

## Sensitivity Analysis of Burgers' Equation with Shocks\*

Qin Li<sup>†</sup>, Jian-Guo Liu<sup>‡</sup>, and Ruiwen Shu<sup>§</sup>

**Abstract.** The generalized polynomial chaos (gPC) method has been extensively used in uncertainty quantification problems where equations contain random variables. For gPC to achieve high accuracy, PDE solutions need to have high regularity in the random space, but this is what hyperbolic type problems cannot provide. We provide a counterargument in this paper and show that even though the solution profile develops singularities in the random space, which destroys the spectral accuracy of gPC, the physical quantities (such as the shock emergence time, the shock location, and the shock strength) are all smooth functions of the uncertainties coming from both initial data and the wave speed. With proper shifting, the solution's polynomial interpolation approximates the real solution accurately, and the error decays as the order of the polynomial increases. Therefore this work provides a new perspective to “quantify uncertainties” and significantly improves the accuracy of the gPC method with a slight reformulation. We use the Burgers' equation as an example for thorough analysis, and the analysis could be extended to general conservation laws with convex fluxes.

**Key words.** sensitivity analysis, hyperbolic conservation law, uncertainty quantification

**AMS subject classifications.** 65M70, 35L67

**DOI.** 10.1137/18M1211763

**1. Introduction.** Hyperbolic conservation laws describe many important physics balance laws such as conservation of mass, momentum, and energy. They describe various important continuum physics including wave propagation and wave interactions. It is a very classical mathematical subject that has a long tradition tracing back to Euler. In all these studies, the equations are deterministic, with prescribed boundary and initial conditions. Typically there are parameters in the equations that are simply predetermined using constitutive laws. However, from a realistic point of view, uncertainties are generic, in the sense that the initial/boundary conditions and equation parameters usually come from experiments and therefore inevitably have measurement error. If the initial/boundary conditions or the constitutive laws are uncertain and inaccurate, is the solution affected dramatically by such uncertainties? And how does one quantify the influence of the uncertainties on the solution?

\*Received by the editors September 5, 2018; accepted for publication (in revised form) September 28, 2020; published electronically December 3, 2020.

<https://doi.org/10.1137/18M1211763>

**Funding:** The work of the first author was partially supported by National Science Foundation grants DMS-1619778, DMS-1740707, and DMS-1107291: RNMS KI-Net. The work of the second author was partially supported by National Science Foundation grants DMS-1812573 and DMS-1107444: RNMS KI-Net. The work of the third author was partially supported by National Science Foundation grant DMS-1107291: RNMS KI-Net.

<sup>†</sup>Mathematics Department, University of Wisconsin-Madison, Madison, WI 53706 USA ([qinli@math.wisc.edu](mailto:qinli@math.wisc.edu)).

<sup>‡</sup>Department of Mathematics and Department of Physics, Duke University, Durham, NC 27708 USA ([jliu@phy.duke.edu](mailto:jliu@phy.duke.edu)).

<sup>§</sup>Department of Mathematics, University of Maryland, College Park, College Park, MD 20742 USA ([rshu@cscamm.umd.edu](mailto:rshu@cscamm.umd.edu)).

One particular example is from ocean science, in which scientists need to determine the arrival time of a tsunami at a particular location (the land, for example). Such quantity is affected by the time and location of the explosion (an underwater earthquake, or underwater landslides or volcanoes), the undersea topography, the strength of wind, and many others. In practice, we only have limited information about them, and mathematically it is natural to model the unknowns as uncertain parameters in the equations. It is then a mathematical question to understand how the solution behaves as the parameters change the value, and to assess the associated sensitivities.

Suppose we use the one-dimensional shallow water wave equation to model a tsunami:

$$(1.1) \quad \begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + \frac{1}{2}h^2\right) = 0, \end{cases}$$

where  $x$  and  $t$  are space and time coordinates,  $h$  is the depth of water, and  $u$  is the velocity of the seawater. The explosion that triggers the tsunami is typically modeled by a shock profile in the initial data  $(h_{\text{in}}, u_{\text{in}})$ . The data is certainly unknown but we can assume the initial data depends on a random variable  $Z : \Omega \rightarrow \mathbb{R}^d$  that lives in a probability space  $(\Omega, \mathcal{B}, \mathbb{P})$ . The joint probability density function of the random vector  $z = Z(\omega) \in \mathbb{R}^d$  is denoted as  $\pi(z)$ . There are many physical quantities that are of interest. One example is the arrival time of the tsunami to  $x_0$ , the location of the land, denoted by  $t^\sharp$ . The ultimate goal is to predict

$$\mathbb{E}(t^\sharp) = \int t^\sharp(z) \pi(z) \, dz,$$

the expected arrival time (assuming the random variable  $z$  in the initial data correctly describes the explosion), and

$$\text{var}(t^\sharp) = \int \left(t^\sharp(z) - \mathbb{E}(t^\sharp)\right)^2 \pi(z) \, dz,$$

the variance that quantifies the reliability of the prediction.

It is a standard procedure to simplify the model (1.1) using the two Riemann invariants  $u \pm 2\sqrt{h}$  [14]. The equations now read

$$\partial_t \left(u \pm 2\sqrt{h}\right) + \left(u \pm \sqrt{h}\right) \partial_x \left(u \pm 2\sqrt{h}\right) = 0.$$

Suppose one cares about one Riemann invariant  $v = u + 2\sqrt{h}$ , and setting  $(u - 2\sqrt{h})|_{t=0} = c(z)$  to have  $z$ -dependence, (1.1) is then reduced to

$$\partial_t v + \left(\frac{3v}{4} + \frac{c}{4}\right) \partial_x v = 0,$$

which can be reformulated into the form of the Burgers' equation

$$(1.2) \quad \partial_t u + \partial_x \left(\frac{\alpha(z)}{2} u^2\right) = 0, \quad u(t=0, x, z) = u_{\text{in}}(x, z),$$

whose wave speed  $\alpha(z)u$  and the initial condition vary according to the initial condition of (1.1). Our goal then is to study  $u$ , or some physical quantities derived from  $u$ , such as  $t^\sharp$ 's dependence on  $z$ .

There are several conventional ways of computing the solution to equations with unknown parameters. Among them, the generalized polynomial chaos (gPC) method has been quite popular in recent several years. To a large extent, it can be regarded as a spectral type method applied onto the  $z$ -space. Although having the same way of representing functions by orthogonal polynomial expansions, the gPC method has many variations (such as gPC-stochastic Galerkin, gPC-stochastic collocation (gPC-SC), gPC-sparse grid, etc.) [9, 24, 25, 18]. Take the gPC-SC method, for example: a few sample points  $\{z_j\}_{j=1}^N$  are preselected according to the probability distribution of  $\pi(z) dz$  (typically one uses collocation points), and with these  $z_j$  fixed, the equations are deterministic and can be easily computed by existing numerical methods. Upon getting the solutions at these preset sample points, the solutions of the equation on other configurations of  $z$  are then interpolated in a polynomial way. The interpolation is regarded as an accurate surrogate to the true solution.

The method gains its popularity largely due to its “spectral” nature: in many cases it gives spectral convergence, which is faster than most other methods. But it also inherits the strong requirement on the data: for the spectral convergence to be valid, regularity of the solution has to be justified—only when the to-be-interpolated functions are shown to be smooth can one prove the fast convergence (depending on the regularity). For “lucky” cases like elliptic or parabolic equations, one can show such regularity, as was done in [3, 2, 26], but it is not the case for most hyperbolic conservation law equations [4]. On the contrary, it is a well-known fact that the Burgers' equation, the toy model equation in scalar conservation laws, develops, in finite time, singular points (or shocks as they are termed) even with  $C^\infty$  initial data, and such singularity in  $x$  will naturally result in the development of discontinuity in  $z$ . More importantly, such irregularity is generic. Because of this, although the convergence of gPC expansions may still be guaranteed [7], there is no hope to declare any results on the spectral accuracy of gPC methods, which is based on that of polynomial expansions [10, 23].

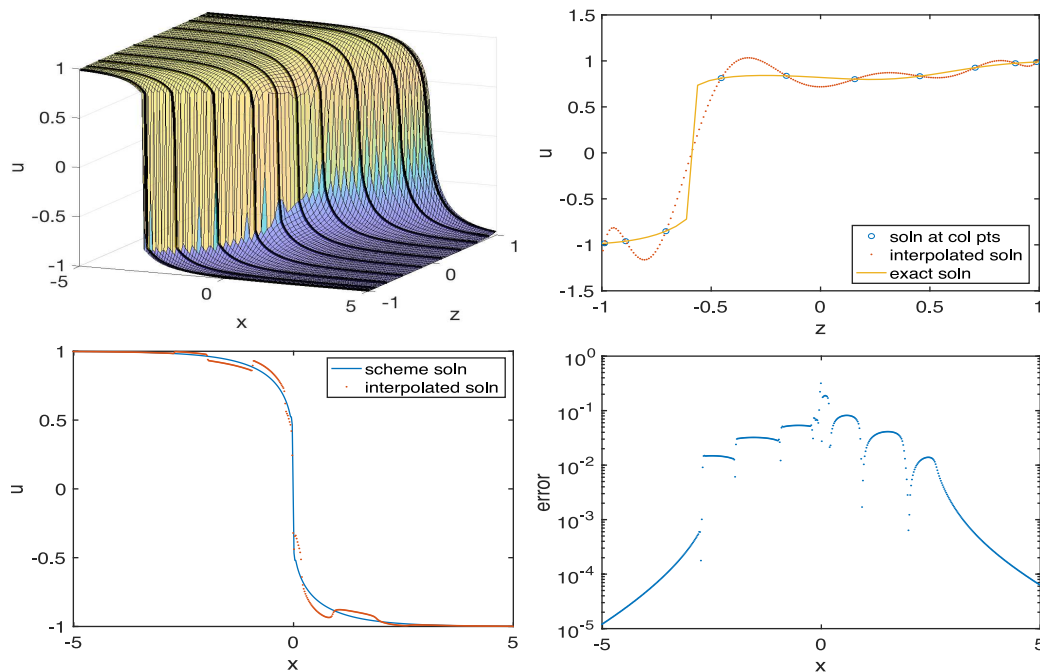
We present some preliminary computation in Figure 1, which shows our numerical results of the Burgers' equation with random initial data,

$$(1.3) \quad u_{\text{in}}(x, z) = v(x, z) - 0.2(v(x, z) + 0.5)(1 - v(x, z))^2 z \quad \text{with} \quad v(x, z) = \frac{1 - e^{x-3z^3}}{1 + e^{x-3z^3}},$$

using the gPC-SC method. In what follows, we sometimes omit the variable  $z$  in functions if it is clear from the content. We have assumed the random variable  $z \in [-1, 1]$  satisfies the Chebyshev distribution. According to the gPC-SC method, the Chebyshev quadrature points are selected,

$$(1.4) \quad z_j = \cos\left(\frac{2N_z + 1 - 2j}{2N_z}\pi\right), \quad j = 1, \dots, N_z = 10,$$

and the equation is numerically solved at each  $z_j$ . We run the experiment up to  $T = 2.2$  and collect  $u(t = 2.2, x, z_j)$  for all  $z_j$  (fifth order WENO scheme with the third order strong stability preserving Runge–Kutta in time used to minimize discretization error from time



**Figure 1.** Top left: the numerical scheme solutions at sample points  $\{z_j\}$ ; top right: the result by direct polynomial interpolation at  $t = 2.2$ ,  $x = -0.5$ ; bottom left: comparing the direct interpolation solution (dots) with the numerical scheme solution (line) at  $z_0 = 0.234$ ; bottom right: error (difference between the two solutions in the bottom left picture).

and space). The solution to all other  $z \in [-1, 1]$  is then interpolated using a ninth order polynomial. We clearly see the spurious oscillations in Figure 1, where we compared the interpolated solution of  $u(t = 2.2, x, z_0 = 0.234)$  with the true solution as a function of  $x$ , and the interpolated solution of  $u(t = 2.2, x = -0.5, z)$  with the true solution as a function of  $z$ .

Such loss of spectral accuracy is expected: due to the Gibbs phenomenon, the spectral method, albeit performing well in  $L_2$ , drastically fails in  $L_\infty$ , in the sense that it gives inaccurate interpolations to discontinuous functions like this one, and the  $L_\infty$  error does not decay when one increases the order of the interpolation. This poor numerical performance was noted in earlier works of [10, 23, 4, 20, 22], and largely for this reason, the results obtained using the gPC type methods are regarded as unsatisfactory in the hyperbolic conservation laws setup.

However, we provide a counterargument in the current paper. In particular, we will be dealing with the inaccuracy in the  $L_\infty$  sense. It is based on a simple observation: even though the solution  $u(t, x, z)$  may be discontinuous in  $z$  for all  $x$ , making the interpolation severely inaccurate, the physical quantities that “practically” matter are still smooth functions in  $z$ -space. They do not contain jump discontinuities in the random space, and we view them as being insensitive to the random perturbation in the initial/boundary conditions. Such physical quantities include the shock location, the shock strength, the shock emerging time, and the arrival time of the shock at certain locations. With these quantities identified, we can perform suitable “shifting” of the solutions, which permits pointwise accuracy, meaning the gPC interpolation can produce accurate approximations even in the  $L_\infty$  sense.

To be more precise, for the one-shock solutions, denote (as plotted in Figure 2)

1.  $t^*(z)$ , the shock appearing time;
2.  $x^c(t, z)$ , the shock location that moves with respect to time for  $t \geq t^*$ ;
3.  $u_1(t, z) - u_2(t, z)$ , the shock strength (for  $t \geq t^*(z)$ ), with

$$(1.5) \quad u_1(t, z) = \lim_{x \rightarrow x^c(t, z)^-} u(t, x, z) \quad \text{and} \quad u_2(t, z) = \lim_{x \rightarrow x^c(t, z)^+} u(t, x, z)$$

being the upper and lower boundaries of the shock;

4.  $t^\sharp(z) = \inf\{t : x^c(t, z) \geq x_0\}$ , the shock arriving time.<sup>1</sup> Here  $x_0$  denotes the predetermined location of land.

We repeat the previous example and find numerical evidence that shows these physical quantities are indeed smooth functions in  $z$ , seen in Figure 3.

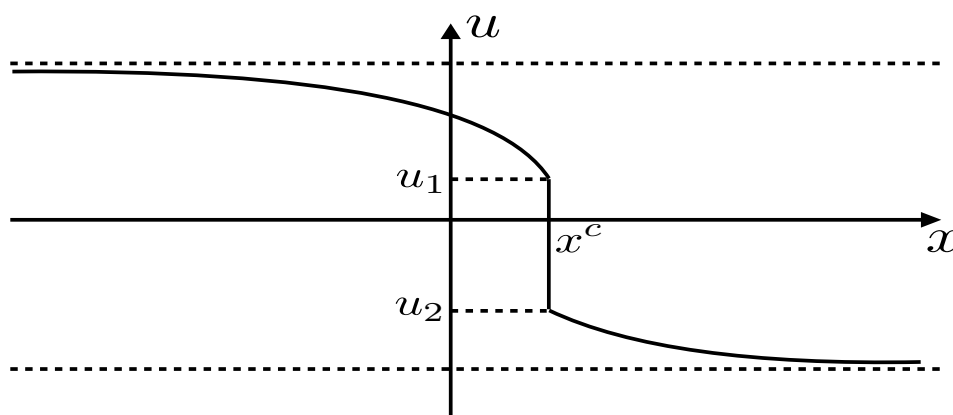


Figure 2. A demonstration of the quantities  $u_1, u_2, x^c$ , in the case of one-shock solution.

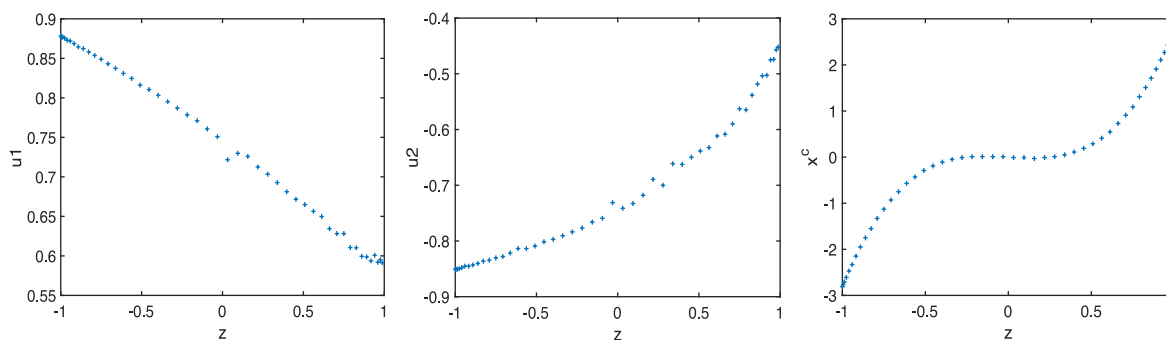


Figure 3. Left to right:  $u_1, u_2, x^c$  as functions of  $z$ . The small zigzags in plots are from numerical error. Numerically to identify these quantities, we look for the biggest jump between two adjacent grid points. This procedure brings some error and is not easily eliminated by adopting another numerical solver for the conservation law.

<sup>1</sup>In case  $x^c(t) \geq x_0$  is never satisfied,  $t^\sharp$  is understood as infinity.

The aim of the current paper is to mathematically prove this observation. Our main claim in this paper is that even the solution to the scalar hyperbolic conservation law varies drastically with respect to the random inputs in initial condition or the equation parameters, the physical quantities are insensitive to them, or more precisely, we have as follows.

**Theorem 1.1 (formal statement).** *Let  $u(t, x, z)$  be the solution to the Burgers' equation (1.2) a parameter  $z$  representing uncertainty. Assume the initial data  $u_{in}(x, z)$  is smooth and satisfies a set of conditions (to be made precise later) so that only one shock will appear for each  $z$ , and assume  $\alpha$  is smooth on  $z$ . Then*

1. *the shock appearing time  $t^*(z)$  depends smoothly on  $z$ ;*
2. *the shock location  $x^c(t, z)$  depends smoothly on  $z$ ;*
3. *the shock strength  $u_1(t, z) - u_2(t, z)$  depends smoothly on  $z$ ;*
4. *the shock arriving time  $t^\sharp(z)$ , depends smoothly on  $z$ .*

As a direct corollary of this theorem, we also find that if the solutions are “shifted correctly,” the shifted solution  $\tilde{u}(t, x, z)$  becomes a smooth function in  $z$  for every  $t$  and  $x$ , granting the accuracy to the gPC type method. This could be summarized as follows.

**Theorem 1.2 (formal statement).** *Define the shifted solution,*

$$\tilde{u}(t, x, z) = u(t + t^*(z), x + x^c(t + t^*(z), z), z),$$

*so that shocks are aligned for all  $z$  to the same emerging time and location; then with the same assumptions as Theorem 1.1,  $\tilde{u}$  is smooth in  $z$  away from the set  $\{(t, x, z) : x = 0\}$  for  $t > 0$ .*

There are many groups of researchers working on similar topics. In [11, 12], the authors adopted the patchwise low-rank studies. In [15, 16], the authors proved the wellposedness of entropy solution when randomness is present in initial data and flux, and  $L_1$  contraction is used for estimating the error from the interpolation method. In [20], the authors gave a very detailed analysis on the shock location of the Burgers' equation with Heaviside function as the initial data. A similar approach was taken in [22], where the author presents very powerful numerical evidence that demonstrates the shifting indeed “saves” the regularities of the solutions. Another approach proposed in [17] and the references therein is to explore Monte Carlo (and multilevel Monte Carlo) methods. In [19], the authors introduce entropic variables and expand the solution as polynomials of the new variable with the understanding that it is smoother when represented by the new variable. In [5] the authors derive the reduced order equations for the one-point and two-point probability density functions of the solution field and design algorithms to compute the statistical properties of the random shock wave accordingly. Other approaches include [1, 8], where authors employed the so-termed truncate-encode framework, and [6], where a kinetic formulation is utilized. The approach we are taking is in line with [22] and [20] but we emphasize giving a quantitative mathematical justification in general cases. Our approach is also closely related to the stochastic transformation proposed in [13] for the viscous Burgers' equation with random forcing.

We remark that Theorems 1.1 and 1.2 also hold for general scalar conservation laws with smooth convex flux functions. This means that the smoothness in  $z$  of the physical quantities (shock appearing time, shock location, shock strength, etc.) is a generic fact, and the smoothness in  $z$  of the solution profile can always be recovered by shifting correctly.

We emphasize before finishing the introduction that the goal of the paper is not to justify the use of the gPC method on Burgers' equation, or hyperbolic conservation laws in general, but to bring one more aspect to understand the shock structure in wave-like equations when uncertainties are present. In fact, almost all numerical methods somewhat rely on the regularity of some to-be-computed quantities, and the result obtained in this paper serves as a justification for these algorithms applied in  $z$ -space.

The rest of this paper is organized as follows. In section 2 we introduce some notation and state the precise quantitative versions of Theorems 1.1 and 1.2; in section 3 we focus on the deterministic case and prepare some necessary tools for analyzing  $u_1$  and  $u_2$ ; these tools are crucial in section 4, where we prove Theorems 1.1 and 1.2. We also extend the results to treat conservation laws with general convex fluxes in section 5. Some proofs are tedious and not essential to the main context, and they are left to the appendix.

**2. Notation and precise statement of main results.** There is a large variety of solution behavior of conservation laws, and we restrict ourselves to the class of smooth initial data such that only one shock is developed for  $t > 0$ . Mathematically, we have as follows.

*Assumption 1.* Denote  $u_{in}(x, z)$  the initial data. We require  $u_{in}(x, z)$  to be smooth in  $(x, z)$  and  $u_{in}(\cdot, z)$ , as a function of  $x$ , to satisfy the following:

- $u_{in}$  monotonically decreases in  $x$ , i.e.,  $u'_{in}(x) < 0$  for all  $x$ , and  $\lim_{x \rightarrow \pm\infty} u_{in}(x) = u_{\pm}$ ; here  $u_+ < u_-$  are constants independent of  $z$ ;
- $u_{in}$  has a unique inflection point  $(x^*, u^*)$ , meaning  $u_{in}(x^*) = u^*$  and  $u''_{in}(x^*) = 0$ ;
- $u'''_{in}(x^*) > 0$ ;
- $\alpha(z)$  has uniform bounds:  $0 < \alpha_0 < \alpha(z) < \alpha_1$ .

Under these assumptions, we restate Theorem 1.1 rigorously and quantitatively.

**Theorem 2.1.** *Let  $u(t, x, z)$  be the solution to the Burgers' equation (1.2) with uncertainty. Assume that the initial data  $u_{in}(x, z)$  satisfies Assumption 1 and that  $u_+ + \delta \leq u^*(z) \leq u_- - \delta$  for all  $z$ , with  $\delta > 0$ . Then with  $\alpha(z) > 0$  being smooth in  $z$ , one has the following:*

1. The shock appearing time is given by

$$(2.1) \quad t^*(z) = -\frac{1}{\alpha(z)u'_{in}(x^*(z), z)},$$

where  $x^*(z)$  is as in Assumption 1. It follows that  $t^*(z)$  depends smoothly on  $z$ .

2. The shock location  $x^c(t, z)$  (defined for  $t \geq t^*(z)$ ) depends smoothly on  $z$  and satisfies the estimate

$$(2.2) \quad \partial_z^k x^c = \mathcal{O}\left((t - t^*)^{\min\{3/2-k, 0\}}\right)$$

for  $t - t^*$  small enough.

3. The shock strength  $u_1(t, z) - u_2(t, z)$  depends smoothly on  $z$  and satisfies the estimate

$$(2.3) \quad \partial_z^k (u_1 - u_2) = \mathcal{O}\left((t - t^*)^{1/2-k}\right)$$

for  $t - t^*$  small enough.



4. Let  $x_0$  be large enough so that  $x_0 > \sup_z x^c(t^*(z), z)$ . Assume

$$(2.4) \quad \partial_t x^c(t^\sharp, z) \neq 0$$

and that  $t^\sharp < \infty$ ; then  $t^\sharp$  depends smoothly on  $z$ .

*Remark 1.* We now comment that assumption (2.4) is not restrictive. In fact, as will be seen in section 3,  $\partial_t x^c$  has the explicit expression (3.5), and thus (2.4) can be checked explicitly. In the same section, we can also derive that in long time,

$$(2.5) \quad \lim_{t \rightarrow \infty} \partial_t x^c(t, z) = \alpha \frac{u_- + u_+}{2},$$

and thus  $u_- + u_+ > 0$  automatically leads to (2.4) for large  $t$ . This exactly corresponds to the realistic case when a tsunami forms far away from the land and takes a long time to propagate to the land.

Theorem 1.2 is a corollary of the theorem above, and in the precise form it states as follows.

**Theorem 2.2.** *Under the same assumptions as Theorem 2.1, the translated solution*

$$(2.6) \quad \tilde{u}(t, x, z) = u(t + t^*(z), x + x^c(t + t^*(z), z), z)$$

is smooth in  $z$  away from the set  $\{(t, x, z) : x = 0, t = 0\}$  and has the estimate

$$(2.7) \quad \left| \partial_z^k \tilde{u}(t, x, z) \right| = \mathcal{O}(|x|^{1-2k} t^{\min\{3/2-k, 0\}})$$

if  $t > 0$  is small enough.

This theorem implies that a proper shifting of the solution could eliminate the irregular jumps in the solution, and this would allow the spectral type method such as gPC to apply well. In particular, using the same example as in section 1, assuming the random variable  $z \in [-1, 1]$  satisfying the Chebyshev distribution, we denote

$$(2.8) \quad u^N(t, x, z) = \sum_{j=0}^N \tilde{u}(t, x, z_j) \ell_j(z),$$

where  $z_j$  are the Chebyshev quadrature points defined in (1.4), and  $\ell_j$  are the corresponding Lagrange polynomials in the  $z$  domain. Then we have the following theorem. We note that if  $z$  satisfies another distribution, similar techniques can still be applied with  $z_j$  shifted accordingly. For the conciseness of the statement here we stick to this particular kind of random variable.

**Theorem 2.3.** *Assume the random variable  $z \in [-1, 1]$  satisfying the Chebyshev distribution. Under the same assumptions as Theorem 2.1, the error of the interpolated solution  $u^N$  can be estimated by*

$$(2.9) \quad \left| \tilde{u}(t, x, z) - u^N(t, x, z) \right| \leq \frac{C(m) |x|^{-1-2m} t^{1/2-m}}{N^m} \quad \forall z \in [-1, 1], \quad t > 0,$$



for any  $m \geq 1$ , i.e., it has  $m$ th order accuracy away from the shock location and the shock appearing time.

Furthermore, if we use  $\mathbb{E}(u^N)$ ,  $\text{var}(u^N)$  to approximate the mean and variance of  $\tilde{u}$ , and assuming that  $z \in [-1, 1]$ , then we have the error estimate

$$(2.10) \quad |\mathbb{E}(\tilde{u}) - \mathbb{E}(u^N)| \leq \frac{C(m)|x|^{-1-2m}t^{1/2-m}}{N^m},$$

$$(2.11) \quad |\text{var}(\tilde{u}) - \text{var}(u^N)| \leq \frac{C(m)(1 + \min\{\|\tilde{u} - u^N\|_{L^\infty_z}, N^2\})|x|^{-1-2m}t^{1/2-m}}{N^m}.$$

We remark that all three estimates in Theorem 2.3 are pointwise in  $t$  and  $x$ . They deteriorate as  $|x|$  or  $t$  gets small, i.e., the location is close to the shock or the time is close to the shock appearing time.

Note that the min in the estimate of  $\text{var}(\tilde{u}) - \text{var}(u^N)$  in (2.11) is necessary. For fixed  $x$ , as  $N$  goes to infinity, according to (2.9),  $\tilde{u}(t, x, z) - u^N(t, x, z)$  shrinks to zero, but on the other hand, for fixed  $N$  and  $x \sim 0$  (close to the shock), the difference between  $\tilde{u}$  and  $u^N$  could be significant and we use  $N^2$  as the bound there.

We note that Theorems 2.1 and 2.2 only state the smoothness in  $z$ -space regarding  $z$  as an external unknown parameter. It was not until Theorem 2.3 where we incorporate the statistical behavior, and thus  $\pi(z)$ , the distribution is needed.

Proofs for Theorems 2.1 and 2.2 heavily depend on the delicate analysis of  $u_1$  and  $u_2$ , while Theorem 2.3 immediately follows the regularity results in Theorem 2.2. Since the dynamics of  $u_{1,2}$  are so crucial, for a clear presentation, we devote section 3 to developing the necessary tools for the equation in the deterministic setting. These results will be used in later sections, where an energy estimate is used for showing the two main theorems.

**3. Burgers' equation—deterministic case.** In this section we will mainly focus on the shock behavior, and the main tool is the hodograph transform. The reformulation is performed in section 3.1 and the local-in-time shock behavior is presented in Theorem 3.1 in section 3.2.

**3.1. Reformulation of the Burgers' equation.** The monotonicity assumption of  $u$  on  $x$  makes the application of the hodograph transform possible: by flipping  $x, u$  coordinates we can study the evolution of  $x(u)$  in time. Denote  $x = x(t, u)$  the inverse function of  $u(t, x)$  for all  $t$ ; then the domain is  $(t, u) \in [0, \infty) \times (u_+, u_-)$ .

**3.1.1. Before the formation of the shock at time  $t^*$ .**  $x(u)$  is the coordinate that sits on the  $u$ -level set. Since the Burgers' equation has its wave propagating with speed  $\alpha u$ , then before  $t^*$ , we have

$$(3.1) \quad \begin{cases} \partial_t x(t, u) = \alpha u, & u \in (u_+, u_-), \\ x(t = 0, u) = x_{\text{in}}(u). \end{cases}$$

Assumption 1 on  $u_{\text{in}}$  could be translated to assumption on  $x_{\text{in}}$ :

- $x'_{\text{in}}(u) < 0$ ;
- $x_{\text{in}}$  has one unique inflection point at  $(u^*, x^*)$ ;
- $(x_{\text{in}})'''(u^*) < 0$ .

Taking  $\partial_u$  of (3.1), we also have

$$(3.2) \quad \partial_t \partial_u x(t, u) = \alpha \quad \Rightarrow \quad \partial_u x(t, u) = x'_{\text{in}}(u) + \alpha t = -f(u) + \alpha t,$$

where we used the notation

$$f(u) = -x'_{\text{in}}(u).$$

Combined with the property  $x'_{\text{in}}(u) < 0$ , we have  $\partial_u x(t, u) \leq 0$  for  $t < t^*$  where

$$(3.3) \quad t^* = \min_u \left( -\frac{1}{\alpha} x'_{\text{in}}(u) \right) = -\frac{1}{\alpha} x'_{\text{in}}(u^*)$$

is the earliest time for a shock to emerge. Such a shock appears at  $(u^*, x^*)$ .

**3.1.2. After the formation of the shock at  $t^* = -\frac{1}{\alpha} x'_{\text{in}}(u^*)$ .** The strong solution to (1.2) breaks down, and (3.1) no longer correctly characterizes the solution behavior. As the weak formulation and the entropy condition are used to replace the strong form to characterize  $u(t, x)$  on the  $(x, u)$ -plane, a different set of equations is needed for  $x(t, u)$  on the  $(x, u)$ -plane: to do that we first utilize the Rankine–Hugoniot condition. Denote  $u_1(t)$  and  $u_2(t)$  as the top and the bottom of the shock; then the shock speed is

$$(3.4) \quad s = \alpha \frac{u_1(t)^2/2 - u_2(t)^2/2}{u_1(t) - u_2(t)} = \alpha \frac{u_1(t) + u_2(t)}{2},$$

meaning the shock location  $x^c(t)$  satisfies the ODE

$$(3.5) \quad \frac{d}{dt} x^c = \alpha \frac{u_1(t) + u_2(t)}{2} \quad \text{with} \quad x^c(t^*) = x^*.$$

Notice that under the hodograph transform, the shock in  $u(t, x)$  becomes a flat region in  $x(t, u)$ , ranging from  $u_2$  to  $u_1$ , and has height  $x^c$ .

The ODE system for  $u_{1,2}(t)$  can also be derived, as seen in Figure 4. If one focuses on the neighborhood of  $u_1$ , the flat region propagates in the vertical direction with speed  $s$ , while  $x$  for  $u > u_1$  propagates in the horizontal direction with a faster speed  $u_1 > s$ . These coordinates that are supposed to travel faster then get absorbed into the flat region, widening it (around  $u_1$ ) by

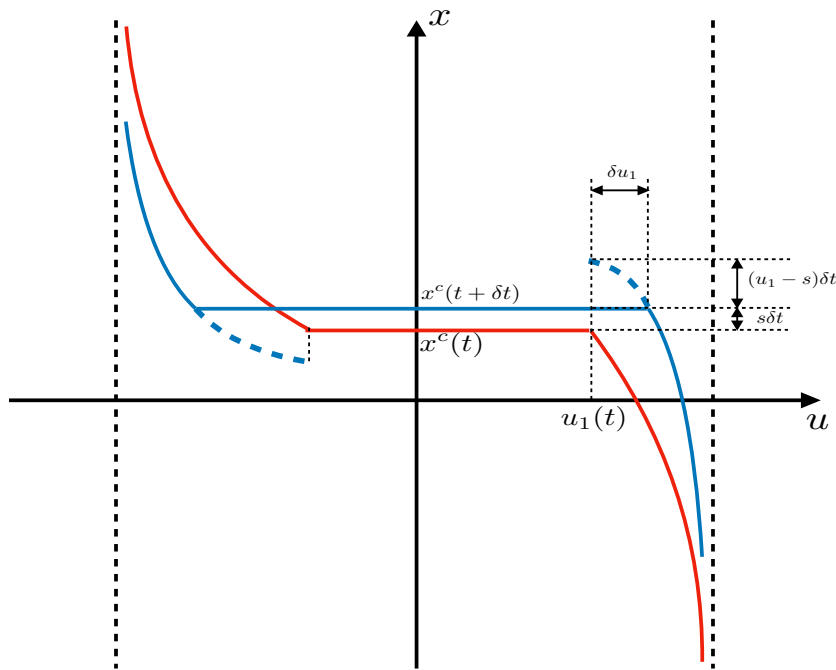
$$(3.6) \quad \delta u_1 = -\partial_x u(t, x)|_{u=u_1} \delta x = -\partial_x u(t, x)|_{u=u_1} \alpha \frac{u_1(t) - u_2(t)}{2} \delta t = -\frac{\alpha}{\partial_u x(t, u_1)} \frac{u_1(t) - u_2(t)}{2} \delta t,$$

where

$$\delta x = (\alpha u_1 - s) \delta t = \alpha \frac{u_1(t) - u_2(t)}{2} \delta t$$

presents the “overshoot” before entropy condition is applied to “cut” the multivalued solution. Considering (3.2) and conducting the same analysis for  $u_2$ , one has

$$(3.7) \quad \frac{d}{dt} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} F_1(u_1, u_2) \\ F_2(u_1, u_2) \end{pmatrix} = \frac{\alpha}{2} (u_1 - u_2) \begin{pmatrix} \frac{1}{f(u_1) - \alpha t} \\ \frac{-1}{f(u_2) - \alpha t} \end{pmatrix}$$



**Figure 4.** Evolution of  $u_1, u_2$  in the hodograph-transformed picture. Here the red curve is the solution  $x(t, u)$  at some time spot, and the solid blue curve is the solution after a small time period  $\delta t$ . The dashed blue curve is the dynamics of (3.1).

with the initial condition

$$(3.8) \quad u_1(t^*) = u_2(t^*) = u^*.$$

In the equation  $F_{1,2}$  denotes the forcing terms for  $u_{1,2}$ , respectively.

*Remark 2.* Some comments are in line:

- According to (3.2) and the monotonicity of  $\partial_u x$ , there holds  $-f(u_{1,2}) + \alpha t = x'_{in}(u_{1,2}) + \alpha t \leq 0$ . Combined with (3.7), it is shown that  $u_1(t)$  monotonically increases in time and  $u_2(t)$  monotonically decreases in time, meaning

$$(3.9) \quad u_1 \geq u^*, \quad u_2 \leq u^*, \quad f(u_{1,2}) - \alpha t \geq 0.$$

- The system (3.7) is self-consistent. This means the information in the shock is fully represented by  $u_1$  and  $u_2$ . The general profile of  $x(u)$  is irrelevant.

**3.1.3. Summary.** To summarize the reformulation, the Burgers' equation, when written on the  $(x, u)$ -plane, becomes

$$(3.10) \quad \begin{cases} t < t^* = -\frac{1}{\alpha} x'_{in}(u^*) : & \text{equation (3.1),} \\ t > t^* : & \begin{cases} \text{equation (3.1) with } u \in (u_+, u_2) \cup (u_1, u_-), \\ \text{equation (3.5) with } u \in (u_2, u_1), \end{cases} \end{cases}$$

where  $u_1$  and  $u_2$  are the shock locations satisfying the ODE system (3.7).

**3.2. Shock behavior for small time.** The shock behavior is fully described by (3.7), which we study in depth in this section. To start, we first shift the coordinates and the time frame so that<sup>2</sup>  $u^* = 0$  and  $t^* = 0$ . Since  $t^* = -\frac{1}{\alpha}x'_{\text{in}}(u^*)$ ,  $u_1(t^*) = u_2(t^*) = u^*$ , one has

$$(3.11) \quad u_1(0) = u_2(0) = 0 \quad \text{and} \quad f(0) = 0.$$

Physically it means the flat region starts forming at  $t = 0$ ,  $u = 0$ .

Assumption 1 on  $u_{\text{in}}$  is now formulated as the following.

*Assumption 2.* Denote  $f(u) = -x'_{\text{in}}(u)$ ; then we have the following:

- $f(u) \geq 0$ . This follows from item 2 of Assumption 1 and the fact that  $x'_i(u) = \frac{1}{u'_{\text{in}}(x)}$ .
- $f'(u) = 0$  at only one point  $u^*$ . To see this, one first differentiates  $x'_i(u) = \frac{1}{u'_{\text{in}}(x)}$  to obtain  $x''_i(u) = -\frac{u''_{\text{in}}(x)}{(u'_{\text{in}}(x))^3}$ , and then notices item 3 of Assumption 1.
- $f'(u) < 0$  when  $u < 0$ ,  $f'(u) > 0$  when  $u > 0$ , and  $f''(0) > 0$ . The sign of  $f'$  can be seen by  $x''_i(u) = -\frac{u''_{\text{in}}(x)}{(u'_{\text{in}}(x))^3}$  and the signs of  $u'_{\text{in}}(x)$ ,  $u''_{\text{in}}(x)$ . The sign of  $f''(0)$  can be seen by differentiating  $x''_i$  and evaluating at  $u = x = 0$  to see  $x'''_{\text{in}}(0) = -\frac{u'''_{\text{in}}(0)}{(u'_{\text{in}}(0))^4}$ , and combining with item 4 of Assumption 1.

*Remark 3.* These assumptions, when combined with (3.11), indicate that around  $u = 0$ ,  $f(u)$  behaves like a quadratic function

$$(3.12) \quad f(u) \sim au^2 \quad \text{with} \quad a = \frac{1}{2}f''(0) > 0.$$

With this simpler quadratic form, in small time,  $u$  is also small. The ODE system (3.7) then gets simplified,

$$(3.13) \quad \begin{cases} \frac{du_1}{dt} = \frac{\alpha}{2}(u_1 - u_2)\frac{1}{au_1^2 - \alpha t}, \\ \frac{du_2}{dt} = -\frac{\alpha}{2}(u_1 - u_2)\frac{1}{au_2^2 - \alpha t}, \end{cases}$$

and one has the explicit solution:

$$(3.14) \quad u_1 = -u_2 = \left(\frac{3\alpha}{a}t\right)^{1/2}.$$

This result implies that  $u_1$  and  $u_2$  approximately grow in time with the power  $\frac{1}{2}$ . Generally,  $f(u)$  is not a quadratic function, but one can still use two quadratic functions with different  $a$  to sandwich the solution for a  $t^{1/2}$  growth rate.

We now state our theorem.

**Theorem 3.1.** *Under Assumption 2 on  $f$ , assume  $u_{1,2}(t)$  solve the ODE system (3.7) with initial condition (3.11) and satisfy  $u_1 > 0$ ,  $u_2 < 0$ ,  $f(u_{1,2}) - \alpha t > 0$  for  $t > 0$ . Then, for any  $\epsilon > 0$ , there holds*

<sup>2</sup>This assumption implies that  $x_{\text{in}}(u) \sim u^3$  for  $u$  close to zero. In other words,  $u_{\text{in}}(x) \sim (x - x^*)^{1/3}$  for  $x$  close to  $x^*$ . This is known as the shock formulation profile for the Burgers' equation.

$$(3.15) \quad \left(\frac{3\alpha}{a+\epsilon}t\right)^{1/2} \leq u_1 \leq \left(\frac{3\alpha}{a-\epsilon}t\right)^{1/2}, \quad -\left(\frac{3\alpha}{a-\epsilon}t\right)^{1/2} \leq u_2 \leq -\left(\frac{3\alpha}{a+\epsilon}t\right)^{1/2},$$

for  $t$  small enough, with  $a$  defined in (3.12).

Note that the wellposedness of the system is not discussed in the theorem. In fact, away from the initial time, the forcing terms are Lipschitz, making the proof of the wellposedness standard, which we leave to Appendix A.1. To prove the small time behavior of  $u_{1,2}$ , we first start with an ODE (analyzed in Lemma 3.2), and then utilize the symmetry condition (Lemma 3.3) for a solution to the ODE system (3.7) when  $f$  is a quadratic function. The monotonicity (Lemma 3.4) is then applied to sandwich the solution to the problem in which  $f$  is not quadratic. Without loss of generality,  $\alpha$  is set to be 1 below.

**Lemma 3.2.** *Under Assumption 2, and assuming  $u(t) > 0$  and  $f(u) > t$  for all  $t > 0$ , the ODE*

$$(3.16) \quad \frac{du}{dt} = \frac{u}{f(u) - t}, \quad u(0) = 0,$$

has a unique solution given by the implicit function

$$(3.17) \quad t = \frac{1}{u} \int_0^u f(s) ds.$$

**Remark 4.** We note that we do not have the wellposedness if we remove the condition  $f(u) > t$  and  $u(t) > 0$ . In fact,  $u = 0$  for all  $t > 0$  is also a solution. The extra condition allows us to obtain the uniqueness for all  $t$ .

**Proof.** The condition  $f(u) > t$  excludes the possibility of  $u(t) = 0$  for  $t > 0$ , and thus  $\frac{du}{dt} \neq 0$ , and one can write  $t = t(u)$ . Then  $t(u)$  satisfies

$$(3.18) \quad \frac{dt}{du} = \frac{f(u) - t}{u}, \quad t(0) = 0,$$

which is a linear ODE, and has the general solution

$$(3.19) \quad t = \frac{1}{u} \left( \int_0^u f(s) ds + C \right),$$

away from  $u = 0$ . Since  $f(s) \sim as^2$  for small  $s$ , it is clear that  $\lim_{u \rightarrow 0} \frac{1}{u} \left( \int_0^u f(s) ds + C \right) = 0$  holds only when  $C = 0$ . This means (3.17) gives the only solution to (3.16) satisfying the assumptions. ■

**Lemma 3.3.** *If  $f(u) = f(-u)$  is a symmetric function, then  $(u, -u)$  solves (3.7) if  $u$  solves (3.16).*

The proof is rather straightforward and we omit it.

**Lemma 3.4.** *Let  $(u_1, u_2)$  solve (3.7) with initial condition (3.11), and let  $(v, -v)$  solve (3.7) with initial condition (3.11) where  $f$  is replaced by an even function  $g$  ( $g(u) = g(-u)$ ). Then for small  $t$ ,*

- if  $f(u) < g(u)$  for all  $u$ , then  $u_1 \geq v, u_2 \leq -v$ ;
- if  $f(u) > g(u)$  for all  $u$ , then  $u_1 \leq v, u_2 \geq -v$ .

*Proof.* Again, we use the shooting method to prove this lemma. We prove the first statement by contradiction. Suppose it is not true; then there exists a  $t_0 > 0$  small enough such that  $u_1(t_0) < v(t_0)$  or  $u_2(t_0) > -v(t_0)$ . Without loss of generality, we assume the former case and that  $-u_2(t_0) \geq u_1(t_0)$ . From the previous lemma, there exists a  $g$ -solution  $(v_1, -v_1)$  with  $v_1(t_0) > u_1(t_0)$ , and this solution hits  $g(v_1) - t = 0$  line before  $t = 0$ , meaning there is  $t_1 > 0$  such that

$$f\left(v_1 - \int_{t_1}^{t_0} F_1^g(v_1, -v_1) ds\right) = t_1.$$

Here  $F_1^g(u_1, u_2) := \frac{1}{2}(u_1 - u_2) \frac{1}{g(u_1) - t}$  is the forcing term for  $v_1$  defined by  $g$ .

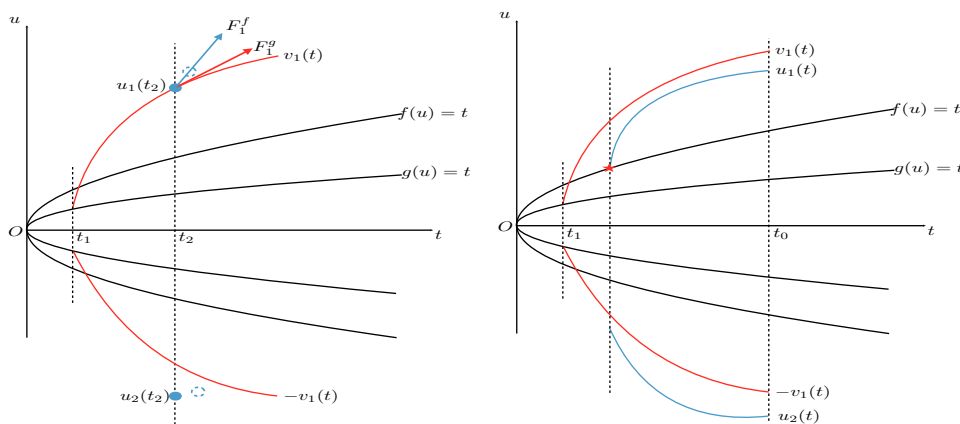
We claim that there is no  $t \leq t_0$  such that  $u_1(t) \geq v_1(t), u_2(t) \leq -v_1(t)$ . In fact, this is not true at  $t = t_0$ . Suppose  $t_2$  is the largest time such that this holds; then without loss of generality, we assume that  $u_1(t_2) = v_1(t_2)$ . Then we have

$$\begin{aligned} F_1^f(u_1(t_2), u_2(t_2)) &:= \frac{1}{2}(u_1(t_2) - u_2(t_2)) \frac{1}{f(u_1(t_2)) - t_2} \geq \frac{1}{2}(v_1(t_2) + v_1(t_2)) \frac{1}{f(v_1(t_2)) - t_2} \\ &> \frac{1}{2}(v_1(t_2) + v_1(t_2)) \frac{1}{g(v_1(t_2)) - t_2} = F_1^g(v_1(t_2), -v_1(t_2)). \end{aligned}$$

This contradicts the choice of  $t_2$ . See Figure 5 (left) for an illustration.

This claim contradicts the fact that  $(u_1, u_2)$  can be continued to the time  $t_1$ , since at this time,  $f(u_1) - t < g(u_1) - t \leq g(v_1) - t = 0$  or  $f(u_2) - t < g(u_2) - t \leq g(-v_1) - t = 0$ . Thus the first statement is proved. See Figure 5 (right) for an illustration. The second statement can be proved similarly. ■

We now are ready to show Theorem 3.1.



**Figure 5.** Proof of Lemma 3.4. Left: obtaining a contradiction at time  $t_2$ . The dashed circle is the approximate position of  $u_{1,2}$  at time slightly larger than  $t_2$ , which contradicts the choice of  $t_2$ . Right: obtaining the final contradiction.  $u_1$  or  $u_2$  must touch the curve  $f(u) = t$  at some time larger than  $t_1$  (the star in the picture).

*Proof.* Locally at  $t = 0$ ,  $f(u) \sim au^2$  and thus we approximate (3.7) by

$$(3.20) \quad \frac{du_1}{dt} = \frac{1}{2}(u_1 - u_2) \frac{1}{au_1^2 - t}, \quad \frac{du_2}{dt} = -\frac{1}{2}(u_1 - u_2) \frac{1}{au_2^2 - t},$$

which has solution

$$u_1 = -u_2 = \left(\frac{3}{a}t\right)^{1/2}.$$

To estimate the solution of (3.7) near  $t = 0$ , we let  $\epsilon < a$ , and find  $\delta$  such that

$$(3.21) \quad |f(u) - au^2| < \epsilon u^2 \quad (\forall |u| < \delta);$$

then according to Lemma 3.4, for small  $t$  (small enough so that  $u_1^- < \delta$  and  $-u_2^- < \delta$ )

$$(3.22) \quad u_1^+ \leq u_1 \leq u_1^-, u_2^- \leq u_2 \leq u_2^+,$$

where  $u_i^\pm$  is the solution of (3.20) with  $a$  replaced by  $a \pm \epsilon$ , and we conclude the theorem. ■

*Remark 5.* We note that the forcing terms in (3.7) are Lipschitz away from  $f(u) = \alpha t$ , and the system is automatically wellposed there. The main difficulty lies in the “small time” regime where  $f(u^*) = \alpha t^*$ .

**4. Smoothness in  $z$ -space.** As we discussed in the introduction, there are two sources of uncertainties in the Burgers equation (1.2): the initial condition  $u_{in}(x, z)$  and the traveling speed of the wave  $\alpha(z)$ . In Theorem 2.1 we claim that the physical quantities such as  $t^*$  (the shock emerging time),  $t^\#$  (the time for the shock hitting the land), and  $x^c$  (the shock location) are smooth functions of  $z$ , and in Theorem 2.2, we claim that with proper shifting, the solution profile depends on  $z$  smoothly as well. These two theorems are proved in sections 4.1 and 4.2, respectively.

**4.1. Smoothness of physical quantities.** The main goal in this subsection is to prove Theorem 2.1, which states that the physical quantities smoothly depend on  $z$ . We start by proving item 1 of Theorem 2.1.

*Proof of item 1 of Theorem 2.1.* In the deterministic case, we have shown that

$$t^* = -\frac{1}{\alpha u'_{in}(x^*)},$$

and thus to show the regularity of  $t^*$  on  $z$  amounts to showing the regularity of  $x^*$  on  $z$ , since  $\alpha$  and  $u_{in}$  are assumed to be smooth in  $z$ . Take the first derivative, for example:

$$(4.1) \quad \partial_z t^* = \frac{1}{(\alpha u'_{in}(x^*))^2} \partial_z (\alpha u'_{in}(x^*)) = \frac{\partial_z \alpha u'_{in}(t_*) + \alpha \partial_z u'_{in}(x^*, z) + \alpha u''_{in}(x^*, z) \partial_z x^*}{(\alpha u'_{in}(x^*))^2}.$$

$\partial_z u_{in}$  and  $\partial_z \alpha$  are known to be bounded quantities, and by definition

$$u''_{in}(x^*(z), z) = 0,$$

which gives  $|\partial_z t^*| < C$ , meaning  $t^*$  is Lipschitz continuous in  $z$ . Higher derivatives can be analyzed in a similar way except that one also needs to analyze  $\partial_z^k x^*$ . It is a bounded quantity as well and we show it for  $k = 1$ . Since

$$(4.2) \quad u''_{in}(x^*(z), z) = 0 \quad \Rightarrow \quad u'''_{in}(x^*(z), z) \partial_z x^* + \partial_z u''_{in}(x^*(z), z) = 0,$$



which gives

$$(4.3) \quad \partial_z x^* = -\frac{\partial_z u_{\text{in}}''(x^*(z), z)}{u_{\text{in}}'''(x^*(z), z)}. \quad \blacksquare$$

Proving items 2 and 3 in Theorem 2.1 requires more delicate analysis and we leave it to the next subsection. Item 4, however, is a direct corollary of 2.

*Proof of item 4 of Theorem 2.1, assuming item 2.* According to the definition,

$$(4.4) \quad t^\sharp = \inf\{t : x^c(t) \geq x_0\} \quad \Rightarrow \quad x^c(t^\sharp) = x_0,$$

meaning  $t^\sharp = t^\sharp(x_0)$  is the inverse function of  $x^c(t)$  evaluated at  $x_0$ . According to item 2 in Theorem 2.1,  $x^c(t, z)$  is smooth in  $z$ ; then taking the  $z$ -derivative on (4.4) gives

$$\partial_z x^c(t^\sharp, z) + \partial_t x^c(t^\sharp, z) \partial_z t^\sharp = 0,$$

which shows that  $|\partial_z t^\sharp| < \infty$  by the assumption that  $\partial_t x^c(t^\sharp, z) \neq 0$ . Higher order  $z$ -derivatives can be handled in the same way.  $\blacksquare$

We now concentrate on showing items 2 and 3 of Theorem 2.1, which state the smooth dependence of  $x^c$  and  $u_1 - u_2$  on  $z$ . We divide the proof into two parts: we will first prove the smoothness assuming all the initial shocks are generated at  $t^* = 0$  and  $u^* = 0$ , meaning  $u_{1,2}(t = 0, z) = 0$  for all  $z$ ; we then shift  $(t^*, u^*)$  to accommodate the general situation stated in Theorem 2.1. The first part of the proof is summarized in Propositions 4.1 and 4.2, and the second part of the proof follows.

**Proposition 4.1.** *Consider (3.7) with initial condition (3.11) and  $\alpha = 1$ . Suppose the initial profile represented by  $-x'_{\text{in}}(u) = f(u)$  has smooth  $z$ -dependence, i.e.,  $f(u; z) \in C^\infty(\mathbb{R}_u, \mathbb{R}_z)$ ; then for  $t$  small enough,*

(1) *the  $z$ -derivatives of  $u_1, u_2$  satisfy the estimate*

$$\partial_z u_{1,2} = \mathcal{O}\left(t^{1/2}\right),$$

(2) *the higher  $z$ -derivatives of  $u_1, u_2$  satisfy*

$$\partial_z^k u_{1,2} = \mathcal{O}\left(t^{1/2}\right),$$

(3) *the higher  $(z, t)$ -derivatives in time satisfy*

$$\partial_z^k \partial_t^{k'} u_{1,2} = \mathcal{O}\left(t^{1/2-k'}\right).$$

*Proof.* To obtain the regularity in the  $z$  direction, one basically needs to take  $z$ -derivatives and show the bounds. We start with the first order derivative of (3.7) to show item 1:

$$\begin{aligned} \frac{d\partial_z u_1}{dt} &= \frac{1}{2}(\partial_z u_1 - \partial_z u_2) \frac{1}{f(u_1) - t} - \frac{1}{2}(u_1 - u_2) \frac{1}{(f(u_1) - t)^2} (f'(u_1) \partial_z u_1 + \partial_z f(u_1)), \\ \frac{d\partial_z u_2}{dt} &= -\frac{1}{2}(\partial_z u_1 - \partial_z u_2) \frac{1}{f(u_2) - t} + \frac{1}{2}(u_1 - u_2) \frac{1}{(f(u_2) - t)^2} (f'(u_2) \partial_z u_2 + \partial_z f(u_2)). \end{aligned}$$

In a compact form, it becomes

$$(4.5) \quad \begin{cases} \frac{d\partial_z u_1}{dt} = A_{11}\partial_z u_1 + A_{12}\partial_z u_2 + S_1, \\ \frac{d\partial_z u_2}{dt} = A_{21}\partial_z u_1 + A_{22}\partial_z u_2 + S_2, \end{cases}$$

with initial data

$$\partial_z u_1(0) = \partial_z u_2(0) = 0.$$

Here  $A$  terms are the linear terms and  $S_{1,2}$  are sources. The terms in  $A$  and  $S$  can be estimated using the results in Theorem 3.1, which states that  $u_{1,2}(t) \approx \pm(3a^{-1}t)^{1/2}$ , so

$$A_{11} = \frac{1}{2} \frac{1}{f(u_1) - t} - \frac{1}{2}(u_1 - u_2) \frac{1}{(f(u_1) - t)^2} f'(u_1) \approx -\frac{5}{4t},$$

and similarly

$$A_{12} = -\frac{1}{2} \frac{1}{f(u_1) - t} \approx -\frac{1}{4t}, \quad A_{21} \approx -\frac{1}{4t}, \quad A_{22} \approx -\frac{5}{4t}.$$

Here we use the notation  $A(t) \approx B(t)$  to mean  $\lim_{t \rightarrow 0^+} \frac{A(t)}{B(t)} = 1$ . Noting that  $f(u; z) \sim a(z)u^2$  (see (3.12)),  $\partial_z f \sim \partial_z a u^2$ , and thus

$$S_1 = -\frac{1}{2}(u_1 - u_2) \frac{1}{(f(u_1) - t)^2} \partial_z f(u_1) = \mathcal{O}(t^{-1/2}), \quad S_2 = \mathcal{O}(t^{-1/2}).$$

Then we perform the standard energy estimate of  $L^2$  type for (4.5) by multiplying it on both sides with  $\partial_z u_{1,2}$  to have

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} (\partial_z u_1)^2 &\leq -\left(\frac{5}{4} - \epsilon\right) \frac{1}{t} (\partial_z u_1)^2 + \left(\frac{1}{4} + \epsilon\right) \frac{1}{t} |\partial_z u_1 \partial_z u_2| + Ct^{-1/2} |\partial_z u_1|, \\ \frac{1}{2} \frac{d}{dt} (\partial_z u_2)^2 &\leq -\left(\frac{5}{4} - \epsilon\right) \frac{1}{t} (\partial_z u_2)^2 + \left(\frac{1}{4} + \epsilon\right) \frac{1}{t} |\partial_z u_1 \partial_z u_2| + Ct^{-1/2} |\partial_z u_2|. \end{aligned}$$

Adding the two inequalities and use the fact that

$$|\partial_z u_1 \partial_z u_2| \leq \frac{1}{2} ((\partial_z u_1)^2 + (\partial_z u_2)^2) \quad \text{and} \quad t^{-1/2} |\partial_z u_1| \leq \epsilon_1 \frac{1}{t} (\partial_z u_1)^2 + \frac{1}{4\epsilon_1} \quad \forall \epsilon_1 > 0,$$

one gets

$$\frac{1}{2} \frac{d}{dt} ((\partial_z u_1)^2 + (\partial_z u_2)^2) \leq -\left(\frac{5}{4} - \epsilon - \left(\frac{1}{4} + \epsilon\right) - C\epsilon_1\right) \frac{1}{t} ((\partial_z u_1)^2 + (\partial_z u_2)^2) + \frac{C}{2\epsilon_1}.$$

Choosing  $\epsilon = 1/4$ ,  $\epsilon_1 = 1/(2C)$ , one gets

$$\frac{d}{dt} ((\partial_z u_1)^2 + (\partial_z u_2)^2) \leq 2C^2,$$

which finishes the proof of item 1.

Extending it to higher derivatives requires mathematical induction. We assume  $\partial_z^j u_{1,2} = \mathcal{O}(t^{1/2})$  holds true for all  $j < k$ , and we now show it for the  $k$ th derivative as well. Taking the  $k$ th derivative on  $z$  we have

$$(4.6) \quad \begin{cases} \frac{d\partial_z^k u_1}{dt} = A_{11}\partial_z^k u_1 + A_{12}\partial_z^k u_2 + S_1^k, \\ \frac{d\partial_z^k u_2}{dt} = A_{21}\partial_z^k u_1 + A_{22}\partial_z^k u_2 + S_2^k, \end{cases}$$

where  $A_{mn}$  has the same definition as in (4.5) ( $m, n = 1, 2$ ). The source term  $S_1^k$  is, however, much more complicated,

$$(4.7) \quad S_1^k = \sum \frac{c\partial_z^{r_1}(u_1 - u_2)}{(f(u_1) - t)^{1+r_2}} \prod_{j=1}^{r_2} \left( \partial_z^{r_{3,j}} \partial_u^{r_{4,j}} f(u_1) \prod_{l=1}^{r_{4,j}} \partial_z^{r_{5,j,l}} u_1 \right),$$

where  $c$  is a constant depending on the summation indices, and the indices in the summation satisfy the relation

$$r_1 + \sum_{j=1}^{r_2} \left( r_{3,j} + \sum_{l=1}^{r_{4,j}} r_{5,j,l} \right) = k, \quad r_1 \leq k-1, \quad r_{5,j,l} \leq k-1.$$

Noting that

$$f(u_1) \sim u_1^2 = \mathcal{O}(t), \quad f'(u_1) \sim u_1 = \mathcal{O}(t^{1/2}), \quad \text{and} \quad \partial_u^r f(u_1) = \mathcal{O}(1), \quad r \geq 2,$$

and as a result,  $\partial_u^r f(u_1) \lesssim \mathcal{O}(t^{1-r/2})$ ,  $r \geq 0$ . In view of  $\partial_z^l u_1 = \mathcal{O}(t^{1/2})$  for  $l \leq k-1$ , the order of the term (4.7) is (in term of the power of  $t$ )

$$\frac{1}{2} - (1 + r_2) + \sum_{j=1}^{r_2} \left( \left(1 - \frac{r_{4,j}}{2}\right) + \frac{r_{4,j}}{2} \right) = -\frac{1}{2}.$$

The term  $S_2^k$  can be analyzed in the same way. Using the energy estimate we conclude with the result. Item 3 is obtained using the induction argument as well, and we leave the proof to the appendix. ■

**Proposition 4.2.** *With the same assumptions as in Proposition 4.1, one has*

$$\partial_z^k \partial_t^{k'} x^c = \chi_{(k'=0)} \partial_z^k x^*(z) + \mathcal{O}\left(t^{3/2-k'}\right).$$

*Proof.* It follows easily from checking (3.5), which holds true on  $t > t^* = 0$  with  $x^c(0) = x^*$ . Integrating (3.5) in  $t$  we get

$$(4.8) \quad x^c(t, z) = x^*(z) + \int_0^t \frac{u_1(s, z) + u_2(s, z)}{2} ds.$$

Taking its  $(z, t)$ -derivative, we have

$$(4.9) \quad \partial_z^k x^c(t, z) = \partial_z^k x^*(z) + \int_0^t \frac{\partial_z^k u_1(s, z) + \partial_z^k u_2(s, z)}{2} ds$$

when  $k' = 0$  and

$$(4.10) \quad \partial_z^k \partial_t^{k'} x^c = \partial_t^{k'} \left( \int_0^t \frac{\partial_z^k u_1(s, z) + \partial_z^k u_2(s, z)}{2} ds \right) = \frac{\partial_z^k \partial_t^{k'-1} u_1(t, z) + \partial_z^k \partial_t^{k'-1} u_2(t, z)}{2}$$

when  $k' > 0$ . These combined with Proposition 4.1 give the conclusion. ■

With these preparations for equations with special initial data  $u^* = t^* = 0$ , we are ready to perform shifting for proving Theorem 2.1.

*Proof of items 2 and 3 of Theorem 2.1.* We translate  $u_1, u_2$  to enforce the initial condition (3.11) by defining  $\bar{u}_1, \bar{u}_2$  as

$$(4.11) \quad \bar{u}_{1,2}(t, z) = u_{1,2} \left( \frac{1}{\alpha(z)} t + t^*(z), z \right) - u^*(z).$$

$\bar{u}_{1,2}(t, z)$  then satisfies the same system (3.7) with initial condition (3.11) and  $\alpha = 1$ .  $f$ , however, is also shifted:

$$(4.12) \quad \bar{f}(\bar{u}, z) = f(\bar{u} + u^*(z)) - \alpha(z)t^*(z).$$

It is clear that  $\bar{f}$  is smooth and satisfies Assumption 2. By the assumption that there exists  $\delta > 0$  such that  $u_+ + \delta \leq u^*(z) \leq u_- - \delta$  for all  $z$ , we have  $\bar{f}$  well-defined for  $u \in [-\delta, \delta]$ .

According to Proposition 4.1,

$$(4.13) \quad \partial_z^k \partial_t^{k'} \bar{u}_{1,2} = \mathcal{O} \left( t^{1/2-k'} \right),$$

and considering  $u_{1,2}(t, z) = \bar{u}_{1,2}(\alpha(z)(t - t^*(z)), z) + u^*(z)$ , and taking the smoothness of  $\alpha(z)$ ,  $t^*(z)$ , and  $u^*(z)$  into account, we obtain the estimate

$$(4.14) \quad \partial_z^k u_{1,2} = \mathcal{O} \left( (t - t^*)^{1/2-k} \right),$$

which implies item 3 of Theorem 2.1. To estimate  $x^c$ , we change  $u_{1,2}$  to  $\bar{u}_{1,2}$  as well:

$$\begin{aligned} x^c(t, z) &= x^*(z) + \alpha(z) \int_{t^*(z)}^t \frac{u_1(s, z) + u_2(s, z)}{2} ds \\ &= x^*(z) + \alpha(z) \int_{t^*(z)}^t \frac{\bar{u}_1(\alpha(s - t^*(z)), z) + \bar{u}_2(\alpha(s - t^*(z)), z)}{2} + u^*(z) ds \\ &= x^*(z) + \frac{1}{2} \int_0^{\alpha(z)(t - t^*(z))} [\bar{u}_1(s, z) + \bar{u}_2(s, z)] ds + \alpha(z)u^*(z)(t - t^*(z)) \\ &= \bar{x}^c(\alpha(z)(t - t^*(z)), z) + \alpha(z)u^*(z)(t - t^*(z)). \end{aligned}$$

Taking its  $z$ -derivative up to order  $k$  on both sides and using Proposition 4.2 together with the smoothness of  $\alpha(z)$ ,  $t^*(z)$ , and  $u^*(z)$ , we conclude the theorem:

$$\partial_z^k x^c(t, z) = \partial_z^k x^*(z) + \mathcal{O} \left( t^{3/2-k} \right) + \mathcal{O}(1) = \mathcal{O} \left( t^{\min\{3/2-k, 0\}} \right). \quad \blacksquare$$

**4.2. Smoothness of the shifted solution profile.** This is to mathematically justify Theorem 2.2, which guarantees the  $z$ -regularity of the shifted solution  $\tilde{u}$ .

*Proof of Theorem 2.2.* Recall  $\tilde{u}$  in (2.6) and we take its  $k$ th derivative in  $z$ :

$$(4.15) \quad \partial_z^k \tilde{u} = \sum c \partial_t^{r_1} \partial_x^{r_2} \partial_z^{r_3} u \prod_{j=1}^{r_2} \partial_z^{r_{4,j}} t^*(z) \prod_{j=1}^{r_2} \partial_z^{r_{5,j}} (x^c(t + t^*(z), z)),$$

where the indices satisfy

$$(4.16) \quad r_3 + \sum_{j=1}^{r_1} r_{4,j} + \sum_{j=1}^{r_2} r_{5,j} = k.$$

Then we further expand

$$(4.17) \quad \partial_z^{r_{5,j}} (x^c(t + t^*(z), z)) = \sum c \partial_t^{r_{6,j}} \partial_z^{r_{7,j}} x^c \prod_{l=1}^{r_{6,j}} \partial_z^{r_{8,j,l}} t^*(z) \quad \text{with} \quad r_{7,j} + \sum_{l=1}^{r_{6,j}} r_{8,j,l} = r_{5,j}.$$

The second factor in (4.15) is  $\partial_z^{r_{4,j}} t^*(z)$  which is of  $\mathcal{O}(1)$ . To deal with the third factor of (4.15), we use (4.17) and only evaluate  $\partial_t^{r_{6,j}} \partial_z^{r_{7,j}} x^c$  with  $r_{6,j} + r_{7,j} \leq r_{5,j}$ . According to Proposition 4.2 we have

$$(4.18) \quad \partial_t^{r_{6,j}} \partial_z^{r_{7,j}} x^c(t + t^*(z), z) = \mathcal{O} \left( t^{\min\{3/2 - r_{6,j}, 0\}} \right).$$

The first factor  $\partial_t^{r_1} \partial_x^{r_2} \partial_z^{r_3} u$  is more complicated. To do that we take the derivative  $\partial_t^{r_1} \partial_x^{r_2} \partial_z^{r_3}$  of  $x(t, u(t, x, z), z) = x$  for

$$\partial_u x \partial_t^{r_1} \partial_x^{r_2} \partial_z^{r_3} u + \sum c \partial_t^{r'_1} \partial_u^{r'_2} \partial_z^{r'_3} x \prod_{j=1}^{r'_2} \partial_t^{r'_{4,j}} \partial_x^{r'_{5,j}} \partial_z^{r'_{6,j}} u = \chi_{(r_1=r_3=0, r_2=1)}$$

with

$$r'_1 + \sum_{j=1}^{r'_2} r'_{4,j} = r_1, \quad r'_3 + \sum_{j=1}^{r'_2} r'_{6,j} = r_3, \quad \sum_{j=1}^{r'_2} r'_{5,j} = r_2.$$

Since  $x(t, u, z) = x_{\text{in}}(u, z) + ut$  is away from the shock, all the derivatives of  $x$  are of  $\mathcal{O}(1)$ . Thus one can show by induction on  $r_1 + r_2 + r_3$  that

$$(4.19) \quad \partial_t^{r_1} \partial_x^{r_2} \partial_z^{r_3} u = \mathcal{O} \left( |\partial_x u|^{2(r_1+r_2+r_3)-1} \right),$$

where  $|\partial_x u|$  is supposed to be large near shock. We then claim that

$$(4.20) \quad |\partial_x u(t + t^*(z), x + x^c(t + t^*(z), z), z)| \leq \frac{2}{|x|}.$$

In fact, suppose  $x > 0$ , then

$$u(t + t^*, x + x^c, z) - u(t + t^*, x^c, z) = \int_{x^c}^{x+x^c} \partial_x u(t + t^*, y, z) dy \leq x \partial_x u(t + t^*, x + x^c, z),$$

where the inequality holds because  $u(t, x, z)$  is convex in  $x$  for  $x > x^c(t, z)$ , and thus  $\partial_x u(t, y, z) \leq \partial_x u(t, x + x^c, z)$ . Taking absolute value and using  $\partial_x u < 0$ , we get

$$\begin{aligned} |x \partial_x u(t + t^*, x + x^c, z)| &\leq |u(t + t^*, x + x^c, z) - u(t + t^*, x^c, z)| \\ &\leq |u(t + t^*, x + x^c, z)| + |u(t + t^*, x^c, z)| \leq 2, \end{aligned}$$

which leads to (4.20). The case  $x < 0$  is similar.

In conclusion, plugging (4.19) and (4.18) into (4.15), and using the fact that  $r_1 + r_2 + r_3 \leq k$  and  $\sum_{j=1}^{r_2} r_{6,j} \leq \sum_{j=1}^{r_2} r_{5,j} \leq k$ , we conclude the theorem. ■

Theorem 2.3 immediately follows from the following proposition. It is a standard result from approximation theory and we leave the proof to the appendix.

**Proposition 4.3.** *Let  $f = f(z) \in C^{m+1}(-1, 1)$ . Then the polynomial interpolation (2.8) has  $m$ th order accuracy:*

$$(4.21) \quad |f(z) - f^N(z)| \leq \frac{C(m) \|\partial_z^{m+1} f\|_{L^\infty}}{N^m} \quad \forall z \in [-1, 1]$$

for  $N \geq 2m$ . Furthermore, if  $\pi(z)$  is supported on  $[-1, 1]$ , then we have the error estimate

$$(4.22) \quad |\mathbb{E}(f) - \mathbb{E}(f^N)| \leq \frac{C(m) \|\partial_z^{m+1} f\|_{L^\infty}}{N^m},$$

$$(4.23) \quad |\text{var}(f) - \text{var}(f^N)| \leq \frac{C(m) (\min\{\|f - f^N\|_{L^\infty}, N^2 \|f\|_{L^\infty}\} + \|f\|_{L^\infty}) \|\partial_z^{m+1} f\|_{L^\infty}}{N^m}.$$

**5. General scalar conservation laws with convex fluxes.** All the results for the Burgers' equation can be extended to study general scalar conservation laws with convex flux term. The proof itself is tedious but contains little novelty and thus we only outline the strategies. In the general cases, the equation reads

$$(5.1) \quad \begin{cases} \partial_t u + \partial_x F(u) = 0, \\ \lim_{x \rightarrow \pm\infty} u_{\text{in}}(x) = \mp 1, \end{cases}$$

where the flux function  $F$ , deterministic, is smooth and strictly convex. We also assume that the initial data is decreasing, and therefore the inverse function  $x(u)$  is well-defined on  $(-1, 1)$ . Note that the domain can be generalized to treat  $(u_+, u_-)$ . The derivation is the same: one flips  $x - u$  coordinates and derives the equation for  $x(u)$  as a function of  $u$ . The reformulation allows us to obtain an explicit expression for the dynamics of the physical quantities such as  $t^*$ ,  $t^\sharp$ , and  $x^c$ . As done for the Burgers' equation we first reformulate the equation, obtain the ODE system, and study its dependence on the unknown variable  $z$ .

### 5.1. Reformulation of the equation.

**5.1.1. Before shocking emergence.** As done for the Burgers' equation, upon flipping  $x$  and  $u$ , one writes the dynamics of  $x(u)$  as

$$(5.2) \quad \begin{cases} \partial_t x(t, u) = F'(u), & u \in (-1, 1), \\ x_{\text{in}}(u) = x(t = 0, u). \end{cases}$$

Convex  $F$  gives the increasing  $F'$ . We denote the inverse function  $G$ :

$$(5.3) \quad G(F'(u)) = F'(G(u)) = u.$$

Plugging it back into (5.2) and denoting  $y(t, u) = x(t, G(u))$ , we have

$$\begin{cases} \partial_t y(t, u) = \partial_t x(t, G(u)) = u, & u \in (-1, 1), \\ y_{\text{in}}(u) = x(t = 0, G(u)) = x_{\text{in}}(G(u)). \end{cases}$$

As was done for the Burgers' equation, we take one more derivative on  $u$  and obtain

$$\begin{cases} \partial_t \partial_u y(t, u) = 1, & u \in (-1, 1), \\ y'_{\text{in}}(u) = x'_{\text{in}}(G(u))G'(u). \end{cases}$$

Therefore

$$\partial_u y = y'_{\text{in}}(u) + t,$$

and equivalently, the earliest shock appears at  $t^* = -\min y'_{\text{in}}(u)$  and we assume there is one and only one and set it as

$$t^* = -\min y'_{\text{in}}(u) = -y'_{\text{in}}(F'(u^*)).$$

**5.1.2. After the emergence of the shock.** Once the shock appears, on the  $u-x$  plane, a “flat” region appears. We denote  $u_1$  and  $u_2$  the top and the bottom of the shock point, then between  $(u_2, u_1)$  the solution is a constant, which moves horizontally with speed:

$$s = \frac{F(u_1) - F(u_2)}{u_1 - u_2},$$

meaning

$$(5.4) \quad \frac{d}{dt} x^c = \frac{F(u_1) - F(u_2)}{u_1 - u_2} \quad \text{with} \quad x^c(t^*) = x^*,$$

where  $x^c$  denotes the shock location. With the same derivation as in section 3.1, one has

$$(5.5) \quad \begin{cases} \frac{du_1}{dt} = F_1(u_1, u_2) = \left( F'(u_1) - \frac{F(u_1) - F(u_2)}{u_1 - u_2} \right) (f(u_1) - F''(u_1)t)^{-1}, \\ \frac{du_2}{dt} = F_2(u_1, u_2) = - \left( \frac{F(u_1) - F(u_2)}{u_1 - u_2} - F'(u_2) \right) (f(u_2) - F''(u_2)t)^{-1} \end{cases}$$

with initial condition  $u_1(t^*) = u_2(t^*) = u^*$ . Here we denote  $f(u) = -x'_{\text{in}}(u)$ . Considering  $F'$  is an increasing function, we see that

$$\frac{du_1}{dt} > 0 > \frac{du_2}{dt}.$$



**5.1.3. Summary.** To summarize the reformulation, in the general convex flux case, when written on the  $x(u)$  plane,  $x$  satisfies the equation

$$(5.6) \quad \begin{cases} t < t^* = -y'_{\text{in}}(u^*) : & \text{equation (5.2)}, \\ t > t^* : & \begin{cases} \text{equation (5.2)} & \text{with } u \in (-1, u_2) \cup (u_1, 1), \\ \text{equation (5.4)} & \text{with } u \in (u_2, u_1), \end{cases} \end{cases}$$

with  $u_1$  and  $u_2$  being the shock locations satisfying the ODE system (5.5).

**5.2. Shock behavior in small time (general flux).** Assume a shock emerges at  $t^* = 0, u^* = 0$ ; then to understand the short time behavior of the ODE system is equivalent to understanding the forcing terms in (5.5). Near  $u_{1,2} = 0$  we can approximate

$$(5.7) \quad \begin{aligned} F'(u_1) - \frac{F(u_1) - F(u_2)}{u_1 - u_2} &= F'(0) + F''(0)u_1 - \frac{F'(0)(u_1 - u_2) + \frac{1}{2}F''(0)(u_1^2 - u_2^2)}{u_1 - u_2} + \mathcal{O}(u_1^2, u_2^2, u_1u_2) \\ &= \frac{1}{2}F''(0)(u_1 - u_2) + \mathcal{O}(u_1^2, u_2^2, u_1u_2), \end{aligned}$$

and thus in the leading order,

$$(5.8) \quad \begin{cases} \frac{du_1}{dt} = \frac{1}{2}F''(0)(u_1 - u_2)(au_1^2 - F''(0)t)^{-1}, \\ \frac{du_2}{dt} = -\frac{1}{2}F''(0)(u_1 - u_2)(au_2^2 - F''(0)t)^{-1}, \end{cases}$$

where  $a = \frac{a_1 F''(0)^2}{G'(F''(0))}$  is a positive number. For small time, the solution is explicit:

$$(5.9) \quad u_1 = -u_2 = (ct)^{1/2}, \quad c = \frac{3F''(0)}{a}.$$

**5.3. Regularities in the random space.** Studying the solution's regularity in the random space is the same as the analysis carried out in section 4. Due to the complexity of the formula, we present only the first derivative in  $z$  of (5.5). One takes the first derivation of (5.5):

$$(5.10) \quad \begin{aligned} \frac{d\partial_z u_1}{dt} &= \left[ F'_1 \partial_z u_1 - \frac{F'_1 \partial_z u_1 - F'_2 \partial_z u_2}{u_1 - u_2} + \frac{(F_1 - F_2)(\partial_z u_1 - \partial_z u_2)}{(u_1 - u_2)^2} \right] (f_1 - F'_1 t)^{-1} \\ &\quad - \left[ F'_1 - \frac{F_1 - F_2}{u_1 - u_2} \right] (f_1 - F'_1 t)^{-2} (f'_1 \partial_z u_1 - F'''_1 \partial_z u_1 t + \partial_z f_1), \\ \frac{d\partial_z u_2}{dt} &= - \left[ \frac{F'_1 \partial_z u_1 - F'_2 \partial_z u_2}{u_1 - u_2} - \frac{(F_1 - F_2)(\partial_z u_1 - \partial_z u_2)}{(u_1 - u_2)^2} - F''_1 \partial_z u_1 \right] (f_2 - F'_2 t)^{-1} \\ &\quad + \left[ \frac{F_1 - F_2}{u_1 - u_2} - F'_2 \right] (f_2 - F'_2 t)^{-2} (f'_2 \partial_z u_2 - F'''_2 \partial_z u_2 t + \partial_z f_2), \end{aligned}$$

where we have used  $\gamma_{1,2}$  to denote  $\gamma(u_1)$  or  $\gamma(u_2)$ , respectively, for all quantities. In a compact form, it reads

$$\begin{cases} \frac{d\partial_z u_1}{dt} = A_{11} \partial_z u_1 + A_{12} \partial_z u_2 + S_1, \\ \frac{d\partial_z u_2}{dt} = A_{21} \partial_z u_1 + A_{22} \partial_z u_2 + S_2, \end{cases} \quad \text{with } \partial_z u_{1,2}(0) = 0.$$

In the equation,

$$\begin{aligned} A_{11} &= \left[ F_1'' - \frac{F_1'}{u_1 - u_2} + \frac{F_1 - F_2}{(u_1 - u_2)^2} \right] (f_1 - F_1''t)^{-1} - \left[ F_1' - \frac{F_1 - F_2}{u_1 - u_2} \right] (f_1 - F_1''t)^{-2} (f_1' - F_1'''t), \\ A_{12} &= \left[ \frac{F_2'}{u_1 - u_2} - \frac{F_1 - F_2}{(u_1 - u_2)^2} \right] (f_1 - F_1''t)^{-1}, \\ S_1 &= - \left[ F_1' - \frac{F_1 - F_2}{u_1 - u_2} \right] (f_1 - F_1''t)^{-2} \partial_z f_1. \end{aligned}$$

To analyze the term  $A_{11}$ , we note that

$$\begin{aligned} (5.11) \quad F_1'' &\approx F''(0), \quad -\frac{1}{2}F''(0) \approx -\frac{F_1'}{u_1 - u_2} + \frac{F_1 - F_2}{(u_1 - u_2)^2}, \\ f(u_1) - F''(u_1)t &\approx (ac - F''(0))t = 2F''(0)t, \\ f'(u_1) - F'''(u_1)t &\approx 2a(ct)^{1/2}, \end{aligned}$$

which allows us to bound

$$A_{11} \approx \frac{1}{2}F''(0) \cdot (2F''(0)t)^{-1} - 2actF''(0)(2F''(0)t)^{-2} = -\frac{5}{4}t^{-1}.$$

Similarly one has

$$A_{22} \approx -\frac{1}{4}t^{-1} \quad \text{and} \quad S_1 = \mathcal{O}(t^{-1/2}).$$

All together,

$$\frac{d}{dt} [(\partial_z u_1)^2 + (\partial_z u_2)^2] \leq C,$$

and the  $H_1(dz)$  norm of  $u_{1,2}$  grows no more than a rate of  $\mathcal{O}(t^{1/2})$ .

**6. Conclusion.** Uncertainty quantification for hyperbolic conservation laws is considered a very challenging task due to the intrinsic discontinuities in the solution in both physical and random spaces. Such discontinuities in the solution profile prevent the gPC type methods from being effective. We give a counterargument in this paper, and we demonstrate, under some mild assumptions on the initial condition, that

1. there exist physical observables depending smoothly on external randomness;
2. with proper shifts of the solution in time and space, the entire solution profile also smoothly depends on the external randomness.

We have to emphasize that the main goal of the paper is not to justify the gPC method's use on hyperbolic systems, but rather to provide a new perspective: for wave-like equations with randomness, the solution profile may not be the right "quantity of interests" to evaluate, and a slight change (the proper shifts) could regularize the problem significantly.

**Appendix A. Supplementary proofs.** For the completeness of the paper we include the proofs with tedious calculation here.

**A.1. Wellposedness of the ODE system (3.7).** We show the wellposedness of the ODE system (3.7). In fact, the two forcing terms  $F_1$  and  $F_2$  in (3.7) are Lipschitz continuous on  $u_1$  and  $u_2$  if  $f(u_{1,2}) - \alpha t$  are away from 0, and the lemma below shows that they keep being Lipschitz as long as  $f(u_{1,2}) - \alpha t > 0$ .

**Lemma A.1.** *Assume  $u_{1,2}(t)$  solves (3.7) with  $f(u_{1,2}(t_1)) - \alpha t_1 > 0$  for some  $t_1$ . Then there exists  $c > 0$  such that  $f(u_{1,2}(t)) - \alpha t > c$  for all  $t > t_1$ .*

*Proof.* Using (3.7), we obtain

$$(A.1) \quad \frac{d}{dt}(f(u_1) - \alpha t) = f'(u_1) \frac{d}{dt} u_1 - \alpha = \frac{\alpha}{2}(u_1 - u_2) f'(u_1) \frac{1}{f(u_1) - \alpha t} - \alpha.$$

Since  $u_1$  is increasing,  $u_2$  is decreasing, and  $f'(u) > 0$  is increasing for  $u > u^*$ , one has

$$(A.2) \quad \frac{1}{2}(u_1 - u_2) f'(u_1) \geq \left[ \frac{1}{2}(u_1 - u_2) f'(u_1) \right] \Big|_{t=t_1} =: c_1 > 0 \quad \forall t \geq t_1.$$

According to the ODE (A.1),

$$(A.3) \quad f(u_1) - \alpha t \geq \min\{c_1/2, (f(u_1) - \alpha t)_{t=t_1}\} =: c > 0 \quad \forall t \geq t_1.$$

In fact, if at any time  $t$  one has  $0 < f(u_1) - \alpha t < \frac{2}{3}c_1$ , then one has  $\frac{d}{dt}(f(u_1) - \alpha t) \geq \alpha(c_1 \frac{3/2}{c_1} - 1) = \alpha/2$ . Therefore starting from  $t = t_1$ ,  $f(u_1) - \alpha t$  keeps increasing unless it becomes larger than  $\frac{2}{3}c_1$ . This implies (A.3).

The proof for  $f(u_2) - t$  is similar. ■

**A.2. Proof for item 3 in Proposition 4.1.** The proof is moved here merely because of the highly involved calculation. The idea still follows that for the rest of the proposition.

*Proof.* We use induction on  $(k, k')$ . Since we already have the cases  $(k, 0)$  (Proposition 4.1), we may assume that all cases  $(j, j')$  with  $j < k$  and  $j = k, j' \leq k'$  are already proved, and then we prove the case  $(k, k' + 1)$ . Taking the  $k'$ th  $t$ -derivative of (4.6) gives

$$(A.4) \quad \begin{aligned} \partial_z^k \partial_t^{k'+1} u_1 &= S_1^{k,k'} := \partial_t^{k'} (A_{11} \partial_z^k u_1 + A_{12} \partial_z^k u_2 + S_1^k), \\ \partial_z^k \partial_t^{k'+1} u_2 &= S_2^{k,k'} := \partial_t^{k'} (A_{21} \partial_z^k u_1 + A_{22} \partial_z^k u_2 + S_2^k). \end{aligned}$$

Notice that every term in  $(A_{11} \partial_z^k u_1 + A_{12} \partial_z^k u_2 + S_1^k)$  is of the form (4.7), with possibly  $r_1 = k$  or  $r_{5,j,l} = k$ . Taking  $\partial_t^{k'}$  of (4.7) gives terms of the form

$$(A.5) \quad \frac{c \partial_z^{r_1} \partial_t^{r'_1} (u_1 - u_2)}{(f(u_1) - t)^{1+r_2+r'_2}} \prod_{j=1}^{r'_2} \partial_t^{r'_{3,j}} (f(u_1) - t) \prod_{j=1}^{r_2} \left( \partial_z^{r_{3,j}} \partial_u^{r_{4,j}+r'_{4,j}} f(u_1) \prod_{l=1}^{r'_{4,j}} \partial_t^{r'_{5,j,l}} u_1 \prod_{l=1}^{r_{4,j}} \partial_z^{r_{5,j,l}} \partial_t^{r'_{6,j,l}} u_1 \right)$$

with

$$(A.6) \quad r'_1 + \sum_{j=1}^{r'_2} r'_{3,j} + \sum_{j=1}^{r_2} \left( \sum_{l=1}^{r'_{4,j}} r'_{5,j,l} + \sum_{l=1}^{r_{4,j}} r'_{6,j,l} \right) = k'.$$

Here we can further write

$$(A.7) \quad \partial_t^{r'_{3,j}}(f(u_1) - t) = \sum \partial_u^{r'_{7,j}} f(u_1) \prod_{l=1}^{r'_{7,j}} \partial_t^{r'_{8,j,l}} u_1 - \chi_{(r'_{3,j}=1)}, \quad \sum_{l=1}^{r'_{7,j}} r'_{8,j,l} = r'_{3,j},$$

where  $\chi_{(r'_{3,j}=1)}$  means 1 when  $r'_{3,j} = 1$  and 0 otherwise.

Therefore, the power of  $t$  of the term (A.5) is

$$(A.8) \quad \begin{aligned} & \frac{1}{2} - r'_1 - (1 + r_2 + r'_2) + \sum_{j=1}^{r'_2} \left( 1 - \frac{r'_{7,j}}{2} + \sum_{l=1}^{r'_{7,j}} \left( \frac{1}{2} - r'_{8,j,l} \right) \right) \\ & + \sum_{j=1}^{r_2} \left( 1 - \frac{r_{4,j} + r'_{4,j}}{2} + \sum_{l=1}^{r'_{4,j}} \left( \frac{1}{2} - r'_{5,j,l} \right) + \sum_{l=1}^{r_{4,j}} \left( \frac{1}{2} - r'_{6,j,l} \right) \right) \\ & = \frac{1}{2} - r'_1 - (1 + r_2 + r'_2) + \sum_{j=1}^{r'_2} \left( 1 - \sum_{l=1}^{r'_{7,j}} r'_{8,j,l} \right) + \sum_{j=1}^{r_2} \left( 1 - \sum_{l=1}^{r'_{4,j}} r'_{5,j,l} - \sum_{l=1}^{r_{4,j}} r'_{6,j,l} \right) \\ & = -\frac{1}{2} - r'_1 - \sum_{j=1}^{r'_2} r'_{3,j} - \sum_{j=1}^{r_2} \left( \sum_{l=1}^{r'_{4,j}} r'_{5,j,l} + \sum_{l=1}^{r_{4,j}} r'_{6,j,l} \right) = -\frac{1}{2} - k' = \frac{1}{2} - (k' + 1). \end{aligned}$$

Notice that in the case  $r'_{3,j} = 1$  the term  $\sum \partial_u^{r'_{7,j}} f(u_1) \prod_{l=1}^{r'_{7,j}} \partial_t^{r'_{8,j,l}} u_1$  has  $t$  power 0, the same as the term  $\chi_{(r'_{3,j}=1)} = 1$ , thus the latter can be ignored. This finishes the induction for  $(k, k' + 1)$ . ■

**A.3. Proof of Proposition 4.3.** We first state a classical result from approximation theory [21, Theorem 7.2].

**Lemma A.2.** *For an integer  $m \geq 1$ , let  $f = f(z)$  and its derivatives through  $f^{(m-1)}$  be absolutely continuous on  $[-1, 1]$  and suppose the  $m$ th derivative  $\partial_z^m f$  is of bounded variation  $V$ . Then for any  $n > m$ , its Chebyshev interpolants (2.8) satisfy*

$$(A.9) \quad \|f - f^N\|_{L^\infty} \leq \frac{4V}{\pi m(N - m)^m}.$$

With this, we can show the following.

*Proof of Proposition 4.3.* (4.21) follows from Lemma A.2 by noticing that  $V \leq 2\|\partial_z^{m+1} f\|_{L^\infty}$  and  $N - m \geq N/2$  if  $N \geq 2m$ . To see (4.22), we use

$$\begin{aligned} |\mathbb{E}(f) - \mathbb{E}(f^N)| &= \left| \int (f^N(z) - f(z)) \pi(z) dz \right| \\ &\leq \frac{C(m)\|\partial_z^{m+1} f\|_{L^\infty}}{N^m} \int \pi(z) dz = \frac{C(m)\|\partial_z^{m+1} f\|_{L^\infty}}{N^m}. \end{aligned}$$

To show (4.23) is similar:

$$\begin{aligned}
 \text{(A.10)} \quad & |\text{var}(f) - \text{var}(f^N)| \\
 & \leq \int |f^N(z)^2 - f(z)^2| \pi(z) \, dz + |(\mathbb{E}(f^N))^2 - (\mathbb{E}(f))^2| \\
 & \leq \|f^N + f\|_{L^\infty} \|f^N - f\|_{L^\infty} + |\mathbb{E}(f^N) + \mathbb{E}(f)| \cdot |\mathbb{E}(f^N) - \mathbb{E}(f)| \\
 & \leq (2\|f\|_{L^\infty} + \|f^N - f\|_{L^\infty}) \|f^N - f\|_{L^\infty} + (2\|f\|_{L^\infty} + \|f^N - f\|_{L^\infty}) |\mathbb{E}(f^N) - \mathbb{E}(f)| \\
 & \leq \frac{C(m)(\|f^N - f\|_{L^\infty} + \|f\|_{L^\infty}) \|\partial_z^{m+1} f\|_{L^\infty}}{N^m},
 \end{aligned}$$

where in the third inequality we used

$$\text{(A.11)} \quad \|f^N + f\|_{L^\infty} \leq \|f^N - f\|_{L^\infty} + 2\|f\|_{L^\infty},$$

and in the last inequality we used (4.21). To finally obtain (4.23), we define the piecewise linear function  $f_1(z)$  by

$$\text{(A.12)} \quad f_1(z) = f(z_j) + \frac{f(z_{j+1}) - f(z_j)}{z_{j+1} - z_j} (z - z_j), \quad \text{for } z_j \leq z < z_{j+1},$$

so that  $f_1$  is absolutely continuous and satisfies  $f_1(z_j) = f(z_j)$  for every  $j$ . Thus its Chebyshev interpolant is also  $f^N$ . Since

$$\text{(A.13)} \quad f'_1(z) = \frac{f(z_{j+1}) - f(z_j)}{z_{j+1} - z_j}, \quad \text{for } z_j \leq z < z_{j+1},$$

is a piecewise constant function, whose total variation is

$$\begin{aligned}
 \text{(A.14)} \quad V &= \sum_{j=1}^{N-1} \left| \frac{f(z_{j+1}) - f(z_j)}{z_{j+1} - z_j} - \frac{f(z_j) - f(z_{j-1})}{z_j - z_{j-1}} \right| \leq 2 \sum_{j=1}^{N-1} \left| \frac{f(z_{j+1}) - f(z_j)}{z_{j+1} - z_j} \right| \\
 &\leq 4N \|f\|_{L^\infty} \sum_{j=1}^{N-1} \frac{1}{z_{j+1} - z_j} \leq C \|f\|_{L^\infty} N^3,
 \end{aligned}$$

where we used  $z_{j+1} - z_j \geq \frac{C}{N^2}$ , which is easily checked<sup>3</sup> by using the mean value theorem for the function  $\cos z$ , then Lemma A.2 for  $f_1$  with  $m = 1$  gives

$$\text{(A.15)} \quad \|f^N\|_{L^\infty} \leq C \|f\|_{L^\infty} N^2,$$

and thus

$$\text{(A.16)} \quad \|f^N - f\|_{L^\infty} \leq C \|f\|_{L^\infty} N^2.$$

This combined with (A.10) gives (4.23). ■

---

<sup>3</sup>In fact,  $z_{j+1} - z_j = \cos(\frac{2N+1-2(j+1)}{2N} \pi) - \cos(\frac{2N+1-2j}{2N} \pi) = \sin \xi \cdot (\frac{2N+1-2j}{2N} \pi - \frac{2N+1-2(j+1)}{2N} \pi) = \sin \xi \cdot \frac{\pi}{N}$ , where  $\xi \in (\frac{2N+1-2(j+1)}{2N} \pi, \frac{2N+1-2j}{2N} \pi) \subset (\frac{1}{2N} \pi, \frac{2N-1}{2N} \pi)$ . Therefore  $\sin \xi \geq \frac{2}{\pi} \cdot \frac{1}{2N} \pi = \frac{1}{N}$ , and we get  $z_{j+1} - z_j \geq \frac{\pi}{N^2}$ .

## REFERENCES

- [1] R. ABGRALL, P. CONGEDO, AND G. GERACI, *A one-time truncate and encode multiresolution stochastic framework*, *J. Comput. Phys.*, 257 (2014), pp. 19–56.
- [2] I. BABUŠKA, F. NOBILE, AND R. TEMPONE, *A stochastic collocation method for elliptic partial differential equations with random input data*, *SIAM J. Numer. Anal.*, 45 (2007), pp. 1005–1034.
- [3] I. BABUŠKA, R. TEMPONE, AND G. E. ZOURARIS, *Galerkin finite element approximations of stochastic elliptic partial differential equations*, *SIAM J. Numer. Anal.*, 42 (2004), pp. 800–825.
- [4] M. BRANICKI AND A. J. MAJDA, *Fundamental limitations of polynomial chaos for uncertainty quantification in systems with intermittent instabilities*, *Commun. Math. Sci.*, 11 (2013), pp. 55–103.
- [5] H. CHO, D. VENTURI, AND G. E. KARNIADAKIS, *Statistical analysis and simulation of random shocks in stochastic Burgers equation*, *Proc. A*, 470 (2014), 20140080.
- [6] B. DESPRES AND B. PERTHAME, *Uncertainty propagation: Intrusive kinetic formulations of scalar conservation laws*, *SIAM/ASA J. Uncertain. Quantif.*, 4 (2016), pp. 980–1013.
- [7] O. G. ERNST, A. MUGLER, H.-J. STARKLOFF, AND E. ULLMANN, *On the convergence of generalized polynomial chaos expansions*, *ESAIM Math. Model. Numer. Anal.*, 46 (2012), pp. 317–339.
- [8] G. GERACI, P. M. CONGEDO, R. ABGRALL, AND G. IACCARINO, *A novel weakly-intrusive non-linear multiresolution framework for uncertainty quantification in hyperbolic partial differential equations*, *J. Sci. Comput.*, 66 (2016), pp. 358–405.
- [9] R. G. GHANEM AND R. M. KRUGER, *Numerical solution of spectral stochastic finite element systems*, *Comput. Methods Appl. Mech. Engrg.*, 129 (1996), pp. 289–303.
- [10] M. D. GUNZBURGER, C. G. WEBSTER, AND G. ZHANG, *Stochastic finite element methods for partial differential equations with random input data*, *Acta Numer.*, 23 (2014), pp. 521–650.
- [11] T. Y. HOU, Q. LI, AND P. ZHANG, *Exploring the locally low dimensional structure in solving random elliptic PDEs*, *Multiscale Model. Simul.*, 15 (2017), pp. 661–695.
- [12] T. Y. HOU, Q. LI, AND P. ZHANG, *A sparse decomposition of low rank symmetric positive semidefinite matrices*, *Multiscale Model. Simul.*, 15 (2017), pp. 410–444.
- [13] T. Y. HOU, W. LUO, B. ROZOVSKII, AND H.-M. ZHOU, *Wiener chaos expansions and numerical solutions of randomly forced equations of fluid mechanics*, *J. Comput. Phys.*, 216 (2006), pp. 687–706.
- [14] L. LANDAU AND E. LIFSHITZ, *Fluid Mechanics: Course of Theoretical Physics*, Amsterdam, Elsevier Science, 1959.
- [15] S. MISHRA, N. H. RISEBRO, C. SCHWAB, AND S. TOKAREVA, *Numerical solution of scalar conservation laws with random flux functions*, *SIAM/ASA J. Uncertain. Quantif.*, 4 (2016), pp. 552–591.
- [16] S. MISHRA AND C. SCHWAB, *Sparse tensor multi-level Monte Carlo finite volume methods for hyperbolic conservation laws with random initial data*, *Math. Comp.*, 81 (2012), pp. 1979–2018.
- [17] S. MISHRA AND C. SCHWAB, *Monte-Carlo finite-volume methods in uncertainty quantification for hyperbolic conservation laws*, in *Uncertainty Quantification for Hyperbolic and Kinetic Equations*, S. Jin and L. Pareschi, eds., Springer, Cham, 2017, pp. 231–277.
- [18] 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 (2008), pp. 2309–2345, <https://doi.org/10.1137/060663660>.
- [19] G. POETTE, B. DESPRES, AND D. LUCOR, *Uncertainty quantification for systems of conservation laws*, *J. Comput. Phys.*, 228 (2009), pp. 2443–2467.
- [20] C. SCHWAB AND S. TOKAREVA, *High order approximation of probabilistic shock profiles in hyperbolic conservation laws with uncertain initial data*, *ESAIM Math. Model. Numer. Anal.*, 47 (2013), pp. 807–835, <https://doi.org/10.1051/m2an/2012060>.
- [21] L. N. TREFETHEN, *Approximation Theory and Approximation Practice*, SIAM, Philadelphia, 2012.
- [22] G. WELPER, *Interpolation of functions with parameter dependent jumps by transformed snapshots*, *SIAM J. Sci. Comput.*, 39 (2017), pp. A1225–A1250.
- [23] D. XIU, *Numerical Methods for Stochastic Computations: A Spectral Method Approach*, Princeton University Press, Princeton, NJ, 2010.
- [24] D. XIU AND J. S. HESTHAVEN, *High-order collocation methods for differential equations with random inputs*, *SIAM J. Sci. Comput.*, 27 (2005), pp. 1118–1139.

- [25] D. XIU AND G. E. KARNIADAKIS, *Modeling uncertainty in flow simulations via generalized polynomial chaos*, J. Comput. Phys., 187 (2003), pp. 137–167.
- [26] G. ZHANG AND M. GUNZBURGER, *Error analysis of a stochastic collocation method for parabolic partial differential equations with random input data*, SIAM J. Numer. Anal., 50 (2012), pp. 1922–1940.