

GEOMETRIC QUASILINEARIZATION FRAMEWORK FOR ANALYSIS AND DESIGN OF BOUND-PRESERVING SCHEMES

KAILIANG WU* AND CHI-WANG SHU†

Abstract. Solutions to many partial differential equations satisfy certain bounds or constraints. For example, the density and pressure are positive for equations of fluid dynamics, and in the relativistic case the fluid velocity is upper bounded by the speed of light, etc. As widely realized, it is crucial to develop bound-preserving numerical methods that preserve such intrinsic constraints. Exploring provably bound-preserving schemes has attracted much attention and is actively studied in recent years. This is however still a challenging task for many systems especially those involving nonlinear constraints.

Based on some key insights from geometry, we systematically propose an innovative and general framework, referred to as geometric quasilinearization (GQL), which paves a new effective way for studying bound-preserving problems with nonlinear constraints. The essential idea of GQL is to *equivalently* transfer all nonlinear constraints into *linear* ones, through properly introducing some free auxiliary variables. We establish the fundamental principle and general theory of GQL via the geometric properties of convex regions, and propose three simple effective methods for constructing GQL. We apply the GQL approach to a variety of partial differential equations, and demonstrate its effectiveness and remarkable advantages for studying bound-preserving schemes, by diverse challenging examples and applications which cannot be easily handled by direct or traditional approaches.

Key words. Geometric quasilinearization, nonlinear constraints, bound-preserving numerical schemes, time-dependent PDE systems, convex invariant regions, hyperbolic conservation laws

AMS subject classifications. 65M08, 65M60, 65M12, 65M06, 35L65

1. Introduction. Solutions to many partial differential equations (PDEs) satisfy certain algebraic constraints, which are usually derived from some (physical) bound principles, for example, the positivity of density and pressure. Consider such time-dependent PDE systems in a general form

$$\partial_t \mathbf{u} + \mathcal{L}(\mathbf{u}) = \mathbf{0}, \quad \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}), \quad (1.1)$$

where \mathcal{L} denotes the differential operator associated with the spatial coordinates \mathbf{x} , and suppose the system (1.1) is defined in a bounded domain with suitable boundary conditions. An important class of such systems, which we are particularly interested in, are the hyperbolic conservation laws:

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) = \mathbf{0}, \quad \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}), \quad (1.2)$$

and other related hyperbolic or convection dominated equations.

Assume that the algebraic constraints (bound principles) can be expressed by either the positivity or the non-negativity of several (linear or nonlinear) functions of \mathbf{u} as

$$g_i(\mathbf{u}) > 0 \quad \forall i \in \mathbb{I}, \quad g_i(\mathbf{u}) \geq 0 \quad \forall i \in \widehat{\mathbb{I}}, \quad (1.3)$$

where $\mathbb{I} \cup \widehat{\mathbb{I}} = \{1, \dots, I\}$ with the positive integer I denoting the total number of the constraints. In other words, the evolved variables $\mathbf{u} = (u_1, \dots, u_N)^\top$ belong to the admissible state set:

$$G = \left\{ \mathbf{u} \in \mathbb{R}^N : g_i(\mathbf{u}) > 0 \quad \forall i \in \mathbb{I}, \quad g_i(\mathbf{u}) \geq 0 \quad \forall i \in \widehat{\mathbb{I}} \right\}. \quad (1.4)$$

Throughout this paper, we assume G is convex, which is valid for many physical systems (several typical examples will be given in [section 2](#)). It is worth noting that the functions $\{g_i(\mathbf{u}), 1 \leq i \leq I\}$ are *not* necessarily concave (and *not* required to be concave in this paper). Moreover, we assume that G is an *invariant region* for the exact solution of the system (1.1), namely,

- If $\mathbf{u}(\mathbf{x}, 0) \in G$ for all \mathbf{x} , then $\mathbf{u}(\mathbf{x}, t) \in G$ for all \mathbf{x} and $t > 0$.

A basic goal behind the design of numerical methods solving (1.1) is that they can inherit as much as possible the intrinsic properties of the system (1.1). The constraints (1.3) and the associated invariant region G carry important properties of the exact solution. It is natural and meaningful to explore bound-preserving schemes that keep the numerical solutions within the region G :

*Department of Mathematics, Southern University of Science and Technology, Shenzhen, Guangdong 518055, China (wukl@sustech.edu.cn). K. Wu is supported in part by NSFC grant 12171227.

†Division of Applied Mathematics, Brown University, Providence, RI 02912, USA (chi-wang.shu@brown.edu). C.-W. Shu is supported in part by NSF grant DMS-2010107 and AFOSR grant FA9550-20-1-0055.

- If $\mathbf{u}_h(\cdot, t_0) \in G$, then $\mathbf{u}_h(\cdot, t_n) \in G$ for all $n \in \mathbb{N}$,

where $\mathbf{u}_h(\cdot, t_n)$ denotes the numerical solutions at n th time level. In fact, preserving such constraints is not only necessary for physical significance, but also very crucial for theoretical analysis and numerical stability. If any of the intrinsic physical constraints (1.3) are violated numerically, the PDE system (1.1) and its discrete equations may become ill-posed outside the physical regimes. For example, when negative density and/or pressure are produced in numerically solving the compressible Euler equations, the key hyperbolicity of the system would be lost. As a result, failure to preserve such physically relevant constraints may cause serious numerical problems, for example, nonlinear instability, nonphysical solutions or phenomena, blowups of the code, etc. Therefore, it is significant and highly desirable to develop bound-preserving schemes.

In the past decades, the exploration of bound-preserving high-order numerical methods has attracted extensive attention and is actively studied, especially for hyperbolic and convection dominated equations (e.g. [40, 67, 68, 69, 71, 72, 63, 23, 51, 55]), and recently for some other types of time-dependent PDEs (e.g. [44, 8, 15, 25, 32]). For example, a general framework was established in [67, 68] for constructing bound-preserving high-order finite volume and discontinuous Galerkin schemes for scalar conservation laws and compressible Euler equations. A key step in this framework is to look for high-order schemes that have a provable “weak” bound-preserving property keeping the cell averages of the numerical solutions in the region G . Once such a property is proven, a simple scaling limiter can be used to enforce the constraints for the numerical solutions at any specified points [67, 68, 71]. The idea of this methodology has been applied to many other hyperbolic or convection dominated systems; see, for example, [60, 69, 9, 43, 10, 7, 41, 65, 66, 58, 24, 13, 59]. Another bound-preserving framework [63, 23, 33] is built on flux-correction limiters, which modify any high-order numerical fluxes to enforce the constraints by combining a provably bound-preserving (lower-order) numerical flux as the building block. This approach has also been applied to various physical systems (cf. [11, 12, 61, 56, 62, 50]). Recently, continuous finite element approximations with convex limiting were developed in [21, 22, 20] to preserve invariant regions for hyperbolic equations. Thorough reviews on bound-preserving efforts can be found in the survey articles [64, 45].

Yet, due to the lack of a general theory, how to rigorously analyze or prove whether a numerical scheme is genuinely bound-preserving remains a challenging task. Despite the success of the limiter-based frameworks (cf. [67, 68, 63, 23]) in constructing high-order bound-preserving schemes, the validity of those limiters is actually based on some (weak or lower-order) bound-preserving properties of the cell-average schemes and/or of the numerical fluxes as the key building blocks. Proving such properties is therefore necessary, but often very difficult [45, 51, 55]. To illustrate the challenges, we suppose that a numerical scheme for (1.1) may be written as

$$\mathbf{u}_j^{n+1} = \mathcal{E}_h(\mathbf{u}_{j-k}^n, \mathbf{u}_{j-k+1}^n, \dots, \mathbf{u}_j^n, \dots, \mathbf{u}_{j+s-1}^n, \mathbf{u}_{j+s}^n), \quad (1.5)$$

where \mathcal{E}_h is the discretization operator, the superscripts on \mathbf{u} denote the time levels, and the subscripts on \mathbf{u} indicate the indexes of the spatial grid or nodal points. The bound-preserving problem for the scheme (1.5) can boil down to answer

$$\text{whether } \mathbf{u}_j^n \in G \quad \forall j \quad \text{implies} \quad \mathbf{u}_j^{n+1} \in G \quad \forall j ?$$

In essence, it is to explore whether or not the range of the high-dimensional function \mathcal{E}_h is always contained in G : $\mathcal{E}_h(G^{s+k+1}) \subseteq G$. For some scalar PDEs with linear constraints, for instance, the scalar conservation laws with the constraints linearly defined by maximum principle, a general approach for bound-preserving analysis and design is to exploit certain monotonicity in schemes; see, e.g., [67, 14, 31]. Yet, for PDE systems especially with nonlinear constraints, there is no unified tool like monotonicity, so that direct and complicated algebraic verification usually has to be performed for each constraint case-by-case for different schemes and different PDEs; see, e.g., [68, 38, 56, 41, 66, 35, 53]. Therefore, the design and analysis of bound-preserving schemes involving nonlinear constraints are highly nontrivial, even for first-order schemes; cf. [49, 2, 48, 39, 26, 34, 36, 57, 51].

Nonlinear constraints widely exist in many physical PDE systems; see several representative examples in section 2. For instance, the physical constraints for solutions of the special relativistic magnetohydrodynamic (MHD) equations (2.15) include: the positivity of density D and thermal pressure p , and the upper bound of fluid velocity field \mathbf{v} by the speed of light c , namely,

$$D > 0, \quad p(\mathbf{u}) > 0, \quad c - \|\mathbf{v}(\mathbf{u})\| > 0, \quad (1.6)$$

where the evolved variables $\mathbf{u} = (D, \mathbf{m}, \mathbf{B}, E)^\top$ with the momentum vector $\mathbf{m} \in \mathbb{R}^3$, the magnetic field $\mathbf{B} \in \mathbb{R}^3$, and the total energy E ; see [Example 2.7](#) and [\[55\]](#) for more details. *The second and third constraints in (1.6) are highly nonlinear with respect to \mathbf{u} , because $p(\mathbf{u})$ and $\mathbf{v}(\mathbf{u})$ cannot be explicitly formulated in terms of \mathbf{u} .* These implicit functions $p(\mathbf{u})$ and $\mathbf{v}(\mathbf{u})$ are often expressed via another implicit function $\hat{\phi}(\mathbf{u})$ as

$$p(\mathbf{u}) = \frac{\Gamma - 1}{\Gamma \Upsilon_{\mathbf{u}}^2(\hat{\phi})} \left(\hat{\phi} - D \Upsilon_{\mathbf{u}}(\hat{\phi}) \right), \quad \mathbf{v}(\mathbf{u}) = \left(\mathbf{m} + (\mathbf{m} \cdot \mathbf{B}) \mathbf{B} / \hat{\phi} \right) / (\hat{\phi} + |\mathbf{B}|^2), \quad (1.7)$$

where $\hat{\phi} = \hat{\phi}(\mathbf{u})$ is implicitly defined by the positive root of the nonlinear function $F(\phi; \mathbf{u}) := \phi - E + \|\mathbf{B}\|^2 - \frac{1}{2} \left(\frac{(\mathbf{m} \cdot \mathbf{B})^2}{\phi^2} + \frac{\|\mathbf{B}\|^2}{\Upsilon_{\mathbf{u}}^2(\phi)} \right) + \frac{\Gamma - 1}{\Gamma} \left(\frac{D}{\Upsilon_{\mathbf{u}}(\phi)} - \frac{\phi}{\Upsilon_{\mathbf{u}}^2(\phi)} \right)$, the constant Γ is the ratio of specific heats, and $\Upsilon_{\mathbf{u}}(\phi) := \left(\frac{\phi^2(\phi + \|\mathbf{B}\|^2)^2 - [\phi^2 \|\mathbf{m}\|^2 + (2\phi + \|\mathbf{B}\|^2)(\mathbf{m} \cdot \mathbf{B})^2]}{\phi^2(\phi + \|\mathbf{B}\|^2)^2} \right)^{-\frac{1}{2}}$. If we substitute a scheme [\(1.5\)](#) into the implicit functions $p(\mathbf{u})$ and $\mathbf{v}(\mathbf{u})$, then evaluating these implicit functions and analytically verifying the nonlinear constraints in [\(1.6\)](#) for the scheme [\(1.5\)](#) are indeed very complicated and difficult (if not impossible).

In this paper we discover that, through properly introducing some extra auxiliary variables *independent* of the system variables \mathbf{u} , nonlinear constraints can be *equivalently represented* by using only *linear* constraints, if the region G is convex. For example, the simple nonlinear constraint

$$g(\mathbf{u}) = u_2 - u_1^2 > 0 \quad (1.8)$$

is exactly equivalent to¹

$$\varphi(\mathbf{u}; \theta_*) := u_2 - 2u_1\theta_* + \theta_*^2 > 0 \quad \forall \theta_* \in \mathbb{R}, \quad (1.9)$$

where the extra parameter θ_* is *independent* of \mathbf{u} and called *free auxiliary variable* in this paper. Clearly, the new constraint [\(1.9\)](#) becomes linear² with respect to \mathbf{u} . As we will show, such equivalent linear representation can be found for general nonlinear constraints, even if the constraints cannot be explicitly formulated. For instance, as it will be shown in [Theorem 4.13](#), the constraints in [\(1.6\)](#) can be equivalently represented as

$$D > 0, \quad \mathbf{u} \cdot \mathbf{n}_* + p_m^* > 0 \quad \forall \mathbf{B}_* \in \mathbb{R}^3 \quad \forall \mathbf{v}_* \in \mathbb{B}_1(\mathbf{0}), \quad (1.10)$$

where $\{\mathbf{B}_*, \mathbf{v}_*\}$ are the free auxiliary variables; the vector \mathbf{n}_* and scalar p_m^* are functions of $\{\mathbf{B}_*, \mathbf{v}_*\}$, defined by [\(4.20\)–\(4.21\)](#); $\mathbb{B}_1(\mathbf{0}) := \{\mathbf{x} \in \mathbb{R}^3 : \|\mathbf{x}\| < 1\}$. Note that the equivalent constraints in [\(1.10\)](#) are all linear with respect to \mathbf{u} . Benefited from such linearity, this novel equivalent form [\(1.10\)](#) has significant advantages over the original form [\(1.6\)](#) in designing and analytically analyzing the bound-preserving schemes [\[55\]](#). Several important questions naturally arise: Are there any intrinsic mechanisms behind such an equivalent linear representation? What is the condition for its existence? In general, how to find or construct it?

The aim of this article is to establish a universal framework, termed as geometric quasilinearization (GQL), for constructing equivalent linear representations for general nonlinear constraints. It will be based on some key insights from geometry to understand a convex region G . The GQL framework would shed new light on challenging bound-preserving problems involving nonlinear constraints. The novelty and significance of the proposed GQL framework include:

- A distinctive innovation of GQL lies in a novel geometric point of view on the nonlinear algebraic constraints and the convex invariant region G .
- Through introducing some extra free auxiliary variables, this framework provides a simple yet unified approach to derive the equivalent linear representation (termed as GQL representation) for a general convex region G .
- GQL offers a highly effective approach for bound-preserving analysis and design for problems with nonlinear constraints.
- The GQL representations have simple formulations and are very easy to construct. We will propose three effective methods for constructing GQL.

¹The equivalence of [\(1.8\)](#) and [\(1.9\)](#) can be easily proven by $\min_{\theta_* \in \mathbb{R}} \varphi(\mathbf{u}; \theta_*) = g(\mathbf{u})$.

²This paper broadly uses the word ‘‘linear’’, which means ‘‘affine’’ for functions or constraints with respect to \mathbf{u} .

The idea of GQL is motivated from a series of our recent works on seeking bound-preserving schemes for the (single-component) compressible MHD systems [51, 57, 53, 54, 55]. For the invariant region of the ideal MHD equations, its equivalent linear representation was first established by technical algebraic manipulations [51]. Such a representation played crucial roles in obtaining the first rigorous positivity-preserving analysis of numerical schemes for the ideal MHD system [51], and also in designing the provably positivity-preserving multidimensional MHD schemes [53, 54, 55]. The success of the GQL idea in these special cases strongly encourages us to explore its essential mechanisms and universal framework for general systems.

Our efforts in this article include:

- We interpret, from a geometric viewpoint, the fundamental principle behind the GQL representations for general nonlinear algebraic constraints.
- We establish the universal GQL framework and its mathematical theory.
- We propose three simple effective methods for constructing GQL representations using extra free auxiliary variables in exchange for linearity. As examples, the GQL representations are derived for the invariant regions of various physical systems.
- We illustrate the GQL methodology and related techniques for nonlinear bound-preserving analysis and design, demonstrating its effectiveness and remarkable advantages, by diverse challenging applications which cannot be easily handled by direct or traditional approaches.

We emphasize that GQL has no restriction on the specific forms of the equations (1.1). This makes the framework applicable to general time-dependent PDE systems that possess convex invariant regions with nonlinear constraints.

The paper is organized as follows. Section 2 presents several examples of physical PDE systems with convex invariant regions and nonlinear constraints. Section 3 explores the fundamental principle and general theory for the GQL framework. We propose in section 4 three simple effective methods for constructing GQL representations, along with extensive examples. Section 5 illustrates the GQL approach for bound-preserving analysis. In section 6 we apply the GQL approach to design bound-preserving schemes for the multicomponent MHD system, and further demonstrate its powerful capabilities in addressing challenging bound-preserving problems that could not be coped with by direct or traditional approaches. Several experimental results are given in section 7 to verify the performance of the bound-preserving schemes developed via GQL. The conclusions follow in section 8. Throughout this paper, we will use $\text{cl}(G)$, $\text{int}(G)$, and ∂G to denote the closure, the interior, and the boundary of a region G , respectively. We employ $\|\mathbf{a}\|$ to denote the 2-norm of vector \mathbf{a} . We use $\mathbf{a} \cdot \mathbf{b}$ to denote the inner product of two vectors \mathbf{a} and \mathbf{b} , and $\mathbf{a} \otimes \mathbf{b}$ to denote the outer product, i.e., in index notation, $(\mathbf{a} \otimes \mathbf{b})_{ij} = a_i b_j$.

2. Examples of PDE systems with nonlinear constraints. In this section, we present several examples of physical PDE systems involving nonlinear algebraic constraints. For convenience, the ideal equation of state $p = (\Gamma - 1)\rho e$ is used to close the systems in Examples 2.1, 2.2, 2.4, 2.6, and 2.7, with p denoting the thermal pressure, ρ the (rest-mass) density, e the specific internal energy, and the constant $\Gamma > 1$ denoting the ratio of specific heats. For the relativistic models in Examples 2.3, 2.4, and 2.7, normalized units are employed such that the speed of light $c = 1$.

Example 2.1 (Euler System). Consider the 1D compressible Euler equations [68]

$$\partial_t \mathbf{u} + \partial_x \mathbf{f}(\mathbf{u}) = \mathbf{0}, \quad \mathbf{u} = \begin{pmatrix} \rho \\ m \\ E \end{pmatrix}, \quad \mathbf{f}(\mathbf{u}) = \begin{pmatrix} m \\ mv + p \\ (E + p)v \end{pmatrix}, \quad (2.1)$$

where ρ , m , $v = m/\rho$, and p denote the fluid density, momentum, velocity, and pressure, respectively. The quantity $E = \rho e + \frac{1}{2}\rho v^2$ is the total energy, with e being the specific internal energy. For this system, the density ρ and the internal energy ρe are positive, namely, \mathbf{u} should stay in the region

$$G = \left\{ \mathbf{u} = (\rho, m, E)^\top \in \mathbb{R}^3 : \rho > 0, g(\mathbf{u}) := E - \frac{m^2}{2\rho} > 0 \right\}, \quad (2.2)$$

which is a convex invariant region of the system (2.1). If we further consider Tadmor's minimum entropy principle [47], $S(\mathbf{u}) \geq S_{\min} := \min_{\mathbf{x}} S(\mathbf{u}_0(\mathbf{x}))$, for the specific entropy $S = p\rho^{-\Gamma}$, then we obtain another convex invariant region

$$\tilde{G} = \left\{ \mathbf{u} = (\rho, m, E)^\top \in \mathbb{R}^3 : \rho > 0, \tilde{g}(\mathbf{u}) \geq 0 \right\} \quad (2.3)$$

with

$$\tilde{g}(\mathbf{u}) := S(\mathbf{u}) - S_{min} = \frac{\Gamma - 1}{\rho^\Gamma} \left(E - \frac{m^2}{2\rho} \right) - S_{min}.$$

The readers are referred to [68, 70] for proofs of the convexity of G and \tilde{G} . Convex invariant regions for the 2D and 3D Euler systems are analogous and omitted here.

Example 2.2 (Navier–Stokes System). Consider the 1D dimensionless compressible Navier–Stokes equations (see, for example, [66]):

$$\partial_t \mathbf{u} + \partial_x \mathbf{f}(\mathbf{u}) = \frac{\eta}{\text{Re}} \partial_{xx} \mathbf{r}(\mathbf{u}), \quad \mathbf{r}(\mathbf{u}) = \begin{pmatrix} 0 \\ v \\ \frac{v^2}{2} + \frac{\Gamma}{\text{Pr}} e \end{pmatrix}, \quad (2.4)$$

where $\{\eta, \text{Re}, \text{Pr}\}$ are positive constants, and the definitions of \mathbf{u} and $\mathbf{f}(\mathbf{u})$ are the same as [Example 2.1](#). Both sets in (2.2) and (2.3) are also invariant regions for system (2.4).

Example 2.3 (M1 Model of Radiative Transfer). For the solutions of the gray M1 moment system of radiative transfer (see, for example, [38, 3]), a convex invariant region is

$$G = \{ \mathbf{u} = (E_r, \mathcal{F}_r)^\top \in \mathbb{R}^4 : g(\mathbf{u}) := E_r - \|\mathcal{F}_r\| \geq 0 \}, \quad (2.5)$$

where E_r is the radiation energy, and \mathcal{F}_r is the radiation energy flux.

Example 2.4 (Relativistic Hydrodynamic System). Consider the 1D governing equations of the special relativistic hydrodynamics (RHD) [56, 41]:

$$\partial_t \mathbf{u} + \partial_x \mathbf{f}(\mathbf{u}) = \mathbf{0}, \quad \mathbf{u} = \begin{pmatrix} D \\ m \\ E \end{pmatrix}, \quad \mathbf{f}(\mathbf{u}) = \begin{pmatrix} Dv \\ mv + p \\ m \end{pmatrix} \quad (2.6)$$

with the density $D = \rho\gamma$, the momentum $m = \rho h\gamma^2 v$, the energy $E = \rho h\gamma^2 - p$. Here, ρ , v , p , and $\gamma = (1 - v^2)^{-\frac{1}{2}}$ denote the rest-mass density, velocity, pressure, and Lorentz factor, respectively. The quantity $h = 1 + e + p/\rho$ represents the specific enthalpy, with e being the specific internal energy. For this system, the density and the pressure are positive, and the magnitude of v must be smaller than the speed of light ($c = 1$). These physical constraints define the invariant region

$$G = \{ \mathbf{u} \in \mathbb{R}^3 : D > 0, p(\mathbf{u}) > 0, 1 - |v(\mathbf{u})| > 0 \}. \quad (2.7)$$

It was proven in [56] that the region G is convex and can be equivalently represented as

$$G = \{ \mathbf{u} \in \mathbb{R}^3 : D > 0, g(\mathbf{u}) := E - \sqrt{D^2 + m^2} > 0 \}. \quad (2.8)$$

As shown in [52], the minimum entropy principle $S(\mathbf{u}) \geq S_{min}$ also holds for the RHD system (2.6), yielding another invariant region

$$\tilde{G} = \{ \mathbf{u} \in \mathbb{R}^3 : D > 0, g(\mathbf{u}) > 0, \tilde{g}(\mathbf{u}) \geq 0 \}, \quad (2.9)$$

where $\tilde{g}(\mathbf{u}) := p(\mathbf{u})(\rho(\mathbf{u}))^{-\Gamma} - S_{min}$ is a highly nonlinear implicit function. In the RHD case, the functions $p(\mathbf{u})$ and $\rho(\mathbf{u})$ cannot be explicitly expressed in terms of \mathbf{u} . Specifically, $p(\mathbf{u})$ is implicitly defined by the positive root of the nonlinear function $F(p; \mathbf{u}) := \frac{m^2}{E+p} + D(1 - \frac{m^2}{(E+p)^2})^{\frac{1}{2}} + \frac{p}{\Gamma-1} - E$, and then $\rho(\mathbf{u}) = D\sqrt{1 - m^2/(E + p(\mathbf{u}))^2}$.

Example 2.5 (Ten-Moment Gaussian Closure System). In 2D, this system [35, 36] reads

$$\partial_t \mathbf{u} + \partial_x \mathbf{f}_1(\mathbf{u}) + \partial_y \mathbf{f}_2(\mathbf{u}) = \mathbf{0}, \quad (2.10)$$

$$\mathbf{u} = \begin{pmatrix} \rho \\ m_1 \\ m_2 \\ E_{11} \\ E_{12} \\ E_{22} \end{pmatrix}, \quad \mathbf{f}_j(\mathbf{u}) = \begin{pmatrix} m_j \\ m_1 v_j + p_{1j} \\ m_2 v_j + p_{2j} \\ E_{11} v_j + p_{1j} v_1 \\ E_{12} v_j + \frac{1}{2}(p_{1j} v_2 + p_{2j} v_1) \\ E_{22} v_j + p_{2j} v_2 \end{pmatrix}, \quad j = 1, 2.$$

Here ρ , $\mathbf{m} = (m_1, m_2)$, $\mathbf{v} = \mathbf{m}/\rho$, $\mathbf{E} = (E_{ij})_{1 \leq i, j \leq 2}$, and $\mathbf{p} = (p_{ij})_{1 \leq i, j \leq 2}$ are respectively the density, momentum vector, velocity, symmetric energy tensor, and symmetric anisotropic pressure tensor. The system (2.10) is closed by $\mathbf{p} = 2\mathbf{E} - \rho\mathbf{v} \otimes \mathbf{v}$. For this system, the density ρ is positive, and the pressure tensor \mathbf{p} is positive-definite, namely, the evolved variables \mathbf{u} should belong to the following invariant region

$$G = \left\{ \mathbf{u} \in \mathbb{R}^6 : \rho > 0, \mathbf{E} - \frac{\mathbf{m} \otimes \mathbf{m}}{2\rho} \text{ is positive-definite} \right\} \quad (2.11)$$

$$= \left\{ \mathbf{u} \in \mathbb{R}^6 : \rho > 0, \mathbf{z}^\top \left(\mathbf{E} - \frac{\mathbf{m} \otimes \mathbf{m}}{2\rho} \right) \mathbf{z} > 0 \quad \forall \mathbf{z} \in \mathbb{R}^2 \setminus \{\mathbf{0}\} \right\}. \quad (2.12)$$

Example 2.6 (Ideal MHD System). This system [51, 53] can be written as

$$\partial_t \begin{pmatrix} \rho \\ \mathbf{m} \\ \mathbf{B} \\ E \end{pmatrix} + \nabla \cdot \begin{pmatrix} \mathbf{m} \\ \mathbf{m} \otimes \mathbf{v} - \mathbf{B} \otimes \mathbf{B} + p_{tot} \mathbf{I} \\ \mathbf{v} \otimes \mathbf{B} - \mathbf{B} \otimes \mathbf{v} \\ (E + p_{tot}) \mathbf{v} - (\mathbf{v} \cdot \mathbf{B}) \mathbf{B} \end{pmatrix} = \mathbf{0} \quad (2.13)$$

with ρ being the density, \mathbf{m} the momentum vector, $\mathbf{v} = \mathbf{m}/\rho$ the velocity, $E = \rho e + \frac{1}{2}(\rho \|\mathbf{v}\|^2 + \|\mathbf{B}\|^2)$ denoting the total energy, $p_{tot} = p + \frac{1}{2} \|\mathbf{B}\|^2$ being the total pressure, p the thermal pressure, and \mathbf{B} the magnetic field which satisfies the extra divergence-free condition $\nabla \cdot \mathbf{B} = 0$. For this system, the density ρ and the internal energy ρe are positive, namely, \mathbf{u} should stay in the invariant region

$$G = \left\{ \mathbf{u} = (\rho, \mathbf{m}, \mathbf{B}, E)^\top \in \mathbb{R}^8 : \rho > 0, g(\mathbf{u}) := E - \frac{\|\mathbf{m}\|^2}{2\rho} - \frac{\|\mathbf{B}\|^2}{2} > 0 \right\}. \quad (2.14)$$

Example 2.7 (Relativistic MHD System). This system [55] takes the form of

$$\partial_t \begin{pmatrix} D \\ \mathbf{m} \\ \mathbf{B} \\ E \end{pmatrix} + \nabla \cdot \begin{pmatrix} D\mathbf{v} \\ \mathbf{m} \otimes \mathbf{v} - \mathbf{B} \otimes (\gamma^{-2} \mathbf{B} + (\mathbf{v} \cdot \mathbf{B}) \mathbf{v}) + p_{tot} \mathbf{I} \\ \mathbf{v} \otimes \mathbf{B} - \mathbf{B} \otimes \mathbf{v} \\ \mathbf{m} \end{pmatrix} = \mathbf{0} \quad (2.15)$$

with the mass density $D = \rho\gamma$, the momentum vector $\mathbf{m} = (\rho h \gamma^2 + \|\mathbf{B}\|^2) \mathbf{v} - (\mathbf{v} \cdot \mathbf{B}) \mathbf{B}$, the energy $E = \rho h \gamma^2 - p_{tot} + \|\mathbf{B}\|^2$, and the magnetic field \mathbf{B} satisfies $\nabla \cdot \mathbf{B} = 0$ as the ideal MHD case. The total pressure p_{tot} consists of the magnetic pressure $p_m := \frac{1}{2} (\gamma^{-2} \|\mathbf{B}\|^2 + (\mathbf{v} \cdot \mathbf{B})^2)$ and the thermal pressure p . Analogously to Example 2.4, the quantities ρ , \mathbf{v} , h , and $\gamma = (1 - \|\mathbf{v}\|^2)^{-\frac{1}{2}}$ are respectively the rest-mass density, velocity, specific enthalpy, and Lorentz factor. The positivity of density and pressure as well as the subluminal constraint $\|\mathbf{v}\| < c = 1$ constitute the invariant region

$$G = \{ \mathbf{u} = (D, \mathbf{m}, \mathbf{B}, E)^\top \in \mathbb{R}^8 : D > 0, p(\mathbf{u}) > 0, 1 - \|\mathbf{v}(\mathbf{u})\| > 0 \}, \quad (2.16)$$

where $p(\mathbf{u})$ and $\mathbf{v}(\mathbf{u})$ are highly nonlinear and cannot be explicitly formulated, as discussed in (1.7).

3. Framework and theory of geometric quasilinearization. This section establishes the universal GQL framework, with the geometric insights into understanding the fundamental principle behind the GQL representations.

Let $G \subset \mathbb{R}^N$ be an invariant region or admissible state set of a physical system. Assume that G can be formulated into the general form (1.4). For notational convenience, we represent G as

$$G = \{ \mathbf{u} \in \mathbb{R}^N : g_i(\mathbf{u}) \succ 0, 1 \leq i \leq I \}, \quad (3.1)$$

where the symbol “ \succ ” denotes “ $>$ ” if $i \in \mathbb{I}$, or “ \geq ” if $i \in \widehat{\mathbb{I}}$. Let $G_L = \{ \mathbf{u} \in \mathbb{R}^N : g_i(\mathbf{u}) \succ 0 \forall i \in \mathbb{I}_L \}$ be the region formed by all the linear constraints in G , i.e., the function g_i is linear for $i \in \mathbb{I}_L$. If $\mathbb{I}_L = \emptyset$, then we define $G_L = \mathbb{R}^N$.

We consider the nontrivial case that at least one of the functions $\{g_i(\mathbf{u})\}$ is nonlinear, namely, $G \subset G_L$ and $G \neq G_L$. The goal of our GQL methodology is to use some extra free auxiliary variables in exchange for linearity, and more precisely, is to equivalently represent G by using only *linear* constraints with the help of free auxiliary variables.

DEFINITION 3.1. We say a set G_* is an equivalent linear representation (termed as GQL representation) of the region G , if $G_* = G$ and G_* takes the form

$$G_* = \{ \mathbf{u} \in \mathbb{R}^N : \varphi_i(\mathbf{u}; \boldsymbol{\theta}_{i*}) > 0 \quad \forall \boldsymbol{\theta}_{i*} \in \Theta_i, 1 \leq i \leq I \}, \quad (3.2)$$

where the functions $\{\varphi_i\}$ are all **linear** (affine) with respect to \mathbf{u} ; the parameters $\boldsymbol{\theta}_{i*}$ are independent of \mathbf{u} and stand for the (possible) extra free auxiliary variables with Θ_i denoting their ranges.

Based on Definition 3.1, we immediately have:

THEOREM 3.2. Assume that a set G_* is of the form (3.2) with φ_i being linear with respect to \mathbf{u} and satisfying

$$\min_{\boldsymbol{\theta}_{i*} \in \Theta_i} \varphi_i(\mathbf{u}; \boldsymbol{\theta}_{i*}) = \lambda_i(\mathbf{u})g_i(\mathbf{u}) \quad (3.3)$$

with $\lambda_i(\mathbf{u}) > 0$ for all $\mathbf{u} \in G_L$. Then $G_* = G$, and G_* is the GQL representation of G .

Remark 3.3. For $i \in \mathbb{I}_L$, the function $g_i(\mathbf{u})$ is already linear, thus we can simply take $\varphi_i(\mathbf{u}; \boldsymbol{\theta}_{i*}) = g_i(\mathbf{u})$, without free auxiliary variable $\boldsymbol{\theta}_{i*}$ in this case. That is, all the linear constraints remain unchanged in the GQL representation.

Theorem 3.2 points out a way to seek the GQL representation, namely, by constructing linear functions $\{\varphi_i\}$ such that (3.3) holds. We have used this approach in [51] to establish the GQL representation of the invariant region (2.14) for the ideal MHD equations. However, this constructive approach needs some empirical observations or trial-and-error procedures, as Theorem 3.2 does not provide any insight on how to find the qualified $\{\varphi_i\}$. In the following, we explore a simpler yet universal approach from the geometric point of view.

Given that $\{\varphi_i\}$ in (3.2) are all linear with respect to \mathbf{u} , the set G_* is always convex. This means if the region G has GQL representation (3.2), then G must also be convex. Hence we should make the following basic (minimal) assumption.

ASSUMPTION 3.4. The invariant region G is convex, and $\text{int}(G) \neq \emptyset$.

This basic assumption is valid for many physical systems including all those introduced in section 2. Again, we emphasize that the functions $\{g_i(\mathbf{u})\}$ are *not* necessarily concave.

3.1. A heuristic example. Before deriving the general theory, let us look at an example to gain some insight, which inspires us to achieve the GQL framework.

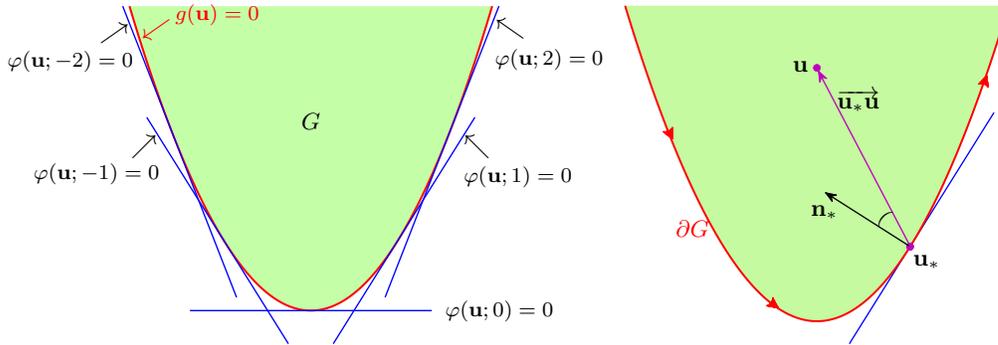


Fig. 1: Illustrations for Example 3.5.

Example 3.5. Consider the simple example mentioned in (1.8)–(1.9), i.e., $G = \{ \mathbf{u} = (u_1, u_2)^\top \in \mathbb{R}^2 : g(\mathbf{u}) = u_2 - u_1^2 > 0 \}$. According to Theorem 3.2, the GQL representation of G is

$$G_* = \{ \mathbf{u} = (u_1, u_2)^\top \in \mathbb{R}^2 : \varphi(\mathbf{u}; \theta_*) = u_2 - 2u_1\theta_* + \theta_*^2 > 0 \quad \forall \theta_* \in \mathbb{R} \}. \quad (3.4)$$

As such, we gain the linearity by introducing the extra free auxiliary variable θ_* . To understand the intrinsic mechanisms, we draw the graph of the region G and its boundary curve $\partial G = \{ \mathbf{u} : g(\mathbf{u}) = 0 \}$ on the u_1 – u_2 plane in Figure 1. We also plot the graphs of $\{ \mathbf{u} : \varphi(\mathbf{u}; \theta_*) = 0 \}$ for several special values of $\theta_* \in \{ \pm 2, \pm 1, 0 \}$ in the left subfigure of Figure 1. It is observed that all the lines $\{ \mathbf{u} : \varphi(\mathbf{u}; \theta_*) = 0 \}$ are tangent to the parabolic curve ∂G , which exactly forms an envelope of the tangent lines.

Let $\mathbf{u}_* = (\theta_*, \theta_*^2)^\top$ denote an arbitrary point on ∂G . One can verify that $\mathbf{n}_* = (-2\theta_*, 1)^\top$ is an inward-pointing normal vector of ∂G at \mathbf{u}_* , and

$$\varphi(\mathbf{u}; \theta_*) = \mathbf{u} \cdot \mathbf{n}_* - \mathbf{u}_* \cdot \mathbf{n}_* = \overrightarrow{\mathbf{u}_* \mathbf{u}} \cdot \mathbf{n}_* > 0 \quad \forall \mathbf{u} \in G.$$

Imagine we are walking along the boundary ∂G in the direction shown in the right subfigure of [Figure 1](#), then the region G always lie entirely on the left side of the tangent lines, namely, the angle between the two vectors $\overrightarrow{\mathbf{u}_* \mathbf{u}}$ and \mathbf{n}_* is always less than 90° for all $\mathbf{u} \in G$ and all $\mathbf{u}_* \in \partial G$. This intuitively interprets the GQL representation [\(3.4\)](#) from the geometric viewpoint.

3.2. Concepts from geometry and convex sets. Let us recall some concepts and results from theory of geometry and convex analysis [\[28, 42, 19\]](#).

A hyperplane in \mathbb{R}^N is a plane of dimension $N - 1$. Let $\mathbf{n}_* \neq \mathbf{0}$ denote a normal vector of a hyperplane H , and let \mathbf{u}_* be a point on H . Then H can be expressed as $H = \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = 0\}$, and it divides \mathbb{R}^N into two halfspaces: $H^+ = \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* \geq 0\}$ and $H^- = \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* \leq 0\}$.

DEFINITION 3.6 (Supporting Hyperplane and Halfspace). *The hyperplane $H = \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = 0\}$ through $\mathbf{u}_* \in \partial G$ is called a supporting hyperplane to G at \mathbf{u}_* , if G lies in one of the two closed halfspaces determined by H . Furthermore, if the normal vector \mathbf{n}_* points towards G , then the closed halfspace containing G is $H^+ = \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* \geq 0\}$ and is called a closed supporting halfspace to G . See [Figure 2](#).*

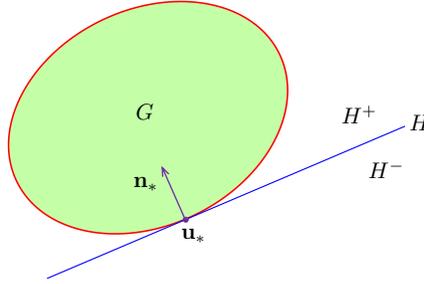


Fig. 2: Supporting hyperplane and halfspace.

THEOREM 3.7 (Supporting Hyperplane Theorem [\[28\]](#)). *If G is a convex set and $\text{int}(G) \neq \emptyset$, then for any $\mathbf{u}_* \in \partial G$, there exists a supporting hyperplane to G at \mathbf{u}_* .*

Remark 3.8. If the boundary ∂G is smooth at a point \mathbf{u}_* , then the supporting hyperplane to G at \mathbf{u}_* is unique and coincide with the tangent [\[42, 19\]](#).

3.3. GQL framework. We are now in the position to establish the GQL framework.

3.3.1. A special case. Inspired by [Example 3.5](#), we first consider a special case that G is either open or closed with differentiable boundary. The general case will be discussed in [subsection 3.3.2](#).

THEOREM 3.9. *Suppose that [Assumption 3.4](#) holds, the region G is either open or closed, and ∂G is differentiable. Then G has the following GQL representation:*

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* \succ 0 \quad \forall \mathbf{u}_* \in \partial G \right\}, \quad (3.5)$$

where the symbol “ \succ ” is taken as “ $>$ ” if G is open, or as “ \geq ” if G is closed, and \mathbf{n}_* is only dependent on \mathbf{u}_* and denotes an inward-pointing normal vector of ∂G at \mathbf{u}_* .

The proof of [Theorem 3.9](#) is presented in [Appendix A](#). Following the proof, one can further extend the above result to any closed convex region G , whose boundary is typically not everywhere smooth so that the supporting hyperplanes at each nonsmooth boundary point are not unique. Let $\mathcal{N}(\mathbf{u}_*)$ denote the set of the inward-pointing unit normal vectors of all the supporting hyperplanes to G at $\mathbf{u}_* \in \partial G$. Then one can prove that

$$G = \left\{ \mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n} \geq 0 \quad \forall \mathbf{n} \in \mathcal{N}(\mathbf{u}_*), \quad \forall \mathbf{u}_* \in \partial G \right\}. \quad (3.6)$$

This means any closed convex region is the intersection of all its closed supporting halfspaces [28]. However, the representation (3.6) is *not* applicable to a general convex region that is neither closed nor open (e.g. the invariant regions in (2.3) and (2.9)). Moreover, the representation (3.6) requires the information of *all* the supporting hyperplanes at each nonsmooth boundary point, which can be difficult to explicitly formulate or verify, so that (3.6) is not desirable for bound-preserving study. A practical GQL representation for more general regions will be derived in subsection 3.3.2.

3.3.2. General case. Consider a general convex region G that may be *not necessarily* open or closed and its boundary may be not everywhere smooth. Note that the boundary of a convex region can be partitioned into several pieces, each of which can be *locally* represented as the graph of a convex function (with respect to a suitable supporting hyperplane). Recall that any convex function is locally Lipschitz continuous and twice differentiable almost everywhere, according to the classical theorems of Rademacher and Alexandrov (cf. [37]). Based on these facts and for convenience, we make a considerably mild assumption on the convex invariant region G . We assume that the boundary of G is piecewise C^1 , and without loss of generality, for each $i \in \{1, \dots, I\}$, the function $g_i(\mathbf{u})$ in (3.1) is C^1 at any points on

$$\mathcal{S}_i := \partial G \cap \partial G_i, \quad \text{with } G_i := \{\mathbf{u} \in \mathbb{R}^N : g_i(\mathbf{u}) \succ 0\},$$

where $\{\mathcal{S}_i\}$ are C^1 hypersurfaces in \mathbb{R}^N and constitute the smooth pieces of ∂G , i.e., $\partial G = \cup_{1 \leq i \leq I} \mathcal{S}_i$. Notice that in general, \mathcal{S}_i may not equal ∂G_i , the region G_i may be not convex, and G may be neither open nor closed; see an example in Figure 3. These make our following discussions nontrivial.

We remark that $G_i = \{\mathbf{u} : g_i(\mathbf{u}) \geq 0\}$ is closed for $i \in \hat{\mathbb{I}}$, and $G_i = \{\mathbf{u} : g_i(\mathbf{u}) > 0\}$ is open for $i \in \mathbb{I}$. Since for each $i \in \mathbb{I}$, the set G_i is not necessarily convex, there is a possibility that G may not be entirely contained in an *open* supporting halfspace at $\mathbf{u}_* \in \mathcal{S}_i$. This issue is avoid if the open region $\cap_{i \in \mathbb{I}} G_i$ is convex, which is satisfied by all the examples in section 2 and implies that

$$G \cap (\cup_{\mathbf{u}_* \in \mathcal{S}_i} \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} = 0\}) = \emptyset \quad \forall i \in \mathbb{I}, \quad (3.7)$$

where \mathbf{n}_{i*} is an inward-pointing normal vector of \mathcal{S}_i at \mathbf{u}_* .

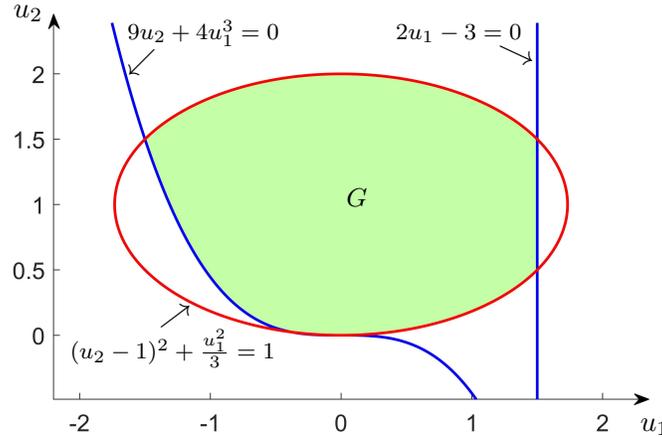


Fig. 3: A convex region G involving nonlinear constraints. $G = \{\mathbf{u} \in \mathbb{R}^2 : g_1(\mathbf{u}) \geq 0, g_2(\mathbf{u}) \geq 0, g_3(\mathbf{u}) > 0\}$ with $g_1(\mathbf{u}) = 3 - 2u_1$, $g_2(\mathbf{u}) = 9u_2 + 4u_1^3$, and $g_3(\mathbf{u}) = 1 - u_1^2/3 - (u_2 - 1)^2$.

THEOREM 3.10. *Suppose that Assumption 3.4 holds, condition (3.7) is satisfied when $\mathbb{I} \neq \emptyset$, and the boundary of G is piecewise C^1 . Then the region G has the following GQL representation:*

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} \succ 0 \quad \forall \mathbf{u}_* \in \mathcal{S}_i, 1 \leq i \leq I \right\}, \quad (3.8)$$

where the symbol “ \succ ” is taken as “ $>$ ” if $i \in \mathbb{I}$, or as “ \geq ” if $i \in \hat{\mathbb{I}}$; the nonzero vector \mathbf{n}_{i*} denotes an inward-pointing normal vector of \mathcal{S}_i at \mathbf{u}_* .

Proof. The proof is divided into three steps.

(i) **Prove that** $G \subseteq G_*$. Let $\partial G =: \widetilde{\partial G} \cup \widehat{\partial G}$ with $\widetilde{\partial G}$ denoting the set of smooth boundary points and $\widehat{\partial G}$ the set of nonsmooth boundary points. For any $\mathbf{u}_* \in \widetilde{\partial G} \cap \mathcal{S}_i$, the hyperplane $(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} = 0$ supports the region G , implying that

$$G \subseteq \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} \geq 0\} \quad \forall \mathbf{u}_* \in \widetilde{\partial G} \cap \mathcal{S}_i, 1 \leq i \leq I. \quad (3.9)$$

Next, we consider an arbitrary nonsmooth boundary point $\mathbf{u}_* \in \widehat{\partial G} \cap \mathcal{S}_i$. There exists a sequence of smooth boundary points $\{\mathbf{u}_*^{(j)}\}_{j \in \mathbb{N}} \subset \widetilde{\partial G} \cap \mathcal{S}_i$ such that $\lim_{j \rightarrow \infty} \mathbf{u}_*^{(j)} = \mathbf{u}_*$. For every $\mathbf{u}_*^{(j)}$, it follows from (3.9) that

$$(\mathbf{u} - \mathbf{u}_*^{(j)}) \cdot \mathbf{n}_{i, \mathbf{u}_*^{(j)}} \geq 0 \quad \forall \mathbf{u} \in G, \quad (3.10)$$

where $\mathbf{n}_{i, \mathbf{u}_*^{(j)}}$ is the inward-pointing normal vector of \mathcal{S}_i at $\mathbf{u}_*^{(j)}$ satisfying $\lim_{j \rightarrow \infty} \mathbf{n}_{i, \mathbf{u}_*^{(j)}} = \mathbf{n}_{i*}$. Taking the limit $j \rightarrow +\infty$ in (3.10) gives

$$(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} \geq 0 \quad \forall \mathbf{u} \in G \quad \forall \mathbf{u}_* \in \widehat{\partial G} \cap \mathcal{S}_i, 1 \leq i \leq I,$$

which along with (3.9) yields

$$G \subseteq \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} \geq 0 \quad \forall \mathbf{u}_* \in \mathcal{S}_i, 1 \leq i \leq I\}. \quad (3.11)$$

Based on (3.7), we then conclude that $G \subseteq G_*$.

(ii) **Prove that** $G_* \subseteq \text{cl}(G)$ by contradiction. Assume that $G_* \not\subseteq \text{cl}(G)$, namely, there exists $\mathbf{u}_0 \in G_*$ but $\mathbf{u}_0 \notin \text{cl}(G)$. According to the theory of convex optimization [5], the minimum of the convex function $\zeta(\mathbf{u}) := \frac{1}{2} \|\mathbf{u} - \mathbf{u}_0\|^2$ over the closed convex region $\text{cl}(G)$ is attained at certain boundary point $\mathbf{u}_{*0} \in \partial G$. In other words, \mathbf{u}_{*0} is a solution to the following optimization problem

$$\begin{aligned} & \underset{\mathbf{u} \in \text{cl}(G)}{\text{minimize}} \quad \zeta(\mathbf{u}) \\ & \text{subject to} \quad -g_i(\mathbf{u}) < 0 \quad \forall i \in \mathbb{I}; \quad -g_i(\mathbf{u}) \leq 0 \quad \forall i \in \widehat{\mathbb{I}}. \end{aligned} \quad (3.12)$$

Since the function $-g_i(\mathbf{u})$ is not necessarily convex, the problem (3.12) is generally not the standard form of convex optimization. Note that the condition $\text{int}(G) \neq \emptyset$ ensures the Slater condition [5, 4] is satisfied. The Karush–Kuhn–Tucker (KKT) conditions [5, 4] tell us that there exist $\{\lambda_0, \lambda_1, \dots, \lambda_I\}$ such that

$$0 = \nabla \zeta(\mathbf{u}_{*0}) - \sum_{i=1}^I \lambda_i \nabla g_i(\mathbf{u}_{*0}), \quad (3.13)$$

$$0 = \lambda_i g_i(\mathbf{u}_{*0}), \quad 1 \leq i \leq I, \quad (3.14)$$

$$\lambda_i \geq 0, \quad 0 \leq i \leq I. \quad (3.15)$$

Define $\mathbb{I}_+ := \{1 \leq i \leq I : \lambda_i > 0\}$. Obviously $\mathbb{I}_+ \neq \emptyset$; otherwise $\lambda_i = 0$ for all $1 \leq i \leq I$, so that $\mathbf{u}_{*0} - \mathbf{u}_0 = \nabla \zeta(\mathbf{u}_{*0}) = \mathbf{0}$ which leads to the contradiction $\partial G \ni \mathbf{u}_{*0} = \mathbf{u}_0 \notin \text{cl}(G)$. This also implies $\mathbf{u}_{*0} \neq \mathbf{u}_0$. Let \mathbf{n}_{i*0} be the inward-pointing normal vector of \mathcal{S}_i at \mathbf{u}_{*0} . Since there exist $\mu_i \geq 0$ such that $\nabla g_i(\mathbf{u}_{*0}) = \mu_i \mathbf{n}_{i*0}$, condition (3.13) can be rewritten as

$$\mathbf{u}_{*0} - \mathbf{u}_0 = \sum_{i \in \mathbb{I}_+} \lambda_i \mu_i \mathbf{n}_{i*0}. \quad (3.16)$$

Thanks to (3.14), we obtain $g_i(\mathbf{u}_{*0}) = 0$ for all $i \in \mathbb{I}_+$, which along with $\mathbf{u}_{*0} \in \partial G$ leads to

$$\mathbf{u}_{*0} \in \mathcal{S}_i = \partial G_i \cap \partial G \quad \forall i \in \mathbb{I}_+.$$

Because $\mathbf{u}_0 \in G_*$, we then have $(\mathbf{u}_0 - \mathbf{u}_{*0}) \cdot \mathbf{n}_{i*0} > 0$ for all $i \in \mathbb{I}_+$. This, together with (3.16) and $\mathbf{u}_{*0} \neq \mathbf{u}_0$, leads to a contradiction:

$$\begin{aligned} 0 &> -\|\mathbf{u}_0 - \mathbf{u}_{*0}\|_2^2 = (\mathbf{u}_0 - \mathbf{u}_{*0}) \cdot (\mathbf{u}_{*0} - \mathbf{u}_0) \\ &= (\mathbf{u}_0 - \mathbf{u}_{*0}) \cdot \left(\sum_{i \in \mathbb{I}_+} \lambda_i \mu_i \mathbf{n}_{i*0} \right) = \sum_{i \in \mathbb{I}_+} \lambda_i \mu_i (\mathbf{u}_0 - \mathbf{u}_{*0}) \cdot \mathbf{n}_{i*0} \geq 0. \end{aligned}$$

Thus the assumption $G_* \not\subseteq \text{cl}(G)$ is incorrect. We have $G_* \subseteq \text{cl}(G)$.

(iii) Prove that $G_* \subseteq G$. If $\mathbb{I} = \emptyset$, then G is a closed region and $G = \text{cl}(G)$. We immediately obtain $G_* \subseteq G$ from step (ii) of this proof. In the following, we focus on $\mathbb{I} \neq \emptyset$ and prove $G_* \subseteq G$ by contradiction. Assume that there exists $\mathbf{u}_0 \in G_*$ but $\mathbf{u}_0 \notin G$. Because we have already shown $G_* \subseteq \text{cl}(G)$ in step (ii) of this proof, we then get $\mathbf{u}_0 \in \text{cl}(G) \setminus G = \partial G$. Note that $\mathbf{u}_0 \in G_*$ implies

$$(\mathbf{u}_0 - \mathbf{u}_*) \cdot \mathbf{n}_{i*} > 0 \quad \forall \mathbf{u}_* \in \mathcal{S}_i, \quad \forall i \in \mathbb{I},$$

which leads to $\mathbf{u}_0 \notin \mathcal{S}_i = \partial G_i \cap \partial G$ for all $i \in \mathbb{I}$. It follows that $\mathbf{u}_0 \notin \partial G_i$ for all $i \in \mathbb{I}$. Note for $i \in \mathbb{I}$, one has $\mathbf{u}_0 \in \text{cl}(G) \subseteq \text{cl}(G_i)$, which gives

$$\mathbf{u}_0 \in \text{cl}(G_i) \setminus \partial G_i = G_i \quad \forall i \in \mathbb{I}. \quad (3.17)$$

On the other hand, $\mathbf{u}_0 \in \text{cl}(G) \subseteq \bigcap_{i \in \widehat{\mathbb{I}}} G_i$, which along with (3.17) implies $\mathbf{u}_0 \in (\bigcap_{i \in \mathbb{I}} G_i) \cap (\bigcap_{i \in \widehat{\mathbb{I}}} G_i) = G$. This contradicts the assumption that $\mathbf{u}_0 \notin G$. Hence the assumption is incorrect, and we have $G_* \subseteq G$.

Combining the conclusions proven in steps (i) and (iii) gives $G = G_*$. \square

Remark 3.11. If we replace \mathcal{S}_i with $\mathcal{S}_i \cap \partial \widetilde{G}$ for $i \in \widehat{\mathbb{I}}$ in (3.8), [Theorem 3.10](#) remains valid, because for $i \in \widehat{\mathbb{I}}$ we have $\{\mathbf{u} : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} \geq 0 \forall \mathbf{u}_* \in \mathcal{S}_i\} = \{\mathbf{u} : (\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_{i*} \geq 0 \forall \mathbf{u}_* \in \mathcal{S}_i \cap \partial \widetilde{G}\}$.

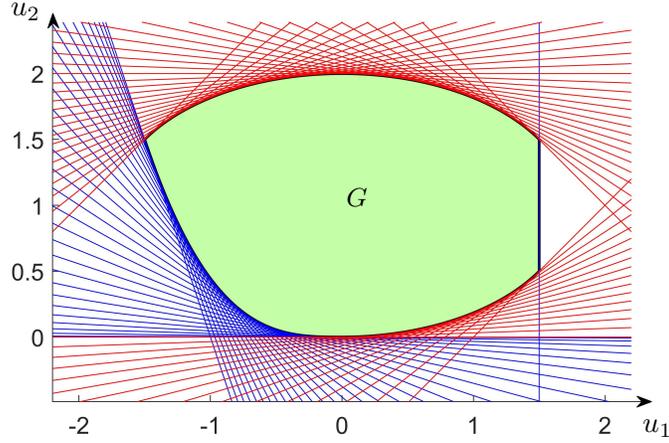


Fig. 4: Illustration of the GQL representation for the convex region G given in [Figure 3](#). The blue (resp. red) lines correspond to closed (resp. open) supporting halfspaces.

Remark 3.12. An illustration of the GQL representation (3.8) is shown in [Figure 4](#). Different from (3.6), the GQL representation (3.8) involves only at most N rather than *all* the supporting halfspaces at each nonsmooth “junction” point. This makes the GQL representation (3.8) easier to formulate or construct. Besides, [Theorem 3.10](#) does not require G to be closed or open.

Remark 3.13 (Significance of GQL). Compared to the original form (3.1) of the invariant region G with nonlinear constraints, its equivalent GQL representation G_* in (3.8) is described with only linear constraints. Such linearity gives the GQL representation some significant advantages over the original form (3.1) in analyzing and designing bound-preserving schemes; see [sections 5](#) and [6](#).

4. Construction of geometric quasilinearization. With three methods and several examples, this section discusses how to construct GQL for convex invariant regions.

4.1. Methods for constructing GQL representations. Based on [Theorems 3.2](#) and [3.10](#), we introduce three simple effective methods for constructing the GQL representation of G .

4.1.1. Gradient-based method. The first method is based on the following result, which is a direct consequence of [Theorem 3.10](#).

THEOREM 4.1. Assume that the hypotheses of [Theorem 3.10](#) hold and

$$\nabla g_i(\mathbf{u}_*) \neq \mathbf{0} \quad \forall \mathbf{u}_* \in \mathcal{S}_i, \quad 1 \leq i \leq I, \quad (4.1)$$

then the invariant region G is exactly equivalent to

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^N : \varphi_i(\mathbf{u}; \mathbf{u}_*) \succ 0 \quad \forall \mathbf{u}_* \in \mathcal{S}_i, 1 \leq i \leq I \right\}, \quad (4.2)$$

where the function φ_i is linear with respect to \mathbf{u} , defined by

$$\varphi_i(\mathbf{u}; \mathbf{u}_*) := (\mathbf{u} - \mathbf{u}_*) \cdot \nabla g_i(\mathbf{u}_*). \quad (4.3)$$

Theorem 4.1 says if $\{\nabla g_i\}$ are computable and satisfy (4.1), then we can directly obtain the GQL representation in the form (4.2) with (4.3).

In some cases, it is, however, difficult to calculate the gradients of nonlinear functions $\{g_i\}$, e.g., the implicit functions in (2.9) and (2.16). This motivates us to propose the following cross-product method based on suitable parametrization of the hypersurface \mathcal{S}_i . The use of parametrization can also help to reduce or decouple the free auxiliary variables, which is highly desirable for bound-preserving applications; see the examples in subsection 4.2 and Remark 4.5.

4.1.2. Cross-product method. Assume that for each i the hypersurface \mathcal{S}_i has the following parametric expression

$$\mathcal{S}_i = \left\{ \mathbf{u}_* = \mathbf{U}_i(\boldsymbol{\theta}_{i*}) : \boldsymbol{\theta}_{i*} \in \Theta_i \subseteq \mathbb{R}^{N-1} \right\}, \quad (4.4)$$

where \mathbf{U}_i is a C^1 vector function defined on the parameter domain Θ_i with \mathcal{S}_i being the function range. Denote $\theta_{i*}^{(k)}$ as the k th component of $\boldsymbol{\theta}_{i*}$. For each i , we define

$$\boldsymbol{\tau}_{i,k}(\boldsymbol{\theta}_{i*}) := \frac{\partial \mathbf{U}_i}{\partial \theta_{i*}^{(k)}}, \quad 1 \leq k \leq N-1.$$

The vectors $\{\boldsymbol{\tau}_{i,k}(\boldsymbol{\theta}_{i*}) : 1 \leq k \leq N-1\}$ are $(N-1)$ tangent vectors of the hypersurface \mathcal{S}_i and generate its local tangent space at \mathbf{u}_* . Then, the normal vector of \mathcal{S}_i at \mathbf{u}_* can be constructed using the $(N-1)$ -ary analogue of the cross product (cf. [46, Pages 83–85]) in \mathbb{R}^N :

$$\mathbf{n}_{i*} = \delta_{i*} \bigwedge_{k=1}^{N-1} \boldsymbol{\tau}_{i,k}(\boldsymbol{\theta}_{i*}) := \delta_{i*} \boldsymbol{\tau}_{i,1}(\boldsymbol{\theta}_{i*}) \times \boldsymbol{\tau}_{i,2}(\boldsymbol{\theta}_{i*}) \times \cdots \times \boldsymbol{\tau}_{i,N-1}(\boldsymbol{\theta}_{i*}),$$

where δ_{i*} is a nonzero factor which may be used to simplify the final formula or/and to adjust the sign such that \mathbf{n}_{i*} is directed towards the interior of G .

As a direct consequence of **Theorem 3.10**, the following result holds.

THEOREM 4.2. *Suppose the hypotheses of **Theorem 3.10** hold and*

$$\bigwedge_{k=1}^{N-1} \boldsymbol{\tau}_{i,k}(\boldsymbol{\theta}_{i*}) \neq \mathbf{0} \quad \forall \boldsymbol{\theta}_{i*} \in \Theta_i, 1 \leq i \leq I,$$

then the region G is exactly equivalent to

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^N : \varphi_i(\mathbf{u}; \boldsymbol{\theta}_{i*}) \succ 0 \quad \forall \boldsymbol{\theta}_{i*} \in \Theta_i, 1 \leq i \leq I \right\}, \quad (4.5)$$

where the function φ_i is linear with respect to \mathbf{u} , defined by

$$\varphi_i(\mathbf{u}; \boldsymbol{\theta}_{i*}) := (\mathbf{u} - \mathbf{U}_i(\boldsymbol{\theta}_{i*})) \cdot \left(\delta_{i*} \bigwedge_{k=1}^{N-1} \boldsymbol{\tau}_{i,k}(\boldsymbol{\theta}_{i*}) \right). \quad (4.6)$$

Remark 4.3. In many cases, there exists a natural (usually physics-based) parametrization of the hypersurface \mathcal{S}_i , typically with the primitive quantities as parametric variables; see the examples in subsection 4.2. The advantages of using the parametric form (4.4) in the GQL representation will become more clear in those examples and the bound-preserving applications in sections 5 and 6.

4.1.3. Constructive method. For completeness, we also summarize the constructive approach and its variant as our third method. Recall that **Theorem 3.2** has told us: if we can construct linear functions $\varphi_i(\mathbf{u}; \boldsymbol{\theta}_{i*})$, $1 \leq i \leq I$, such that (3.3) holds, then the GQL representation of G is (3.2). The constructive approach does not require the assumptions in **Theorems 4.1** and **4.2**, but often needs some empirical trial-and-error techniques to find the qualified $\{\varphi_i\}$. In practice, one can use the proposed three methods in a hybrid way: first formally formulate $\{\varphi_i\}$ via either (4.3) or (4.6) and then verify (3.3). Such a hybrid approach is efficient, as it may exempt the assumptions in **Theorems 4.1** and **4.2** and also avoid the trial-and-error procedure.

4.2. Examples of GQL representations. We give several examples for constructing GQL representations of convex invariant regions.

Example 1: Euler and Navier–Stokes systems.

THEOREM 4.4. *For the 1D Euler and Navier–Stokes systems, the GQL representation of the invariant region G in (2.2) is given by*

$$G_* = \left\{ \mathbf{u} = (\rho, m, E)^\top : \rho > 0, \quad \varphi(\mathbf{u}; v_*) > 0 \quad \forall v_* \in \mathbb{R} \right\} \quad (4.7)$$

with $\varphi(\mathbf{u}; v_*) := E - mv_* + \rho \frac{v_*^2}{2}$ being linearly dependent on \mathbf{u} .

Proof. We respectively use the three methods proposed in subsection 4.1 to derive the GQL representation for this example. Note the first constraint in (2.2) is linear.

(i) **Gradient-based method.** For the second constraint in (2.2), the gradient $\nabla g(\mathbf{u}) = \left(\frac{m^2}{2\rho^2}, -\frac{m}{\rho}, 1 \right)^\top$, and the associated boundary hypersurface $\mathcal{S} = \{ \mathbf{u}_* = (\rho_*, m_*, E_*)^\top : \rho_* > 0, g(\mathbf{u}_*) = 0 \}$ can be parameterized as

$$\mathcal{S} = \left\{ \mathbf{u}_* = \left(\rho_*, \rho_* v_*, \frac{\rho_*}{2} v_*^2 \right)^\top : \rho_* > 0, v_* \in \mathbb{R} \right\}. \quad (4.8)$$

For $\mathbf{u}_* \in \mathcal{S}$ and $\mathbf{u} = (\rho, m, E)^\top$, we have

$$(\mathbf{u} - \mathbf{u}_*) \cdot \nabla g(\mathbf{u}_*) = (\rho - \rho_*) \frac{v_*^2}{2} + (m - \rho_* v_*) (-v_*) + E - \frac{\rho_*}{2} v_*^2 = \varphi(\mathbf{u}; v_*). \quad (4.9)$$

By Theorem 4.1, we obtain the GQL representation (4.7) of G .

(ii) **Cross-product method.** Based on the parametrization of \mathcal{S} in (4.8), we can compute the normal vector of \mathcal{S} at \mathbf{u}_* by cross product

$$\frac{\partial \mathbf{u}_*}{\partial \rho_*} \times \frac{\partial \mathbf{u}_*}{\partial v_*} = \left(1, v_*, \frac{1}{2} v_*^2 \right)^\top \times (0, \rho_*, \rho_* v_*)^\top = \rho_* \left(\frac{1}{2} v_*^2, -v_*, 1 \right)^\top =: \frac{1}{\delta_*} \mathbf{n}_*,$$

where $\delta_* = 1/\rho_*$ is a nonzero factor. By Theorem 4.2 and $(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = \varphi(\mathbf{u}; v_*)$, we get the GQL representation (4.7).

(iii) **Constructive method.** Observe that

$$\varphi(\mathbf{u}; v_*) = E - mv_* + \rho \frac{v_*^2}{2} = \frac{\rho}{2} \left(v_* - \frac{m}{\rho} \right)^2 + g(\mathbf{u}) \geq g(\mathbf{u}), \quad (4.10)$$

which implies $\min_{v_* \in \mathbb{R}} \varphi(\mathbf{u}; v_*) = g(\mathbf{u})$ for $\rho > 0$. According to Theorem 3.2, we also achieve the GQL representation (4.7). \square

Remark 4.5. Note only one free auxiliary variable v_* explicitly appears in the GQL representation (4.7). This is benefited from the use of parametric form (4.8).

Remark 4.6 (Physical Interpretation of GQL). It seems that the linear function $\varphi(\mathbf{u}; v_*)$ plays an energy-like role from a physical point of view. For the present example, $\varphi(\mathbf{u}; v_*) = \frac{1}{2} \rho (v - v_*)^2 + \rho e$, which represents the total energy in the reference frame moving at a velocity of v_* .

We now utilize the cross-product method to construct the GQL representation of the invariant region \tilde{G} in (2.3), where the minimum entropy principle $S(\mathbf{u}) := p\rho^{-\Gamma} \geq S_{min}$ is also included.

THEOREM 4.7. *For the 1D Euler and Navier–Stokes systems, the GQL representation of the invariant region \tilde{G} in (2.3) is given by*

$$\tilde{G}_* = \left\{ \mathbf{u} = (\rho, m, E)^\top : \rho > 0, \quad \tilde{\varphi}(\mathbf{u}; \rho_*, v_*) \geq 0 \quad \forall \rho_* \in \mathbb{R}^+ \quad \forall v_* \in \mathbb{R} \right\} \quad (4.11)$$

with $\tilde{\varphi}(\mathbf{u}; \rho_*, v_*) := \mathbf{u} \cdot \mathbf{n}_* + S_{min} \rho_*^\Gamma$ and $\mathbf{n}_* := \left(\frac{v_*^2}{2} - \frac{S_{min} \Gamma \rho_*^{\Gamma-1}}{\Gamma-1}, -v_*, 1 \right)^\top$.

Proof. We only need to handle the nonlinear constraint $\tilde{g}(\mathbf{u}) > 0$ in (2.3), with the boundary hypersurface $\tilde{\mathcal{S}} := \{ \mathbf{u}_* = (\rho_*, m_*, E_*) : \rho_* > 0, \tilde{g}(\mathbf{u}_*) = 0 \}$. Motivated from the equivalence of $\tilde{g}(\mathbf{u}) = 0$ and $p = S_{min} \rho^\Gamma$, we find a natural physics-based parametrization of $\tilde{\mathcal{S}}$ as

$$\tilde{\mathcal{S}} = \left\{ \mathbf{u}_* = \left(\rho_*, \rho_* v_*, \frac{1}{2} \rho_* v_*^2 + \frac{S_{min} \rho_*^\Gamma}{\Gamma-1} \right)^\top : \rho_* > 0, v_* \in \mathbb{R} \right\}.$$

Then we can derive the normal vector \mathbf{n}_* of \mathcal{S} at \mathbf{u}_* by cross product:

$$\frac{\partial \mathbf{u}_*}{\partial \rho_*} \times \frac{\partial \mathbf{u}_*}{\partial v_*} = \left(1, v_*, \frac{S_{\min} \Gamma \rho_*^{\Gamma-1}}{\Gamma-1} + \frac{v_*^2}{2} \right)^\top \times (0, \rho_*, \rho_* v_*)^\top = \rho_* \mathbf{n}_*.$$

By [Theorem 4.2](#) and $(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = \tilde{\varphi}(\mathbf{u}; \rho_*, v_*)$, we obtain the GQL representation [\(4.11\)](#). \square

Example 2: M1 model of radiative transfer.

THEOREM 4.8. *For the gray M1 moment system of radiative transfer, the GQL representation of the invariant region G in [\(2.5\)](#) is given by*

$$G_* = \{ \mathbf{u} = (E_r, \mathcal{F}_r)^\top \in \mathbb{R}^4 : E_r - \mathcal{F}_r \cdot \boldsymbol{\theta}_* \geq 0 \quad \forall \boldsymbol{\theta}_* \in \mathbb{S}_1(\mathbf{0}) \} \quad (4.12)$$

with $\mathbb{S}_1(\mathbf{0}) := \{ \mathbf{x} \in \mathbb{R}^3 : \|\mathbf{x}\| = 1 \}$ denoting the unit 3D sphere.

Proof. The constructive method is used for this example. The Cauchy–Schwarz inequality yields

$$\varphi(\mathbf{u}; \boldsymbol{\theta}_*) := E_r - \mathcal{F}_r \cdot \boldsymbol{\theta}_* \geq g(\mathbf{u}) \quad \forall \boldsymbol{\theta}_* \in \mathbb{S}_1(\mathbf{0}),$$

where equality holds for $\mathcal{F}_r \neq \mathbf{0}$ with $\boldsymbol{\theta}_* = \mathcal{F}_r / \|\mathcal{F}_r\|$ and for $\mathcal{F}_r = \mathbf{0}$ with any $\boldsymbol{\theta}_*$. Thus, $\min_{\boldsymbol{\theta}_* \in \mathbb{S}_1(\mathbf{0})} \varphi(\mathbf{u}; \boldsymbol{\theta}_*) = g(\mathbf{u})$, and by [Theorem 3.2](#) we get the GQL representation [\(4.12\)](#). \square

Example 3: Relativistic hydrodynamic system.

THEOREM 4.9. *For the 1D RHD system [\(2.6\)](#), the GQL representation of the invariant region G in [\(2.8\)](#) is given by*

$$G_* = \{ \mathbf{u} = (D, m, E)^\top : D > 0, \varphi(\mathbf{u}; v_*) > 0 \quad \forall v_* \in (-1, 1) \} \quad (4.13)$$

with $\varphi(\mathbf{u}; v_*) := E - mv_* - D\sqrt{1-v_*^2}$ being a linear function of \mathbf{u} .

Proof. The first constraint in [\(2.8\)](#) is linear. We deal with the second one by the constructive method. The Cauchy–Schwarz inequality implies

$$\varphi(\mathbf{u}; v_*) \geq E - \sqrt{D^2 + m^2} \sqrt{v_*^2 + \left(\sqrt{1-v_*^2} \right)^2} = E - \sqrt{D^2 + m^2} = g(\mathbf{u}),$$

where equality holds if $v_* = m/\sqrt{D^2 + m^2}$. This means $\min_{v_* \in (-1, 1)} \varphi(\mathbf{u}; v_*) = g(\mathbf{u})$. According to [Theorem 3.2](#), we get the GQL representation [\(4.13\)](#). \square

We now utilize the cross-product method to construct the GQL representation of the invariant region \tilde{G} in [\(2.9\)](#), where the minimum entropy principle is also included as a constraint.

THEOREM 4.10. *For the 1D RHD system [\(2.6\)](#), the GQL representation of the invariant region \tilde{G} in [\(2.9\)](#) is given by*

$$\tilde{G}_* = \left\{ \mathbf{u} = (D, m, E)^\top : \rho > 0, \tilde{\varphi}(\mathbf{u}; \rho_*, v_*) \geq 0 \quad \forall \rho_* \in \mathbb{R}^+ \quad \forall v_* \in (-1, 1) \right\} \quad (4.14)$$

with $\tilde{\varphi}(\mathbf{u}; \rho_*, v_*) := \mathbf{u} \cdot \mathbf{n}_* + S_{\min} \rho_*^\Gamma$ and $\mathbf{n}_* := \left(-\sqrt{1-v_*^2} \left(1 + \frac{S_{\min} \Gamma \rho_*^{\Gamma-1}}{\Gamma-1} \right), -v_*, 1 \right)^\top$.

Proof. We only need to tackle the second and third constraints in [\(2.9\)](#). For the third constraint $\tilde{g}(\mathbf{u}) \geq 0$, the corresponding boundary hypersurface is $\tilde{\mathcal{S}} := \{ \mathbf{u}_* = (\rho_*, m_*, E_*) : \rho_* > 0, g(\mathbf{u}_*) > 0, \tilde{g}(\mathbf{u}_*) = 0 \}$. Based on the equivalence of $\tilde{g}(\mathbf{u}) = 0$ and $p = S_{\min} \rho^\Gamma$, we obtain a natural physics-based parametrization of $\tilde{\mathcal{S}}$, namely,

$$\tilde{\mathcal{S}} = \left\{ \mathbf{u}_* = \left(\frac{\rho_*}{\sqrt{1-v_*^2}}, \frac{\left(\rho_* + \frac{S_{\min} \Gamma \rho_*^\Gamma}{\Gamma-1} \right) v_*}{1-v_*^2}, \frac{\rho_* + \frac{S_{\min} \Gamma \rho_*^\Gamma}{\Gamma-1}}{1-v_*^2} - S_{\min} \rho_*^\Gamma \right)^\top : \rho_* > 0, v_* \in (-1, 1) \right\}.$$

We can then derive the normal vector \mathbf{n}_* of \mathcal{S} at \mathbf{u}_* by cross product:

$$\frac{\partial \mathbf{u}_*}{\partial \rho_*} \times \frac{\partial \mathbf{u}_*}{\partial v_*} = \frac{1}{\delta_*} \mathbf{n}_*, \quad \text{with } \delta_* := (1-v_*)^{5/2} \left(\rho_* + \frac{S_{\min} \Gamma}{\Gamma-1} \rho_*^\Gamma (1+v_*^2 - \Gamma v_*^2) \right)^{-1}.$$

By [Theorem 4.2](#) and $(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = \tilde{\varphi}(\mathbf{u}; \rho_*, v_*)$, the GQL representation for $\tilde{g}(\mathbf{u}) \geq 0$ is

$$\tilde{\varphi}(\mathbf{u}; \rho_*, v_*) \geq 0 \quad \forall \rho_* \in \mathbb{R}^+, \quad \forall v_* \in \mathbb{R}. \quad (4.15)$$

Note that $S_{min} > 0$ and

$$g(\mathbf{u}) > g(\mathbf{u}) - \frac{S_{min}}{\Gamma - 1} \left(\frac{D^2}{\sqrt{D^2 + m^2}} \right)^\Gamma = \tilde{\varphi} \left(\mathbf{u}; \frac{D^2}{\sqrt{D^2 + m^2}}, \frac{m}{\sqrt{D^2 + m^2}} \right),$$

which means that [\(4.15\)](#) also implies $g(\mathbf{u}) > 0$ in [\(2.9\)](#). That is, the second and third constraints in [\(2.9\)](#) can be equivalently represented by [\(4.15\)](#). Therefore, we obtain the GQL representation [\(4.14\)](#). \square

Example 4: Ten-moment Gaussian closure system.

THEOREM 4.11. *For the 2D ten-moment Gaussian closure system [\(2.10\)](#), the GQL representation of the invariant region G in [\(2.12\)](#) is given by*

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^6 : \rho > 0, \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) > 0 \quad \forall \mathbf{v}_* \in \mathbb{R}^2 \quad \forall \mathbf{z} \in \mathbb{R}^2 \setminus \{\mathbf{0}\} \right\}, \quad (4.16)$$

where $\mathbf{u} := (\rho, \mathbf{m}, E_{11}, E_{12}, E_{22})^\top$, and the function $\varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*)$ is linear with respect to \mathbf{u} :

$$\varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) := \mathbf{z}^\top \left(\mathbf{E} - \mathbf{m} \otimes \mathbf{v}_* + \rho \frac{\mathbf{v}_* \otimes \mathbf{v}_*}{2} \right) \mathbf{z}. \quad (4.17)$$

Proof. We only need to deal with the nonlinear constraint in [\(2.12\)](#). Note that

$$\varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) = \mathbf{z}^\top \left(\mathbf{E} - \frac{\mathbf{m} \otimes \mathbf{m}}{2\rho} \right) \mathbf{z} + \frac{\rho}{2} \left| \mathbf{z} \cdot \left(\mathbf{v}_* - \frac{\mathbf{m}}{\rho} \right) \right|^2,$$

which implies $\min_{\mathbf{v}_* \in \mathbb{R}^2} \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) = \mathbf{z}^\top \left(\mathbf{E} - \frac{\mathbf{m} \otimes \mathbf{m}}{2\rho} \right) \mathbf{z}$. By [Theorem 3.2](#), we immediately obtain the GQL representation [\(4.16\)](#). \square

Example 5: Ideal MHD system.

THEOREM 4.12. *For the ideal MHD system [\(2.13\)](#), the GQL representation of the invariant region G in [\(2.14\)](#) is given by*

$$G_* = \left\{ \mathbf{u} = (\rho, \mathbf{m}, \mathbf{B}, E)^\top \in \mathbb{R}^8 : \rho > 0, \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) > 0 \quad \forall \mathbf{v}_*, \mathbf{B}_* \in \mathbb{R}^3 \right\} \quad (4.18)$$

with $\varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) := \mathbf{u} \cdot \mathbf{n}_* + \frac{\|\mathbf{B}_*\|^2}{2}$ and $\mathbf{n}_* := \left(\frac{\|\mathbf{v}_*\|^2}{2}, -\mathbf{v}_*, -\mathbf{B}_*, 1 \right)^\top$.

Proof. We use the gradient-based method. For the nonlinear constraint in [\(2.14\)](#), the gradient of $g(\mathbf{u})$ is $\nabla g(\mathbf{u}) = \left(\frac{\|\mathbf{m}\|^2}{2\rho^2}, -\frac{\mathbf{m}}{\rho}, -\mathbf{B}, 1 \right)^\top$, and the corresponding boundary hypersurface is $\mathcal{S} := \{ \mathbf{u}_* = (\rho_*, \mathbf{m}_*, \mathbf{B}_*, E_*)^\top : \rho_* > 0, g(\mathbf{u}_*) = 0 \}$. Based on the equivalence of $g(\mathbf{u}) = 0$ and $p = 0$, we obtain a natural physics-based parametrization of \mathcal{S} , namely,

$$\mathcal{S} = \left\{ \mathbf{u}_* = \left(\rho_*, \rho_* \mathbf{v}_*, \mathbf{B}_*, \frac{1}{2} (\rho_* \|\mathbf{v}_*\|^2 + \|\mathbf{B}_*\|^2) \right)^\top : \rho_* > 0, \mathbf{v}_* \in \mathbb{R}^3, \mathbf{B}_* \in \mathbb{R}^3 \right\}.$$

For $\mathbf{u}_* \in \mathcal{S}$ and $\mathbf{u} = (\rho, \mathbf{m}, \mathbf{B}, E)^\top$, we have $(\mathbf{u} - \mathbf{u}_*) \cdot \nabla g(\mathbf{u}_*) = \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*)$. By [Theorem 4.1](#), we obtain the GQL representation [\(4.18\)](#). \square

Example 6: Relativistic MHD system.

THEOREM 4.13. *For the relativistic MHD system [\(2.15\)](#), the GQL representation of the invariant region G in [\(2.16\)](#) is given by*

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^8 : D > 0, \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) > 0 \quad \forall \mathbf{B}_* \in \mathbb{R}^3 \quad \forall \mathbf{v}_* \in \mathbb{B}_1(\mathbf{0}) \right\}, \quad (4.19)$$

where $\mathbf{u} = (D, \mathbf{m}, \mathbf{B}, E)^\top$, $\mathbb{B}_1(\mathbf{0}) := \{\mathbf{x} \in \mathbb{R}^3 : \|\mathbf{x}\| \leq 1\}$ is a unit 3D ball, and the linear function $\varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) := \mathbf{u} \cdot \mathbf{n}_* + p_m^*$ with

$$p_m^* := \frac{1}{2} \left((1 - \|\mathbf{v}_*\|^2) \|\mathbf{B}_*\|^2 + (\mathbf{v}_* \cdot \mathbf{B}_*)^2 \right), \quad (4.20)$$

$$\mathbf{n}_* := \left(-\sqrt{1 - \|\mathbf{v}_*\|^2}, -\mathbf{v}_*, -(1 - \|\mathbf{v}_*\|^2)\mathbf{B}_* - (\mathbf{v}_* \cdot \mathbf{B}_*)\mathbf{v}_*, 1 \right)^\top. \quad (4.21)$$

Note that p_m^* and \mathbf{n}_* only depend on the free auxiliary variables $(\mathbf{v}_*, \mathbf{B}_*)$.

Proof. As shown in [57], the region G in (2.16) can be equivalently represented as

$$G = \{\mathbf{u} \in \mathbb{R}^8 : D > 0, g_2(\mathbf{u}) > 0, p(\mathbf{u}) > 0\} \quad (4.22)$$

with $g_2(\mathbf{u}) := E - \sqrt{D^2 + \|\mathbf{m}\|^2}$. Although the implicit function $p(\mathbf{u})$ defined in (1.7) can not be explicitly formulated, the corresponding boundary hypersurface $\mathcal{S} := \{\mathbf{u}_* = (D_*, \mathbf{m}_*, \mathbf{B}_*, E_*)^\top : D_* > 0, g_2(\mathbf{u}_*) > 0, p(\mathbf{u}_*) = 0\}$ has an explicit physics-based parameterization:

$$\mathcal{S} = \left\{ \mathbf{u}_* = \left(\rho_* \gamma_*, \rho_* \gamma_*^2 \mathbf{v}_* + \|\mathbf{B}_*\|^2 \mathbf{v}_* - (\mathbf{v}_* \cdot \mathbf{B}_*) \mathbf{B}_*, \mathbf{B}_*, \right. \right. \\ \left. \left. \rho_* \gamma_*^2 + \|\mathbf{B}_*\|^2 - p_m^* \right)^\top : \rho_* > 0, \mathbf{B}_* \in \mathbb{R}^3, \mathbf{v}_* \in \mathbb{B}_1(\mathbf{0}) \right\}$$

with p_m^* defined in (4.20) and $\gamma_* := (1 - \|\mathbf{v}_*\|^2)^{\frac{1}{2}}$. This parameterization is helpful for dealing with the highly nonlinear constraint $p(\mathbf{u}) > 0$ by the cross-product method. For $1 \leq i \leq 3$, denote $\mathbf{e}_i := (\delta_{1i}, \delta_{2i}, \delta_{3i})$ with δ_{ij} being the Kronecker delta. Taking the partial derivatives of \mathbf{u}_* with respect to the parametric variables $\{\rho_*, \mathbf{v}_*, \mathbf{B}_*\}$ gives

$$\begin{aligned} \frac{\partial \mathbf{u}_*}{\partial \rho_*} &= (\gamma_*, \gamma_*^2 \mathbf{v}_*, 0, 0, 0, \gamma_*^2)^\top, \\ \frac{\partial \mathbf{u}_*}{\partial v_{i*}} &= \left(\rho_* \gamma_*^3 v_{i*}, (\rho_* \gamma_*^2 + \|\mathbf{B}_*\|^2) \mathbf{e}_i + 2\rho_* \gamma_*^4 v_{i*} \mathbf{v}_* - B_{i*} \mathbf{B}_*, 0, 0, 0, 2\rho_* \gamma_*^4 v_{i*} - B_{i*} (\mathbf{v}_* \cdot \mathbf{B}_*) + \|\mathbf{B}_*\|^2 v_{i*} \right)^\top, \\ \frac{\partial \mathbf{u}_*}{\partial B_{i*}} &= \left(0, -(\mathbf{v}_* \cdot \mathbf{B}_*) \mathbf{e}_i + 2B_{i*} \mathbf{v}_* - v_{i*} \mathbf{B}_*, \mathbf{e}_i, B_{i*} (1 + \|\mathbf{v}_*\|^2) - v_{i*} (\mathbf{v}_* \cdot \mathbf{B}_*) \right), \quad 1 \leq i \leq 3, \end{aligned}$$

which are all perpendicular to the nonzero vector \mathbf{n}_* defined in (4.21). This means \mathbf{n}_* is parallel to the cross product $\frac{\partial \mathbf{u}_*}{\partial \rho_*} \times \left(\bigwedge_{i=1}^3 \frac{\partial \mathbf{u}_*}{\partial v_{i*}} \right) \times \left(\bigwedge_{i=1}^3 \frac{\partial \mathbf{u}_*}{\partial B_{i*}} \right)$, implying that \mathbf{n}_* is a normal vector of \mathcal{S} at \mathbf{u}_* . It can be verified that \mathbf{n}_* is always directed towards the concave side of \mathcal{S} . By Theorem 4.2 and $(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*)$, we know that the GQL representation for $p(\mathbf{u}) > 0$ is

$$\varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) > 0 \quad \forall \mathbf{B}_* \in \mathbb{R}^3 \quad \forall \mathbf{v}_* \in \mathbb{B}_1(\mathbf{0}). \quad (4.23)$$

If taking $\mathbf{v}_* = \mathbf{m} / \sqrt{D^2 + \|\mathbf{m}\|^2}$ and $\mathbf{B}_* = \mathbf{0}$, we obtain $\varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) = g_2(\mathbf{u})$, which means that (4.23) also implies $g_2(\mathbf{u}) > 0$ in (4.22). In other words, the second and third constraints in (4.22) can be equivalently represented by (4.23). Therefore, we obtain the GQL representation (4.19). \square

5. Geometric quasilinearization for bound-preserving analysis. This section applies the GQL approach to analyze the bound-preserving property of numerical schemes and shows its remarkable advantages over direct and traditional approaches by diverse examples covering different schemes of three PDE systems in one and two dimensions. We only focus on first-order schemes for illustrative purposes, while the GQL approach is readily extensible to high-order schemes. The application of GQL to design high-order bound-preserving scheme will also be explored in section 6 for the multicomponent MHD system, to further demonstrate its capability in addressing challenging bound-preserving problems that could not be coped with by direct approaches.

5.1. Example 1: Euler system. Consider a finite volume scheme

$$\bar{\mathbf{u}}_j^{n+1} = \bar{\mathbf{u}}_j^n - \sigma \left(\hat{\mathbf{f}}_{j+\frac{1}{2}} - \hat{\mathbf{f}}_{j-\frac{1}{2}} \right), \quad (5.1)$$

for solving the 1D Euler system (2.1) on a uniform spatial mesh $\{[x_{j-1/2}, x_{j+1/2}]\}$ with $\sigma := \Delta t / \Delta x$ denoting the ratio of the temporal step-size Δt to the spatial step-size Δx . Here $\bar{\mathbf{u}}_j^n$ is an approximation to the average of $\mathbf{u}(x, t_n)$ on cell $[x_{j-1/2}, x_{j+1/2}]$, and $\hat{\mathbf{f}}_{j+1/2}$ is a numerical flux at $x_{j+1/2}$. For system (2.1), it holds that $\mathbf{f}(\mathbf{u}) = v\mathbf{u} + p(0, 1, v)^\top$, which will be used in the following analysis.

We apply the GQL approach to analyze the bound-preserving property of the scheme (5.1) with the invariant region G defined in (2.2). Thanks to the GQL representation in Theorem 4.4, we have

$$G = G_* = \{\mathbf{u} : \mathbf{u} \cdot \mathbf{e}_1 > 0, \mathbf{u} \cdot \mathbf{n}_* > 0 \quad \forall v_* \in \mathbb{R}\} \quad (5.2)$$

with $\mathbf{e}_1 := (1, 0, 0)^\top$ and $\mathbf{n}_* := (\frac{v_*^2}{2}, -v_*, 1)^\top$. GQL transfers the bound-preserving problem into preserving the positivity of $\mathbf{u} \cdot \mathbf{e}_1$ and $\mathbf{u} \cdot \mathbf{n}_*$, which are all linear with respect to \mathbf{u} and helpful for bound-preserving study.

Example 1.1: Lax–Friedrichs scheme. To clearly illustrate the basic idea, we begin with the simple Lax–Friedrichs scheme with the numerical flux $\hat{\mathbf{f}}_{j+1/2}$ taken as

$$\hat{\mathbf{f}}^{\text{LF}}(\bar{\mathbf{u}}_j^n, \bar{\mathbf{u}}_{j+1}^n) := \frac{1}{2} \left(\mathbf{f}(\bar{\mathbf{u}}_j^n) + \mathbf{f}(\bar{\mathbf{u}}_{j+1}^n) - \alpha_n (\bar{\mathbf{u}}_{j+1}^n - \bar{\mathbf{u}}_j^n) \right), \quad (5.3)$$

where $\alpha_n := \max_j \alpha(\bar{\mathbf{u}}_j^n)$ with $\alpha(\mathbf{u}) := |v| + \sqrt{\Gamma p / \rho}$ being the spectral radius of the Jacobian matrix $\partial \mathbf{f} / \partial \mathbf{u}$. Given that $\bar{\mathbf{u}}_j^n \in G$ for all j , we wish $\bar{\mathbf{u}}_j^{n+1} \in G$.

For respectively $\mathbf{n} = \mathbf{e}_1$ and $\mathbf{n} = \mathbf{n}_*$, thanks to the linearity of $\mathbf{u} \cdot \mathbf{n}$ we obtain

$$\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n} = (1 - \sigma \alpha_n) \bar{\mathbf{u}}_j^n \cdot \mathbf{n} + \frac{\sigma}{2} \left(\alpha_n \bar{\mathbf{u}}_{j+1}^n \cdot \mathbf{n} - \mathbf{f}(\bar{\mathbf{u}}_{j+1}^n) \cdot \mathbf{n} + \alpha_n \bar{\mathbf{u}}_{j-1}^n \cdot \mathbf{n} + \mathbf{f}(\bar{\mathbf{u}}_{j-1}^n) \cdot \mathbf{n} \right).$$

The problem boils down to control the effect of $\mathbf{f}(\bar{\mathbf{u}}_{\pm 1}^n) \cdot \mathbf{n}$ by using the positivity of $\bar{\mathbf{u}}_{\pm 1}^n \cdot \mathbf{n}$. For any $\mathbf{u} \in G$, we have $\mathbf{u} \cdot \mathbf{n} > 0$ and

$$\begin{aligned} \pm \mathbf{f}(\mathbf{u}) \cdot \mathbf{e}_1 &= \pm v (\mathbf{u} \cdot \mathbf{e}_1) < \alpha(\mathbf{u}) \mathbf{u} \cdot \mathbf{e}_1, \\ \pm \mathbf{f}(\mathbf{u}) \cdot \mathbf{n}_* &= \pm v (\mathbf{u} \cdot \mathbf{n}_*) \pm p (v - v_*) \\ &\leq |v| (\mathbf{u} \cdot \mathbf{n}_*) + \left(\frac{1}{2} \rho (v - v_*)^2 + \rho e \right) \frac{p}{\rho \sqrt{2e}} \\ &= \left(|v| + \frac{p}{\rho \sqrt{2e}} \right) \mathbf{u} \cdot \mathbf{n}_* < \alpha(\mathbf{u}) \mathbf{u} \cdot \mathbf{n}_*, \end{aligned}$$

which yield $\alpha_n \bar{\mathbf{u}}_{j\pm 1}^n \cdot \mathbf{n} \mp \mathbf{f}(\bar{\mathbf{u}}_{j\pm 1}^n) \cdot \mathbf{n} > 0$. Thus we obtain $\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n} > (1 - \sigma \alpha_n) \bar{\mathbf{u}}_j^n \cdot \mathbf{n} \geq 0$ provided that $\sigma \alpha_n \leq 1$. This proves that the scheme (5.1) with the Lax–Friedrichs flux (5.3) is bound-preserving under the standard CFL condition $\sigma \alpha_n \leq 1$.

Remark 5.1. As we have seen, unlike the traditional approaches that require substituting the target scheme into the original nonlinear constraint of G in (2.2), the GQL approach skillfully transfers all the constraints into linear ones which can be investigated in a unified way.

Example 1.2: Gas-kinetic scheme. In order to demonstrate the advantages of the GQL approach in bound-preserving analysis, we consider a challenging example—the gas-kinetic scheme with the numerical flux $\hat{\mathbf{f}}_{j+1/2}$ taken as

$$\hat{\mathbf{f}}^{\text{GK}}(\bar{\mathbf{u}}_j^n, \bar{\mathbf{u}}_{j+1}^n) := \mathbf{f}^+(\bar{\mathbf{u}}_j^n) + \mathbf{f}^-(\bar{\mathbf{u}}_{j+1}^n). \quad (5.4)$$

$$\mathbf{f}^\pm(\mathbf{u}) := \int_{\mathbb{R}^\pm} \int_{\mathbb{R}^M} \begin{pmatrix} w \\ w^2 \\ \frac{w}{2}(w^2 + \boldsymbol{\xi}^2) \end{pmatrix} F(w, \boldsymbol{\xi}; \mathbf{u}) d\boldsymbol{\xi} dw, \quad (5.5)$$

where w is the particle velocity, $\boldsymbol{\xi} \in \mathbb{R}^M$ denotes the internal variables whose degrees of freedom $M = (3 - \Gamma) / (\Gamma - 1)$, the equilibrium distribution function F is

$$F(w, \boldsymbol{\xi}; \mathbf{u}) := \rho \left(\frac{\lambda}{\pi} \right)^{\frac{M+1}{2}} e^{-\lambda((w-v)^2 + \|\boldsymbol{\xi}\|^2)} \quad (5.6)$$

with ρ being the fluid density, v being the fluid velocity, and $\lambda = \rho / (2p)$.

In a traditional approach [49], the bound-preserving property of this scheme was studied by: (i) first, evaluating the integration (5.5) as

$$\mathbf{f}^\pm(\mathbf{u}) = \rho \begin{pmatrix} \frac{v}{2} \operatorname{erfc}(\mp \sqrt{\lambda} v) \pm \frac{1}{2} \frac{e^{-\lambda v^2}}{\sqrt{\pi \lambda}} \\ \left(\frac{v^2}{2} + \frac{1}{4\lambda} \right) \operatorname{erfc}(\mp \sqrt{\lambda} v) \pm \frac{v}{2} \frac{e^{-\lambda v^2}}{\sqrt{\pi \lambda}} \\ \left(\frac{v^3}{4} + \frac{M+3}{8\lambda} v \right) \operatorname{erfc}(\mp \sqrt{\lambda} v) \pm \left(\frac{v^2}{4} + \frac{M+2}{8\lambda} \right) \frac{e^{-\lambda v^2}}{\sqrt{\pi \lambda}} \end{pmatrix} \quad (5.7)$$

with $\operatorname{erfc}(x) := \frac{2}{\sqrt{\pi}} \int_x^{+\infty} e^{-w^2} dw$; (ii) then, plugging the numerical flux (5.4) with (5.7) into (5.1) and splitting the scheme (5.1) into two steps; (iii) and finally checking the bound-preserving properties of the split schemes by verifying the original constraints of G in (2.2). For this scheme, verifying the nonlinear constraint in (2.2) are difficult and complicated.

Benefited from its linear feature, the GQL approach is highly effective for this challenging case. For $\mathbf{n} = \mathbf{e}_1$ or $\mathbf{n} = \mathbf{n}_*$, thanks to the linearity of $\mathbf{u} \cdot \mathbf{n}$ we obtain

$$\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n} = \bar{\mathbf{u}}_j^n \cdot \mathbf{n} - \sigma (\mathbf{f}^+(\bar{\mathbf{u}}_j^n) - \mathbf{f}^-(\bar{\mathbf{u}}_j^n)) \cdot \mathbf{n} - \sigma \mathbf{f}^-(\bar{\mathbf{u}}_{j+1}^n) \cdot \mathbf{n} + \sigma \mathbf{f}^+(\bar{\mathbf{u}}_{j-1}^n) \cdot \mathbf{n}.$$

Note that for any $\mathbf{u} \in G$, we have $F(w, \boldsymbol{\xi}; \mathbf{u}) > 0$ and

$$\begin{aligned} \pm \mathbf{f}^\pm(\mathbf{u}) \cdot \mathbf{e}_1 &= \int_{\mathbb{R}^\pm} \int_{\mathbb{R}^M} |w| F(w, \boldsymbol{\xi}; \mathbf{u}) d\boldsymbol{\xi} dw > 0, \\ \pm \mathbf{f}^\pm(\mathbf{u}) \cdot \mathbf{n}_* &= \int_{\mathbb{R}^\pm} \int_{\mathbb{R}^M} \frac{|w|}{2} \left((w - v_*)^2 + \|\boldsymbol{\xi}\|^2 \right) F(w, \boldsymbol{\xi}; \mathbf{u}) d\boldsymbol{\xi} dw > 0. \end{aligned}$$

It follows, for $\mathbf{n} = \mathbf{e}_1$ and $\mathbf{n} = \mathbf{n}_*$ respectively, that

$$\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n} > \bar{\mathbf{u}}_j^n \cdot \mathbf{n} - \sigma (\mathbf{f}^+(\bar{\mathbf{u}}_j^n) - \mathbf{f}^-(\bar{\mathbf{u}}_j^n)) \cdot \mathbf{n}. \quad (5.8)$$

Next, we use the positivity of $\mathbf{u} \cdot \mathbf{n}$ to bound the effect of $(\mathbf{f}^+(\mathbf{u}) - \mathbf{f}^-(\mathbf{u})) \cdot \mathbf{n}$ as follows:

$$\begin{aligned} (\mathbf{f}^+(\mathbf{u}) - \mathbf{f}^-(\mathbf{u})) \cdot \mathbf{e}_1 &= (\mathbf{u} \cdot \mathbf{e}_1) \left(\frac{\lambda}{\pi} \right)^{\frac{1}{2}} \left(\int_{\mathbb{R}} |w| e^{-\lambda(w-v)^2} dw \right) \\ &\leq (\mathbf{u} \cdot \mathbf{e}_1) \left(\frac{\lambda}{\pi} \right)^{\frac{1}{2}} \left(\int_{\mathbb{R}} (|v| + |w-v|) e^{-\lambda(w-v)^2} dw \right) \\ &= (\mathbf{u} \cdot \mathbf{e}_1) \left(|v| + 1/\sqrt{\pi\lambda} \right) < a(\mathbf{u}) \mathbf{u} \cdot \mathbf{e}_1, \\ (\mathbf{f}^+(\mathbf{u}) - \mathbf{f}^-(\mathbf{u})) \cdot \mathbf{n}_* &= \int_{\mathbb{R}} \int_{\mathbb{R}^M} \frac{|w|}{2} \left((w - v_*)^2 + \|\boldsymbol{\xi}\|^2 \right) F(w, \boldsymbol{\xi}; \mathbf{u}) d\boldsymbol{\xi} dw \\ &\leq \int_{\mathbb{R}} \int_{\mathbb{R}^M} \frac{|v| + |w-v|}{2} \left((w - v_*)^2 + \|\boldsymbol{\xi}\|^2 \right) F(w, \boldsymbol{\xi}; \mathbf{u}) d\boldsymbol{\xi} dw \\ &= |v| (\mathbf{u} \cdot \mathbf{n}_*) + \frac{\rho}{2\sqrt{\pi\lambda}} \left((v - v_*)^2 + \frac{M+2}{2\lambda} \right) \\ &\leq |v| (\mathbf{u} \cdot \mathbf{n}_*) + \frac{\rho}{2\sqrt{\pi\lambda}} \left((v - v_*)^2 + \frac{M+1}{2\lambda} \right) \frac{M+2}{M+1} \\ &= \left(|v| + \frac{M+2}{M+1} (\pi\lambda)^{-\frac{1}{2}} \right) (\mathbf{u} \cdot \mathbf{n}_*) < a(\mathbf{u}) (\mathbf{u} \cdot \mathbf{n}_*). \end{aligned}$$

This implies $(\mathbf{f}^+(\bar{\mathbf{u}}_j^n) - \mathbf{f}^-(\bar{\mathbf{u}}_j^n)) \cdot \mathbf{n} < \alpha(\bar{\mathbf{u}}_j^n) \bar{\mathbf{u}}_j^n \cdot \mathbf{n} \leq \alpha_n \bar{\mathbf{u}}_j^n \cdot \mathbf{n}$. It then follows from (5.8) that $\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n} > (1 - \sigma\alpha_n) \bar{\mathbf{u}}_j^n \cdot \mathbf{n} \geq 0$ provided that $\sigma\alpha_n \leq 1$. This proves that the scheme (5.1) with the gas-kinetic flux (5.4) is bound-preserving under the standard CFL condition $\sigma\alpha_n \leq 1$.

Remark 5.2. The linearity of GQL brought by introducing the free auxiliary variable v_* gives remarkable advantages in our above analysis. Because v_* is independent of all the system variables \mathbf{u} , it can freely move cross the integrals. We no longer need to substitute a complicated scheme into the nonlinear function $g(\mathbf{u})$ in (2.2) to verify $g(\mathbf{u}) > 0$. Instead, we work on the simpler but equivalent linear constraint $\mathbf{u} \cdot \mathbf{n}_* > 0$. The interested readers may compare the above analysis based on GQL and the traditional analysis in [49].

5.2. Example 2: Navier–Stokes system. Consider the scheme

$$\bar{\mathbf{u}}_j^{n+1} = \bar{\mathbf{u}}_j^n - \sigma (\hat{\mathbf{f}}_{j+\frac{1}{2}} - \hat{\mathbf{f}}_{j-\frac{1}{2}}) + \frac{\Delta t}{\Delta x^2} \frac{\eta}{\operatorname{Re}} \mathbf{H}_j \quad (5.9)$$

with $\mathbf{H}_j := \mathbf{r}(\bar{\mathbf{u}}_{j+1}^n) - 2\mathbf{r}(\bar{\mathbf{u}}_j^n) + \mathbf{r}(\bar{\mathbf{u}}_{j-1}^n)$, for solving the 1D dimensionless compressible Navier–Stokes equations (2.4). Here $\hat{\mathbf{f}}_{j+1/2}$ is taken as a bound-preserving numerical flux for the 1D Euler system

(2.1), for example, the Lax–Friedrichs flux (5.3) or the gas-kinetic flux (5.4), which satisfy: if $\bar{\mathbf{u}}_j^n \in G$ for all j , then

$$(\hat{\mathbf{f}}_{j+\frac{1}{2}} - \hat{\mathbf{f}}_{j-\frac{1}{2}}) \cdot \mathbf{n} < \alpha_n \bar{\mathbf{u}}_j^n \cdot \mathbf{n} \quad \forall j$$

holds for respectively $\mathbf{n} = \mathbf{e}_1$ and $\mathbf{n} = \mathbf{n}_*$, according to the analysis in subsection 5.1. Thus we have

$$\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n} > (1 - \sigma\alpha_n) \bar{\mathbf{u}}_j^n \cdot \mathbf{n} + \frac{\Delta t}{\Delta x^2} \frac{\eta}{\text{Re}} \mathbf{H}_j \cdot \mathbf{n}. \quad (5.10)$$

Thanks to GQL, we clearly see that the bound-preserving essence is to control the potentially negative term $\mathbf{H}_j \cdot \mathbf{n}$ by the positive term $\bar{\mathbf{u}}_j^n \cdot \mathbf{n}$. Note $\mathbf{H}_j \cdot \mathbf{e}_1 = 0$, thereby $\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{e}_1 > (1 - \sigma\alpha_n) \bar{\mathbf{u}}_j^n \cdot \mathbf{e}_1 \geq 0$ if $\sigma\alpha_n \leq 1$. For any $\mathbf{u} \in G$ and $v_* \in \mathbb{R}$, we have

$$-\frac{v_*^2}{2} < \mathbf{r}(\mathbf{u}) \cdot \mathbf{n}_* = \frac{1}{2}(v - v_*)^2 + \frac{\Gamma}{\text{Pr} \eta} e - \frac{v_*^2}{2} \leq \max \left\{ 1, \frac{\Gamma}{\text{Pr} \eta} \right\} \frac{1}{\rho} (\mathbf{u} \cdot \mathbf{n}_*) - \frac{v_*^2}{2}.$$

This gives

$$\begin{aligned} \mathbf{H}_j \cdot \mathbf{n}_* &= \left(\mathbf{r}(\bar{\mathbf{u}}_{j+1}^n) \cdot \mathbf{n}_* + \mathbf{r}(\bar{\mathbf{u}}_{j-1}^n) \cdot \mathbf{n}_* \right) - 2\mathbf{r}(\bar{\mathbf{u}}_j^n) \cdot \mathbf{n}_* \\ &\geq \left(-\frac{v_*^2}{2} - \frac{v_*^2}{2} \right) - 2 \left(\max \left\{ 1, \frac{\Gamma}{\text{Pr} \eta} \right\} \frac{1}{\rho_j^n} (\bar{\mathbf{u}}_j^n \cdot \mathbf{n}_*) - \frac{v_*^2}{2} \right) \\ &= -\frac{2}{\rho_j^n} \max \left\{ 1, \frac{\Gamma}{\text{Pr} \eta} \right\} (\bar{\mathbf{u}}_j^n \cdot \mathbf{n}_*). \end{aligned}$$

It then follows from (5.10) that

$$\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n}_* > (1 - \sigma\alpha_n) \bar{\mathbf{u}}_j^n \cdot \mathbf{n}_* - \frac{\Delta t}{\Delta x^2} \frac{\eta}{\text{Re}} \frac{2}{\rho_j^n} \max \left\{ 1, \frac{\Gamma}{\text{Pr} \eta} \right\} (\bar{\mathbf{u}}_j^n \cdot \mathbf{n}_*).$$

We then immediately have $\bar{\mathbf{u}}_j^{n+1} \cdot \mathbf{n}_* > 0$, provided that

$$\alpha_n \frac{\Delta t}{\Delta x} + \frac{\Delta t}{\Delta x^2} \frac{2}{\rho_j^n \text{Re}} \max \left\{ \eta, \frac{\Gamma}{\text{Pr}} \right\} \leq 1. \quad (5.11)$$

In conclusion, the scheme (5.9) is bound-preserving under condition (5.11).

Remark 5.3. A standard approach for handling bound-preserving problems with multiple terms (e.g., convection term and diffusion term [66], or convection term and source term [69]) is based on decomposing the schemes into a convex combination of some subterms, and then enforcing all the subterms in G . This may lead to stricter conditions on the time step-size Δt . Since the linear feature of GQL has already naturally incorporated the convexity of G into the GQL representation, technical convex decomposition is not necessary in the GQL approach.

5.3. Example 3: Ten-moment Gaussian closure system. Consider the scheme

$$\bar{\mathbf{u}}_{ij}^{n+1} = \bar{\mathbf{u}}_{ij}^n - \sigma_1 (\hat{\mathbf{f}}_{1,i+\frac{1}{2},j} - \hat{\mathbf{f}}_{1,i-\frac{1}{2},j}) - \sigma_2 (\hat{\mathbf{f}}_{2,i,j+\frac{1}{2}} - \hat{\mathbf{f}}_{2,i,j-\frac{1}{2}}), \quad (5.12)$$

for solving the 2D Gaussian closure equations (2.10) on a uniform Cartesian mesh $\{[x_{i-1/2}, x_{i+1/2}] \times [y_{j-1/2}, y_{j+1/2}]\}$, with $\sigma_1 = \frac{\Delta t}{\Delta x}$, $\sigma_2 = \frac{\Delta t}{\Delta y}$. Here $\bar{\mathbf{u}}_{ij}^n$ denotes an approximation to the average of $\mathbf{u}(x, y, t_n)$ on each cell, and the Lax-Friedrichs numerical fluxes are considered, i.e.

$$\hat{\mathbf{f}}_{1,i+1/2,j} = \hat{\mathbf{f}}_1^{\text{LF}}(\bar{\mathbf{u}}_{ij}^n, \bar{\mathbf{u}}_{i+1,j}^n), \quad \hat{\mathbf{f}}_{2,i,j+1/2} = \hat{\mathbf{f}}_2^{\text{LF}}(\bar{\mathbf{u}}_{ij}^n, \bar{\mathbf{u}}_{i,j+1}^n), \quad (5.13)$$

$$\hat{\mathbf{f}}_\ell^{\text{LF}}(\mathbf{u}^L, \mathbf{u}^R) := \frac{1}{2} \left(\mathbf{f}_\ell(\mathbf{u}^L) + \mathbf{f}_\ell(\mathbf{u}^R) - \alpha_{\ell,n}(\mathbf{u}^R - \mathbf{u}^L) \right), \quad \ell = 1, 2, \quad (5.14)$$

where $\alpha_{\ell,n} = \max_{ij} \alpha_\ell(\bar{\mathbf{u}}_{ij}^n)$, and $\alpha_\ell(\mathbf{u}) := |v_\ell| + \sqrt{p_{\ell\ell}/\rho}$.

In the original form (2.11) of G , the second constraint is the positive definiteness of a matrix $\mathbf{E} - \frac{\mathbf{m} \otimes \mathbf{m}}{2\rho}$ which nonlinearly depends on \mathbf{u} . This leads to the challenges in the bound-preserving study. Thanks to Theorem 4.11, the invariant region G is equivalently represented as

$$G_* = \left\{ \mathbf{u} \in \mathbb{R}^6 : \mathbf{u} \cdot \mathbf{e}_1 > 0, \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) > 0 \quad \forall \mathbf{v}_* \in \mathbb{R}^2 \quad \forall \mathbf{z} \in \mathbb{R}^2 \setminus \{\mathbf{0}\} \right\}, \quad (5.15)$$

where $\mathbf{e}_1 := (1, 0, \dots, 0)^\top$ and the linear function $\varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*)$ is defined by (4.17).

We apply the GQL approach to investigate the bound-preserving property of the scheme (5.12) with (5.13). Similar to the Euler system, for any $\mathbf{u} \in G$ we have $\mathbf{f}_\ell(\mathbf{u}) \cdot \mathbf{e}_1 = v_\ell(\mathbf{u} \cdot \mathbf{e}_1)$ and

$$\pm \mathbf{f}_\ell(\mathbf{u}) \cdot \mathbf{e}_1 \leq |v_\ell|(\mathbf{u} \cdot \mathbf{e}_1) < \alpha_\ell(\mathbf{u})(\mathbf{u} \cdot \mathbf{e}_1),$$

which gives $\bar{\mathbf{u}}_{ij}^{n+1} \cdot \mathbf{e}_1 > 0$ under the CFL condition $\sigma_1 \alpha_{1,n} + \sigma_2 \alpha_{2,n} < 1$. In the following, we focus on the second constraint in (5.15). Thanks to the linearity of $\varphi(\cdot; \mathbf{z}, \mathbf{v}_*)$, we obtain

$$\varphi(\mathbf{f}_1(\mathbf{u}); \mathbf{z}, \mathbf{v}_*) = v_1 \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) + [\mathbf{z} \cdot (\mathbf{v} - \mathbf{v}_*)] (\mathbf{p}_1 \cdot \mathbf{z}) \quad (5.16)$$

with the vector $\mathbf{p}_1 := (p_{11}, p_{12})^\top$. For any $\mathbf{u} \in G$, using the AM-GM inequality gives

$$\begin{aligned} \left| [\mathbf{z} \cdot (\mathbf{v} - \mathbf{v}_*)] (\mathbf{p}_1 \cdot \mathbf{z}) \right| &\leq \frac{1}{2} \sqrt{\rho p_{11}} |\mathbf{z} \cdot (\mathbf{v} - \mathbf{v}_*)|^2 + \frac{1}{2\sqrt{\rho p_{11}}} |\mathbf{p}_1 \cdot \mathbf{z}|^2 \\ &= \sqrt{\frac{p_{11}}{\rho}} \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) - \frac{z_2^2 \det(\mathbf{p})}{2\sqrt{\rho p_{11}}} \leq \sqrt{\frac{p_{11}}{\rho}} \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*), \end{aligned}$$

which together with the identity (5.16) yields

$$\pm \varphi(\mathbf{f}_1(\mathbf{u}); \mathbf{z}, \mathbf{v}_*) \leq \left(|v_1| + \sqrt{p_{11}/\rho} \right) \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*) = \alpha_1(\mathbf{u}) \varphi(\mathbf{u}; \mathbf{z}, \mathbf{v}_*). \quad (5.17)$$

Using the linearity of $\varphi(\cdot; \mathbf{z}, \mathbf{v}_*)$ again and (5.17), we obtain

$$\begin{aligned} \varphi\left(\hat{\mathbf{f}}_{1,i+\frac{1}{2},j} - \hat{\mathbf{f}}_{1,i-\frac{1}{2},j}; \mathbf{z}, \mathbf{v}_*\right) &= \frac{1}{2} [\varphi(\mathbf{f}_1(\bar{\mathbf{u}}_{i+1,j}^n); \mathbf{z}, \mathbf{v}_*) - \alpha_{1,n} \varphi(\bar{\mathbf{u}}_{i+1,j}^n; \mathbf{z}, \mathbf{v}_*)] \\ &\quad + \frac{1}{2} [-\varphi(\mathbf{f}_1(\bar{\mathbf{u}}_{i-1,j}^n); \mathbf{z}, \mathbf{v}_*) - \alpha_{1,n} \varphi(\bar{\mathbf{u}}_{i-1,j}^n; \mathbf{z}, \mathbf{v}_*)] \\ &\quad + \alpha_{1,n} \varphi(\bar{\mathbf{u}}_{ij}^n; \mathbf{z}, \mathbf{v}_*) \leq \alpha_{1,n} \varphi(\bar{\mathbf{u}}_{ij}^n; \mathbf{z}, \mathbf{v}_*). \end{aligned}$$

Similarly, we have $\varphi(\hat{\mathbf{f}}_{2,i,j+\frac{1}{2}} - \hat{\mathbf{f}}_{2,i,j-\frac{1}{2}}; \mathbf{z}, \mathbf{v}_*) \leq \alpha_{2,n} \varphi(\bar{\mathbf{u}}_{ij}^n; \mathbf{z}, \mathbf{v}_*)$. It then follows that

$$\varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{z}, \mathbf{v}_*) \geq (1 - \sigma_1 \alpha_{1,n} - \sigma_2 \alpha_{2,n}) \varphi(\bar{\mathbf{u}}_{ij}^n; \mathbf{z}, \mathbf{v}_*) > 0, \quad (5.18)$$

under the CFL condition $\sigma_1 \alpha_{1,n} + \sigma_2 \alpha_{2,n} < 1$. This, along with $\bar{\mathbf{u}}_{ij}^{n+1} \cdot \mathbf{e}_1 > 0$, implies $\bar{\mathbf{u}}_{ij}^{n+1} \in G_* = G$ and the bound-preserving property of the scheme (5.12) with (5.13).

6. Application of GQL to design bound-preserving schemes for multicomponent MHD. This section applies the GQL approach to develop bound-preserving high-order finite volume and discontinuous Galerkin schemes for the multicomponent MHD system. We mainly focus on the 2D case, while our discussions are extensible to the 3D case. The 2D multicomponent compressible MHD system for a ideal fluid mixture with N_c components can be written as

$$\partial_t \mathbf{u} + \partial_x \mathbf{f}_1(\mathbf{u}) + \partial_y \mathbf{f}_2(\mathbf{u}) = \mathbf{0}, \quad (6.1a)$$

$$\mathbf{u} = \begin{pmatrix} \rho \mathbf{Y} \\ \rho \\ \mathbf{m} \\ \mathbf{B} \\ E \end{pmatrix}, \quad \mathbf{f}_\ell(\mathbf{u}) = \begin{pmatrix} \rho \mathbf{Y} v_\ell \\ \rho v_\ell \\ \mathbf{m} v_\ell - \mathbf{B} B_\ell + p_{tot} \mathbf{e}_\ell \\ \mathbf{B} v_\ell - \mathbf{v} B_\ell \\ v_\ell (E + p_{tot}) - B_\ell (\mathbf{v} \cdot \mathbf{B}) \end{pmatrix}, \quad \ell = 1, 2, \quad (6.1b)$$

along with the extra divergence-free condition on the magnetic field \mathbf{B} :

$$\nabla \cdot \mathbf{B} := \partial_x B_1 + \partial_y B_2 = 0. \quad (6.2)$$

In (6.1b), ρ denotes the total density, $\mathbf{m} = \rho \mathbf{v}$ is the momentum with \mathbf{v} being the fluid velocity, $\mathbf{Y} = (Y_1, \dots, Y_{n_c-1})^\top$ denotes the mass fractions of the first $(n_c - 1)$ components, the mass fraction of the n_c th component is $Y_{n_c} := 1 - \sum_{k=1}^{n_c-1} Y_k$, and $p_{tot} = p + \frac{\|\mathbf{B}\|^2}{2}$ is the total pressure with the thermal pressure p calculated by

$$p = (\Gamma(\mathbf{u}) - 1) \left(E - \frac{\|\mathbf{m}\|^2}{2\rho} - \frac{\|\mathbf{B}\|^2}{2} \right), \quad \Gamma(\mathbf{u}) := \frac{\sum_{k=1}^{n_c} \Gamma_k C_{v_k} Y_k}{\sum_{k=1}^{n_c} C_{v_k} Y_k}, \quad (6.3)$$

where $C_{v_k} > 0$ and $\Gamma_k > 1$ respectively denote the heat capacity at constant volume and the ratio of specific heats for species k .

6.1. GQL representation of invariant region. For the system (6.1), the total density ρ and the thermal pressure p are all positive, and the mass fractions $\{Y_k\}_{k=1}^{n_c}$ are between 0 and 1. These constraints constitute the following invariant region

$$G = \{\mathbf{u} \in \mathbb{R}^{n_c+7} : 0 \leq Y_k \leq 1, 1 \leq k \leq n_c, \rho > 0, p(\mathbf{u}) > 0\} \quad (6.4)$$

with $p(\mathbf{u})$ is a highly nonlinear function defined by (6.3). Due to the strong nonlinearity and the underlying connections between the bound-preserving and divergence-free properties, the design and analysis of bound-preserving schemes for system (6.1) are highly challenging.

Following the GQL framework, the convex region G in (6.4) can be equivalently represented as

$$G_* = \{\mathbf{u} \in \mathbb{R}^{n_c+7} : \mathbf{u} \cdot \mathbf{e}_k \geq 0, 0 \leq k < n_c, \mathbf{u} \cdot \mathbf{e}_{n_c} > 0, \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) > 0 \forall \mathbf{v}_*, \mathbf{B}_* \in \mathbb{R}^3\}, \quad (6.5)$$

where $\mathbf{e}_0 := \mathbf{e}_{n_c} - \sum_{k=1}^{n_c-1} \mathbf{e}_k$, the vector \mathbf{e}_k for $k \geq 1$ has a 1 in the k th component and zeros elsewhere, and $\varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) := \mathbf{u} \cdot \mathbf{n}_* + \frac{\|\mathbf{B}_*\|^2}{2}$ with $\mathbf{n}_* = (\mathbf{0}_{n_c-1}, \frac{\|\mathbf{v}_*\|^2}{2}, -\mathbf{v}_*, -\mathbf{B}_*, 1)^\top$. In the following, we will derive bound-preserving schemes for (6.1) based on the GQL representation (6.5). *The GQL approach will not only help overcome the difficulties arising from the nonlinearity, but also play a crucial role in establishing the key relations between the bound-preserving property and a discrete divergence-free (DDF) condition on the numerical magnetic field.*

6.2. GQL bridges bound-preserving property and DDF condition. We focus on the Euler forward method for time discretization, while all our discussions are directly extensible to high-order strong-stability-preserving time discretizations [18] which are formally convex combinations of Euler forward. Consider the finite volume methods and the scheme of the cell averages of the discontinuous Galerkin method, which can be written into a unified form as

$$\bar{\mathbf{u}}_{ij}^{n+1} = \bar{\mathbf{u}}_{ij}^n - \sigma_1 (\hat{\mathbf{f}}_{1,i+\frac{1}{2},j} - \hat{\mathbf{f}}_{1,i-\frac{1}{2},j}) - \sigma_2 (\hat{\mathbf{f}}_{2,i,j+\frac{1}{2}} - \hat{\mathbf{f}}_{2,i,j-\frac{1}{2}}), \quad (6.6)$$

for solving (6.1) on a uniform Cartesian mesh $\{\mathcal{I}_{ij} := [x_{i-1/2}, x_{i+1/2}] \times [y_{j-1/2}, y_{j+1/2}]\}$, with $\sigma_1 = \frac{\Delta t}{\Delta x}$ and $\sigma_2 = \frac{\Delta t}{\Delta y}$. Here $\bar{\mathbf{u}}_{ij}^n$ denotes the approximate cell average of $\mathbf{u}(x, y, t_n)$ on \mathcal{I}_{ij} . For a $(K+1)$ th-order accurate scheme, in each cell \mathcal{I}_{ij} a polynomial vector of degree K , denoted by $\mathbf{U}_{ij}^n(x, y)$, is also constructed as the approximate solution, which is either the reconstructed polynomial solution in a finite volume scheme or the discontinuous Galerkin polynomial solution. Denote $\{\omega_q, x_i^{(q)}\}_{q=1}^Q$ and $\{\omega_q, y_j^{(q)}\}_{q=1}^Q$ as the Gauss quadrature weights and nodes in $[x_{i-1/2}, x_{i+1/2}]$ and $[y_{j-1/2}, y_{j+1/2}]$, respectively. Let $\mathbf{u}_{i\mp\frac{1}{2},j}^{\pm,q} = \mathbf{U}_{ij}^n(x_{i\mp\frac{1}{2}}, y_j^{(q)})$, $\mathbf{u}_{i,j\mp\frac{1}{2}}^{q,\pm} = \mathbf{U}_{ij}^n(x_i^{(q)}, y_{j\mp\frac{1}{2}})$. The numerical fluxes in (6.6) are then given by

$$\hat{\mathbf{f}}_{1,i+\frac{1}{2},j} = \sum_{q=1}^Q \omega_q \hat{\mathbf{f}}_1^{\text{LF}}(\mathbf{u}_{i+\frac{1}{2},j}^{-,q}, \mathbf{u}_{i+\frac{1}{2},j}^{+,q}), \quad \hat{\mathbf{f}}_{2,i,j+\frac{1}{2}} = \sum_{q=1}^Q \omega_q \hat{\mathbf{f}}_2^{\text{LF}}(\mathbf{u}_{i,j+\frac{1}{2}}^{q,-}, \mathbf{u}_{i,j+\frac{1}{2}}^{q,+}), \quad (6.7)$$

where $\hat{\mathbf{f}}_\ell^{\text{LF}}(\cdot, \cdot)$ is taken as the Lax-Friedrichs flux (5.14) with the numerical viscosity parameters

$$\alpha_{1,n} \geq \max_{i,j,\mu} \hat{\alpha}_1(\mathbf{u}_{i+\frac{1}{2},j}^{\mp,q}, \mathbf{u}_{i-\frac{1}{2},j}^{\pm,q}), \quad \alpha_{2,n} \geq \max_{i,j,q} \hat{\alpha}_2(\mathbf{u}_{i,j+\frac{1}{2}}^{q,\mp}, \mathbf{u}_{i,j-\frac{1}{2}}^{q,\pm}). \quad (6.8)$$

Here $\hat{\alpha}_\ell(\mathbf{u}, \tilde{\mathbf{u}}) = \max\{|v_\ell| + \mathcal{C}_\ell, |\tilde{v}_\ell| + \tilde{\mathcal{C}}_\ell, \frac{|\sqrt{\rho}v_\ell + \sqrt{\tilde{\rho}}\tilde{v}_\ell|}{\sqrt{\rho} + \sqrt{\tilde{\rho}}}\} + \max\{\mathcal{C}_\ell, \tilde{\mathcal{C}}_\ell\} + \frac{\|\mathbf{B} - \tilde{\mathbf{B}}\|}{\sqrt{\rho} + \sqrt{\tilde{\rho}}}$, $\ell = 1, 2$, and \mathcal{C}_1 and \mathcal{C}_2 are the fast magneto-acoustic speeds in the x - and y -directions, respectively.

Seeking a condition for the scheme (6.6) to be bound-preserving is very challenging, due to the complexity of the system (6.1) and the region (6.4) as well as the intrinsic relations between the bound-preserving property and the DDF condition, On one hand, it is very difficult to establish such relations, since the bound-preserving property is an *algebraic* property while the DDF condition is a discrete *differential* property. In fact, their relations remained unclear for a long time, until the recent work [51] on the single-component MHD case. On the other hand, the DDF condition strongly couples the states $\{\mathbf{u}_{i\mp\frac{1}{2},j}^{\pm,q}, \mathbf{u}_{i,j\mp\frac{1}{2}}^{q,\pm}\}$, making the traditional or standard analysis approaches (which typically rely on decomposing high-order or/and multidimensional schemes into convex combinations of first-order 1D schemes [67, 68, 71]) *inapplicable* to the present case.

First, let us consider the first-order scheme to gain some insights. In this case, the polynomial degree $K = 0$ so that $\mathbf{U}_{ij}^n(x, y) \equiv \bar{\mathbf{u}}_{ij}^n$ for all $(x, y) \in \mathcal{I}_{ij}$, and we can reformulate the scheme (6.6) as

$$\bar{\mathbf{u}}_{ij}^{n+1} = (1 - \lambda)\bar{\mathbf{u}}_{ij}^n + \sigma_1\alpha_{1,n}\mathbf{\Pi}_1 + \sigma_2\alpha_{2,n}\mathbf{\Pi}_2, \quad (6.9)$$

with $\lambda := \sigma_1\alpha_{1,n} + \sigma_2\alpha_{2,n}$, and

$$\mathbf{\Pi}_1 = \frac{1}{2} \left(\bar{\mathbf{u}}_{i+1,j}^n - \frac{\mathbf{f}_1(\bar{\mathbf{u}}_{i+1,j}^n)}{\alpha_{1,n}} + \bar{\mathbf{u}}_{i-1,j}^n + \frac{\mathbf{f}_1(\bar{\mathbf{u}}_{i-1,j}^n)}{\alpha_{1,n}} \right), \mathbf{\Pi}_2 = \frac{1}{2} \left(\bar{\mathbf{u}}_{i,j+1}^n - \frac{\mathbf{f}_2(\bar{\mathbf{u}}_{i,j+1}^n)}{\alpha_{2,n}} + \bar{\mathbf{u}}_{i,j-1}^n + \frac{\mathbf{f}_2(\bar{\mathbf{u}}_{i,j-1}^n)}{\alpha_{2,n}} \right).$$

THEOREM 6.1. *If $\bar{\mathbf{u}}_{ij}^n \in G$ for all i and j , then, under the CFL condition $\lambda \leq 1$, the solution $\bar{\mathbf{u}}_{ij}^{n+1}$ of (6.9) satisfies*

$$\bar{\mathbf{u}}_{ij}^{n+1} \cdot \mathbf{e}_k \geq 0, \quad 0 \leq k < n_c, \quad \bar{\mathbf{u}}_{ij}^{n+1} \cdot \mathbf{e}_{n_c} > 0, \quad (6.10)$$

$$\varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{v}_*, \mathbf{B}_*) > -\Delta t(\mathbf{v}_* \cdot \mathbf{B}_*) \operatorname{div}_{ij} \bar{\mathbf{B}} \quad \forall \mathbf{v}_*, \mathbf{B}_* \in \mathbb{R}^3, \quad (6.11)$$

where $\operatorname{div}_{ij} \bar{\mathbf{B}} := \frac{\bar{B}_{1,i+1,j}^n - \bar{B}_{1,i-1,j}^n}{2\Delta x} + \frac{\bar{B}_{2,i,j+1}^n - \bar{B}_{2,i,j-1}^n}{2\Delta y}$ is a discrete divergence. Furthermore, if the states $\{\bar{\mathbf{u}}_{ij}^n\}$ satisfy the DDF condition $\operatorname{div}_{ij} \bar{\mathbf{B}} = 0$, then (6.10)–(6.11) imply $\bar{\mathbf{u}}_{ij}^{n+1} \in G_* = G$.

Proof. For $0 \leq k < n_c$ and any $\mathbf{u} \in \{\bar{\mathbf{u}}_{ij}^n\}$, we have $\pm \mathbf{f}_\ell(\mathbf{u}) \cdot \mathbf{e}_k = \pm v_\ell(\mathbf{u} \cdot \mathbf{e}_k) \leq \alpha_{\ell,n}(\mathbf{u} \cdot \mathbf{e}_k)$, which implies $\mathbf{\Pi}_\ell \cdot \mathbf{e}_k \geq 0$. Similarly, $\mathbf{\Pi}_\ell \cdot \mathbf{e}_{n_c} > 0$. These lead to (6.10). Following [51, Lemma 2.6], we can derive that

$$\varphi(\mathbf{\Pi}_1; \mathbf{v}_*, \mathbf{B}_*) > \frac{\mathbf{v}_* \cdot \mathbf{B}_*}{2\alpha_{1,n}} (\bar{B}_{1,i-1,j}^n - \bar{B}_{1,i+1,j}^n), \quad \varphi(\mathbf{\Pi}_2; \mathbf{v}_*, \mathbf{B}_*) > \frac{\mathbf{v}_* \cdot \mathbf{B}_*}{2\alpha_{2,n}} (\bar{B}_{2,i,j-1}^n - \bar{B}_{2,i,j+1}^n).$$

Thanks to the linearity of $\varphi(\cdot; \mathbf{v}_*, \mathbf{B}_*)$, it then follows from (6.9) that

$$\begin{aligned} \varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{v}_*, \mathbf{B}_*) &= (1 - \lambda)\varphi(\bar{\mathbf{u}}_{ij}^n; \mathbf{v}_*, \mathbf{B}_*) + \sigma_1\alpha_{1,n}\varphi(\mathbf{\Pi}_1; \mathbf{v}_*, \mathbf{B}_*) + \sigma_2\alpha_{2,n}\varphi(\mathbf{\Pi}_2; \mathbf{v}_*, \mathbf{B}_*) \\ &> (1 - \lambda)\varphi(\bar{\mathbf{u}}_{ij}^n; \mathbf{v}_*, \mathbf{B}_*) - \Delta t(\mathbf{v}_* \cdot \mathbf{B}_*) \operatorname{div}_{ij} \bar{\mathbf{B}}, \end{aligned}$$

which yields (6.11) under the CFL condition $\lambda \leq 1$. \square

Theorem 6.1 shows the connection between the bound-preserving property and a DDF condition, which is bridged by (6.11) with the help of the free auxiliary variables $\{\mathbf{v}_*, \mathbf{B}_*\}$ in the GQL representation (6.5). This demonstrates the essential importance of the GQL approach in establishing this connection and its significant advantages for bound-preserving analysis and design.

Now, we use the GQL approach to explore bound-preserving high-order schemes with $K \geq 1$. Denote $\{\hat{x}_i^{(\beta)}\}_{\beta=1}^L$ and $\{\hat{y}_j^{(\beta)}\}_{\beta=1}^L$ as the Gauss–Lobatto quadrature points in $[x_{i-1/2}, x_{i+1/2}]$ and $[y_{j-1/2}, y_{j+1/2}]$, respectively, and $\{\hat{\omega}_\beta\}_{\beta=1}^L$ as the weights, with $L = \lceil \frac{K+3}{2} \rceil$. Similar to **Theorem 6.1** and [51, Theorem 4.7], the following result can be derived with the proof omitted here.

THEOREM 6.2. *If, for all i and j , $\bar{\mathbf{u}}_{ij}^n \in G$ and the polynomial vector $\mathbf{U}_{ij}^n(x, y)$ satisfies*

$$\mathbf{U}_{ij}^n(\hat{x}_i^{(\beta)}, y_j^{(q)}), \mathbf{U}_{ij}^n(x_i^{(q)}, \hat{y}_j^{(\beta)}) \in G \quad \forall \beta, q, \quad (6.12)$$

then, the solution $\bar{\mathbf{u}}_{ij}^{n+1}$ of the scheme (6.6) satisfies

$$\varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{v}_*, \mathbf{B}_*) > 2(\hat{\omega}_1 - \lambda)\varphi(\mathbf{\Pi}; \mathbf{v}_*, \mathbf{B}_*) - \Delta t(\mathbf{v}_* \cdot \mathbf{B}_*) \operatorname{div}_{ij} \mathbf{B} \quad (6.13)$$

with $\mathbf{\Pi} := \frac{1}{2\lambda} \sum_q \omega_q [\sigma_1\alpha_{1,n}(\mathbf{u}_{i+\frac{1}{2},j}^{-,q} + \mathbf{u}_{i-\frac{1}{2},j}^{+,q}) + \sigma_2\alpha_{2,n}(\mathbf{u}_{i,j+\frac{1}{2}}^{q,-} + \mathbf{u}_{i,j-\frac{1}{2}}^{q,+})] \in G$. Furthermore, under the CFL condition $\lambda \leq \hat{\omega}_1$, we have

$$\bar{\mathbf{u}}_{ij}^{n+1} \cdot \mathbf{e}_k \geq 0, \quad 0 \leq k < n_c, \quad \bar{\mathbf{u}}_{ij}^{n+1} \cdot \mathbf{e}_{n_c} > 0, \quad (6.14)$$

$$\varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{v}_*, \mathbf{B}_*) > -\Delta t(\mathbf{v}_* \cdot \mathbf{B}_*) \operatorname{div}_{ij} \mathbf{B} \quad \forall \mathbf{v}_*, \mathbf{B}_* \in \mathbb{R}^3, \quad (6.15)$$

where the discrete divergence is defined as $\operatorname{div}_{ij} \mathbf{B} := \frac{1}{2} (\operatorname{div}_{ij}^- \mathbf{B} + \operatorname{div}_{ij}^+ \mathbf{B})$ with

$$\operatorname{div}_{ij}^\mp \mathbf{B} := \frac{1}{\Delta x} \sum_{q=1}^Q \omega_q (B_{1,i+\frac{1}{2},j}^{\mp,q} - B_{1,i-\frac{1}{2},j}^{\pm,q}) + \frac{1}{\Delta y} \sum_{q=1}^Q \omega_q (B_{2,i,j+\frac{1}{2}}^{q,\mp} - B_{2,i,j-\frac{1}{2}}^{q,\pm}).$$

Remark 6.3. The condition (6.12) in Theorem 6.2 is a standard condition for bound-preserving finite volume and discontinuous Galerkin schemes (see [67, 68]). This condition can be easily enforced by a simple scaling limiter (see Appendix B). Thanks to the GQL representation (6.5), we conclude from (6.14)–(6.15) that, in order to ensure $\bar{\mathbf{u}}_{ij}^{n+1} \in G_* = G$, a DDF condition

$$\operatorname{div}_{ij} \mathbf{B} := \frac{1}{2} (\operatorname{div}_{ij}^- \mathbf{B} + \operatorname{div}_{ij}^+ \mathbf{B}) = 0 \quad (6.16)$$

is also required. Unfortunately, the high-order schemes (6.6) do not preserve the DDF condition (6.16), which depends on the numerical solution information from adjacent cells. Although a few globally divergence-free techniques (e.g. [30, 16, 6]) were developed and can enforce the condition (6.16), the local scaling limiter for (6.12) will destroy the globally divergence-free property. Notice that the locally divergence-free technique (e.g. [29]) is compatible with the local scaling limiter, but can only guarantee $\operatorname{div}_{ij}^- \mathbf{B} = 0$. In subsection 6.3, we will use the GQL approach to explore how to eliminate the effect of the remaining part $\operatorname{div}_{ij}^+ \mathbf{B}$ by properly modifying the scheme (6.6).

6.3. Seek high-order provably bound-preserving schemes via GQL. We have established the relations between the bound-preserving and divergence-free properties at the numerical level. Interestingly, at the continuous level, bound preservation is also closely related to the divergence-free condition (6.2): If condition (6.2) is slightly violated, then even the exact solution of system (6.1) may not stay in G ; see [53] for a discussion which is also valid for system (6.1). To address this issue, we consider a modified formulation of the multicomponent MHD equations

$$\partial_t \mathbf{u} + \partial_x \mathbf{f}_1(\mathbf{u}) + \partial_y \mathbf{f}_2(\mathbf{u}) + (\nabla \cdot \mathbf{B}) \mathbf{S}(\mathbf{u}) = \mathbf{0} \quad (6.17)$$

by adding an extra source term to (6.1a) with $\mathbf{S}(\mathbf{u}) = (\mathbf{0}_{nc}, \mathbf{B}, \mathbf{v}, \mathbf{v} \cdot \mathbf{B})^\top$. Such a formulation was first proposed by Godunov [17] for the purpose of entropy symmetrization in the single-component MHD case. Notice that, for divergence-free initial conditions, the exact solutions of the modified form (6.17) and the standard form (6.1) are the same. However, if the divergence-free condition (6.2) is violated, the extra source term in the modified form (6.17) becomes beneficial and helps keep the exact solutions always in G ; see [54] for an analysis which also works for system (6.17). This finding motivates us to explore bound-preserving schemes based on suitable discretization of the modified form (6.17). Thus we consider

$$\bar{\mathbf{u}}_{ij}^{n+1} = \bar{\mathbf{u}}_{ij}^n - \sigma_1 (\hat{\mathbf{f}}_{1,i+\frac{1}{2},j} - \hat{\mathbf{f}}_{1,i-\frac{1}{2},j}) - \sigma_2 (\hat{\mathbf{f}}_{2,i,j+\frac{1}{2}} - \hat{\mathbf{f}}_{2,i,j-\frac{1}{2}}) - \hat{\mathbf{S}}_{ij} \quad (6.18)$$

by adding a properly discretized source term $\hat{\mathbf{S}}_{ij}$ into the standard finite volume or discontinuous Galerkin schemes (6.6). As discussed in Remark 6.3, we can adopt a locally divergence-free technique for the magnetic components of $\mathbf{U}_{ij}^n(x, y)$ such that $\operatorname{div}_{ij}^- \mathbf{B} = 0$. This gives $2\operatorname{div}_{ij} \mathbf{B} = \operatorname{div}_{ij}^+ \mathbf{B} = \operatorname{div}_{ij}^+ \mathbf{B} - \operatorname{div}_{ij}^- \mathbf{B}$, thereby leading to

$$\operatorname{div}_{ij} \mathbf{B} = \frac{1}{2\Delta x} \sum_{q=1}^Q \omega_q \left(\llbracket B_1 \rrbracket_{i+\frac{1}{2},j}^q + \llbracket B_1 \rrbracket_{i-\frac{1}{2},j}^q \right) + \frac{1}{2\Delta y} \sum_{q=1}^Q \omega_q \left(\llbracket B_2 \rrbracket_{i,j+\frac{1}{2}}^q + \llbracket B_2 \rrbracket_{i,j-\frac{1}{2}}^q \right), \quad (6.19)$$

where $\llbracket B_1 \rrbracket_{i+\frac{1}{2},j}^q = B_{1,i+\frac{1}{2},j}^{+,q} - B_{1,i+\frac{1}{2},j}^{-,q}$ and $\llbracket B_2 \rrbracket_{i,j+\frac{1}{2}}^q = B_{2,i,j+\frac{1}{2}}^{q,-} - B_{2,i,j+\frac{1}{2}}^{q,-}$ are the jumps of the normal magnetic component across the cell interface. Using the GQL approach with the linearity of $\varphi(\cdot; \mathbf{v}_*, \mathbf{B}_*)$ and the estimate (6.13) under the hypothesis of Theorem 6.2, we obtain

$$\varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{v}_*, \mathbf{B}_*) > 2(\hat{\omega}_1 - \lambda) \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*) - \left[\Delta t (\mathbf{v}_* \cdot \mathbf{B}_*) \operatorname{div}_{ij} \mathbf{B} + \hat{\mathbf{S}}_{ij} \cdot \mathbf{n}_* \right]. \quad (6.20)$$

Then the key is to carefully design $\hat{\mathbf{S}}_{ij}$ to exactly offset the effect of $\operatorname{div}_{ij} \mathbf{B}$ in (6.20), so that the resulting schemes (6.18) become bound-preserving. Observing that for any $b \in \mathbb{R}$ and any $\mathbf{u} \in G$,

$$b(\mathbf{v}_* \cdot \mathbf{B}_* + \mathbf{S}(\mathbf{u}) \cdot \mathbf{n}_*) = b(\mathbf{v} - \mathbf{v}_*) \cdot (\mathbf{B} - \mathbf{B}_*) \leq |b| \rho^{1/2} \varphi(\mathbf{u}; \mathbf{v}_*, \mathbf{B}_*), \quad (6.21)$$

we devise

$$\begin{aligned} \hat{\mathbf{S}}_{ij} &= \frac{\sigma_1}{2} \sum_{q=1}^Q \omega_q \left[\llbracket B_1 \rrbracket_{i+\frac{1}{2},j}^q \mathbf{S}(\mathbf{u}_{i+\frac{1}{2},j}^{-,q}) + \llbracket B_1 \rrbracket_{i-\frac{1}{2},j}^q \mathbf{S}(\mathbf{u}_{i-\frac{1}{2},j}^{+,q}) \right] \\ &\quad + \frac{\sigma_2}{2} \sum_{q=1}^Q \omega_q \left[\llbracket B_2 \rrbracket_{i,j+\frac{1}{2}}^q \mathbf{S}(\mathbf{u}_{i,j+\frac{1}{2}}^{q,-}) + \llbracket B_2 \rrbracket_{i,j-\frac{1}{2}}^q \mathbf{S}(\mathbf{u}_{i,j-\frac{1}{2}}^{q,+}) \right], \end{aligned} \quad (6.22)$$

such that the last term in (6.20) satisfies

$$\begin{aligned}
& \Delta t(\mathbf{v}_* \cdot \mathbf{B}_*) \operatorname{div}_{ij} \mathbf{B} + \widehat{\mathbf{S}}_{ij} \cdot \mathbf{n}_* \\
&= \frac{\sigma_1}{2} \sum_{q=1}^Q \omega_q \left[\llbracket B_1 \rrbracket_{i+\frac{1}{2},j}^q \left(\mathbf{v}_* \cdot \mathbf{B}_* + \mathbf{S}(\mathbf{u}_{i+\frac{1}{2},j}^{-,q}) \cdot \mathbf{n}_* \right) + \llbracket B_1 \rrbracket_{i-\frac{1}{2},j}^q \left(\mathbf{v}_* \cdot \mathbf{B}_* + \mathbf{S}(\mathbf{u}_{i-\frac{1}{2},j}^{+,q}) \cdot \mathbf{n}_* \right) \right] \\
&+ \frac{\sigma_2}{2} \sum_{q=1}^Q \omega_q \left[\llbracket B_2 \rrbracket_{i,j+\frac{1}{2}}^q \left(\mathbf{v}_* \cdot \mathbf{B}_* + \mathbf{S}(\mathbf{u}_{i,j+\frac{1}{2}}^{q,-}) \cdot \mathbf{n}_* \right) + \llbracket B_2 \rrbracket_{i,j-\frac{1}{2}}^q \left(\mathbf{v}_* \cdot \mathbf{B}_* + \mathbf{S}(\mathbf{u}_{i,j-\frac{1}{2}}^{q,+}) \cdot \mathbf{n}_* \right) \right] \\
&\leq \varepsilon \lambda \varphi(\mathbf{\Pi}; \mathbf{v}_*, \mathbf{B}_*), \tag{6.23}
\end{aligned}$$

where we use (6.19) in the equality and (6.21) in the inequality, and $\varepsilon = \max\{\beta_1/\alpha_{1,n}, \beta_2/\alpha_{2,n}\}$ with $\beta_1 = \max_{i,j,q} \{ \llbracket B_1 \rrbracket_{i+\frac{1}{2},j}^q |(\rho_{i+\frac{1}{2},j}^{\pm,q})^{1/2}| \}$ and $\beta_2 = \max_{i,j,q} \{ \llbracket B_2 \rrbracket_{i,j+\frac{1}{2}}^q |(\rho_{i,j+\frac{1}{2}}^{q,\pm})^{1/2}| \}$. Combining (6.20) with (6.23), we obtain

$$\varphi(\bar{\mathbf{u}}_{ij}^{n+1}; \mathbf{v}_*, \mathbf{B}_*) > 2(\widehat{\omega}_1 - \lambda - \varepsilon \lambda) \varphi(\mathbf{\Pi}; \mathbf{v}_*, \mathbf{B}_*) \geq 0,$$

under the CFL condition $(1 + \varepsilon)\lambda \leq \widehat{\omega}_1$. Notice that the first n_c components of $\widehat{\mathbf{S}}_{ij}$ are zeros, which implies (6.14) in Theorem 6.2 also holds for the modified schemes (6.18). In summary, we obtain:

THEOREM 6.4. *If for all i and j , $\bar{\mathbf{u}}_{ij}^n \in G$ and the polynomial vector $\mathbf{U}_{ij}^n(x, y)$ satisfies (6.12) and $\operatorname{div}_{ij}^- \mathbf{B} = 0$, then, under the CFL condition $(1 + \varepsilon)\lambda \leq \widehat{\omega}_1$, the solution $\bar{\mathbf{u}}_{ij}^{n+1}$ of (6.18) is always preserved in G_* .*

Theorem 6.4 indicates that, if we use the scaling limiter in Appendix B to enforce (6.12) and a locally divergence-free technique to ensure $\operatorname{div}_{ij}^- \mathbf{B} = 0$, then the schemes (6.18) with (6.22) are bound-preserving. The bounds are also preserved if a high-order strong-stability-preserving time discretization [18] is used to replace the Euler forward method.

7. Experimental results. This section gives two highly demanding numerical examples to further demonstrate our theoretical analysis as well as the robustness and effectiveness of the bound-preserving schemes designed via GQL in subsection 6.3 for the 2D multicomponent MHD. We use the proposed bound-preserving third-order locally divergence-free discontinuous Galerkin method for spatial discretization. As the tests involve strong discontinuities, the locally divergence-free WENO limiter [73] is also employed in some trouble cells adaptively detected by the indicator of [27]. The third-order strong-stability-preserving Runge-Kutta method [18] is adopted for time discretization, with the CFL number set as 0.15.

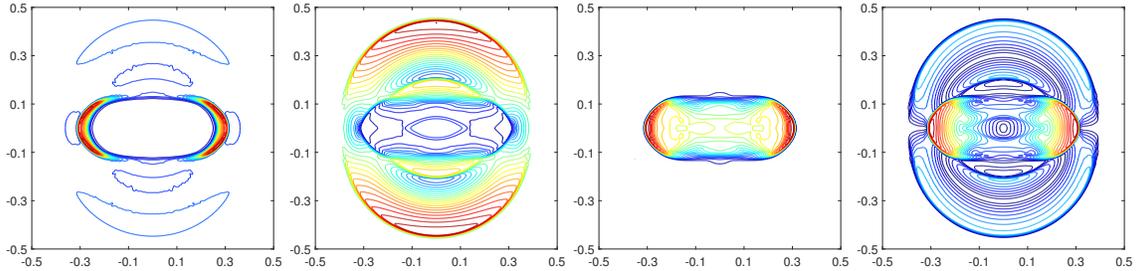


Fig. 5: The contour plots of ρ , p_m , p , and $\|\mathbf{v}\|$ (from left to right) for the blast problem at $t = 0.01$.

Example 7.1 (Blast problem). This test simulates a benchmark MHD problem in the domain $[-0.5, 0.5]^2$ with outflow boundary conditions. The setup is similar to that in [1] except for a fluid mixture with $n_c = 2$, $C_{v_1} = 2.42$, $C_{v_2} = 0.72$, $\Gamma_1 = 5/3$, and $\Gamma_2 = 1.4$. Initially, the fluid is stationary, with $(\rho, p, Y_1, Y_2) = (1, 1000, 1, 0)$ in the explosion region ($x^2 + y^2 \leq 0.01$) and $(1, 0.1, 0, 1)$ in the ambient region ($x^2 + y^2 > 0.01$). The magnetic field \mathbf{B} is initialized as $(100/\sqrt{4\pi}, 0, 0)$. Due to the large jump in p and the strong magnetic field, negative numerical p can be easily produced and often cause failure of the numerical simulations. Figure 5 presents the contour plots of the density ρ , the magnetic pressure $p_m = \frac{1}{2}\|\mathbf{B}\|^2$, the thermal pressure p , and the velocity magnitude $\|\mathbf{v}\|$ computed by the proposed bound-preserving discontinuous Galerkin method with 400×400

uniform cells. We observe the flow structures are well captured, and our method is highly robust and always preserves the bound principles (6.4) in the whole simulation.

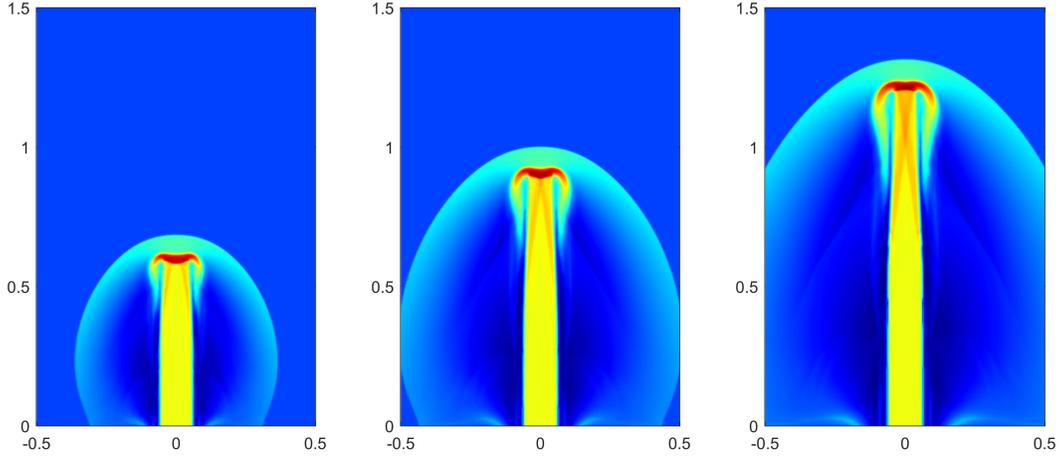


Fig. 6: The plots of $\log(\rho)$ for the jet problem. From left to right: $t = 0.001, 0.0015$ and 0.002 .

Example 7.2 (Astrophysical jet). This test simulates a high-speed MHD jet flow in the domain $[-0.5, 0.5] \times [0, 1.5]$ with $n_c = 2$, $C_{v_1} = 0.72$, $C_{v_2} = 2.42$, $\Gamma_1 = 1.4$, and $\Gamma_2 = 5/3$. The domain is initially filled with static fluid with $(\rho, p, Y_1, Y_2) = (0.14, 1, 0, 1)$. The inflow jet condition is fixed on boundary $\{|x| < 0.05, y = 0\}$ with $(\rho, p, Y_1, Y_2) = (1.4, 1, 1, 0)$ and $\mathbf{v} = (0, 800, 0)$, while the outflow conditions are specified on the other boundaries. There is a strong magnetic field \mathbf{B} initialized as $(0, \sqrt{4000}, 0)$, which makes this test more challenging. Our simulation is based on the proposed bound-preserving method with 200×600 uniform cells in $[0, 0.5] \times [0, 1.5]$. The numerical results are shown Figure 6. The flow pattern is captured with high resolution and similar to the single-component MHD case reported in [53, 54]. In such an extreme test, our bound-preserving method exhibits good robustness. However, if the proposed scaling limiter is not used to enforce (6.12), or if the locally divergence-free technique is not employed to ensure $\text{div}_{ij}^- \mathbf{B} = 0$, or if the proposed source term (6.22) is dropped, the resulting method even with the WENO limiter is not bound-preserving and would fail quickly due to nonphysical numerical solutions out of the bounds. This confirms our theoretical analyses and the importance of the proposed conditions and techniques.

8. Conclusions. We have systematically proposed a novel and general framework, called geometric quasilinearization (GQL), for studying bound-preserving problems with nonlinear constraints. GQL skillfully transfers all nonlinear constraints into linear ones, via properly introducing some free auxiliary variables independent of the system variables. We have established the fundamental principle and general theory of GQL, and provided three simple methods for constructing GQL representations. The GQL approach equivalently casts the nonlinear bound-preserving problems into preserving the positivity of linear functions, thereby opening up a new effective way for bound-preserving study. Several examples have been provided to demonstrate the effectiveness and advantages of the GQL approach in addressing nonlinear bound-preserving problems that are highly challenging and could not be easily handled by direct or traditional approaches. Besides the examples in this paper, recently the GQL approach also achieved successes in finding (high-order) bound-preserving schemes for several complicated PDE systems in [51, 53, 57, 52, 54, 55].

As the proposed GQL framework is not restricted to the specific forms of the PDEs, it applies to general time-dependent PDE systems that possess convex invariant regions with nonlinear constraints. Moreover, it can be used in conjunction with the well-developed limiters in [67, 68, 63, 23] to design high-order bound-preserving schemes. It can be expected the GQL approach will be useful for addressing more challenging bound-preserving problems for a variety of PDEs in the future.

Appendix A. Proof of Theorem 3.9. The proof is divided into two steps.

(i) Prove that $G \subseteq G_*$. For any $\mathbf{u}_* \in \partial G$, the hyperplane $(\mathbf{u} - \mathbf{u}_*) \cdot \mathbf{n}_* = 0$ supports the convex region G at \mathbf{u}_* . Thus we have

$$G \subseteq \{\mathbf{u} \in \mathbb{R}^N : (\mathbf{u} - \mathbf{u}_*, \mathbf{n}_*) \geq 0 \quad \forall \mathbf{u}_* \in \partial G\}. \quad (\text{A.1})$$

If G is closed, then (A.1) means $G \subseteq G_*$. Next, we assume G is open and show $G \subseteq G_*$ by contradiction. Assume that

$$\text{there exists } \mathbf{u}_0 \in G \text{ but } \mathbf{u}_0 \notin G_*. \quad (\text{A.2})$$

Then, according to (A.1), there exists $\mathbf{u}_* \in \partial G$ such that $(\mathbf{u}_0 - \mathbf{u}_*) \cdot \mathbf{n}_* = 0$. Since G is open, there exists $\delta > 0$ such that $\Omega_\delta := \{\mathbf{u} \in \mathbb{R}^N : \|\mathbf{u} - \mathbf{u}_0\| < \delta\} \subset G$. We take $\mathbf{u}_\delta := \mathbf{u}_0 - \frac{\delta}{2\|\mathbf{n}_*\|} \mathbf{n}_* \in \Omega_\delta$. Then $\mathbf{u}_\delta \in G$. However, using $(\mathbf{u}_0 - \mathbf{u}_*) \cdot \mathbf{n}_* = 0$ gives

$$(\mathbf{u}_\delta - \mathbf{u}_*) \cdot \mathbf{n}_* = \left(\mathbf{u}_0 - \mathbf{u}_* - \frac{\delta}{2\|\mathbf{n}_*\|} \mathbf{n}_* \right) \cdot \mathbf{n}_* = -\frac{\delta}{2} \|\mathbf{n}_*\| < 0,$$

which contradicts (A.1) and $\mathbf{u}_\delta \in G$. Thus the assumption (A.2) is incorrect, and we have $G \subseteq G_*$.

(ii) Prove that $G_* \subseteq G$. We first show that $G_* \subseteq \text{cl}(G)$ by contradiction. Assume that

$$\text{there exists } \mathbf{u}_0 \in G_* \text{ but } \mathbf{u}_0 \notin \text{cl}(G). \quad (\text{A.3})$$

According to the theory of convex optimization [5], the minimum of the convex function $\zeta(\mathbf{u}) := \|\mathbf{u} - \mathbf{u}_0\|^2$ over the closed convex region $\text{cl}(G)$ is attained at certain boundary point $\mathbf{u}_* \in \partial G$. Let $\hat{\mathbf{u}}$ be an arbitrary interior point of G . Thanks to the convexity $\text{cl}(G)$, one has $\mathbf{u}_\lambda := \lambda \hat{\mathbf{u}} + (1 - \lambda) \mathbf{u}_* \in \text{cl}(G)$ for any $\lambda \in [0, 1]$. We then know that the quadratic function

$$\hat{\zeta}(\lambda) := \zeta(\mathbf{u}_\lambda) = \lambda^2 \|\hat{\mathbf{u}} - \mathbf{u}_*\|^2 + 2\lambda (\hat{\mathbf{u}} - \mathbf{u}_*) \cdot (\mathbf{u}_* - \mathbf{u}_0) + \|\mathbf{u}_* - \mathbf{u}_0\|^2$$

attains its minimum over $[0, 1]$ at $\lambda = 0$. This implies $(\hat{\mathbf{u}} - \mathbf{u}_*) \cdot (\mathbf{u}_* - \mathbf{u}_0) \geq 0$, for an arbitrary interior point $\hat{\mathbf{u}}$ of G . Thus $\text{int}(G) \subseteq \{\mathbf{u} : (\mathbf{u} - \mathbf{u}_*) \cdot (\mathbf{u}_* - \mathbf{u}_0) \geq 0\} =: H_*^+$, where H_*^+ is a closed halfspace. It follows that H_*^+ is a supporting halfspace to G , and $\mathbf{u}_* - \mathbf{u}_0$ is an inward-pointing normal vector of G at \mathbf{u}_* . Because ∂G is smooth, there exists $\mu > 0$ such that $\mathbf{n}_* = \mu(\mathbf{u}_* - \mathbf{u}_0)$, which implies $(\mathbf{u}_0 - \mathbf{u}_*) \cdot \mathbf{n}_* = -\mu \|\mathbf{u}_0 - \mathbf{u}_*\|^2 < 0$. This contradicts the assumption $\mathbf{u}_0 \in G_*$. Thus the assumption (A.3) is incorrect, and we have $G_* \subseteq \text{cl}(G)$. If G is closed, then we obtain $G_* \subseteq G$. If G is open, then $\partial G \cap G_* = \emptyset$, which along with $G_* \subseteq \text{cl}(G)$ yields $G_* \subseteq G$.

In summary, we have $G = G_*$, and the proof is completed.

Appendix B. A simple scaling limiter to enforce (6.12). The condition (6.12) is not always automatically satisfied by the polynomial vector $\mathbf{U}_{ij}^n(x, y)$ of the high-order schemes. If this happens, the following limiter is used to modify $\mathbf{U}_{ij}^n(x, y)$ into $\tilde{\mathbf{U}}_{ij}^n(x, y)$ such that $\tilde{\mathbf{U}}_{ij}^n(x, y)$ satisfies (6.12). Define $\mathbb{Q}_{ij} = \{(\hat{x}_i^{(\beta)}, \hat{y}_j^{(q)}), (x_i^{(q)}, \hat{y}_j^{(\beta)}) \mid \forall \beta, q\}$ as the set of all the points involved in (6.12). Since the limiter is performed separately for each cell, the subscripts ij and superscript n of all quantities are omitted below for convenience. First, modify the density as

$$\hat{\rho}(x, y) = \bar{\rho} + \theta_1 (\rho(x, y) - \bar{\rho}), \quad \theta_1 := (\bar{\rho} - \epsilon_1) / (\bar{\rho} - \min_{(x, y) \in \mathbb{Q}_{ij}} \rho(x, y)),$$

where ϵ_1 is a small positive number and may be taken as $\min\{10^{-13}, \bar{\rho}\}$. Define $\mathbb{S}_k = \{(x, y) \in \mathbb{Q}_{ij} : \rho Y_k(x, y) \leq 0\}$. Then, modify the mass fractions [13] as

$$\widehat{\rho Y}_k(x, y) = \rho Y_k(x, y) + \theta_2 \left(\frac{\widehat{\rho Y}_k}{\bar{\rho}} \hat{\rho}(x, y) - \rho Y_k(x, y) \right), \quad 1 \leq k \leq n_c - 1,$$

where $\theta_2 = \max_{1 \leq k \leq n_c} \max_{(x, y) \in \mathbb{S}_k} \left\{ \frac{-\rho Y_k(x, y)}{\widehat{\rho Y}_k \hat{\rho}(x, y) / \bar{\rho} - \rho Y_k(x, y)} \right\}$ with $\rho Y_{n_c} = \hat{\rho} - \sum_{k=1}^{n_c-1} \rho Y_k$. Denote $\hat{\mathbf{U}} = (\widehat{\rho \mathbf{Y}}, \hat{\rho}, \mathbf{m}, \mathbf{B}, E)^\top$. Finally, modify $\hat{\mathbf{U}}$ to enforce the positivity of $g(\mathbf{U}) = E - \frac{1}{2}(\|\mathbf{m}\|^2 / \rho + \|\mathbf{B}\|^2)$ by

$$\tilde{\mathbf{U}}(x, y) = \bar{\mathbf{U}} + \theta_3 (\hat{\mathbf{U}}(x, y) - \bar{\mathbf{U}}), \quad \theta_3 := (g(\bar{\mathbf{U}}) - \epsilon_2) / (g(\bar{\mathbf{U}}) - \min_{(x, y) \in \mathbb{Q}_{ij}} g(\hat{\mathbf{U}}(x, y))),$$

where ϵ_2 is a small positive number and may be taken as $\min\{10^{-13}, g(\bar{\mathbf{U}})\}$. Note that the pressure function $p(\mathbf{U})$ in (6.3) is generally not concave so we use the concave function $g(\mathbf{U})$ instead of $p(\mathbf{U})$. It can be verified that the limited solution $\tilde{\mathbf{U}}(x, y) \in G$ for all $(x, y) \in \mathbb{Q}_{ij}$ and its cell average equals $\bar{\mathbf{U}}$. Such type of limiters do not lose the high-order accuracy, as demonstrated in [67, 68, 66].

- [1] D. S. BALSARA AND D. SPICER, *A staggered mesh algorithm using high order Godunov fluxes to ensure solenoidal magnetic fields in magnetohydrodynamic simulations*, J. Comput. Phys., 149 (1999), pp. 270–292.
- [2] P. BATTEN, N. CLARKE, C. LAMBERT, AND D. M. CAUSON, *On the choice of wavespeeds for the HLLC Riemann solver*, SIAM J. Sci. Comput., 18 (1997), pp. 1553–1570.
- [3] C. BERTHON, P. CHARRIER, AND B. DUBROCA, *An HLLC scheme to solve the M1 model of radiative transfer in two space dimensions*, J. Sci. Comput., 31 (2007), pp. 347–389.
- [4] J. BORWEIN AND A. S. LEWIS, *Convex Analysis and Nonlinear Optimization: Theory and Examples*, Springer Science & Business Media, 2010.
- [5] S. BOYD, S. P. BOYD, AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, 2004.
- [6] P. CHANDRASHEKAR, *A global divergence conforming DG method for hyperbolic conservation laws with divergence constraint*, J. Sci. Comput., 79 (2019), pp. 79–102.
- [7] J. CHENG AND C.-W. SHU, *Positivity-preserving Lagrangian scheme for multi-material compressible flow*, J. Comput. Phys., 257 (2014), pp. 143–168.
- [8] Q. CHENG AND J. SHEN, *Global constraints preserving scalar auxiliary variable schemes for gradient flows*, SIAM J. Sci. Comput., 42 (2020), pp. A2489–A2513.
- [9] Y. CHENG, I. GAMBA, AND J. PROFT, *Positivity-preserving discontinuous Galerkin schemes for linear Vlasov-Boltzmann transport equations*, Math. Comp., 81 (2012), pp. 153–190.
- [10] Y. CHENG, F. LI, J. QIU, AND L. XU, *Positivity-preserving DG and central DG methods for ideal MHD equations*, J. Comput. Phys., 238 (2013), pp. 255–280.
- [11] A. J. CHRISTLIEB, Y. LIU, Q. TANG, AND Z. XU, *High order parametrized maximum-principle-preserving and positivity-preserving WENO schemes on unstructured meshes*, J. Comput. Phys., 281 (2015), pp. 334–351.
- [12] A. J. CHRISTLIEB, Y. LIU, Q. TANG, AND Z. XU, *Positivity-preserving finite difference weighted ENO schemes with constrained transport for ideal magnetohydrodynamic equations*, SIAM J. Sci. Comput., 37 (2015), pp. A1825–A1845.
- [13] J. DU, C. WANG, C. QIAN, AND Y. YANG, *High-order bound-preserving discontinuous Galerkin methods for stiff multispecies detonation*, SIAM J. Sci. Comput., 41 (2019), pp. B250–B273.
- [14] Q. DU, Z. HUANG, AND P. G. LEFLOCH, *Nonlocal conservation laws. a new class of monotonicity-preserving models*, SIAM J. Numer. Anal., 55 (2017), pp. 2465–2489.
- [15] Q. DU, L. JU, X. LI, AND Z. QIAO, *Maximum bound principles for a class of semilinear parabolic equations and exponential time-differencing schemes*, SIAM Review, 63 (2021), pp. 317–359.
- [16] P. FU, F. LI, AND Y. XU, *Globally divergence-free discontinuous Galerkin methods for ideal magnetohydrodynamic equations*, J. Sci. Comput., 77 (2018), pp. 1621–1659.
- [17] S. K. GODUNOV, *Symmetric form of the equations of magnetohydrodynamics*, Numerical Methods for Mechanics of Continuum Medium, 1 (1972), pp. 26–34.
- [18] S. GOTTLIEB, C.-W. SHU, AND E. TADMOR, *Strong stability-preserving high-order time discretization methods*, SIAM Review, 43 (2001), pp. 89–112.
- [19] P. M. GRUBER, *Convex and Discrete Geometry*, vol. 336, Springer Science & Business Media, 2007.
- [20] J.-L. GUERMOND, M. NAZAROV, B. POPOV, AND I. TOMAS, *Second-order invariant domain preserving approximation of the Euler equations using convex limiting*, SIAM J. Sci. Comput., 40 (2018), pp. A3211–A3239.
- [21] J.-L. GUERMOND AND B. POPOV, *Invariant domains and first-order continuous finite element approximation for hyperbolic systems*, SIAM J. Numer. Anal., 54 (2016), pp. 2466–2489.
- [22] J.-L. GUERMOND AND B. POPOV, *Invariant domains and second-order continuous finite element approximation for scalar conservation equations*, SIAM J. Numer. Anal., 55 (2017), pp. 3120–3146.
- [23] X. Y. HU, N. A. ADAMS, AND C.-W. SHU, *Positivity-preserving method for high-order conservative schemes solving compressible Euler equations*, J. Comput. Phys., 242 (2013), pp. 169–180.
- [24] Y. JIANG AND H. LIU, *Invariant-region-preserving DG methods for multi-dimensional hyperbolic conservation law systems, with an application to compressible Euler equations*, J. Comput. Phys., 373 (2018), pp. 385–409.
- [25] L. JU, X. LI, Z. QIAO, AND J. YANG, *Maximum bound principle preserving integrating factor Runge–Kutta methods for semilinear parabolic equations*, J. Comput. Phys., 439 (2021), p. 110405.
- [26] B. KHOBALATTE AND B. PERTHAME, *Maximum principle on the entropy and second-order kinetic schemes*, Math. Comp., 62 (1994), pp. 119–131.
- [27] L. KRIVODONOVA, J. XIN, J.-F. REMACLE, N. CHEVAUGEON, AND J. E. FLAHERTY, *Shock detection and limiting with discontinuous Galerkin methods for hyperbolic conservation laws*, Appl. Numer. Math., 48 (2004), pp. 323–338.
- [28] I. E. LEONARD AND J. E. LEWIS, *Geometry of Convex Sets*, John Wiley & Sons, Hoboken, New Jersey, 2015.
- [29] F. LI AND C.-W. SHU, *Locally divergence-free discontinuous Galerkin methods for MHD equations*, J. Sci. Comput., 22 (2005), pp. 413–442.
- [30] F. LI, L. XU, AND S. YAKOVLEV, *Central discontinuous Galerkin methods for ideal MHD equations with the exactly divergence-free magnetic field*, J. Comput. Phys., 230 (2011), pp. 4828–4847.
- [31] H. LI AND X. ZHANG, *On the monotonicity and discrete maximum principle of the finite difference implementation of C^0 - Q^2 finite element method*, Numer. Math., 145 (2020), pp. 437–472.
- [32] J. LI, X. LI, L. JU, AND X. FENG, *Stabilized integrating factor Runge–Kutta method and unconditional preservation of maximum bound principle*, SIAM J. Sci. Comput., 43 (2021), pp. A1780–A1802.
- [33] C. LIANG AND Z. XU, *Parametrized maximum principle preserving flux limiters for high order schemes solving multi-dimensional scalar hyperbolic conservation laws*, J. Sci. Comput., 58 (2014), pp. 41–60.
- [34] D. LING, J. DUAN, AND H. TANG, *Physical-constraints-preserving Lagrangian finite volume schemes for one-and two-dimensional special relativistic hydrodynamics*, J. Comput. Phys., 396 (2019), pp. 507–543.
- [35] A. K. MEENA, H. KUMAR, AND P. CHANDRASHEKAR, *Positivity-preserving high-order discontinuous Galerkin schemes for ten-moment Gaussian closure equations*, J. Comput. Phys., 339 (2017), pp. 370–395.
- [36] A. K. MEENA, R. KUMAR, AND P. CHANDRASHEKAR, *Positivity-preserving finite difference WENO scheme for ten-moment equations with source term*, J. Sci. Comput., 82 (2020), pp. 1–37.
- [37] C. NICULESCU AND L.-E. PERSSON, *Convex Functions and Their Applications*, Springer, 2006.
- [38] E. OLBRANT, C. D. HAUCK, AND M. FRANK, *A realizability-preserving discontinuous Galerkin method for the*

- M1 model of radiative transfer*, J. Comput. Phys., 231 (2012), pp. 5612–5639.
- [39] B. PERTHAME, *Second-order Boltzmann schemes for compressible Euler equations in one and two space dimensions*, SIAM J. Numer. Anal., 29 (1992), pp. 1–19.
- [40] B. PERTHAME AND C.-W. SHU, *On positivity preserving finite volume schemes for Euler equations*, Numer. Math., 73 (1996), pp. 119–130.
- [41] T. QIN, C.-W. SHU, AND Y. YANG, *Bound-preserving discontinuous Galerkin methods for relativistic hydrodynamics*, J. Comput. Phys., 315 (2016), pp. 323–347.
- [42] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, 2015.
- [43] J. A. ROSSMANITH AND D. C. SEAL, *A positivity-preserving high-order semi-Lagrangian discontinuous Galerkin scheme for the Vlasov–Poisson equations*, J. Comput. Phys., 230 (2011), pp. 6203–6232.
- [44] J. SHEN AND J. XU, *Unconditionally bound preserving and energy dissipative schemes for a class of Keller–Segel equations*, SIAM J. Numer. Anal., 58 (2020), pp. 1674–1695.
- [45] C.-W. SHU, *Bound-preserving high-order schemes for hyperbolic equations: Survey and recent developments*, in Theory, Numerics and Applications of Hyperbolic Problems II, C. Klingenberg and M. Westdickenberg, eds., Cham, 2018, Springer International Publishing, pp. 591–603.
- [46] M. SPIVAK, *Calculus on Manifolds*, W. A. Benjamin, New York, 1965.
- [47] E. TADMOR, *A minimum entropy principle in the gas dynamics equations*, Appl. Numer. Math., 2 (1986), pp. 211–219.
- [48] H.-Z. TANG AND K. XU, *Positivity-preserving analysis of explicit and implicit Lax–Friedrichs schemes for compressible Euler equations*, J. Sci. Comput., 15 (2000), pp. 19–28.
- [49] T. TANG AND K. XU, *Gas-kinetic schemes for the compressible Euler equations: positivity-preserving analysis*, Z. Angew. Math. Phys., 50 (1999), pp. 258–281.
- [50] K. WU, *Design of provably physical-constraint-preserving methods for general relativistic hydrodynamics*, Phys. Rev. D, 95 (2017), 103001.
- [51] K. WU, *Positivity-preserving analysis of numerical schemes for ideal magnetohydrodynamics*, SIAM J. Numer. Anal., 56 (2018), pp. 2124–2147.
- [52] K. WU, *Minimum principle on specific entropy and high-order accurate invariant region preserving numerical methods for relativistic hydrodynamics*, SIAM J. Sci. Comput., in press (2021).
- [53] K. WU AND C.-W. SHU, *A provably positive discontinuous Galerkin method for multidimensional ideal magnetohydrodynamics*, SIAM J. Sci. Comput., 40 (2018), pp. B1302–B1329.
- [54] K. WU AND C.-W. SHU, *Provably positive high-order schemes for ideal magnetohydrodynamics: analysis on general meshes*, Numer. Math., 142 (2019), pp. 995–1047.
- [55] K. WU AND C.-W. SHU, *Provably physical-constraint-preserving discontinuous Galerkin methods for multidimensional relativistic MHD equations*, Numer. Math., 148 (2021), pp. 699–741.
- [56] K. WU AND H. TANG, *High-order accurate physical-constraints-preserving finite difference WENO schemes for special relativistic hydrodynamics*, J. Comput. Phys., 298 (2015), pp. 539–564.
- [57] K. WU AND H. TANG, *Admissible states and physical-constraints-preserving schemes for relativistic magnetohydrodynamic equations*, Math. Models Methods Appl. Sci., 27 (2017), pp. 1871–1928.
- [58] K. WU AND H. TANG, *Physical-constraint-preserving central discontinuous Galerkin methods for special relativistic hydrodynamics with a general equation of state*, Astrophys. J. Suppl. Ser., 228 (2017), 3.
- [59] K. WU AND Y. XING, *Uniformly high-order structure-preserving discontinuous Galerkin methods for Euler equations with gravitation: Positivity and well-balancedness*, SIAM J. Sci. Comput., 43 (2021), pp. A472–A510.
- [60] Y. XING, X. ZHANG, AND C.-W. SHU, *Positivity-preserving high order well-balanced discontinuous Galerkin methods for the shallow water equations*, Adv. Water Resour., 33 (2010), pp. 1476–1493.
- [61] T. XIONG, J.-M. QIU, AND Z. XU, *High order maximum-principle-preserving discontinuous Galerkin method for convection-diffusion equations*, SIAM J. Sci. Comput., 37 (2015), pp. A583–A608.
- [62] T. XIONG, J.-M. QIU, AND Z. XU, *Parametrized positivity preserving flux limiters for the high order finite difference WENO scheme solving compressible Euler equations*, J. Sci. Comput., 67 (2016), pp. 1066–1088.
- [63] Z. XU, *Parametrized maximum principle preserving flux limiters for high order schemes solving hyperbolic conservation laws: one-dimensional scalar problem*, Math. Comp., 83 (2014), pp. 2213–2238.
- [64] Z. XU AND X. ZHANG, *Bound-preserving high-order schemes*, in Handbook of Numerical Analysis, vol. 18, Elsevier, 2017, pp. 81–102.
- [65] D. YUAN, J. CHENG, AND C.-W. SHU, *High order positivity-preserving discontinuous Galerkin methods for radiative transfer equations*, SIAM J. Sci. Comput., 38 (2016), pp. A2987–A3019.
- [66] X. ZHANG, *On positivity-preserving high order discontinuous Galerkin schemes for compressible Navier-Stokes equations*, J. Comput. Phys., 328 (2017), pp. 301–343.
- [67] X. ZHANG AND C.-W. SHU, *On maximum-principle-satisfying high order schemes for scalar conservation laws*, J. Comput. Phys., 229 (2010), pp. 3091–3120.
- [68] X. ZHANG AND C.-W. SHU, *On positivity-preserving high order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes*, J. Comput. Phys., 229 (2010), pp. 8918–8934.
- [69] X. ZHANG AND C.-W. SHU, *Positivity-preserving high order discontinuous Galerkin schemes for compressible Euler equations with source terms*, J. Comput. Phys., 230 (2011), pp. 1238–1248.
- [70] X. ZHANG AND C.-W. SHU, *A minimum entropy principle of high order schemes for gas dynamics equations*, Numer. Math., 121 (2012), pp. 545–563.
- [71] X. ZHANG, Y. XIA, AND C.-W. SHU, *Maximum-principle-satisfying and positivity-preserving high order discontinuous Galerkin schemes for conservation laws on triangular meshes*, J. Sci. Comput., 50 (2012), pp. 29–62.
- [72] Y. ZHANG, X. ZHANG, AND C.-W. SHU, *Maximum-principle-satisfying second order discontinuous Galerkin schemes for convection-diffusion equations on triangular meshes*, J. Comput. Phys., 234 (2013), pp. 295–316.
- [73] J. ZHAO AND H. TANG, *Runge-Kutta discontinuous Galerkin methods for the special relativistic magnetohydrodynamics*, J. Comput. Phys., 343 (2017), pp. 33–72.