\Lambda}(\hat{\omega})$  is real as well. We calculate

$$\begin{aligned} D_{\hat{y}}\hat{f}(\hat{y}) &= \mathcal{Q}(\hat{\omega}) \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ a^2\mathbb{1} - D_{\hat{\nu}}^2(\hat{\omega}) & 2D_{\hat{\nu}}(\hat{\omega}) \end{pmatrix} \mathcal{Q}^{-1}(\hat{\omega}) = [\mathcal{Q}(\hat{\omega})\hat{T}(\hat{\omega})]\hat{\Lambda}(\hat{\omega})[\mathcal{Q}(\hat{\omega})\hat{T}(\hat{\omega})]^{-1} \\ \text{for } \hat{T}(\hat{\omega}) &= \begin{pmatrix} \mathbb{1} & \mathbb{1} \\ \hat{\Lambda}^+(\hat{\omega}) & \hat{\Lambda}^-(\hat{\omega}) \end{pmatrix} \quad \text{and} \quad \hat{T}^{-1}(\hat{\omega}) = \frac{1}{2a} \begin{pmatrix} -\hat{\Lambda}^-(\hat{\omega}) & \mathbb{1} \\ \hat{\Lambda}^+(\hat{\omega}) & -\mathbb{1} \end{pmatrix}. \end{aligned}$$

□

#### 4. Discussion of Assumptions

The stochastic Galerkin formulation  $(\hat{\mathcal{R}}_K)$  is proven hyperbolic only in a local domain around given states  $\bar{y}$ , where the matrix  $\mathcal{P}(\bar{\alpha})$  is assumed to be strictly positive definite. For example the state  $\bar{y}$  could be the initial value. We first argue that the assumption (A2) is reasonable for initial states and then, we discuss a possible loss of hyperbolicity.

##### 4.1. Initial Values

An obvious requirement in the deterministic case is a strictly positive density. For the strong formulations  $(\mathcal{C}(\xi))$  and  $(\mathcal{R}(\xi))$  we require the conditions  $\rho(\xi) > 0$   $\mathbb{P}$ -a.s. and  $\alpha(\xi) > 0$   $\mathbb{P}$ -a.s. A straight forward analogue for the stochastic systems  $(\mathcal{C}_K(\xi))$  and  $(\mathcal{R}_K(\xi))$  with truncated gPC expansion would be

$$\mathcal{G}_K[\rho](\xi) > 0 \text{ } \mathbb{P}\text{-a.s.} \quad \text{and} \quad \mathcal{G}_K[\alpha](\xi) > 0 \text{ } \mathbb{P}\text{-a.s.} \quad (\text{A3})$$

Indeed, assumption (A3) is used to guarantee hyperbolicity of general quasilinear forms in [18]. It is also used for an operator splitting based method [19, Remark 4.1] and moderate stochastic variations are assumed for the Roe variable transform in [24, Lemma 1]. *Initial values* that satisfy assumption (A3) can be described for example by Legendre polynomials, by a Wiener-Haar expansion and by a  $\beta$ -distribution, which is proposed in [35, 30] as a “truncated Gaussian model”, since negative densities do not occur due to truncation errors of the series expansion. An analysis of the log-normal and the reflected Gaussian distribution is found in [36].

Note that assumption (A3) implies the assumption (A2). Indeed, define the random variables

$$Z_i(\xi) := \phi_i(\xi) \mathcal{G}_K[\alpha](\xi)^{1/2} \text{ } \mathbb{P}\text{-a.s.} \quad (18)$$

which are well-defined provided that (A3) holds. Hence, the  $i, j$ -th entry  $\mathcal{P}_{i,j}(\hat{\alpha})$  satisfies

$$\mathcal{P}_{i,j}(\hat{\alpha}) = \sum_{\ell=0}^K \hat{\alpha}_\ell \left\langle \phi_\ell, \phi_i \phi_j \right\rangle_{\mathbb{P}} = \left\langle \sum_{\ell=0}^K \hat{\alpha}_\ell \phi_\ell, \phi_i \phi_j \right\rangle_{\mathbb{P}} = \left\langle \mathcal{G}_K[\alpha], \phi_i \phi_j \right\rangle_{\mathbb{P}} = \mathbb{E} \left[ Z_i Z_j \right]. \quad (19)$$

Thus, the matrix  $\mathcal{P}(\hat{\alpha})$  is a covariance matrix, i.e. it is positive semidefinite and symmetric. Since the projection  $\mathcal{G}_K[\alpha]^{1/2}$  is assumed to be  $\mathbb{P}$ -a.s. strictly positive and since  $\phi_i(\xi)$  are basis functions, we deduce even strict positive definiteness of  $\mathcal{P}(\hat{\alpha})$  as in [18, Th. 2.1]. We conclude that assumption (A2) is a relaxation of (A3). Therefore, condition (A2) does *not restrict additionally* the choice of initial values.

**Remark 2.** Condition (A3) is *violated* for example by Gaussian distributed germs with Hermite polynomials. Assumption (A2) is a relaxation of (A3): If the sequence  $\{\hat{\alpha}_\ell\}_{\ell=0}^\infty$  decays fast enough or, in intuitive terms, if the variance given by  $\hat{\alpha}_1, \dots, \hat{\alpha}_K$  is sufficiently small compared to the expected value  $\mathbb{E}[\mathcal{G}_K[\alpha]] = \hat{\alpha}_0$ , then the condition

$$\sigma_{\min} \left[ \sum_{\ell=1}^K \hat{\alpha}_\ell \mathcal{M}_\ell \right] > -\hat{\alpha}_0 \quad (20)$$

implies assumption (A2), where  $\sigma_{\min}$  denotes the smallest eigenvalue. The inequality (20) can be satisfied also for unbounded gPC expansions.

Although stochastic densities that are unbounded from below are not physically relevant for fluid flows, an expansion with Hermite polynomials allows to consider the normal distribution or more general Gaussian processes using a Karhunen-Loève transform. These processes are commonly used in the engineering community to characterize stochastic inputs. There may be interest in such an expansion, because when reconstructing the density from the gPC expansion of the Roe variable  $\alpha$  the density would be positive, i.e.  $\rho(\xi) \approx \mathcal{G}_K[\alpha](\xi)^2 \geq 0$ . The product  $\mathcal{G}_K[\alpha](\xi)^2$ , however, lives in  $\mathcal{S}_{2K}$  and the projection onto  $\mathcal{S}_K$ , i.e. the pseudo-spectral product  $\hat{\mathcal{G}}[\alpha, \alpha](\xi)$  is no more guaranteed positive. Therefore, we have introduced the weaker assumption (A2).

#### 4.2. Possible Loss of Hyperbolicity

Deterministic Euler and shallow water equations have distinct eigenvalues and genuinely nonlinear or linearly degenerate characteristic fields. Furthermore, the systems are endowed with an entropy and entropy flux pair [37]. Then, an entropy solution exists for any fixed time domain as long as the total variation of initial values is sufficiently small [38, Th. 7.1].

It is an open problem if also the system  $(\hat{\mathcal{R}}_K)$  admits an entropy. The assumption of genuine nonlinearity can be weakened [39, 40, 38], but the eigenvalues of stochastic Galerkin formulations coincide in general. Therefore, we expect that the existence of a solution cannot be guaranteed for any fixed time domain. Without the introduction of Roe variables from [24] there would be no existence and uniqueness result, since the system  $(\hat{\mathcal{C}}_K)$  is in general *not hyperbolic* [33]. Furthermore, positive definiteness of the matrix  $\mathcal{P}(\hat{\rho})$  is also a problem for the system  $(\hat{\mathcal{C}}_K)$ .

Theorem 3.2 guarantees a — possibly small — domain  $\bar{\mathbb{A}}_K$  around initial states, where the system remains hyperbolic. However, it is not guaranteed that the solution exists for any time domain. We can verify the hyperbolicity of system  $(\hat{\mathcal{R}}_K)$  at least numerically in a simple way, as described in the following section.

### 5. Numerical Discretization of Isothermal Euler Equations

We describe how to discretize the intrusive formulation  $(\hat{\mathcal{R}}_K)$ .

#### 5.1. Flux Function for gPC Modes

Using Gaussian quadrature with  $n := \lceil \frac{3}{2}(K+1) \rceil$  points  $x_k$  and weights  $\hat{w}(x_k)$ , the norms  $\|\phi_\ell\|_{\mathbb{P}}$  and tensors  $\mathcal{M}_\ell$  can be calculated exactly by

$$\langle \phi_\ell, \phi_i \phi_j \rangle_{\mathbb{P}} = \int \phi_\ell(\xi) \phi_i(\xi) \phi_j(\xi) d\mathbb{P} = \sum_{k=1}^n \phi_\ell(x_k) \phi_i(x_k) \phi_j(x_k) \hat{w}(x_k).$$

We have shown in equation (19) that assumption  $\mathcal{G}_K[\alpha] > 0$   $\mathbb{P}$ -a.s. implies (A2). Writing

$$\mathcal{P}(\hat{\alpha})_{i,j}(t, x) = \sum_{k=1}^n \mathcal{G}_K[\alpha](t, x; x_k) \phi_i(x_k) \phi_j(x_k) \hat{w}(x_k) \quad (21)$$

we see that it suffices to assume

$$\mathcal{G}_K[\alpha](t, x; x_k) > 0 \text{ for all Gauss quadrature points } \{x_1, \dots, x_n\}. \quad (\text{A4})$$

Thus, we can verify hyperbolicity easily by testing condition (A4) at *finitely* many points.

### 5.2. Computing the Roe Variables

Because of the transform in Roe variables  $\hat{\omega} = \hat{\mathcal{Y}}^{-1}(\hat{y})$  a nonlinear system  $\mathcal{P}(\hat{\alpha})\hat{\alpha} = \hat{\rho}$  and a linear system  $\mathcal{P}(\hat{\alpha})\hat{\beta} = \hat{q}$  must be solved. We solve the equations  $\tilde{F}(\hat{\alpha}) := \mathcal{P}(\hat{\alpha})\hat{\alpha} - \hat{\rho} = 0$  using Newton's method. The recursion reads

$$\begin{aligned}\hat{\alpha}_{\text{new}} &= \hat{\alpha}_{\text{old}} - [\mathbf{D}_{\hat{\alpha}} \tilde{F}]^{-1}(\hat{\alpha}_{\text{old}}) \tilde{F}(\hat{\alpha}_{\text{old}}) = \hat{\alpha}_{\text{old}} - \frac{1}{2} \mathcal{P}^{-1}(\hat{\alpha}_{\text{old}}) [\mathcal{P}(\hat{\alpha}_{\text{old}}) \hat{\alpha}_{\text{old}} - \hat{\rho}] \\ &= \frac{1}{2} [\hat{\alpha}_{\text{old}} + \mathcal{P}^{-1}(\hat{\alpha}_{\text{old}}) \hat{\rho}].\end{aligned}$$

Under the positive definiteness assumption the linear systems of equations can be solved efficiently using e.g. the Cholesky decomposition of  $\mathcal{P}(\hat{\alpha})$ .

### 5.3. Numerical Flux Function

Roe suggested in [26] to use a flux function of the form

$$\hat{\mathbf{F}}(\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r) = \frac{1}{2} [\hat{f}(\hat{\mathbf{y}}_{\ell}) + \hat{f}(\hat{\mathbf{y}}_r)] + \frac{1}{2} |\hat{A}(\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r)| (\hat{\mathbf{y}}_{\ell} - \hat{\mathbf{y}}_r), \quad (\text{RoeFlux})$$

where  $\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r$  denote cell averages. The **Roe matrix**  $\hat{A}$  with absolute value  $|\hat{A}| := \hat{\mathcal{X}} |\hat{\mathcal{D}}| \hat{\mathcal{X}}^{-1}$  must have the **Roe properties**:

$$(\text{Roe 1}) \quad \hat{A}(\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r) (\hat{\mathbf{y}}_{\ell} - \hat{\mathbf{y}}_r) = \hat{f}(\hat{\mathbf{y}}_{\ell}) - \hat{f}(\hat{\mathbf{y}}_r)$$

$$(\text{Roe 2}) \quad \hat{A}(\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r) \text{ is diagonalizable with real eigenvalues } \hat{\mathcal{D}}(\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r)$$

$$(\text{Roe 3}) \quad \hat{A}(\hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r) \rightarrow \mathbf{D}_{\bar{\mathbf{y}}} \hat{f}(\bar{\mathbf{y}}) \text{ smoothly as } \hat{\mathbf{y}}_{\ell}, \hat{\mathbf{y}}_r \rightarrow \bar{\mathbf{y}}$$

It is shown in [27, Sec. 14.2.4] for the deterministic case that the Jacobian evaluated at an averaged velocity satisfies the Roe properties. Therefore, the Roe matrix should reduce for  $K = 0$  to

$$\hat{A}(\mathbf{y}_{\ell}, \mathbf{y}_r) := \begin{pmatrix} 0 & 1 \\ a^2 - \bar{u}^2(\mathbf{y}_{\ell}, \mathbf{y}_r) & 2\bar{u}(\mathbf{y}_{\ell}, \mathbf{y}_r) \end{pmatrix} \quad \text{with} \quad \bar{u}(\mathbf{y}_{\ell}, \mathbf{y}_r) := \frac{\sqrt{\rho_{\ell}} u(\mathbf{y}_{\ell}) + \sqrt{\rho_r} u(\mathbf{y}_r)}{\sqrt{\rho_{\ell}} + \sqrt{\rho_r}}. \quad (22)$$

The Roe matrix (22) depends only on the Roe variables, so we define  $\bar{A}(\omega_{\ell}, \omega_r) := \hat{A}(\mathcal{Y}(\omega_{\ell}), \mathcal{Y}(\omega_r))$  in the deterministic case. We extend this definition and write  $\bar{A}(\hat{\omega}_{\ell}, \hat{\omega}_r)$  with  $\hat{\omega}_{\ell, r} := \hat{\mathcal{Y}}^{-1}(\hat{\mathbf{y}}_{\ell, r})$  also for the stochastic case.

**Theorem 5.1** (Roe Matrix). *Properties (Roe 1) – (Roe 3) are satisfied by the matrix*

$$\bar{A}(\hat{\omega}_{\ell}, \hat{\omega}_r) := \begin{pmatrix} \mathbb{O} & \mathbb{1} \\ a^2 - \mathcal{P}_{\hat{\nu}}^2(\bar{\omega}) & 2\mathcal{P}_{\hat{\nu}}(\bar{\omega}) \end{pmatrix} \quad \text{evaluated at} \quad \bar{\omega} := \frac{\hat{\omega}_{\ell} + \hat{\omega}_r}{2} \quad (23)$$

if the positive definiteness assumption (A2) holds.

*Proof.* Due to  $\mathcal{P}(\bar{\alpha}) = \frac{1}{2} [\mathcal{P}(\hat{\alpha}_\ell) + \mathcal{P}(\hat{\alpha}_r)]$  we get

$$\begin{aligned} D_{\bar{\omega}} \hat{\mathcal{Y}}(\bar{\omega}) &= \begin{pmatrix} 2\mathcal{P}(\bar{\alpha}) & \mathbb{O} \\ \mathcal{P}(\bar{\beta}) & \mathcal{P}(\bar{\alpha}) \end{pmatrix} = \frac{1}{2} [D_{\hat{\omega}_\ell} \hat{\mathcal{Y}}(\hat{\omega}_\ell) + D_{\hat{\omega}_r} \hat{\mathcal{Y}}(\hat{\omega}_r)], & \frac{D_{\bar{\omega}} \hat{\mathcal{Y}}(\bar{\omega}) \bar{\omega}}{2} &= \hat{\mathbf{y}}, \\ D_{\bar{\omega}} \hat{F}(\bar{\omega}) &= \begin{pmatrix} \mathcal{P}(\bar{\beta}) & \mathcal{P}(\bar{\alpha}) \\ 2a^2 \mathcal{P}(\bar{\alpha}) & 2\mathcal{P}(\bar{\beta}) \end{pmatrix} = \frac{1}{2} [D_{\hat{\omega}_\ell} \hat{F}(\hat{\omega}_\ell) + D_{\hat{\omega}_r} \hat{F}(\hat{\omega}_r)], & \frac{D_{\bar{\omega}} \hat{F}(\bar{\omega}) \bar{\omega}}{2} &= \hat{f}(\hat{\mathbf{y}}) \end{aligned}$$

such that  $\bar{A}(\hat{\omega}_\ell, \hat{\omega}_r) = D_{\bar{\omega}} \hat{F}(\bar{\omega}) [D_{\bar{\omega}} \hat{\mathcal{Y}}]^{-1}(\bar{\omega})$  and

$$\begin{aligned} D_{\hat{\omega}_\ell} \hat{\mathcal{Y}}(\hat{\omega}_\ell) \hat{\omega}_r &= \begin{pmatrix} 2\mathcal{P}(\hat{\alpha}_\ell) \hat{\alpha}_r \\ \mathcal{P}(\hat{\beta}_\ell) \hat{\alpha}_r + \mathcal{P}(\hat{\alpha}_\ell) \hat{\beta}_r \end{pmatrix} = D_{\hat{\omega}_r} \hat{\mathcal{Y}}(\hat{\omega}_r) \hat{\omega}_\ell \\ \Rightarrow D_{\bar{\omega}} \hat{\mathcal{Y}}(\bar{\omega}) (\hat{\omega}_\ell - \hat{\omega}_r) &= \frac{1}{2} [D_{\hat{\omega}_\ell} \hat{\mathcal{Y}}(\hat{\omega}_\ell) \hat{\omega}_\ell - D_{\hat{\omega}_r} \hat{\mathcal{Y}}(\hat{\omega}_r) \hat{\omega}_r] = \hat{\mathbf{y}}_\ell - \hat{\mathbf{y}}_r, \end{aligned} \quad (24)$$

$$\begin{aligned} D_{\hat{\omega}_\ell} \hat{F}(\hat{\omega}_\ell) \hat{\omega}_r &= \begin{pmatrix} \mathcal{P}(\hat{\beta}_\ell) \hat{\alpha}_r + \mathcal{P}(\hat{\alpha}_\ell) \hat{\beta}_r \\ 2\mathcal{P}(\hat{\beta}_\ell) \hat{\beta}_r + 2a^2 \mathcal{P}(\hat{\alpha}_\ell) \hat{\alpha}_r \end{pmatrix} = D_{\hat{\omega}_r} \hat{F}(\hat{\omega}_r) \hat{\omega}_\ell \\ \Rightarrow D_{\bar{\omega}} \hat{F}(\bar{\omega}) (\hat{\omega}_\ell - \hat{\omega}_r) &= \frac{1}{2} [D_{\hat{\omega}_\ell} \hat{F}(\hat{\omega}_\ell) \hat{\omega}_\ell - D_{\hat{\omega}_r} \hat{F}(\hat{\omega}_r) \hat{\omega}_r] = \hat{f}(\hat{\mathbf{y}}_\ell) - \hat{f}(\hat{\mathbf{y}}_r). \end{aligned} \quad (25)$$

Then, we obtain the following properties:

(Roe 1) Since (A2) holds for  $\mathcal{P}(\bar{\alpha})$ , the inverse  $[D_{\bar{\omega}} \hat{\mathcal{Y}}]^{-1}(\bar{\omega})$  exists. Equations (24) and (25) yield

$$\bar{A}(\hat{\omega}_\ell, \hat{\omega}_r) (\hat{\mathbf{y}}_\ell - \hat{\mathbf{y}}_r) = D_{\bar{\omega}} \hat{F}(\bar{\omega}) [D_{\bar{\omega}} \hat{\mathcal{Y}}]^{-1}(\bar{\omega}) (\hat{\mathbf{y}}_\ell - \hat{\mathbf{y}}_r) = \hat{f}(\hat{\mathbf{y}}_\ell) - \hat{f}(\hat{\mathbf{y}}_r).$$

(Roe 2) This follows from Corollary 2 due to  $\bar{A}(\hat{\omega}_\ell, \hat{\omega}_r) = D_{\bar{\mathbf{y}}} \hat{f}(\bar{\mathbf{y}})$  for  $\bar{\mathbf{y}} := \hat{\mathcal{Y}}(\bar{\omega})$ .

(Roe 3) It follows directly from the structure of  $\bar{A}(\hat{\omega}_\ell, \hat{\omega}_r)$  and  $\bar{A}(\bar{\omega}, \bar{\omega}) = D_{\bar{\mathbf{y}}} \hat{f}(\bar{\mathbf{y}})$ .

□

A derivation of a Roe matrix for the full Euler equations has additional difficulties. We refer the interested reader to [24], where property (Roe 3) is formulated in terms of the flux function in Roe variables, i.e.  $\bar{A}(\hat{\omega}_\ell, \hat{\omega}_r) \rightarrow D_{\bar{\omega}} \hat{F}(\bar{\omega})$  for  $\hat{\omega}_\ell, \hat{\omega}_r \rightarrow \bar{\omega}$ . Then, the Roe properties are satisfied for the Wiener-Haar expansion and for linear multiwavelets. The Roe matrix (23) satisfies the Roe properties for any gPC expansions and in the case  $K = 0$  it reduces to the deterministic case (22).

#### 5.4. Eigendecomposition of the Jacobian

If there is an eigenvalue decomposition of the form  $\mathcal{P}(\hat{\alpha}) = V D_{\mathcal{P}}(\hat{\alpha}) V^T$ , the eigenvalue decomposition of the matrix  $\mathcal{P}_\nu(\hat{\omega})$  is efficient, since eigenvectors are constant. The matrix  $\mathcal{P}_\nu(\hat{\omega})$  in Corollary 2, however, is nonsymmetric for general gPC expansions. Thus, its eigenvalue decomposition causes a computational overhead. This is the computational price that has to be paid to guarantee hyperbolicity for arbitrary gPC expansions.

If the flow direction is unchanged and we have  $q(\xi) \rho^{-1/2}(\xi) \approx \mathcal{G}_K[\beta](\xi) > 0$   $\mathbb{P}$ -a.s., then additionally the matrix  $\mathcal{P}(\hat{\beta})$  is strictly positive definite. Computational cost can be reduced by considering the similar matrix  $\mathcal{P}_2(\hat{\omega})$ , which has the same eigenvalues as the matrix  $\mathcal{P}_\nu(\hat{\omega})$ . Lemma 3.3 states that the similar matrix  $\mathcal{P}_2(\hat{\omega})$  is strictly positive definite if and only if  $\mathcal{P}(\hat{\beta})$  is strictly positive



definite. Thus, the cheaper Cholesky decomposition can be applied. This case can also be checked by condition (A4).

### 5.5. Eigenvalue Estimate

In [18] the estimate for the spectral radius  $\sigma_{\max}$  of the Jacobian

$$\sigma_{\max}[\mathbf{D}_{\hat{y}}\hat{f}(\hat{y}(t, x))] \leq \max_{\xi} \left\{ \sigma_{\max}[\mathbf{D}_y f(y(t, x; \xi))] \right\} \quad (26)$$

based on the random Jacobian of the underlying random conservation law is proven. It gives an estimate for general quasilinear forms with no restriction on a particular basis. It shows that the spectrum of the projected system is within the range of the random spectrum of the underlying random system. This illustrates that the intrusive formulation is weaker than the strong formulation ( $\mathcal{C}(\xi)$ ). This estimate is not meaningful for distributions with unbounded support, since the right hand side may not be bounded. Furthermore, it is computationally expensive, since an optimization problem has to be solved. Therefore, we modify the proof of [18, Th. 2.2] in the following Corollary to state a cheaper estimate under the stronger assumption (A4).

**Corollary 3** (Eigenvalue Estimate). *Under assumption (A4) and with Gauss nodes  $x_1, \dots, x_n$  the spectral radius of the Jacobian  $\mathbf{D}_{\hat{y}}\hat{f}(\hat{y})$  is estimated by*

$$\begin{aligned} \sigma_{\max}[\mathbf{D}_{\hat{y}}\hat{f}(\hat{y})] &\leq D_{\max}(\hat{\omega}) + a \quad \text{with} \\ D_{\max}(\hat{\omega}) &:= \max_{k=1, \dots, n} \left\{ \left\| \left[ \sum_{\ell=0}^K \hat{\alpha}_{\ell} \phi_{\ell}(x_k) \right]^{-1} \left[ \sum_{\ell=0}^K \hat{\beta}_{\ell} \phi_{\ell}(x_k) \right] \right\| \right\} < \infty. \end{aligned}$$

*Proof.* Assumption (A4) guarantees  $D_{\max}(\hat{\omega}) < \infty$  and

$$\sum_{\ell=0}^K [D_{\max}(\hat{\omega}) \hat{\alpha}_{\ell} \pm \hat{\beta}_{\ell}] \phi_{\ell}(x_k) \geq 0 \quad \text{for all } k = 1, \dots, n.$$

Thus, as discussed in (19) and (21), the symmetric matrices  $D_{\max}(\hat{\omega})\mathcal{P}(\hat{\alpha}) \pm \mathcal{P}(\hat{\beta})$  with entries

$$\left( D_{\max}(\hat{\omega})\mathcal{P}(\hat{\alpha}) \pm \mathcal{P}(\hat{\beta}) \right)_{i,j} = \left\langle \sum_{\ell=0}^K [D_{\max} \hat{\alpha}_{\ell} \pm \hat{\beta}_{\ell}] \phi_{\ell}, \phi_i \phi_j \right\rangle_{\mathbb{P}} \quad (27)$$

are positive semidefinite. The matrix  $\mathcal{P}(\hat{\beta})\mathcal{P}^{-1}(\hat{\alpha})$  is diagonalizable and the symmetric square root  $\mathcal{P}^{1/2}(\hat{\alpha})$  exists. Lemma 3.3 yields for the symmetric matrices with entries (27) the equivalence

$$\begin{aligned} &D_{\max}(\hat{\omega})\mathcal{P}(\hat{\alpha}) \pm \mathcal{P}(\hat{\beta}) \\ &= \mathcal{P}^{1/2}(\hat{\alpha}) \left[ D_{\max}(\hat{\omega}) \mathbb{I} \pm \mathcal{P}^{-1/2}(\hat{\alpha}) \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1/2}(\hat{\alpha}) \right] \mathcal{P}^{1/2}(\hat{\alpha}) \quad \text{positive semidefinite} \\ \Leftrightarrow &D_{\max}(\hat{\omega}) \geq \sigma_{\max} \left[ \mathcal{P}^{-1/2}(\hat{\alpha}) \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1/2}(\hat{\alpha}) \right] = \sigma_{\max} \left[ \mathcal{P}(\hat{\beta}) \mathcal{P}^{-1}(\hat{\alpha}) \right] = \sigma_{\max} \left[ D_{\hat{\nu}}(\hat{\omega}) \right]. \end{aligned} \quad (28)$$

Due to Corollary 2 and estimate (28) we obtain

$$\sigma_{\max}[\mathbf{D}_{\hat{y}}\hat{f}(\hat{y})] = \sigma_{\max}[D_{\hat{\nu}}(\hat{\omega})] + a \leq D_{\max}(\hat{\omega}) + a.$$

□

### 5.6. Choice of the Numerical Method

An equidistant space discretization with  $\Delta x > 0$  is used to divide the space interval  $[0, x_{\text{end}}]$  into  $N$  cells such that  $\Delta x N = x_{\text{end}}$  with centers  $x_j := (j + \frac{1}{2})\Delta x$  and edges  $x_{j-1/2} := j\Delta x$ . The discrete time steps are denoted by  $t_k := k\Delta t$  for  $k \in \mathbb{N}_0$ . A variable time discretization  $\Delta t > 0$  is chosen such that the CFL-condition

$$\text{CFL}(\hat{\omega}_j^k) := \max_{j=1,\dots,N} \left\{ \left| \hat{\Lambda}(\hat{\omega}_j^k) \right| \right\} \frac{\Delta t}{\Delta x} < 1 \quad (29)$$

holds. Cell averages at  $t_k$  are approximated by

$$\hat{\omega}_j^k \approx \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} \hat{\omega}(t_k, x) dx \quad \text{and} \quad \hat{y}_j^k \approx \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} \hat{y}(t_k, x) dx.$$

We consider a first order discretization

$$\hat{y}_j^{k+1} = \hat{y}_j^k - \frac{\Delta t}{\Delta x} \left( \hat{F}(\hat{y}_j^k, \hat{y}_{j+1}^k) - \hat{F}(\hat{y}_{j-1}^k, \hat{y}_j^k) \right)$$

with numerical flux function  $\hat{F}(\hat{y}_\ell, \hat{y}_r)$ . We may use the Roe flux ([RoeFlux](#)), which is deduced in Corollary 3, the Lax-Friedrichs or local Lax-Friedrichs flux, i.e.

$$\hat{F}(\hat{y}_\ell, \hat{y}_r) = \frac{1}{2} \left[ \hat{f}(\hat{y}_\ell) + \hat{f}(\hat{y}_r) \right] + \frac{\Delta t}{2\Delta x} (\hat{y}_\ell - \hat{y}_r), \quad (\text{LaxFriedrichs})$$

$$\hat{F}(\hat{y}_\ell, \hat{y}_r) = \frac{1}{2} \left[ \hat{f}(\hat{y}_\ell) + \hat{f}(\hat{y}_r) \right] + \frac{1}{2} \max_{j=\ell, r} \left\{ \sigma \left\{ D_{\hat{y}} \hat{f}(\hat{y}) \Big|_{\hat{y}=\hat{y}_j} \right\} \right\} (\hat{y}_\ell - \hat{y}_r). \quad (\text{LocalLaxFriedrichs})$$

An appropriate choice depends on the application and on numerical cost, which are strongly influenced by the eigenvalue decomposition of the Jacobian ( $D\hat{\mathcal{R}}_K$ ), as discussed in Section 5.4. The local Lax-Friedrichs flux ([LaxFriedrichs](#)) does not need the eigenvalue decomposition of the Jacobian ( $D\hat{\mathcal{R}}_K$ ). The eigenvalue estimate in Corollary 3 is sufficient to satisfy the CFL-condition (29). Thus, this choice would lead to a cheap numerical method. However, the classical Lax-Friedrichs flux is quite dissipative and shocks would not be resolved properly.

For both the local Lax-Friedrichs flux ([LocalLaxFriedrichs](#)) and the Roe flux ([RoeFlux](#)) the eigenvalue decomposition of the Jacobian ( $D\hat{\mathcal{R}}_K$ ) must be determined. It is important to note that the average in the Roe matrix (23) is with respect to the Roe variables  $\hat{\omega}$ . This causes a noteworthy computational overhead, since also the eigenvalue decomposition of the matrix  $\mathcal{P}_\nu(\bar{\omega})$  has to be calculated. Therefore, we choose the time discretization  $\Delta t$  such that the CFL-condition  $\text{CFL}(\bar{\omega}_j^k) < 1$  is satisfied by the Roe average and we do not calculate the eigenvalue decomposition of the Jacobian ( $D\hat{\mathcal{R}}_K$ ). This choice is justified by the properties (Roe 2) and (Roe 3). In this case, computational cost for the numerical flux functions ([LocalLaxFriedrichs](#)) and ([RoeFlux](#)) are similar.

**Remark 3.** An open problem is the use of high order schemes. In principle, our system may be solved numerically with standard methods due to the hyperbolic and conservative formulation.

In previous works [24, 16, 41, 21, 10] Roe-type and MUSCL-schemes have proven successful and adaptivity in the stochastic space decreases computational cost significantly.

For small gPC truncations  $K$  we have also employed a high order Runge-Kutta discontinuous Galerkin (RKDG) method, implemented in the MULTIWAVE [42, 43] software-package. This RKDG scheme uses polynomial elements with the local Lax-Friedrichs flux and the minmod limiter [44]. As time discretization method we have used a strong stability-preserving Runge-Kutta method. The performance is enhanced by a local multi-resolution based grid adaptation [45]. For larger gPC truncations this high order RKDG numerical solver can *violate* the positive definiteness assumption (A2). Probably, a special limiter designed for the intrusive formulation has to be developed that preserves hyperbolicity.

## 6. Numerical Results for Isothermal Euler Equations

We illustrate the hyperbolic systems by comparing the strong and the intrusive formulation. Furthermore, we show that the proposed Roe flux is less dissipative. Therefore, the matter of choice are the numerical flux functions (`LocalLaxFriedrichs`) and (`RoeFlux`) with the relatively high CFL-conditions  $\text{CFL}(\hat{\omega}_j^k) = 0.99$ ,  $\text{CFL}(\bar{\omega}_j^k) = 0.99$  to exclude the possibility that the system or the numerical method seems stable only due to large artificial numerical viscosities. Hyperbolicity of the system ( $\hat{\mathcal{R}}_K$ ) is verified numerically by condition (A4).

### 6.1. Strong Formulation

First, we consider the strong formulation ( $\mathcal{C}(\xi)$ ) and focus on a **shock tube problem**, found e.g. in [27, 46]. For a given left  $\mathbf{y}_\ell$  and right state  $\mathbf{y}_r$  with  $\rho_\ell \geq \rho_r > 0$  and  $\mathbf{q}_\ell = \mathbf{q}_r = 0$ , the solution consists of a rarefaction wave, moving with negative speed, and a shock wave with positive speed. Both waves are connected by an intermediate state  $y_m$ . The entropy solution, satisfying the [46, Lax Entropy Condition], with initial values

$$y(0, x) := \begin{cases} \mathbf{y}_\ell, & x < 0 \\ \mathbf{y}_r, & x > 0 \end{cases} \quad \text{is solved by} \quad y(t, x) = \begin{cases} \mathbf{y}_\ell, & x < t\lambda^-(\mathbf{y}_\ell), \\ y_{\text{rf}}(t, x), & t\lambda^-(\mathbf{y}_\ell) \leq x < t\lambda^-(y_m), \\ y_m, & t\lambda^-(y_m) \leq x < ts, \\ \mathbf{y}_r, & ts < x \end{cases} \quad (30)$$

$$\text{with rarefaction wave} \quad y_{\text{rf}}(t, x) = \rho_\ell \exp\left(\frac{\rho_\ell - x/t}{a}\right) \begin{pmatrix} 1 \\ x/t + a \end{pmatrix},$$

$$\text{intermediate state} \quad (\rho_m, q_m)^T = \mathcal{R}^-(\rho_m; \mathbf{y}_\ell) \quad \text{such that} \quad \mathcal{R}^-(\rho_m; \mathbf{y}_\ell) = \mathcal{S}^+(\rho_m; \mathbf{y}_r),$$

$$\text{and shock speed} \quad s = \frac{\mathbf{q}_r}{\rho_r} + a\left(\frac{\rho_m}{\rho_r}\right)^{1/2}.$$

The integral curve  $\mathcal{R}^-$  and the Hugoniot locus  $\mathcal{S}^+$  are

$$\begin{aligned} \mathcal{R}^-(\cdot; \mathbf{y}) : \mathbb{R}^+ &\rightarrow \mathbb{R}, \quad \theta \mapsto \mathcal{R}^-(\theta; \mathbf{y}) = \left( \theta, \theta \frac{\mathbf{q}}{\rho} - a\theta \ln\left(\frac{\theta}{\rho}\right) \right)^T, \\ \mathcal{S}^+(\cdot; \mathbf{y}) : \mathbb{R}^+ &\rightarrow \mathbb{R}, \quad \theta \mapsto \mathcal{S}^+(\theta; \mathbf{y}) = \left( \theta, \theta \frac{\mathbf{q}}{\rho} + a(\theta - \rho)\left(\frac{\theta}{\rho}\right)^{1/2} \right)^T. \end{aligned}$$

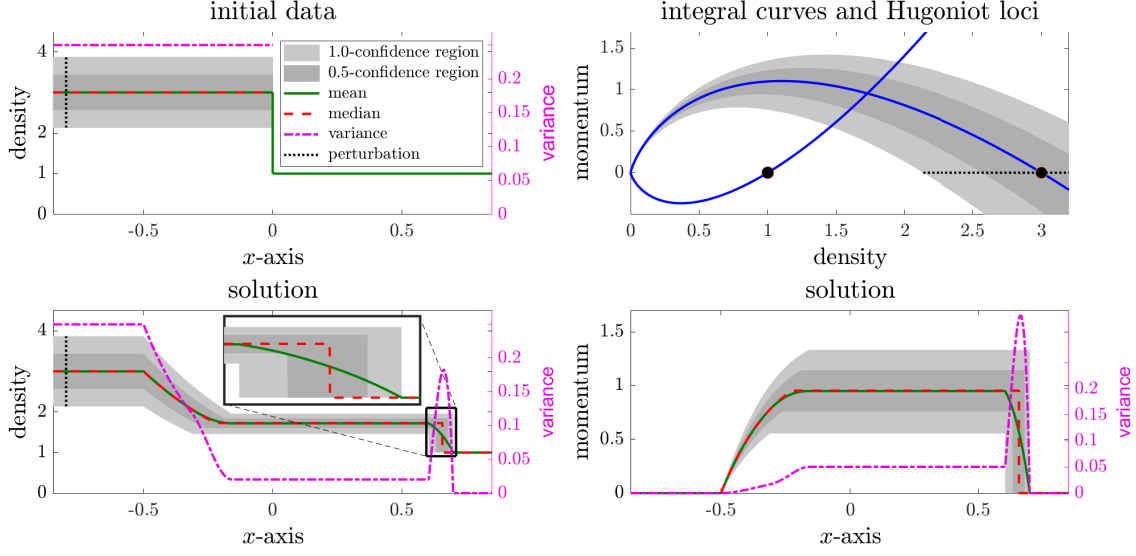


Figure 1: Shock tube problem for  $\mathbf{y}_\ell(\xi) := (3 + 0.5\phi_1(\xi), 0)^\top$ ,  $\xi \sim \mathcal{U}(-1, 1)$  and  $\mathbf{y}_r(\xi) := (1, 0)^\top$ ,  $t = 0.5$

In Figure 1 we consider a Riemann problem with uniformly distributed initial data

$$\boldsymbol{\rho}_\ell(\xi) \sim \mathcal{U}(\mu - \sqrt{3}\sigma, \mu + \sqrt{3}\sigma) \quad \text{such that} \quad \mathbb{E}[\boldsymbol{\rho}_\ell(\xi)] = \mu \quad \text{and} \quad \text{Var}[\boldsymbol{\rho}_\ell(\xi)] = \sigma^2.$$

As illustrated by the 1.0-confidence region, which contains  $\mathbb{P}$ -a.s. all realisations, densities are in fact  $\mathbb{P}$ -a.s. strictly positive. Then, the strong formulation  $(\mathcal{C}(\xi))$  is defined under assumption (A1) as the unique intersection of integral curves and Hugoniot loci. Due to the symmetric perturbation the mean of the initial data coincides with the median.

We distinguish the solution into *quantiles* (median, confidence region) and *moments* (mean, variance). The edges of the 1.0-confidence region are the 0- and 1-quantile. Therefore, *quantiles* have the same smoothness properties as the deterministic solution, i.e. quantiles are discontinuous, too. For *moments*, however, the expectation operator causes a smoothing. Therefore, the mean and variance are smooth even through the shock.

## 6.2. Intrusive Formulation

We consider again the shock tube problem, but now under the intrusive formulation  $(\hat{\mathcal{R}}_K)$  with flux function (17). The Roe flux with  $\Delta x = 2.5 \cdot 10^{-4}$  and gPC expansion  $K = 6$  is illustrated in Figure 2. For all uniform cells we have generated  $10^4$  samples to estimate quantiles with MATLAB. Moments are immediately given by gPC modes. While the rarefaction wave and the intermediate state seem similar to the strong formulation in Figure 1, there is no longer a smooth expected shock.

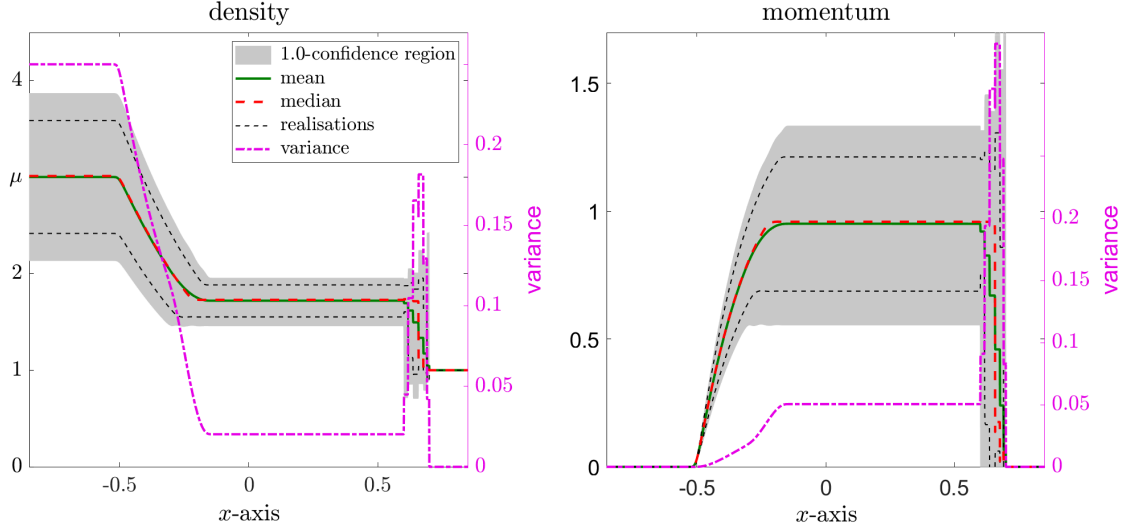


Figure 2: Intrusive formulation ( $\hat{\mathcal{R}}_K$ ) for  $\mathbf{y}_\ell(\xi) := (3 + 0.5\phi_1(\xi), 0)^\top$ ,  $\xi \sim \mathcal{U}(-1, 1)$  and  $\mathbf{y}_r(\xi) := (1, 0)^\top$ ,  $t = 0.5$

Similar results have been observed for Burgers' equation in [29, Sec. 6.2]. We analyze the regularity, which is determined by the gPC expansion, for the intrusive formulation likewise. To investigate the structure of the waves observed in Figure 2, we compare in Figure 3 the solutions for  $K = 0, \dots, 8$ . The deterministic solution, which corresponds to  $K = 0$  is shown in the first subplot. The last subplot shows a standard Monte-Carlo simulation with  $10^5$  samples. The solution corresponding to a particular sample is again given in equation (30). We observe that the shock, moving to the right, forms in the intrusive formulation ( $\hat{\mathcal{R}}_K$ ) a wave structure similar to a wave package. To be detailed, the initial shock splits into several waves all moving with slightly different speeds. Altogether  $2(K + 1)$  waves can emerge from the Riemann problem. Note that when plotting only one component, e.g. the mean, not all waves may be observable. This behaviour is also expected to occur for a larger  $K$ , however on a smaller scale.

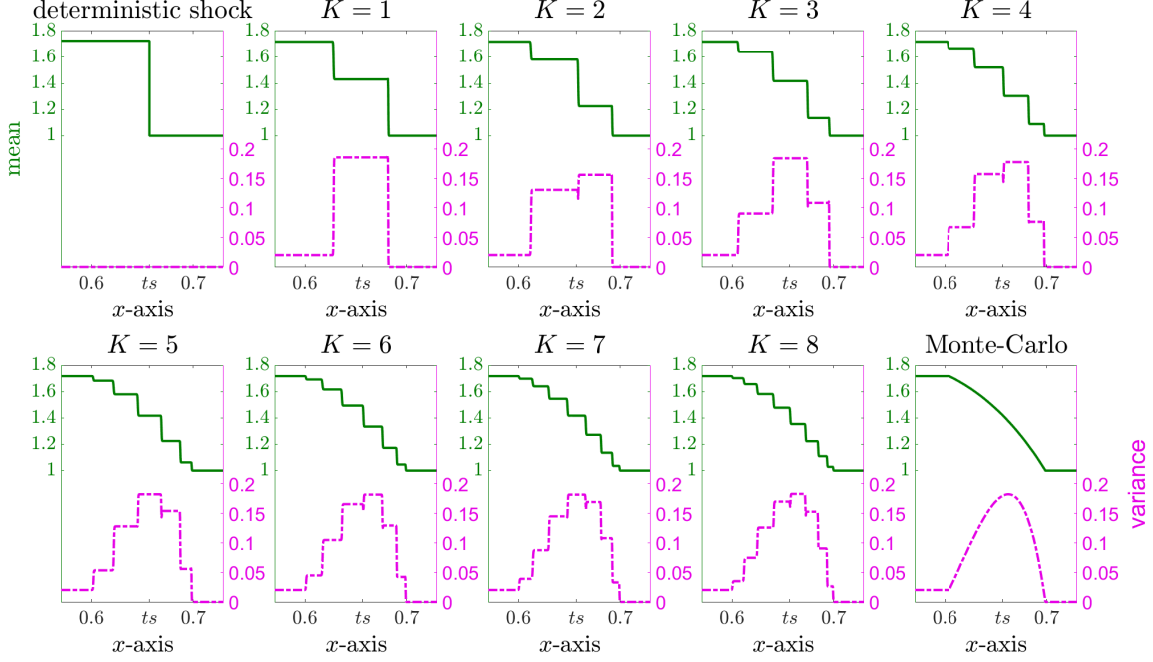


Figure 3: Zoom on shock at  $t = 0.5$  with initial data  $\mathbf{y}_\ell(\xi) := (3 + 0.5\phi_1(\xi), 0)^\top$ ,  $\xi \sim \mathcal{U}(-1, 1)$  and  $\mathbf{y}_r(\xi) := (1, 0)^\top$

The stochastic Galerkin approach aims to minimize the mean squared error (MSE)

$$\text{MSE}[\text{ref}] := \mathbb{E} \left[ \left( \text{ref} - \mathcal{G}_K[\text{ref}] \right)^2 \right],$$

where  $\text{ref}(\xi)$  denotes a reference solution. Additionally, we state errors for the mean and the variance

$$E_K^{(\mathbb{E})}[\text{ref}] := \left| \mathbb{E}[\text{ref}] - \mathbb{E}[\mathcal{G}_K[\text{ref}]] \right|, \quad E_K^{(\mathbb{V})}[\text{ref}] := \left| \text{Var}[\text{ref}] - \text{Var}[\mathcal{G}_K[\text{ref}]] \right|.$$

These errors are approximated for each fixed point in space with a Monte-Carlo method with  $10^5$  samples to obtain the estimates  $\widehat{\text{MSE}}$ ,  $\hat{E}_K^{(\mathbb{E})}$ ,  $\hat{E}_K^{(\mathbb{V})}$ . The Monte-Carlo method causes errors in the order of  $10^{-5}$  to  $10^{-3}$ , but is trustworthy since the reference samples are given according to equation (30). We use the  $L^1$ - and  $L^\infty$ -norms  $\int |\cdot| dx$  and  $\sup_x |\cdot|$  to obtain one value in space for the rarefaction wave for  $x \in (-0.75, 0)$  and the shock for  $x \in (0.55, 0.75)$ .

The rarefaction wave, moving to the left, has a similar shape as in the strong formulation. Indeed, Table 1 hints that the truncation  $K = 2$  is enough to obtain values which are in the magnitude of Monte-Carlo errors. We have already seen in Figure 1 and Figure 3 that a deterministic shock discontinuity may result in a stochastic setting as smooth expected value which is approximated by a wave package. We observe from Table 2 a slow convergence against the reference solution.

rarefaction wave and $L_x^1$ -norm									
gPC truncation: $K$	0	1	2	3	4	5	6	7	8
$\hat{E}_K^{(\mathbb{E})}$	4.02	1.77	1.70	1.71	1.71	1.72	1.72	1.72	1.72
$\hat{E}_K^{(\mathbb{V})}$	108.35	0.53	0.29	0.30	0.29	0.29	0.30	0.29	0.30
$\widehat{\text{MSE}}$	108.46	0.11	0.03	0.02	0.02	0.02	0.02	0.02	0.02

rarefaction wave and $L_x^\infty$ -norm									
gPC truncation: $K$	0	1	2	3	4	5	6	7	8
$\hat{E}_K^{(\mathbb{E})}$	5.97	3.11	3.17	3.19	3.20	3.22	3.23	3.23	3.24
$\hat{E}_K^{(\mathbb{V})}$	25.01	0.67	0.54	0.55	0.55	0.55	0.55	0.55	0.55
$\widehat{\text{MSE}}$	25.01	0.13	0.11	0.11	0.11	0.11	0.11	0.11	0.11

Table 1: Monte-Carlo estimates  $\hat{E}_K^{(\mathbb{E})}[\text{ref}] \approx E_K^{(\mathbb{E})}[\text{ref}]$ ,  $\hat{E}_K^{(\mathbb{V})}[\text{ref}] \approx E_K^{(\mathbb{V})}[\text{ref}]$ ,  $\widehat{\text{MSE}}[\text{ref}] \approx \text{MSE}[\text{ref}]$  for the rarefaction wave  $x \in (-0.75, 0)$  illustrated in Figure 2; units in  $10^{-3}$  for  $L_x^1$ -norm,  $10^{-2}$  for  $L_x^\infty$ -norm

shock wave and $L_x^1$ -norm									
gPC truncation: $K$	0	1	2	3	4	5	6	7	8
$\hat{E}_K^{(\mathbb{E})}$	17.20	8.39	6.41	5.02	4.11	3.43	2.91	2.52	2.22
$\hat{E}_K^{(\mathbb{V})}$	12.63	3.75	2.85	2.18	1.85	1.54	1.31	1.13	1.00
$\widehat{\text{MSE}}$	16.82	5.44	3.98	3.20	2.67	2.28	1.98	1.74	1.55

shock wave and $L_x^\infty$ -norm									
gPC truncation: $K$	0	1	2	3	4	5	6	7	8
$\hat{E}_K^{(\mathbb{E})}$	41.98	23.64	17.41	14.04	10.82	8.92	7.65	6.59	5.85
$\hat{E}_K^{(\mathbb{V})}$	18.22	13.61	8.32	7.16	5.97	4.83	4.04	3.46	3.04
$\widehat{\text{MSE}}$	35.61	16.89	11.17	9.17	7.45	6.20	5.11	4.42	3.98

Table 2: Monte-Carlo estimates  $\hat{E}_K^{(\mathbb{E})}[\text{ref}] \approx E_K^{(\mathbb{E})}[\text{ref}]$ ,  $\hat{E}_K^{(\mathbb{V})}[\text{ref}] \approx E_K^{(\mathbb{V})}[\text{ref}]$ ,  $\widehat{\text{MSE}}[\text{ref}] \approx \text{MSE}[\text{ref}]$  for the shock  $x \in (0.55, 0.75)$  illustrated in Figure 3; units in  $10^{-3}$  for  $L_x^1$ -norm,  $10^{-2}$  for  $L_x^\infty$ -norm

Next, we compare the formulations  $(\hat{\mathcal{C}}_K)$  and  $(\hat{\mathcal{R}}_K)$  and show the gPC modes  $\hat{\rho}(0.5, x)$ . In the left subplot, we revisit the case  $K = 2$  from Figure 3, which is used as reference solution (black line). In the right subplot, we use the initial values  $\hat{\mathbf{y}}_\ell := (3, -0.3, -0.03, 0, 14.7, 9.8)^\top$  and  $\hat{\mathbf{y}}_r := \hat{\mathbf{y}}_\ell/3$ . For these values the system  $(\hat{\mathcal{C}}_K)$  has complex eigenvalues, see [33]. This choice yields realisations in both sub- and supersonic regimes. If the complex part of the eigenvalues of the Jacobian from  $(\hat{\mathcal{C}}_K)$  is neglected, one may apply the local Lax-Friedrichs flux. Still the CFL-condition may not be satisfied.

The possibility of instabilities, when complex eigenvalues occur, is discussed in [17, Sec. 4]. We include the simulation *without* Roe variable transform *for comparison only*.

In contrast, the Roe variable based system ( $\hat{\mathcal{R}}_K$ ) remains hyperbolic. Both formulations are solved with the local Lax-Friedrichs flux (green, blue). If the local Lax-Friedrichs flux is applied, we observe for slow velocities (left subplot), almost no difference between both formulations. In contrast for larger velocities (right subplot), when the influence of the nonlinear term  $q^2/\rho$  is more important, there are differences. In both cases, the Roe flux (RoeFlux) yields a good resolution of all waves. Figure 5 shows a zoom on the results of Figure 4 for a better comparison. Different space discretizations are shown, which confirm advantages of the Roe flux.

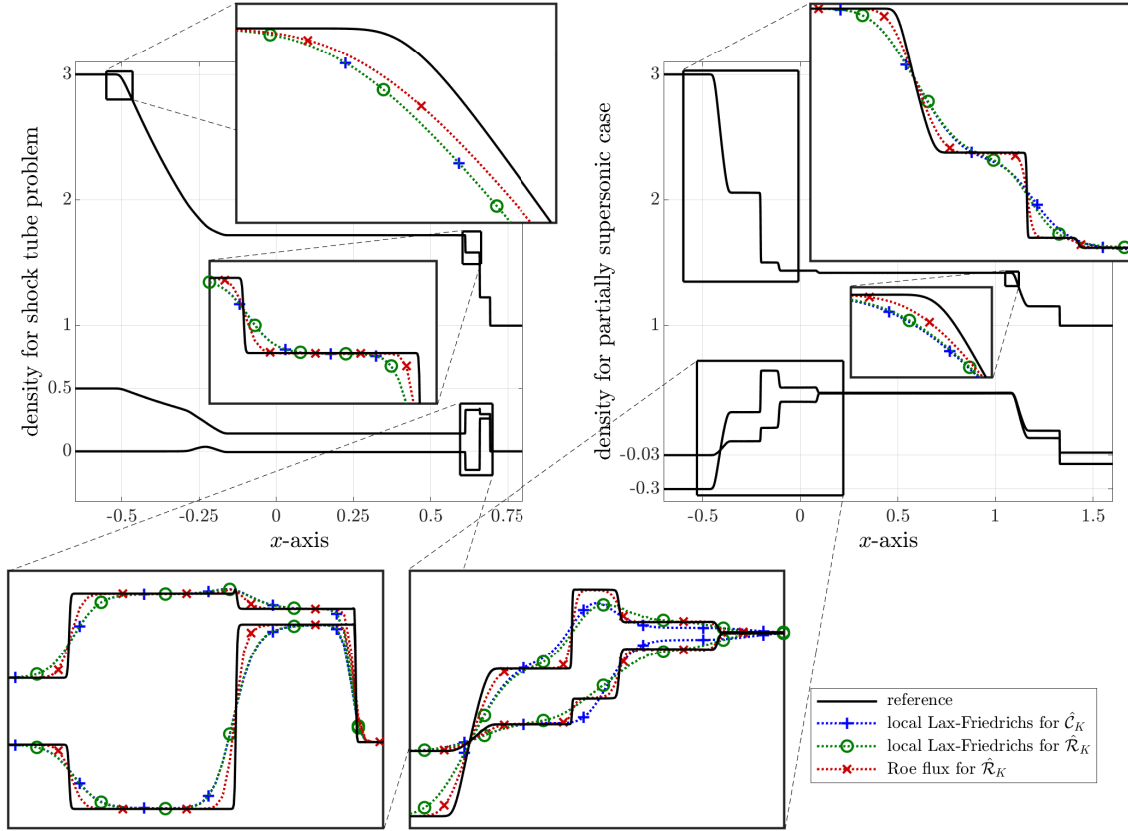


Figure 4: left:  $\hat{\mathbf{y}}_\ell := (3, 0.5, 0, 0, 0, 0)^T$ ,  $\hat{\mathbf{y}}_r := (1, 0, 0, 0, 0, 0)^T$ ; right:  $\hat{\mathbf{y}}_\ell := (3, -0.3, -0.03, 0, 14.7, 9.8)^T$ ;  $\hat{\mathbf{y}}_r := \hat{\mathbf{y}}_\ell/3$ ;  $t = 0.5$



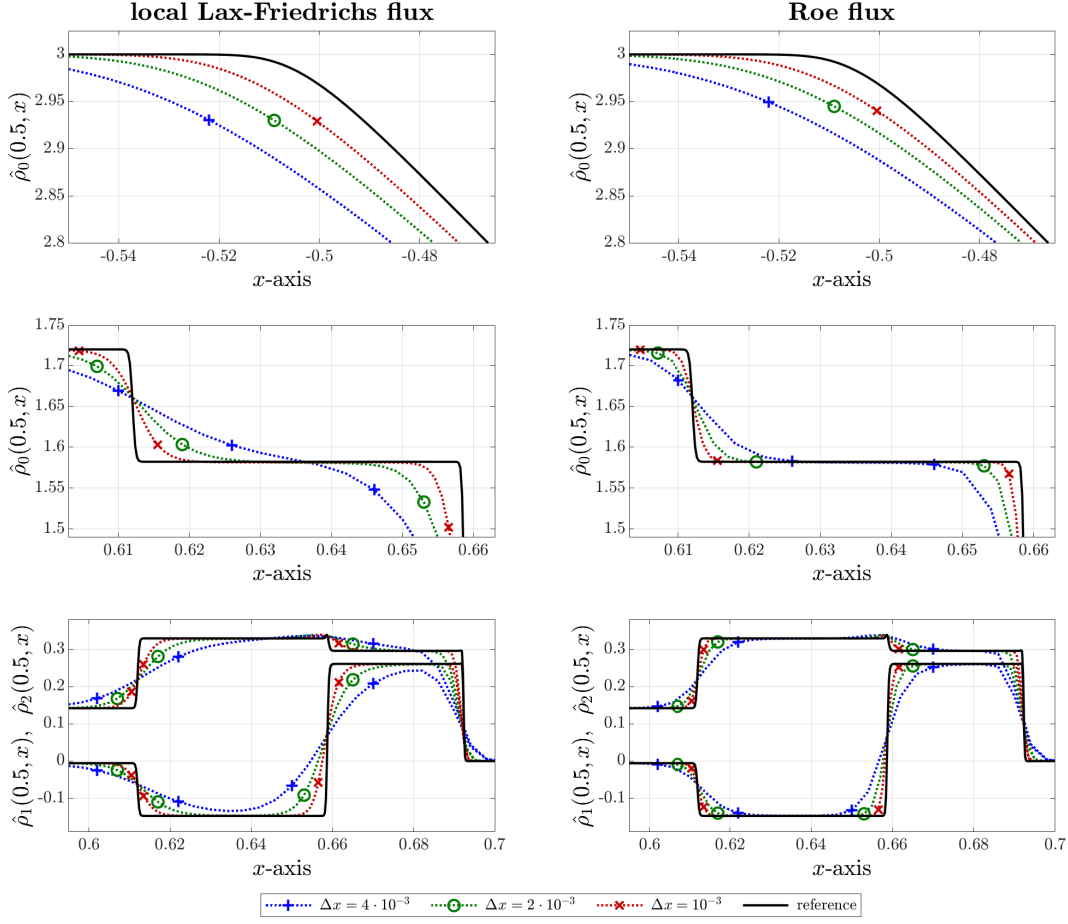


Figure 5: Local Lax-Friedrichs flux versus Roe flux for Riemann problem with  $\mathbf{y}_\ell(\xi) := (3 + 0.5\phi_1(\xi), 0)^\top$  and  $\mathbf{y}_r(\xi) := (1, 0)^\top$  in  $t = 0.5$

## 7. Summary

Stochastic Galerkin formulations for the  $p$ -system have been introduced and conditions that preserve hyperbolicity have been presented. In particular for isothermal Euler equations, the resulting system can be guaranteed hyperbolic for arbitrary polynomial chaos expansions. However, if positive physical quantities are not represented appropriately, hyperbolicity of the stochastic Galerkin formulation is lost. The hyperbolic structure of the presented systems has been illustrated numerically and a Roe flux has been introduced. A first-order solver preserving the appropriate representation of the positive quantities has been implemented.

## Appendix

Table 3 summarizes the notation and eigendecompositions.

	symbol	
	deterministic	stochastic
conserved variables	$y = (\rho, q)^T$	$\hat{y} = (\hat{\rho}, \hat{q})^T$
Roe variables	$\omega = (\alpha, \beta)^T$	$\hat{\omega} = (\hat{\alpha}, \hat{\beta})^T$
Roe transform	$\mathcal{Y}(\omega) = y$	$\hat{\mathcal{Y}}(\hat{\omega}) = \hat{y}$
flux function	$f(y)$	$\hat{f}(\hat{y})$
	$F(\omega) = f(\mathcal{Y}(\omega))$	$\hat{F}(\hat{\omega}) = \hat{f}(\hat{\mathcal{Y}}(\hat{\omega}))$
pressure law	$p(\rho)$	
	$\pi(\alpha) = p(\alpha^2)$	$\hat{\pi}(\hat{\alpha})$
parameters	$a, g > 0$	
admissible set	$\rho > 0$	$\mathbb{A}_K$
Lebesgue space		$\mathbb{L}^2(\Omega, \mathbb{P})$
orthogonal subspace		$\mathcal{S}_K$
gPC basis		$\phi_k(\xi), k = 0, \dots, K$
projection operator		$\mathcal{G}_K[y], \hat{\mathcal{G}}_K[y, y]$
cell averages		$\hat{\mathbf{y}}, \hat{\boldsymbol{\omega}}$
Roe matrix		$\bar{A}(\hat{\boldsymbol{\omega}}_\ell, \hat{\boldsymbol{\omega}}_r)$
numerical flux		$\hat{\mathbf{F}}(\hat{\mathbf{y}}_\ell, \hat{\mathbf{y}}_r)$
Galerkin product		$(*), \mathcal{P}(\hat{\alpha}) = V(\hat{\alpha})D_{\mathcal{P}}(\hat{\alpha})V^T(\hat{\alpha})$
velocity	$u(y) = q/\rho, \nu(\omega) = \beta/\alpha$	$\mathcal{P}_{\hat{\nu}}(\hat{\omega}) = \mathcal{P}(\hat{\beta})\mathcal{P}^{-1}(\hat{\alpha})$
		$\mathcal{P}_2(\hat{\omega}) = \mathcal{P}^{-1/2}(\hat{\alpha})\mathcal{P}(\hat{\beta})\mathcal{P}^{-1/2}(\hat{\alpha})$
eigendecompositions		$\mathcal{P}_{\hat{\nu}}(\hat{\omega}) = Q(\hat{\omega})D_{\hat{\nu}}(\hat{\omega})Q^{-1}(\hat{\omega})$
		$\hat{\pi}(\hat{\alpha}) = VD_{\hat{\pi}}(\hat{\alpha})V^T$
	$D_y f(y) = T(y)\Lambda(y)T^{-1}(y)$	$D_{\hat{y}}\hat{f}(\hat{y}) = [\mathcal{V}\hat{T}(\hat{\omega})]\hat{\Lambda}(\hat{\omega})[\mathcal{V}\hat{T}(\hat{\omega})]^{-1}$
		$D_{\hat{y}}\hat{f}(\hat{y}) = [\mathcal{Q}(\hat{\omega})\hat{T}(\hat{\omega})]\hat{\Lambda}(\hat{\omega})[\mathcal{Q}(\hat{\omega})\hat{T}(\hat{\omega})]^{-1}$
		$\hat{A} = \hat{\mathcal{X}}\hat{\mathcal{D}}\hat{\mathcal{X}}^{-1}$

Table 3: Notation and eigendecompositions

Table 4 summarizes the relation of discussed systems and the corresponding assumptions satisfying

$$(A3) \Rightarrow (A4) \Rightarrow (A2), \quad \text{but} \quad (A1) \not\Rightarrow (A2), \quad (A1) \not\Rightarrow (A3), \quad (A1) \not\Rightarrow (A4).$$

The stochastic systems for conserved ( $\mathcal{C}(\xi)$ ) and Roe variables ( $\mathcal{R}(\xi)$ ) require assumption (A1). The truncated gPC expansion  $\mathcal{C}_K(\xi)$  does not lead to a hyperbolic intrusive formulation. The truncated expansion in terms of Roe variables  $\mathcal{R}_K(\xi)$ , however, results in the hyperbolic system  $\hat{\mathcal{R}}_K$  if assumptions (A2) holds. The stronger assumptions (A3) is redundant. It implies the assumptions (A2) and (A4), the latter one is used for an eigenvalue estimate and for checking hyperbolicity numerically.

system	usage	assumptions for hyperbolicity		reference
$\mathcal{C}(\xi)$	strong formulation	(A1)	$\rho(\xi) > 0$ $\mathbb{P}$ -a.s.	Section 3
$\mathcal{R}(\xi)$			$\alpha(\xi) > 0$ $\mathbb{P}$ -a.s.	
$\mathcal{C}_K(\xi)$	truncated expansion	(A2)	not a hyperbolic system	
$\mathcal{R}_K(\xi)$			$\mathcal{P}(\hat{\alpha})$ strictly positive definite	Theorem 3.2
$\hat{\mathcal{R}}_K$	gPC modes			
		not required		
		(A3)	$\mathcal{G}_K[\rho](\xi) > 0, \mathcal{G}_K[\alpha](\xi) > 0$ $\mathbb{P}$ -a.s.	Section 4
numerical discretization		assumption for eigenvalue estimate		
		(A4)	$\mathcal{G}_K[\alpha](x_k) > 0 \forall k = 1, \dots, n$	Corollary 3

Table 4: Summary of systems and assumptions

## References

- [1] W. Schoutens, Stochastic processes and orthogonal polynomials, Lecture Notes in Statistics, Springer, New York, 2000. [doi:10.1007/978-1-4612-1170-9](https://doi.org/10.1007/978-1-4612-1170-9).
- [2] N. Wiener, The homogeneous chaos, American Journal of Mathematics 60 (4) (1938) 897–936.
- [3] R. H. Cameron, W. T. Martin, The orthogonal development of non-linear functionals in series of Fourier-Hermite functionals, Annals of Mathematics 48 (2) (1947) 385–392.
- [4] R. G. Ghanem, P. D. Spanos, Stochastic finite elements: A Spectral Approach, 1st Edition, Springer, New York, 1991. [doi:10.1007/978-1-4612-3094-6](https://doi.org/10.1007/978-1-4612-3094-6).
- [5] D. Xiu, G. E. Karniadakis, The Wiener-Askey polynomial chaos for stochastic differential equations, SIAM Journal on Scientific Computing 24 (2002) 619–644.
- [6] D. Xiu, J. Shen, Efficient stochastic Galerkin methods for random diffusion equations, Journal of Computational Physics 228 (2) (2009) 266–281.

- [7] M. Eigel, C. J. Gittelsohn, C. Schwab, E. Zander, Adaptive stochastic Galerkin FEM, *Computer Methods in Applied Mechanics and Engineering* 270 (Supplement C) (2014) 247–269.
- [8] R. Shu, J. Hu, S. Jin, A stochastic Galerkin method for the Boltzmann equation with multi-dimensional random inputs using sparse wavelet bases, *Numerical Mathematics: Theory, Methods and Applications* 10 (2) (2017) 465–488.
- [9] J. Hu, S. Jin, A stochastic Galerkin method for the Boltzmann equation with uncertainty, *Journal of Computational Physics* 315 (2016) 150–168.
- [10] S. Jin, L. Pareschi, Uncertainty quantification for hyperbolic and kinetic equations, SEMA SIMAI Springer Series, Springer International Publishing, Cham, Switzerland, 2017. doi:[10.1007/978-3-319-67110-9](https://doi.org/10.1007/978-3-319-67110-9).
- [11] O. P. L. Maître, O. M. Knio, Spectral Methods for uncertainty quantification, 1st Edition, Springer Netherlands, 2010. doi:[10.1007/978-90-481-3520-2](https://doi.org/10.1007/978-90-481-3520-2).
- [12] Q.-Y. Chen, D. Gottlieb, J. S. Hesthaven, Uncertainty analysis for the steady-state flows in a dual throat nozzle, *Journal of Computational Physics* 204 (2005) 378–398.
- [13] D. Gottlieb, D. Xiu, Galerkin method for wave equations with uncertain coefficients, *Communications in Computational Physics* 3 (2) (2008) 505–518.
- [14] R. Pulch, D. Xiu, Generalised polynomial chaos for a class of linear conservation laws, *Journal of Scientific Computing* 51 (2012) 293–312.
- [15] S. Jin, D. Xiu, X. Zhu, A well-balanced stochastic Galerkin method for scalar hyperbolic balance laws with random inputs, *Journal of Scientific Computing* 67 (2016) 1198–1218.
- [16] J. Tryoen, O. P. L. Maître, M. Ndjinga, A. Ern, Intrusive Galerkin methods with upwinding for uncertain nonlinear hyperbolic systems, *Journal of Computational Physics* 229 (2010) 6485–6511.
- [17] B. Després, G. Poëtte, D. Lucor, Robust uncertainty propagation in systems of conservation laws with the entropy closure method 92 (2013) 105–149.
- [18] K. Wu, H. Tang, D. Xiu, A stochastic Galerkin method for first-order quasilinear hyperbolic systems with uncertainty, *Journal of Computational Physics* 345 (2017) 224–244.
- [19] A. Chertock, S. Jin, A. Kurganov, An operator splitting based stochastic Galerkin method for the one-dimensional compressible Euler equations with uncertainty, [www.math.wisc.edu](http://www.math.wisc.edu), preprint (2015) (2015).
- [20] A. Chertock, S. Jin, A. Kurganov, A well-balanced operator splitting based stochastic Galerkin method for the one-dimensional Saint-Venant system with uncertainty, [www.math.wisc.edu](http://www.math.wisc.edu), preprint (2015) (2015).
- [21] R. Abgrall, P. Congedo, G. Geraci, G. Iaccarino, An adaptive multiresolution semi-intrusive scheme for UQ in compressible fluid problems, *International Journal for Numerical Methods in Fluids* 78 (2015) 595–637.

- [22] L. Schlachter, F. Schneider, A hyperbolicity-preserving stochastic Galerkin approximation for uncertain hyperbolic systems of equations, *Journal of Computational Physics* 375 (2018) 80–98.
- [23] B. Després, G. Poëtte, D. Lucor, Uncertainty quantification for systems of conservation laws, *Journal of Computational Physics* 228 (2009) 2443–2467.
- [24] P. Pettersson, G. Iaccarino, J. Nordström, A stochastic Galerkin method for the Euler equations with Roe variable transformation, *Journal of Computational Physics* 257 (2014) 481–500.
- [25] R. V. Field, M. Grigoriu, On the accuracy of the polynomial chaos approximation, *Probabilistic Engineering Mechanics* 19 (1) (2004) 65–80.
- [26] P. L. Roe, Approximate Riemann solvers, parameter vectors, and difference schemes, *Journal of Computational Physics* 43 (1981) 357–372.
- [27] R. J. Leveque, *Numerical Methods for Conservation Laws*, 2nd Edition, *Lectures in Mathematics*. ETH Zürich, Birkhäuser Basel, 1992. doi:10.1007/978-3-0348-8629-1.
- [28] B. J. Deusschere, H. N. Najm, P. P. Pébay, O. M. Knio, R. G. Ghanem, O. P. L. Maître, Numerical challenges in the use of polynomial chaos representations for stochastic processes, *SIAM Journal on Scientific Computing* 26 (2) (2004) 698–719.
- [29] P. Pettersson, G. Iaccarino, J. Nordström, *Polynomial chaos methods for hyperbolic partial differential equations*, Springer International Publishing, Switzerland, 2015. doi:10.1007/978-3-319-10714-1.
- [30] D. Xiu, *Numerical methods for stochastic computations*, Princeton University Press, Princeton, 2010. doi:978-0691142128.
- [31] T. J. Sullivan, *Introduction to uncertainty quantification*, 1st Edition, *Texts in Applied Mathematics*, Springer, Switzerland, 2015. doi:10.1007/978-3-319-23395-6.
- [32] M. D. Gunzburger, C. G. Webster, G. Zhang, Stochastic finite element methods for partial differential equations with random input data, *Acta Numerica* 23 (2014) 521–650. doi:10.1017/S096249214000075.
- [33] S. Jin, R. Shu, A study of hyperbolicity of kinetic stochastic Galerkin system for the isentropic Euler equations with uncertainty, *Chinese Annals of Mathematics, Series B* 40 (2019) 765–780.
- [34] G. H. Golub, C. F. van Loan, *Matrix Computations*, 3rd Edition, *Johns Hopkins Series in the Mathematical Sciences*, Johns Hopkins University Press, Baltimore, 1996.
- [35] D. Xiu, G. E. Karniadakis, Supersensitivity due to uncertain boundary conditions, *International Journal for Numerical Methods in Engineering* 61 (2004) 2114–2138. doi:10.1002/nme.1152.
- [36] O. G. Ernst, A. Mugler, H.-J. Starkloff, E. Ullmann, On the convergence of generalized polynomial chaos expansions, *ESAIM: Mathematical Modelling and Numerical Analysis* 46 (2012) 317–339.

- [37] C. M. Dafermos, Hyperbolic conservation laws in continuum physics, 3rd Edition, Vol. 325 of A series of comprehensive studies in mathematics, Springer-Verlag Berlin Heidelberg, 2010. [doi:10.1007/978-3-642-04048-1](https://doi.org/10.1007/978-3-642-04048-1).
- [38] A. Bressan, Hyperbolic systems of conservation laws: The one dimensional Cauchy problem, Oxford Lecture Series in Mathematics and its Applications, Oxford University Press, New York, 2005.
- [39] T.-P. Liu, The Riemann problem for general  $2 \times 2$  conservation laws, Transactions of the American Mathematical Society 199 (1974) 89–112.
- [40] T.-P. Liu, The Riemann problem for general systems of conservation laws, Journal of Differential Equations 18 (1) (1975) 218–234.
- [41] J. Tryoen, O. P. L. Maître, M. Ndjinga, A. Ern, Roe solver with entropy corrector for uncertain hyperbolic systems, Journal of Computational and Applied Mathematics 235 (2010) 491–506.
- [42] N. Gerhard, F. Iacono, G. May, S. Müller, R. Schäfer, A high-order discontinuous Galerkin discretization with multiwavelet-based grid adaptation for compressible flows, Journal of Scientific Computing 62 (2015) 25–52.
- [43] N. Gerhard, S. Müller, Adaptive multiresolution discontinuous Galerkin schemes for conservation laws: Multi-dimensional case, Computational and Applied Mathematics 35 (2016) 321–349.
- [44] B. Cockburn, C.-W. Shu, The Runge-Kutta discontinuous Galerkin method for conservation laws V: Multidimensional systems, Journal of Computational Physics 141 (1998) 199–224.
- [45] N. Hovhannisyanyan, S. Müller, R. Schäfer, Adaptive multiresolution discontinuous Galerkin schemes for conservation laws, Mathematics of Computation 83 (2014) 113–151.
- [46] R. J. Leveque, Finite volume methods for hyperbolic problems, 1st Edition, Cambridge Texts in Applied Mathematics, Cambridge University Press, 2002.

## Acknowledgements

This work is supported by DFG HE5386/14,15, BMBF 05M18PAA, DFG-GRK 2326.