

NUMERICAL SIMULATIONS OF SURFACE-QUASI GEOSTROPHIC FLOWS ON PERIODIC DOMAINS *

ANDREA BONITO[†] AND MURTAZO NAZAROV[‡]

Abstract. We propose a novel algorithm for the approximation of surface-quasi geostrophic (SQG) flows modeled by a nonlinear partial differential equation coupling transport and fractional diffusion phenomena. The time discretization consists of an explicit strong-stability-preserving three-stage Runge-Kutta method while a flux-corrected-transport (FCT) method coupled with Dunford-Taylor representations of fractional operators is advocated for the space discretization. Standard continuous piecewise linear finite elements are employed and the algorithm does not have restrictions on the mesh structure nor on the computational domain. In the inviscid case, we show that the resulting scheme satisfies a discrete maximum principle property under a standard CFL condition and observe, in practice, its second-order accuracy in space. The algorithm successfully approximates several benchmarks with sharp transitions and fine structures typical of SQG flows. In addition, theoretical Kolmogorov energy decay rates are observed on a freely decaying atmospheric turbulence simulation.

Key words. Geostrophic flows; Finite element method; Dunford-Taylor integral; Fractional Diffusion; Discrete maximum principle; Nonlinear viscosity; FCT algorithm.

AMS subject classifications. 65M60, 65M12, 35L65, 76U05, 35R11

1. Introduction. The Navier-Stokes system models the behavior of incompressible, adiabatic, inviscid fluids in hydrostatic balance. When, in addition, the fluid is constrained by environmental rotation and stratification, Charney [15] derived in the 1940's a three dimensional quasi-geostrophic model to describe large-scale mid-latitude atmospheric motions and oceanographic motions. Charney's quasi-geostrophic model received much attention, we mention [23, 42, 49, 36, 38, 60] for discussions on its validity.

In a surface quasi-geostrophic (SQG) setting, it is further assumed that the potential vorticity is uniform, see for instance [49, 14, 37, 20, 21, 31]. Consequently, on the half plane above the surface $\mathcal{S} := \{(x_1, x_2, x_3) \in \mathbb{R}^3 : x_3 = 0\}$, the vorticity $\tilde{\psi}(x_1, x_2, x_3, t)$ satisfies

$$(1.1) \quad \Delta \tilde{\psi} = 0, \quad \text{where } x_3 > 0 \quad \text{and} \quad \lim_{z \rightarrow \infty} \tilde{\psi}(x_1, x_2, x_3, t) = 0.$$

On \mathcal{S} , the buoyancy (or potential temperature) is given by $\theta := \partial_{x_3} \tilde{\psi}|_{x_3=0}$.

We restrict our considerations to surface consisting in a rectangular periodic domain $\Omega := (0, \pi)^2$ denoted \mathbb{T}^2 in short. When restricted to \mathbb{T}^2 , (1.1) corresponds to a nonlocal elliptic partial differential equation involving $\theta(x_1, x_2, t)$ and $\psi(x_1, x_2, t) := \tilde{\psi}(x_1, x_2, 0, t)$, namely

$$(-\Delta)^{\frac{1}{2}} \psi = \theta,$$

where $(-\Delta)^{\frac{1}{2}}$ stands for the *spectral* fractional laplacian defined in Section 2.1. The buoyancy is transported on \mathbb{T}^2 along the orthogonal directions to the vorticity gradient

*

Funding: A.B. is partially supported by the NSF Grant DMS-1817691; M.N. is partially supported by Esseen scholarship at Uppsala University.

[†]Department of Mathematics, Texas A&M University, College Station, TX 77843.

[‡]Corresponding author. Department of Information Technology, Uppsala University, SE 75105.

$\mathbf{u} := \nabla^\perp \psi$, where for $v : \mathbb{R}^2 \rightarrow \mathbb{R}$ we set $\nabla^\perp v := \begin{pmatrix} -\partial_{x_2} v \\ \partial_{x_1} v \end{pmatrix}$. In addition, we account for the Ekman pumping effect (friction between vertical thin layers of atmosphere) resulting in the following nonlinear advection-diffusion relation for the buoyancy

$$\partial_t \theta + \mathbf{u} \cdot \nabla \theta + \varkappa (-\Delta)^{\frac{1}{2}} \theta = 0,$$

where $\varkappa \geq 0$ stands for the Ekman pumping coefficient. This coefficient is typically small except on narrow boundary layers touching the fluid boundary [49]. The case $\varkappa = 0$ will be referred to as the inviscid case. A detailed derivation of the SQG system can be found in [61], based on the works [49, 14, 31, 39, 32] and [8, 57].

The above nonlinear system of equations features many aspect of large-scale atmospheric motions. Among them, we list the apparition in finite time of discontinuous temperature - called Frontogenesis - and the conservation (for $\varkappa = 0$) of the kinetic energy and helicity, see Section 4. Whether solutions to the SQG equations can develop singularities is a question which concerned many researchers and global regularity for general data remains an open problem [19]. We refer to [11, 18, 20, 53, 33, 10] for additional information.

Numerical methods are of fundamental importance to assess the behavior of the solutions to the SQG system. Particular attention must be made (i) to reproduce accurately discontinuous profiles while conserving quantities such as the kinetic energy and helicity; (ii) to preserve the discrete maximum principle. The conservation properties and the discrete maximum principle are of great interest in the application and analysis of conservation laws. Existing numerical algorithms for the approximation of the SQG system are based on spectral decompositions of the solution coupled with higher-order exponential filters [19, 20, 21], spectral vanishing viscosities [56], and standard hyperviscosity [31]. These filtering approaches are shown to give impressive results in terms of accurately calculating the kinetic energy and helicity as presented in the above-mentioned references. However, it is not known whether these approaches preserve the discrete maximum principle. Our approach is based on standard finite element discretization with nonlinear stabilization based on the second-order elliptic operator. A particular goal of this work is to construct a high-order scheme that has conservation properties and satisfies the discrete maximum principle in general unstructured triangulations.

We summaries in Section 2 the SQG system along with some of its important properties. In Section 3, we propose to adapt the algorithms proposed in [6, 7, 4, 5] to the present periodic setting and employ a flux corrected transport (FCT) limiting blending a low order scheme satisfying a discrete maximum principle (when $\varkappa = 0$) with a higher order shock-capturing method [26, 30, 25]. Notice that compared to the previous analyses of shock-capturing methods for hyperbolic problems, the hyperbolic flux in our system is non-local making the estimation of the wave speed non standard even in the linear case. The resulting scheme retain the maximum preserving property and is observed in practice to retain the higher order accuracy. At this point it is worth mentioning that there seem to be no mathematical explanation of the higher order properties of FCT algorithm available in the literature. We showcase in Section 4 the need of the FCT algorithm to avoid over-diffusive simulations. In fact, the numerical simulations obtained exhibit sharp resolutions of line discontinuities and fine structures. In addition, we propose numerical simulations of freely decaying turbulence and confirm the predictions of [31, 59, 37, 52, 12] for the decay of the kinetic energy cascade for the inviscid ($\varkappa = 0$) and diffuse ($\varkappa > 0$) SQG system. At large scales we recover the $-\frac{5}{3}$ Kolmogorov rate of decay typical to three dimensional

flows while a -3 Kolmogorov rate of decay, this time typical of two dimensional flows, is observed at small scales.

2. Preliminaries.

2.1. The Spectral Fractional Laplacian on the Torus. To define the fractional laplacian in (2.4), we denote by $\{(\lambda_i, \phi_i)\}_{i=0}^\infty \subset \mathbb{R}_+ \times H^1(\mathbb{T}^2)$ the eigenpairs of the Laplacian on the torus \mathbb{T}^2 . We use the convention $0 = \lambda_0 < \lambda_1 \leq \lambda_2 \leq \dots$ and assume that the ϕ_i 's are orthonormal in $L^2(\mathbb{T}^2)$ and orthogonal in $H^1(\mathbb{T}^2)$.

For $-1 \leq r \leq 1$, the fractional power of the Laplacian is defined for smooth functions $v \in C^\infty(\mathbb{T}^2)$ with vanishing mean value as

$$(2.1) \quad (-\Delta)^r v := \sum_{n=1}^{\infty} \lambda_i^r v_i \phi_i, \quad v_i := \int_{\mathbb{T}^2} v(\mathbf{x}) \phi_i(\mathbf{x}) d\mathbf{x}.$$

The definition of the fractional laplacian (2.1) is extended by density to

$$\mathcal{D}((-\Delta)^s) := \left\{ v \in L^2_\#(\mathbb{T}^2) : \sum_{i=0}^{\infty} \left(\int_{\mathbb{T}^2} v \phi_i \right)^2 \lambda_i^{2s} < \infty \right\},$$

where $L^2_\#(\mathbb{T}^2)$ is the subspace of $L^2(\mathbb{T}^2)$ consisting of vanishing mean value functions.

For latter use, we record the following relation directly following from the definition of the fractional laplacian

$$(2.2) \quad \int_{\mathbb{T}^2} (-\Delta)^{s_1} v (-\Delta)^{s_2} w = \int_{\mathbb{T}^2} (-\Delta)^{r_1} v (-\Delta)^{r_2} w, \quad v, w \in C^\infty(\mathbb{T}^2) \cap L^2_\#(\mathbb{T}^2),$$

for $-1 \leq s_1 \leq r_1 \leq r_2 \leq s_2 \leq 1$ satisfying $s_1 + s_2 = r_1 + r_2$.

2.2. The SQG Equations. We denote by T the final time. The solution to the SQG system is a pair $\theta, \psi : \mathbb{T}^2 \times [0, T] \rightarrow \mathbb{R}$ satisfying

$$(2.3) \quad \partial_t \theta + \mathbf{u} \cdot \nabla \theta + \varkappa (-\Delta)^s \theta = 0, \quad \text{in } \mathbb{T}^2 \times (0, T]$$

and

$$(2.4) \quad (-\Delta)^{\frac{1}{2}} \psi = \theta, \quad \mathbf{u} = \nabla^\perp \psi \quad \text{in } \mathbb{T}^2 \times (0, T].$$

As for e.g. in [18], we introduced a parameter $0 < s < 1$ to include additional mathematical models considered in the literature for the design of the numerical method. However, our numerical experiments focus on the physical SQG system and thus on the critical case $s = \frac{1}{2}$. The system of equations (2.3) and (2.4) is supplemented by the initial and mean value conditions

$$\theta(\cdot, 0) = \theta_0 \quad \text{in } \mathbb{T}^2, \quad \int_{\mathbb{T}^2} \theta = \int_{\mathbb{T}^2} \psi = 0 \quad \text{in } (0, T),$$

where $\theta_0 : \mathbb{T}^2 \rightarrow \mathbb{R}$ is a given initial buoyancy satisfying $\int_{\mathbb{T}^2} \theta_0(\mathbf{x}) = 0$.

From now on we assume that there exists a unique sufficiently smooth solution (θ, ψ) and refer to works cited in the introduction for discussions on the existence and uniqueness of solutions as well as their regularity.

2.3. Kinetic Energy and Helicity. The kinetic energy $\mathcal{K}(\theta)$ and helicity $\mathcal{H}(\theta)$ are defined by

$$(2.5) \quad \mathcal{K}(\theta) := \frac{1}{2} \int_{\mathbb{T}^2} \theta^2(\mathbf{x}, t) d\mathbf{x} \quad \text{and} \quad \mathcal{H}(\theta) := - \int_{\mathbb{T}^2} \psi(\mathbf{x}, t) \theta(\mathbf{x}, t) d\mathbf{x}$$

and are monitored in several numerical experiments in Section 4 to showcase the performances of the proposed algorithm. We also compute in Section 4.5 the Kolmogorov energy cascades for the Kinetic model and validate our turbulence model. Both are conserved quantities when $\varkappa = 0$ and dissipated when $\varkappa > 0$. We make this more precise now.

To obtain an evolution relation for the kinetic energy, we multiply (2.3) by θ and integrate over \mathbb{T}^2 to get

$$(2.6) \quad \frac{d}{dt} \mathcal{K}(\theta) = -\varkappa \int_{\mathbb{T}^2} (-\Delta)^s \theta \theta,$$

where we used the definition $\mathbf{u} = \nabla^\perp \psi$ to deduce that $\operatorname{div} \mathbf{u} = 0$ and so $\int_{\mathbb{T}^2} \mathbf{u} \cdot \nabla \theta \theta = 0$. The integration by parts relation (2.2) applied to the right hand side of (2.6) yields

$$(2.7) \quad \frac{d}{dt} \mathcal{K}(\theta) = -\varkappa \int_{\mathbb{T}^2} |(-\Delta)^{\frac{s}{2}} \theta|^2.$$

We now turn our attention to the helicity. We multiply (2.3) by ψ , integrate over \mathbb{T}^2 and invoke the relation $\psi = (-\Delta)^{-\frac{1}{2}} \theta$ to write

$$(2.8) \quad \int_{\mathbb{T}^2} \partial_t \theta (-\Delta)^{-\frac{1}{2}} \theta + \int_{\mathbb{T}^2} \mathbf{u} \cdot \nabla \theta \psi = -\varkappa \int_{\mathbb{T}^2} (-\Delta)^s \theta (-\Delta)^{-\frac{1}{2}} \theta.$$

We rewrite the above three terms separately. The term involving the velocity $\mathbf{u} = \nabla^\perp \psi$ vanishes in this case as well

$$\int_{\mathbb{T}^2} \mathbf{u} \cdot \nabla \theta \psi = - \int_{\mathbb{T}^2} \mathbf{u} \cdot \nabla \psi \theta = - \int_{\mathbb{T}^2} \nabla^\perp \psi \cdot \nabla \psi \theta = 0.$$

For the left most term in (2.8), we invoke the integration by parts formula (2.2) (twice) and the relation $\psi = (-\Delta)^{-\frac{1}{2}} \theta$ to deduce

$$\int_{\mathbb{T}^2} \partial_t \theta (-\Delta)^{-\frac{1}{2}} \theta = \frac{1}{2} \frac{d}{dt} \int_{\mathbb{T}^2} |(-\Delta)^{-\frac{1}{4}} \theta|^2 = \frac{1}{2} \frac{d}{dt} \int_{\mathbb{T}^2} \psi \theta = -\frac{1}{2} \frac{d}{dt} \mathcal{H}(\theta).$$

Using (2.2) once again for the right hand side of (2.8) yields

$$-\varkappa \int_{\mathbb{T}^2} (-\Delta)^s \theta (-\Delta)^{-\frac{1}{2}} \theta = -\varkappa \int_{\mathbb{T}^2} |(-\Delta)^{\frac{1}{2}(s-\frac{1}{2})} \theta|^2.$$

Gathering the above relations, we obtain

$$(2.9) \quad \frac{1}{2} \frac{d}{dt} \mathcal{H}(\theta) = \varkappa \int_{\mathbb{T}^2} |(-\Delta)^{\frac{1}{2}(s-\frac{1}{2})} \theta|^2.$$

3. Numerical Algorithm.

3.1. The Finite Element Spaces. We propose to use continuous piecewise linear finite elements for the space approximation of the potential temperature θ and stream function ψ . Let $\{\mathcal{T}_h\}_{h>0}$ be a sequence of shape-regular, quasi-uniform and conforming triangulations of \mathbb{T}^2 in the sense of [16], where $h := \min_{K \in \mathcal{T}_h} \text{diam}(K)$ stands for the smallest diameter of all the triangles in \mathcal{T}_h .

To each triangulation \mathcal{T}_h , we associate the spaces of continuous piecewise polynomial

$$(3.1) \quad \mathcal{X}_h := \{v_h \in \mathcal{C}^0(\mathbb{T}^2); \forall K \in \mathcal{T}_h, v_h|_K \in \mathbb{P}_1\}, \quad \mathcal{X}_{h,0} := \mathcal{X}_h \cap L^2_{\#}(\mathbb{T}^2),$$

where \mathbb{P}_1 denotes the space of polynomials of degree at most one and $\mathcal{C}^0(\mathbb{T}^2)$ the space of continuous functions on \mathbb{T}^2 (and therefore 2π -periodic on each variable). We denote by $\{\varphi_1, \dots, \varphi_I\}$ the basis of \mathcal{X}_h made of linear Lagrange finite elements (hat functions) associated with the collection of all the vertices $\{\mathbf{x}_j\}_{j=1}^I$ in the triangulation \mathcal{T}_h (not counting twice the periodic nodes). The index list of basis functions interacting with φ_i , $1 \leq i \leq I$, is denoted by

$$(3.2) \quad \mathcal{I}(i) := \{j \in \{1, \dots, I\} : \text{supp}(\varphi_i) \cap \text{supp}(\varphi_j) \neq \emptyset\}.$$

A mass lumping strategy detailed below will be critical to obtain maximum principle preserving schemes. We denote by

$$(3.3) \quad m_{ij} := \int_{\mathbb{T}^2} \varphi_j \varphi_i \quad \text{and} \quad m_i := \sum_{j \in \mathcal{I}(i)} m_{ij} = \int_{\mathbb{T}^2} \varphi_i$$

the elements of the consistent and lumped mass matrices.

To ease the notations, we will use capital letters to denote finite element approximations and drop the subindex h . For instance, $\Theta \in \mathcal{X}_{h,0}$ will denote the approximation of θ .

3.2. Approximations of the Fractional Laplacian with Periodic Boundary Conditions. Several approaches are available for the approximation of the spectral fractional Laplacian. We refer for instance to the reviews [1] and [40]. In this work, we adapt the algorithms developed in [6, 7, 5], which are based on different Balakrishnan-Dunford-Taylor representations described now. We emphasize that the resulting algorithms consist of the agglomerations of solutions to advection-diffusion problems approximated using a standard continuous piecewise linear finite element space \mathcal{X}_h . Their implementations are therefore straightforward and readily available in standard finite element softwares. Also, the algorithms presented do not suffer any restriction regarding the shape of the computational domain.

3.2.1. Approximations of Negative Powers of Fractional Operators. For $f \in L^2_{\#}(\mathbb{T}^2)$ and $s \in (0, 1)$, we have the following representation

$$(-\Delta)^{-s} f = v := \frac{1}{\pi} \int_{-\infty}^{\infty} e^{(1-s)y} w(y) dy,$$

where $w(y) \in H^1(\mathbb{T}^2) \cap L^2_{\#}(\mathbb{T}^2)$ solves

$$e^y w(y) - \Delta w(y) = f \quad \text{in } \mathbb{T}^2,$$

see e.g. [62].

A sinc quadrature is advocated for the approximation of the integral in y , thereby requiring the values of $w(y_\ell)$ at some selected snapshots $y_\ell \in \mathbb{R}$. The latter are approximated using a standard finite element method for reaction-diffusion problems. Given a spacing parameter $k > 0$ and integer $M \sim k^{-2}$, we have

$$(3.4) \quad (-\Delta)^{-s} f \approx V_k := \frac{1}{\pi} k \sum_{\ell=-M}^M e^{(1-s)y_\ell} W(y_\ell),$$

where $y_\ell := \ell k$, $\ell = -M, \dots, M$, and $\mathcal{X}_h \ni W(y_\ell) \approx w(y_\ell)$ solves

$$e^{y_\ell} \int_{\mathbb{T}^2} W(y_\ell) R + \int_{\mathbb{T}^2} \nabla W(y_\ell) \cdot \nabla R = \int_{\mathbb{T}^2} f R, \quad \forall R \in \mathcal{X}_h.$$

Notice that when $\int_{\mathbb{T}^2} f = 0$, we automatically have $\int_{\mathbb{T}^2} W(y_\ell) = 0$ and thus $V_k \in \mathcal{X}_{h,0}$. We refer to [6, 5] for the convergence analysis of v_h^k towards v . We only point out here that the convergence is exponential in $-1/k$ and optimal in h (depending on the regularity of f and the metric used to measure the error).

3.2.2. Approximations of Positive Powers of Fractional Operators. While (3.4) is sufficient to design a numerical scheme approximating (2.4), the explicit nature of our proposed time stepping scheme (see Section 3) also requires, when $s > 0$, an approximation of

$$\int_{\mathbb{T}^2} (-\Delta)^s V W,$$

for $\frac{1}{2} \leq s < 1$ and $V, W \in \mathcal{X}_h$. This time, we use the representation derived in [4]

$$(3.5) \quad \int_{\mathbb{T}^2} (-\Delta)^s V W = 2 \frac{\sin(\pi s)}{\pi} \int_0^\infty e^{sy} \int_{\mathbb{T}^2} (V + \tilde{V}(y; V)) W dy,$$

which is valid for $0 \leq s \leq 1$, $V, W \in \mathcal{X}_h$ and where the function $\tilde{V} := \tilde{V}(y; V) \in \mathcal{X}_h$ are given by the relation

$$(3.6) \quad \int_{\mathbb{T}^2} \tilde{V} R + e^{-y} \int_{\mathbb{T}^2} \nabla \tilde{V} \cdot \nabla R = - \int_{\mathbb{T}^2} V R, \quad \forall R \in \mathcal{X}_h.$$

As in the previous case, the integration in y is approximated by a sinc quadrature: given $k > 0$ and $M \sim 1/k^2$, we define

$$(3.7) \quad A_{h,k}(V, W) := 2 \frac{\sin(\pi s)}{\pi} k \sum_{\ell=-M}^M e^{sy_\ell} \int_{\mathbb{T}^2} (V + \tilde{V}(y_\ell; V)) W \approx \int_{\mathbb{T}^2} (-\Delta)^s V W.$$

An analysis of this approximation strategy in the more complex case of the integral fractional Laplacian is available in [4]. We do not expand on this further but note for later use that because $\int_{\mathbb{T}} \tilde{V}(y, V) = - \int_{\mathbb{T}^2} V$, we deduce that

$$(3.8) \quad A_{h,k}(V, 1) = 0, \quad \forall V \in \mathcal{X}_h.$$

3.2.3. Approximations of the System Velocity. We now discuss the approximation of the velocity \mathbf{u} in (2.4) for a given approximation $\Theta \in \mathcal{X}_{h,0}$ of $\theta \in L^2_{\#}(\mathbb{T}^2)$. It is performed in two steps. First, we use the approximation of the inverse fractional Laplacian (3.4) to define $\Psi_k \in \mathcal{X}_{h,0}$ as

$$\Psi_k := \frac{1}{\pi} k \sum_{\ell=-M}^M e^{(1-s)y_\ell} W(y_\ell) \approx \psi = (-\Delta)^{-1/2} \theta,$$

where $W(y_\ell) \in \mathcal{X}_{h,0}$ solves

$$e^{y_\ell} \int_{\mathbb{T}^2} W(y_\ell) R + \int_{\mathbb{T}^2} \nabla W(y_\ell) \cdot \nabla R = \int_{\mathbb{T}^2} \Theta R, \quad \forall R \in \mathcal{X}_h.$$

Then, the velocity approximation $\mathbf{U}_k := \mathbf{U}_k(\Theta) \in [\mathcal{X}_h]^2$ is defined as the component-wise Clément interpolant [17], see also [54], of $\nabla^\perp \Psi_k$. Notice that this construction does not guarantee that $\operatorname{div}(\mathbf{U}_k) = 0$. This possible lack of conservation property will be accounted for in the design of the algorithm for the temperature potential equation below, see for instance Lemma 3.1.

3.3. Approximation of the Temperature Potential Equation (2.3). The third order (three stages) Strong Stability Preserving Runge-Kutta (SSP-RK3) method [55] is advocated for the approximation of the time evolution in (2.3). We recall that one step of the SSP-RK3 scheme on an homogeneous equation $\frac{d}{dt}v = f(v)$ consist of computing v^{n+1} from v^n as follows:

$$\begin{aligned} v^{(1)} &:= v^n + \Delta t_{n+1} f(v^n), \\ v^{(2)} &:= \frac{3}{4}v^n + \frac{1}{4}(v^{(1)} + \Delta t_{n+1} f(v^{(1)})), \\ v^{n+1} &:= \frac{1}{3}v^n + \frac{2}{3}(v^{(2)} + \Delta t_{n+1} f(v^{(2)})). \end{aligned}$$

Since SSP-RK3 consists of a linear combination of three forward Euler steps, we restrict the discussion below to the construction of the latter. The finite element method for the space discretization is based on the finite element spaces (3.1) enhanced with adequate integration formulas and vanishing entropy residual viscosity stabilizations. One particular aspect the latter is that it does not rely on a mesh-size function and the approximation of the time derivative in the residual equation is by-passed.

As we shall see, these choices lead to a method satisfying a maximum principle when $\varkappa = 0$ (see Theorems 3.2 and 3.4) while retaining in practice the second order accuracy (see Section 4). This is in contrast with the standard entropy viscosity method.

3.3.1. The Conservative Galerkin Method. The time interval $[0, T]$ is split onto N intervals of variable length Δt_n , $n = 1, \dots, N$ and we set $t_n := \sum_{m=1}^n \Delta t_m$, $n = 0, \dots, N$ to denote the breakpoints of this subdivision. Let $\Theta^0 \in \mathcal{X}_{h,0}$ be an approximation of the initial temperature potential $\theta_0 \in L^2_{\#}(\mathbb{T}^2)$. We compute Θ_k^n , $n = 1, \dots, N$ recursively as detailed now. Given the temperature approximation $\Theta_k^n \in \mathcal{X}_{h,0}$ and the velocity approximation $\mathbf{U}_k^n := \mathbf{U}_k(\Theta_k^n) \in [\mathcal{X}_h]^2$ (see Section 3.2.3), we define $\Theta_k^{n+1} \in \mathcal{X}_h$ as the solution to

$$(3.9) \quad \int_{\mathbb{T}^2} \frac{\Theta_k^{n+1} - \Theta_k^n}{\Delta t_{n+1}} \varphi_i + \int_{\mathbb{T}^2} \mathbf{U}_k^n \cdot \nabla \Theta_k^n \varphi_i + \varkappa A_{h,k}(\Theta_k^n, \varphi_i) = 0, \quad 1 \leq \varphi_i \leq I.$$

In general $\int_{\mathbb{T}^2} \Theta_k^{n+1} \neq 0$ due to non-conservative approximation of the velocity, i.e. $\operatorname{div}(\mathbf{U}_k^n) \neq 0$. To circumvent this issue, we follow [28, Sec. 3.2] and replace the flux $\mathbf{U}_k^n \cdot \nabla \Theta_k^n$ by its linear interpolation $\sum_{j=1}^I \mathbf{u}_j^n \theta_j^n \varphi_j$ with $\mathbf{u}_j^n := \mathbf{U}_k^n(\mathbf{x}_j)$ and $\theta_j^n := \Theta_k^n(\mathbf{x}_j)$. The velocity term in (3.9) is thus approximated by

$$\int_{\mathbb{T}^2} \mathbf{U}_k^n \cdot \nabla \Theta_k^n \varphi_i \approx \sum_{j=1}^I \theta_j^n \mathbf{u}_j^n \cdot \int_{\mathbb{T}^2} \nabla \varphi_j \varphi_i = \sum_{j=1}^I \mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n, \quad i = 1, \dots, I,$$

where we introduced the notation $\mathbf{c}_{ij} := \int_{\mathbb{T}^2} \nabla \varphi_j \varphi_i$. In turn, (3.9) reduces to a system of equations for the coefficient $(\theta_j^{n+1})_{j=1}^I$ of $\Theta_k^{n+1} \in \mathcal{X}_h$, namely

$$(3.10) \quad \sum_{j=1}^I m_{ij} \frac{\theta_j^{n+1} - \theta_j^n}{\Delta t_{n+1}} + \sum_{j=1}^I \mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n + \varkappa \sum_{j=1}^I \theta_j^n A_{h,k}(\varphi_j, \varphi_i) = 0, \quad i = 1, \dots, I.$$

Notice that the relations

$$(3.11) \quad \mathbf{c}_{ij} = -\mathbf{c}_{ji} \quad \text{and} \quad \sum_{j=1}^I \mathbf{c}_{ji} = 0$$

hold. As a consequence,

$$(3.12) \quad \sum_{i=1}^I \sum_{j=1}^I \mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n = \sum_{i=1}^I \sum_{j=1}^I (\mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n - \mathbf{u}_i^n \cdot \mathbf{c}_{ij} \theta_i^n) = 0$$

and we now have, recalling the definition (3.3) of m_{ij} ,

$$\int_{\mathbb{T}^2} \Theta_k^{n+1} = \sum_{i,j=1}^I m_{ij} \theta_j^{n+1} = \sum_{i,j=1}^I m_{ij} \theta_j^n = \int_{\mathbb{T}^2} \Theta_k^n = 0$$

thanks to (3.8) as well. From this, we deduce that $\Theta_k^{n+1} \in \mathcal{X}_{h,0}$.

The numerical scheme (3.10) does not satisfy a maximum principle, it is actually not even stable in $L^\infty(\mathbb{T}^2)$. In the next sections we modify the above scheme and obtain low order method satisfying a maximum principle (or invariant domain property when $\varkappa > 0$).

3.3.2. Low Order Scheme. We follow the approach in [24] modifying (3.10) with appropriate mass lumping quadratures and incorporating a low order (graph) viscosity. When $\varkappa = 0$, we show that the resulting scheme satisfies a discrete maximum principle property

$$\min_{j \in \mathcal{I}(i)} \theta_j^n \leq \theta_i^{n+1} \leq \max_{j \in \mathcal{I}(i)} \theta_j^n.$$

When $\varkappa > 0$, we cannot guaranteed that θ_i^{n+1} strictly lies in $[\min_{j \in \mathcal{I}(i)} \theta_j^n, \max_{j \in \mathcal{I}(i)} \theta_j^n]$ without additional assumptions. We postpone this discussion to Remark 3.3 below.

We start with a mass lumping quadrature formula (see the definitions (3.3)) for the time derivative term

$$\sum_{j=1}^I m_{ij} \frac{\theta_j^{n+1} - \theta_j^n}{\Delta t_{n+1}} \approx m_i \frac{\theta_i^{n+1} - \theta_i^n}{\Delta t_{n+1}}$$

and the diffusion term

$$\sum_{j=1}^I \theta_j^n A_{h,k}(\varphi_j, \varphi_i) \approx m_i A_i(\Theta_k^n),$$

where

$$(3.13) \quad A_i(\Theta_k^n) := 2 \frac{\sin(\pi s)}{\pi} k \sum_{\ell=-M}^M e^{sy_\ell} (\theta_i^n + \tilde{v}_i(y_\ell; \Theta_k^n))$$

and $\tilde{v}_i := \tilde{v}_i(y_\ell; \Theta_k^n)$, $i = 1, \dots, I$, satisfies

$$m_i \tilde{v}_i + e^{-y_\ell} \sum_{j \in I(i)} \tilde{v}_j \int_{\mathbb{T}} \nabla \varphi_j \cdot \nabla \varphi_i = -m_i \theta_i^n;$$

compare with (3.7) and (3.6). Notice that because $A_{h,k}(\Theta_k^n, 1) = 0$, see (3.8), we have

$$(3.14) \quad \sum_{i=1}^I m_i A_i(\Theta_k^n) = 0,$$

instrumental property to preserve the average of the buoyancy after each time step (see Lemma 3.1). It is well known that standard continuous Galerkin methods are not stable in the approximation of first order hyperbolic systems [22, Chapter 5]. We propose here an artificial (vanishing) graph viscosity approach to stabilize our system, see [29, Sec.3.2] and [24, Sec.4.2] and incorporate a diffusing term

$$\sum_{j \in \mathcal{I}(i)} d_{ij}^{L,n} \theta_j^n$$

in (3.10). The coefficients $d_{ij}^{L,n}$ are defined for $i \neq j$ as

$$(3.15) \quad d_{ij}^{L,n} := \max(\lambda_{\max}(\mathbf{n}_{ij}, \theta_i^n, \theta_j^n) \|\mathbf{c}_{ij}\|_{\ell^2}, \lambda_{\max}(\mathbf{n}_{ji}, \theta_j^n, \theta_i^n) \|\mathbf{c}_{ji}\|_{\ell^2}),$$

where the local maximum wave speed are given by

$$(3.16) \quad \lambda_{\max} := \lambda_{\max}(\mathbf{n}_{ij}, \theta_i^n, \theta_j^n) := \max(|\mathbf{u}_i^n \cdot \mathbf{n}_{ij}|, |\mathbf{u}_j^n \cdot \mathbf{n}_{ij}|),$$

with $\mathbf{n}_{ij} := \mathbf{c}_{ij} / \|\mathbf{c}_{ij}\|_{\ell^2}$ and $\|\mathbf{c}_{ij}\|_{\ell^2}$ denotes the Euclidian norm of the vector $\mathbf{c}_{ij} \in \mathbb{R}^2$.

For $i = j$, we set

$$(3.17) \quad d_{ii}^{L,n} := - \sum_{i \neq j \in \mathcal{I}(i)} d_{ij}^{L,n}.$$

For future reference, we record the properties of the artificial viscosity coefficients:

$$(3.18) \quad d_{ij}^{L,n} \geq 0, \quad d_{ij}^{L,n} = d_{ji}^{L,n}, \quad \text{and} \quad \sum_{j \in \mathcal{I}(i)} d_{ij}^{L,n} = \sum_{i \in \mathcal{I}(j)} d_{ij}^{L,n} = 0.$$

We are now in position to define the low order scheme associated with (3.10): for $\Theta_k^n = \sum_{i=1}^I \theta_i^n \varphi_i$ and $\mathbf{U}_k^n = \sum_{i=1}^I \mathbf{u}_i^n \varphi_i$, we determine $\Theta_k^{L,n+1} = \sum_{i=1}^I \theta_i^{L,n+1} \varphi_i$ from the independent and explicit relations

$$(3.19) \quad \begin{aligned} m_i \theta_i^{L,n+1} &= m_i \theta_i^n - \Delta t_{n+1} \sum_{j \in \mathcal{I}(i)} \mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n \\ &\quad - \kappa \Delta t_{n+1} m_i A_i(\Theta_k^n) + \Delta t_{n+1} \sum_{j \in \mathcal{I}(i)} d_{ij}^{L,n} \theta_j^n, \quad 1 \leq i \leq I. \end{aligned}$$

The properties of the above scheme are discussed next. We start with a lemma ensuring that the discrete scheme preserves the average of the buoyancy, critical property to define the velocity (see Section 3.2.3).

LEMMA 3.1. *The low order scheme defined by the relations (3.19) is conservative, i.e.*

$$\int_{\mathbb{T}^2} \Theta_k^{L,n} = \int_{\mathbb{T}^2} \Theta_h^n.$$

Proof. After summing for $i = 1, \dots, I$ the relation (3.19) and using the conservation properties (3.12), (3.18) and (3.14), we realize that

$$\sum_{i=1}^I m_i \theta_i^{L,n+1} = \sum_{i=1}^I m_i \theta_i^n.$$

Hence,

$$\int_{\mathbb{T}^2} \Theta_k^{L,n+1} = \sum_{i=1}^I m_i \theta_i^{L,n+1} = \sum_{i=1}^I m_i \theta_i^n = \int_{\mathbb{T}^2} \Theta_k^n,$$

which is the desired estimate. \square

We now turn our attention to the discrete maximum principle when $\varkappa = 0$. The discrete maximum property requires a CFL type condition to hold, namely there exists a real number $0 < \text{CFL} \leq \frac{1}{2}$ such that the time step are selected (on the fly) to satisfy

$$(3.20) \quad \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} d_{ij}^{L,n} \leq \text{CFL}.$$

Note that $m_i \sim h^2$, $d_{ij}^{L,n} \sim \lambda^n h$, where λ^n is a characteristic (local velocity), and thus the above condition requires that for the computation of Θ_k^{n+1} , the time step Δt_{n+1} is selected so that $\Delta t_{n+1} \lambda^n / h$ is sufficiently small.

THEOREM 3.2 (Discrete maximum principle). *Let us assume that $\varkappa = 0$ and assume that condition (3.20) holds for some $0 < \text{CFL} \leq \frac{1}{2}$. Then the solution of the low order scheme (3.19) satisfies*

$$(3.21) \quad \min_{j \in \mathcal{I}(i)} \theta_j^n \leq \theta_i^{L,n+1} \leq \max_{j \in \mathcal{I}(i)} \theta_j^n$$

for all $i = 1, \dots, I$.

Proof. Using the conservation properties (3.12) and (3.18), we rewrite (3.19) as

$$\theta_i^{L,n+1} = \theta_i^n - \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} (\mathbf{u}_j^n \theta_j^n - \mathbf{u}_i^n \theta_i^n) \cdot \mathbf{c}_{ij} + \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} d_{ij}^{L,n} (\theta_j^n - \theta_i^n)$$

or, rearranging the terms, as

$$\theta_i^{L,n+1} = \theta_i^n \left(1 - \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} (-\mathbf{u}_i^n \cdot \mathbf{c}_{ij} + d_{ij}^{L,n}) \right) + \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} (-\mathbf{u}_j^n \cdot \mathbf{c}_{ij} + d_{ij}^{L,n}) \theta_j^n.$$

We obtain (3.21) by showing that the right hand side of the above relation is a convex combinations of θ_j^n , $j \in \mathcal{I}(i)$. To see this, we first note that the coefficients add-up to 1. Moreover, from the definition (3.15) of the low order viscosity coefficients $d_{ij}^{L,n}$, we have $\mathbf{u}_j^n \cdot \mathbf{c}_{ij} \leq d_{ij}^{L,n}$ and $-\mathbf{u}_i^n \cdot \mathbf{c}_{ij} \leq d_{ij}^{L,n}$. The former guarantees that

the coefficients in front of the θ_j^n are positive. The latter, in conjunction with the assumption $\text{CFL} \leq \frac{1}{2}$, yields

$$1 - \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} (-\mathbf{u}_i^n \cdot \mathbf{c}_{ij} + d_{ij}^{L,n}) \geq 1 - \frac{\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} 2d_{ij}^{L,n} \geq 0$$

and so the coefficient in front of θ_i^n is positive as well. This ends the proof. \square

We conclude this section with a remark concerning the viscous case.

Remark 3.3. We have already mentioned that a maximum principle like (3.21) does not necessarily hold without additional assumption. However, proceeding as in the proof of Theorem 3.2, we can rewrite the low order scheme as

$$\begin{aligned} \Theta_i^{L,n+1} = & \frac{1}{2} \left[\left(1 - \frac{2\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} (-\mathbf{u}_i^n \cdot \mathbf{c}_{ij} + d_{ij}^{L,n}) \right) \theta_i^n \right. \\ & \left. + \frac{2\Delta t_{n+1}}{m_i} \sum_{i \neq j \in \mathcal{I}(i)} (-\mathbf{u}_i^n \cdot \mathbf{c}_{ij} + d_{ij}^{L,n}) \theta_j^n \right] + \frac{1}{2} \left[\theta_i^n - 2\Delta t_{n+1} \sum_{i=1}^I m_i A_i(\Theta_k^n) \right]. \end{aligned}$$

From this we see that if for some real numbers $a < b$ we have $\theta_i^n \in [a, b]$ and $\theta_i^n - 2\Delta t_{n+1} \sum_{i=1}^I m_i A_i(\Theta_k^n) \in [a, b]$, both for $i = 1, \dots, I$, then $\theta_i^{L,n+1} \in [a, b]$ for $i = 1, \dots, I$ whenever

$$\Delta t_{n+1} \leq m_i / (4 \sum_{i \neq j \in \mathcal{I}(i)} d_{ij}^{L,n}).$$

The above condition holds provided $0 < \text{CFL} \leq \frac{1}{4}$. This property is called invariant domain in [29].

Alternatively, if the triangulation satisfies the acute angle condition $\int_{\mathbb{T}^2} \nabla \varphi_i \cdot \nabla \varphi_j < 0$ for $i \neq j$ along with a restriction on the sinc quadrature, then the viscous low order scheme satisfies the maximum principle property (3.21). This finer analysis is out of the scope of this paper and we refer to [2] for additional details. \square

3.3.3. Higher Order Scheme. The construction of the higher order scheme starts again from the Galerkin scheme (3.10) but with an artificial viscosity $d_{ij}^{H,n}$ chosen to vanish at a higher rate (with respect to the meshsize h) than (3.15) and (3.17) used for the low order scheme.

Besides being of higher order, the only restriction needed on the artificial viscosity coefficients $d_{ij}^{H,n}$ is that they satisfy the properties (3.18). In this work, we propose an artificial viscosity proportional to the residual of one entropy of the buoyancy equation. This is commonly referred to as an entropy viscosity method. Such strategy for conservation laws was originally proposed by [27] and was later extended to solve compressible flows [46, 47]. A priori error and stability analysis of the method for some entropy functionals were investigated in [45] and [3].

Given $\Theta_k^n = \sum_{j=1}^I \theta_j^n \varphi_i$ and $\mathbf{U}_k^n = \sum_{j=1}^I \mathbf{u}_j^n \varphi_i$, the higher order scheme consists of finding $\Theta_k^{H,n+1} = \sum_{j=1}^I \theta_j^{H,n+1} \varphi_i$ from the system of equations

$$\begin{aligned} (3.22) \quad \sum_{j=1}^I m_{ij} \theta_j^{H,n+1} = & \sum_{j=1}^I m_{ij} \theta_j^n - \Delta t_{n+1} \sum_{j=1}^I \mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n \\ & - \kappa \Delta t_{n+1} m_i A_i(\Theta_n) + \Delta t_{n+1} \sum_{j=1}^I d_{ij}^{H,n} \theta_j^n, \end{aligned}$$

for $i = 1, \dots, I$ and where the higher order entropy residual viscosity coefficients $d_{ij}^{H,n}$ are yet to be determined. This is the focus of the remaining part of this section but before embarking in this discussion, we point out that unlike for the lower order scheme (3.19), we use the consistent mass matrix for the time derivative term (no mass lumping) to reduce the dispersion error generated by the mass lumping in the low order scheme.

Entropy residuals have been discussed in details in the literature, see for instance [27]. In our particular context, for a given q sufficiently smooth, we define the entropy residual by $\mathcal{R}^n(Q) := \sum_{i=1}^I \mathcal{R}_i^n(q) \varphi_i$ where

$$\mathcal{R}_i^n(q) := \int_{\mathbb{T}^2} \left(\frac{q - \Theta_k^n}{\Delta t_{n+1}} + \mathbf{U}_k^n \cdot \nabla \Theta_k^n + \varkappa \sum_{j=1}^I A_j(\Theta_k^n) \varphi_j \right) \eta'(\Theta_k^n) \varphi_i,$$

where the entropy function η is taken to be $\eta(x) := \frac{1}{2}x^2$. Note that the action of the operator $(-\Delta)^s$ is not well defined on Θ_k^n and is therefore replaced in \mathcal{R}_i^n by $\sum_{j=1}^I A_j(\Theta_k^n) \varphi_j$.

One of the difficulty when using residual based viscosity on dynamical systems is the proper handling of the time derivative and in particular what function q to use. In order to avoid interferences from the time discretization in the computation of the residual, we resort to a novel idea from [25], see also [41]. To motivate the final expression of the residual we (formally) consider the solution θ^G of the following implicit time discretization

$$(3.23) \quad \frac{\theta^G - \theta_j^n}{\Delta t_{n+1}} + \mathbf{U}_k^n \cdot \nabla \theta^G + \varkappa (-\Delta)^s \theta^G = 0.$$

The entropy residual evaluated at $q = \theta^G$ reads

$$\mathcal{R}_i^n(\theta^G) = \int_{\mathbb{T}^2} \left(-\mathbf{U}_k^n \cdot \nabla \theta^G - \varkappa (-\Delta)^s \theta^G + \mathbf{U}_k^n \cdot \nabla \Theta_k^n + \varkappa \sum_{j=1}^I A_j(\Theta_k^n) \varphi_j \right) \eta'(\Theta_k^n) \varphi_i.$$

The above expression is not practical because of the cost in computing θ^G . Instead, one can use the plain Galerkin solution $\Theta_k^{G,n+1} := \sum_{j=1}^I \theta_j^{G,n+1} \varphi_j \in \mathcal{X}_h$ defined as the higher order scheme but without artificial viscosity

$$(3.24) \quad \sum_{j=1}^I m_{ij} \frac{\theta_j^{G,n+1} - \theta_j^n}{\Delta t_{n+1}} = - \sum_{j=1}^I \mathbf{u}_j^n \cdot \mathbf{c}_{ij} \theta_j^n - \varkappa m_i A_i(\Theta_k^n), \quad i = 1, \dots, I,$$

which leads to the final expression for the entropy residual

$$(3.25) \quad \mathcal{R}_i^n := \int_{\mathbb{T}^2} \left(\mathbf{U}_k^n \cdot \nabla (\Theta_k^n - \Theta_k^{G,n+1}) + \varkappa \sum_{j=1}^I (A_j(\Theta_k^n) - A_j(\Theta_k^{G,n+1})) \varphi_j \right) \eta'(\Theta_k^n) \varphi_i.$$

Then, the high order nonlinear viscosity in (3.22) is defined by

$$(3.26) \quad d_{ij}^{H,n} := \min \left(d_{ij}^{L,n}, c_{\text{EV}} \max \left(\frac{\mathcal{R}_i^n}{\eta_i^n}, \frac{\mathcal{R}_j^n}{\eta_j^n} \right) \right),$$

where c_{EV} is the stablization parameter (typically $0.1 \leq c_{\text{EV}} \leq 1$). The normalization coefficient $\widetilde{\eta}_i^n$ in (3.26) are given by

$$(3.27) \quad \widetilde{\eta}_i^n := \max(|\max_{j \in \mathcal{I}(i)} \eta(\Theta_j^n) - \min_{j \in \mathcal{I}(i)} \eta(\Theta_j^n)|, \epsilon |\eta(\Theta_i^n)|),$$

with $\epsilon := 10^{-8}$ (or below the scheme accuracy) is a small safety factor. We refer to Section 4 for a discussion on the effect of c_{EV} and on the normalization.

3.3.4. Flux corrected transport (FCT) limiting. In the above sections, we introduced two methods: a first-order maximum principle (or invariant domain) preserving scheme and a high order nonlinear viscosity scheme. The FCT algorithm below, first introduced by [9], ensures that the high order solution satisfies the discrete maximum principle (or invariant domain property).

We relate the high and low order schemes by subtracting (3.19) from (3.22):

$$\sum_{j \in \mathcal{I}(i)} m_{ij} \theta_j^{H,n+1} = m_i \theta_i^{L,n+1} + \sum_{j \in \mathcal{I}(i)} m_{ij} (\theta_j^n - \theta_i^n) + \Delta t_{n+1} \sum_{j \in \mathcal{I}} (d_{ij}^{H,n} - d_{ij}^{L,n}) \theta_j^n.$$

Note that to derive the above relation, we used the definition $m_i = \sum_{j \in \mathcal{I}(i)} m_{ij}$.

Adding $m_i \theta_i^{H,n+1}$ on both sides of the equation, we get

$$(3.28) \quad \begin{aligned} m_i \theta_i^{H,n+1} &= m_i \theta_i^{L,n+1} + \sum_{j \in \mathcal{I}(i)} m_{ij} (\theta_j^n - \theta_i^n) - \sum_{j \in \mathcal{I}(i)} m_{ij} (\theta_j^{H,n+1} - \theta_i^{H,n+1}) \\ &\quad + \Delta t_{n+1} \sum_{j \in \mathcal{I}(i)} (d_{ij}^{H,n} - d_{ij}^{L,n}) \theta_j^n \\ &=: m_i \theta_i^{L,n+1} + \Delta t_{n+1} \sum_{j \in \mathcal{I}(i)} \mathcal{A}_{ij}. \end{aligned}$$

The coefficient \mathcal{A}_{ij} can be rewritten using the conservative properties (3.18), (3.26) of the artificial diffusions coefficient as

$$(3.29) \quad \begin{aligned} \mathcal{A}_{ij} &= -\frac{m_{ij}}{\Delta t_{n+1}} \left((\theta_j^{H,n+1} - \theta_j^n) - (\theta_i^{H,n+1} - \theta_i^n) \right) \\ &\quad + (d_{ij}^{H,n} - d_{ij}^{L,n}) (\theta_j^n - \theta_i^n). \end{aligned}$$

From this representation, one sees that $\mathcal{A}_{ij} = -\mathcal{A}_{ji}$.

The low order solution as proven earlier preserves the discrete maximum principle ($\varkappa = 0$).

$$\theta_{\min}^n := \min_{j=1,\dots,I} \theta_j^n \leq \theta_i^{L,n+1} \leq \max_{j=1,\dots,I} \theta_j^n =: \theta_{\max}^n.$$

However, the high order solution may violate this maximum principle. The idea of Zalesak [63] is to introduce a limiter matrix of coefficient $\mathcal{L}_{ij} \geq 0$ to guarantee that the high order solution remains within $[\theta_{\min}^n, \theta_{\max}^n]$ while retaining high-order accuracy. To make this more precise, we write

$$\theta_{\min}^n = \theta_i^{L,n+1} + (\theta_{\min}^n - \theta_i^{L,n+1}) = \theta_i^{L,n+1} + \frac{\Delta t_{n+1}}{m_i} Q_i^-,$$

where $Q_i^- := \frac{m_i}{\Delta t_{n+1}} (\theta_i^{\min,n} - \theta_i^{L,n+1})$, $i = 1, \dots, I$. Furthermore, using the notations $P_i^- := \sum_{j \in \mathcal{I}(i)} \min\{0, \mathcal{A}_{ij}\}$ and $R_i^- := \min\left\{1, \frac{Q_i^-}{P_i^-}\right\}$ for $i = 1, \dots, I$, we deduce that

$$(3.30) \quad \theta_{\min}^n \leq \theta_i^{L,n+1} + \frac{\Delta t_{n+1}}{m_i} \sum_{j \in \mathcal{I}(i), \mathcal{A}_{ij} \leq 0} R_i^- \mathcal{A}_{ij}.$$

Similarly, upon defining $Q_i^+ := \frac{m_i}{\Delta t_{n+1}}(\theta_{\max}^n - \theta_i^{L,n+1})$, $P_i^+ := \sum_{j \in \mathcal{I}(i)} \max\{0, \mathcal{A}_{ij}\}$ and $R_i^+ := \min\left\{1, \frac{Q_i^+}{P_i^+}\right\}$, we have

$$(3.31) \quad \theta_{\max}^n \geq \theta_i^{L,n+1} + \frac{\Delta t_{n+1}}{m_i} \sum_{\substack{j \in \mathcal{I}(i) \\ \mathcal{A}_{ij} \geq 0}} R_i^+ \mathcal{A}_{ij}.$$

In view of (3.30) and (3.31), we define

$$(3.32) \quad \mathcal{L}_{ij} := \begin{cases} \min\{R_i^+, R_j^-\}, & \text{if } \mathcal{A}_{ij} \geq 0, \\ \min\{R_i^-, R_j^+\}, & \text{otherwise} \end{cases}$$

and the coefficients of the FCT solution Θ_k^{n+1} are obtained from the relation

$$(3.33) \quad \theta_i^{n+1} = \theta_i^{L,n+1} + \frac{\Delta t_{n+1}}{m_i} \sum_{j \in \mathcal{I}(i)} \mathcal{L}_{ij} \mathcal{A}_{ij};$$

compare to (3.28). We have the following result.

THEOREM 3.4. *The FCT solution $\Theta_k^{n+1} = \sum_{i=1}^I \theta_i^{n+1} \varphi_i$ of (3.33) with $\Theta_k^{L,n+1} = \sum_{i=1}^I \theta_i^{L,n+1} \varphi_i$ satisfies*

$$(3.34) \quad \int_{\mathbb{T}^2} \Theta_k^{n+1} = \int_{\mathbb{T}^2} \Theta_k^{L,n+1}.$$

Furthermore, if for some $a, b \in \mathbb{R}$ the lower order scheme satisfies $\theta_i^{L,n+1} \in [a, b]$ for all $i = 1, \dots, I$ and all $n \geq 1$, then

$$(3.35) \quad \theta_i^{n+1} \in [a, b], \quad \forall i = 1, \dots, I, \quad \forall n \geq 1.$$

Proof. For the conservative property, observe that \mathcal{A}_{ij} is skew-symmetric and \mathcal{L}_{ij} is symmetric. Then multiplying (3.33) by $m_i \varphi_i$ and summing over i we get (3.34). Relation (3.35) follows directly from the definition of \mathcal{L}_{ij} . The proof is complete. \square

We conclude this section with a summary of one Euler step for the approximation of the buoyancy equation.

Algorithm 3.1 Limiting algorithm for potential temperature equation

Input: $\Theta_k^n, U_k^n, \varkappa$ and Δt_{n+1}

Output: Θ_k^{n+1}

- 1: Compute $d_{ij}^{L,n}$ from (3.15) and the low order solution $\theta_i^{L,n+1}$ defined by (3.19);
 - 2: Compute $\Theta_k^{G,n+1}$ defined by (3.24) to construct $d_{ij}^{H,n}$ in (3.26);
 - 3: Compute the higher order solution $\theta_i^{H,n+1}$ defined by (3.22);
 - 4: Compute the matrix \mathcal{A}_{ij} in (3.29) and \mathcal{L}_{ij} in (3.32);
 - 5: Compute Θ_k^{n+1} using (3.33).
-

4. Numerical Illustrations. In this section, we solve several benchmark problems to validate the proposed numerical scheme and present novel insightful simulations. The smooth convection problem in Section 4.1 validates the entropy residual

viscosity model (3.26) and observe that the FCT scheme preserve the high order accuracy of the high order scheme. The discretization of the fractional diffusion operator is investigated in Section 4.2. In Sections 4.3 and 4.4, we perform standard benchmarks for the SQG system: rotating vortices and initial data with saddle structures leading to sharp transitions. We conclude with a turbulence study in Section 4.5 and show that our numerical scheme exhibit the theoretical predictions of the Kolomogorov energy decay rate.

In order to plot the evolution of the buoyancy approximation, we denote by $\Theta_k(t)$ the continuous piecewise time reconstruction defined by $\Theta_k(t)|_{[t_n, t_{n+1}]} := \Theta_k^n + (\Theta_k^{n+1} - \Theta_k^n)(t - t_n)/\Delta t_{n+1}$.

4.1. Smooth convection problem. We consider the inviscid SQG equations, i.e., $\varkappa = 0$, the convection field is given by $\mathbf{u} = (1, 1)^\top$, and initial data is a smooth function defined as

$$\theta_0(x_1, x_2) = \sin x_1 \sin x_2 + \cos x_2.$$

This is a pure convection problem with constant transport illustrating the differences between the low order, the higher order and FCT schemes.

We run the problem on a sequence of meshes until the final time $T = 2\pi$. The time steps is uniform over the entire simulation and chosen so that $\text{CFL} = 0.2$, see (3.20). Note that in view of Theorem 3.4, a larger CFL value than needed for the stabilized schemes is chosen to guarantee the stability of the Galerkin scheme used for comparison purposes. The nonlinear entropy residual parameter in (3.26) is set to $c_{\text{EV}} = 1$.

The results of the numerical simulation are collected in Table 1. The first-row block corresponds to Galerkin solution (3.24), i.e., without any stabilization terms. The second and third-row blocks correspond to the entropy viscosity solution described in Section 3.3.3 and the FCT solution described in Section 3.3.4. We compute the errors at the final time for the $L^1(\mathbb{T}^2)$ -, $L^2(\mathbb{T}^2)$ - and $L^\infty(\mathbb{T}^2)$ -norms for several spacial resolutions (uniform triangulations) along with their associated rate of convergence. We observe a second-order convergence rate in the $L^1(\mathbb{T}^2)$ - and $L^2(\mathbb{T}^2)$ -norms but the entropy viscosity solution and the FCT solution deliver suboptimal rate in the $L^\infty(\mathbb{T}^2)$ -norm. Obtaining optimal rates in the $L^\infty(\mathbb{T}^2)$ -norm for limited solutions is notoriously difficult, see e.g., [30]. However, we report in the fourth and fifth row blocks of Table 1, FCT simulations obtained using a different normalization term

$$(4.1) \quad \widetilde{\eta}_2^n := \max(\widetilde{\eta}_i^n, |\eta(\Theta_i^n)|)$$

and leading to optimal second order convergence rates in all norms. Although not optimal in the maximum norm, we use for the rest of the paper the normalization (3.27) because of its robustness on nonlinear problems.

The values of the kinetic energy and helicity (2.5) are given in the last two columns of Table 1. We see that for the finer meshes the method recovers the kinetic energy and helicity to their reference values of 14.8041 and 26.4241 computed with the initial condition on the finest mesh. Note that the exact value of the kinetic energy is $3\pi^2/2 \approx 14.8044$.

4.2. Smooth fractional diffusion problem. In this section we approximate a purely fractional diffusion problem

$$\partial_t \theta + \frac{1}{1000} (-\Delta)^{\frac{1}{4}} \theta = 0, \quad \text{in } \mathbb{T}^2 \times (0, \pi),$$

	# dofs	L^1	rate	L^2	rate	L^∞	rate	$\mathcal{K}(\theta)$	$\mathcal{H}(\theta)$
Galerkin	100	1.58E+00	—	3.13E-01	—	1.44E-01	—	13.5743	23.8296
	400	3.77E-01	2.06	7.52E-02	2.06	3.35E-02	2.10	14.4839	25.7446
	1600	9.39E-02	2.01	1.87E-02	2.01	8.23E-03	2.02	14.7235	26.2526
	6400	2.35E-02	2.00	4.68E-03	2.00	2.06E-03	2.00	14.7841	26.3816
	25600	5.87E-03	2.00	1.17E-03	2.00	5.14E-04	2.00	14.7993	26.4139
	102400	1.47E-03	2.00	2.93E-04	2.00	1.29E-04	2.00	14.8031	26.4220
	409600	3.67E-04	2.00	7.31E-05	2.00	3.21E-05	2.00	14.8041	26.4241
EV	100	9.26E+00	—	1.78E+00	—	6.04E-01	—	7.43147	17.7222
	400	2.46E+00	1.91	5.36E-01	1.73	2.39E-01	1.34	12.7493	24.1963
	1600	6.89E-01	1.84	1.64E-01	1.71	9.89E-02	1.27	14.3682	25.9460
	6400	1.77E-01	1.96	4.60E-02	1.83	3.94E-02	1.33	14.7046	26.3142
	25600	4.52E-02	1.97	1.26E-02	1.87	1.59E-02	1.31	14.7803	26.3980
	102400	1.16E-02	1.96	3.40E-03	1.89	6.34E-03	1.32	14.7985	26.4181
	409600	2.97E-03	1.96	9.10E-04	1.90	2.53E-03	1.32	14.8029	26.4231
FCT+EV	100	9.31E+00	—	1.79E+00	—	6.07E-01	—	7.41079	17.6890
	400	2.50E+00	1.90	5.40E-01	1.73	2.42E-01	1.33	12.7460	24.1918
	1600	6.98E-01	1.84	1.65E-01	1.72	9.94E-02	1.28	14.3681	25.9458
	6400	1.80E-01	1.96	4.63E-02	1.83	3.95E-02	1.33	14.7046	26.3142
	25600	4.59E-02	1.97	1.27E-02	1.87	1.58E-02	1.32	14.7803	26.3980
	102400	1.18E-02	1.96	3.42E-03	1.89	6.37E-03	1.31	14.7985	26.4181
	409600	3.10E-03	1.92	9.30E-04	1.88	2.55E-03	1.32	14.8029	26.4231
$\widetilde{\text{EV}} + \eta_2^2$	100	8.68E+00	—	1.70E+00	—	5.61E-01	—	7.94752	14.2740
	400	1.73E+00	2.33	3.75E-01	2.18	1.54E-01	1.86	13.4028	23.8887
	1600	3.26E-01	2.41	6.88E-02	2.45	3.10E-02	2.31	14.6079	26.0524
	6400	4.83E-02	2.75	1.03E-02	2.74	5.15E-03	2.59	14.7708	26.3583
	25600	8.05E-03	2.59	1.73E-03	2.57	8.82E-04	2.54	14.7977	26.4111
	102400	1.64E-03	2.30	3.47E-04	2.32	1.68E-04	2.39	14.8029	26.4217
	409600	3.82E-04	2.10	7.86E-05	2.14	3.62E-05	2.21	14.8043	26.4246
$\widetilde{\text{FCT+EV}} + \widetilde{\eta_2^2}$	100	8.81E+00	—	1.73E+00	—	5.67E-01	—	7.96202	14.3007
	400	1.92E+00	2.20	3.95E-01	2.13	1.52E-01	1.90	13.4111	23.9015
	1600	3.76E-01	2.35	7.70E-02	2.36	3.13E-02	2.28	14.6085	26.0534
	6400	6.94E-02	2.44	1.34E-02	2.52	5.42E-03	2.53	14.7709	26.3585
	25600	1.51E-02	2.20	2.78E-03	2.27	1.03E-03	2.39	14.7977	26.4111
	102400	3.52E-03	2.10	6.47E-04	2.10	2.33E-04	2.15	14.8029	26.4217
	409600	1.05E-03	1.75	2.06E-04	1.65	1.96E-04	0.25	14.8043	26.4246

Table 1: Convergence tests on smooth convection problem with CFL = 0.2. Comparison between the Galerkin, Entropy Viscosity (EV), EV with FCT, EV with normalization (4.1) and EV with FCT and normalization (4.1) schemes. Errors at time $t = 2\pi$ in the $L^1(\mathbb{T}^2)$, $L^2(\mathbb{T}^2)$, and $L^\infty(\mathbb{T}^2)$ norms are reported for different spacial resolutions along with the corresponding convergence rates. The last two columns along reports the Kinetic energy and Helicity. The reference values for the kinetic energy and helicity are $\mathcal{K}(\theta) = 14.8041$ and $\mathcal{H}(\theta) = 26.4241$.

supplemented with the initial condition $\theta(x_1, x_2, 0) = e^{-2^{\frac{1}{4}} \kappa t} \sin x_2 \cos x_1$. The latter is a scaled eigenfunction of the Laplacian. Hence, in view of the definition (2.1), the exact solution θ is given by

$$\theta(x_1, x_2, t) = e^{-\frac{t}{1000} 2^{\frac{1}{4}}} \sin x_2 \cos x_1.$$

In Table 2, we report the errors of the Galerkin method (3.24) (with $U_k^n \equiv 0$) evaluated at time $t = \pi$ in the $L^1(\mathbb{T}^2)$ -, $L^2(\mathbb{T}^2)$ - and $L^\infty(\mathbb{T}^2)$ -norms on a sequence of uniformly refined subdivisions. Two different choice of the sinc quadrature parameters are investigated, see (3.7). The time step is set to be $\Delta t = 0.1h$ Second order

convergence rates are observed in all norm when the sinc quadrature parameters are chosen to be $k = 0.8$, $M = 12$. However, the rate of convergence in the $L^\infty(\mathbb{T}^2)$ norm is reduced for the finer set of sinc quadrature parameters. Note that the convergence in the L^∞ -norm is not analyzed in [4]. From now on, the sinc quadrature parameters are set to $k = 0.8$ and $M = 12$.

	# dofs	L^1	rate	L^2	rate	L^∞	rate
$k=0.2, M=62$	100	1.28E+00	—	2.51E-01	—	8.44E-02	—
	400	3.25E-01	1.98	6.36E-02	1.98	2.14E-02	1.98
	1600	8.19E-02	1.99	1.59E-02	2.00	6.08E-03	1.82
	6400	2.07E-02	1.99	4.04E-03	1.98	2.01E-03	1.59
	25600	5.41E-03	1.94	1.08E-03	1.90	9.13E-04	1.14
	102400	1.50E-03	1.85	3.36E-04	1.69	6.97E-04	0.39
$k=0.8, M=12$	100	1.26E+00	—	2.47E-01	—	8.28E-02	—
	400	3.20E-01	1.97	6.27E-02	1.98	2.10E-02	1.98
	1600	7.97E-02	2.01	1.55E-02	2.01	5.43E-03	1.95
	6400	1.88E-02	2.08	3.67E-03	2.08	1.36E-03	2.00
	25600	3.55E-03	2.40	7.27E-04	2.34	3.40E-04	2.00
	102400	1.01E-03	1.81	2.00E-04	1.86	1.21E-04	1.49

Table 2: Effect of the sinc quadrature parameters on a smooth fractional diffusion problem. When the sinc quadrature parameters are chosen too fine, the rate of convergence of the Galerkin finite element approximation deteriorates in $L^\infty(\mathbb{T}^2)$ but not in $L^1(\mathbb{T}^2)$ nor in $L^2(\mathbb{T}^2)$.

4.3. Vortex rotation. When the geometry of the level set of the buoyancy is simple and does not contain a hyperbolic saddle, then the solution to the SQG system (2.3), (2.4) does not exhibit singularities [20] even when $\varkappa = 0$ as chosen in this section. To illustrate this, we follow [56] and consider the initial buoyancy profile

$$\theta_0(x_1, x_2) = e^{-(x_1 - \pi)^2 - 16(x_2 - \pi)^2},$$

which develops into a rotating vertex.

We set $\text{CFL} = 0.4$, see (3.20), and perform the simulations using two different space resolution corresponding to uniform triangulations \mathcal{T}_H with 351×351 and \mathcal{T}_h with 512×512 vertices. We also investigate the effect of the residual entropy viscosity parameter c_{EV} in (3.26) chosen to be either $c_{\text{EV}} = 0.1$ or $c_{\text{EV}} = 0.5$.

The buoyancy at several time snapshots is provided in Figure 1. The columns of Figure 1 correspond to four simulations: first order solution (first column) on \mathcal{T}_H , FCT solutions with $c_{\text{EV}} = 0.5$ on \mathcal{T}_H (second column), $c_{\text{EV}} = 0.1$ on \mathcal{T}_H (third column), and $c_{\text{EV}} = 0.1$ on \mathcal{T}_h (fourth column). We observe the significant improvement in accuracy of the limiting algorithm when comparing with the first order scheme. We also remark that the predictions from all the higher order finite element schemes are comparable to the spectral methods used in [56]. In particular, we see that the vortex grows thin tails which eventually generate small structures (spinning vortices). Furthermore, all simulations exhibit a discrete maximum principle property as predicted by Theorem 3.4, which does not seem to be the case for the simulations provided in [56].

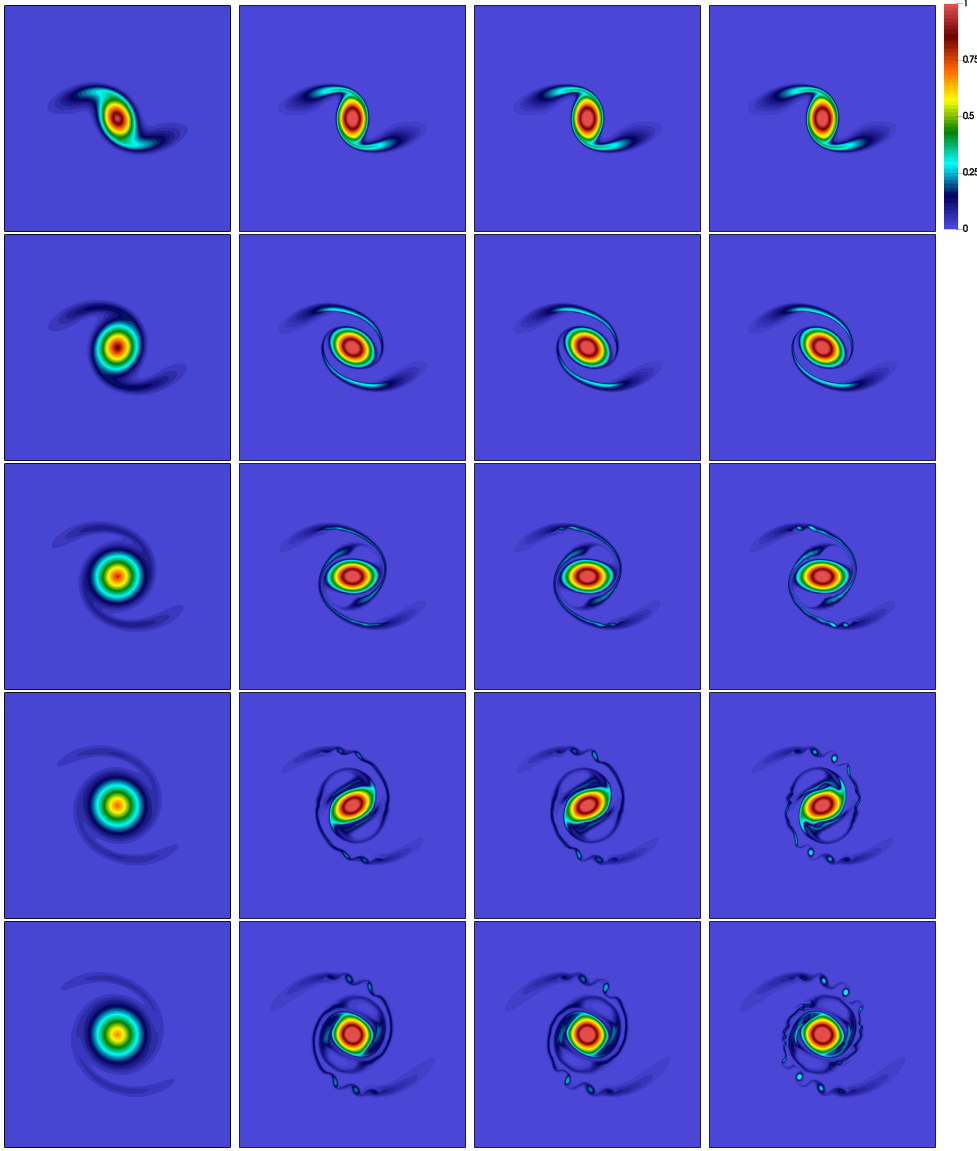


Fig. 1: Single vortex rotation at (by rows) $t = 8$, $t = 16$, $t = 26$, $t = 35$ and $t = 40$; the columns correspond to first order scheme on \mathcal{T}_H (first), and FCT solution with $c_{EV} = 0.5$ on \mathcal{T}_H (second), $c_{EV} = 0.1$ on \mathcal{T}_H (third), and $c_{EV} = 0.1$ on the finer mesh \mathcal{T}_h (fourth). All simulations satisfy a discrete maximum principle property.

In Figure 2, we report the evolution of the kinetic energy and helicity (2.5), which are conserved at the continuous level. These quantities are not conserved by the FCT scheme due to the presence of the artificial viscosity. We also report in Figure 3 (left) the evolution of $\|\nabla\Theta_k(t)\|_{L^\infty(\mathbb{T}^2)}$ to monitor apparition of singularities. The norm of the gradient of the solution oscillates when long vortex filaments develop and

eventually break down to a small scale new vortices (between $t = 10$ and $t = 40$).

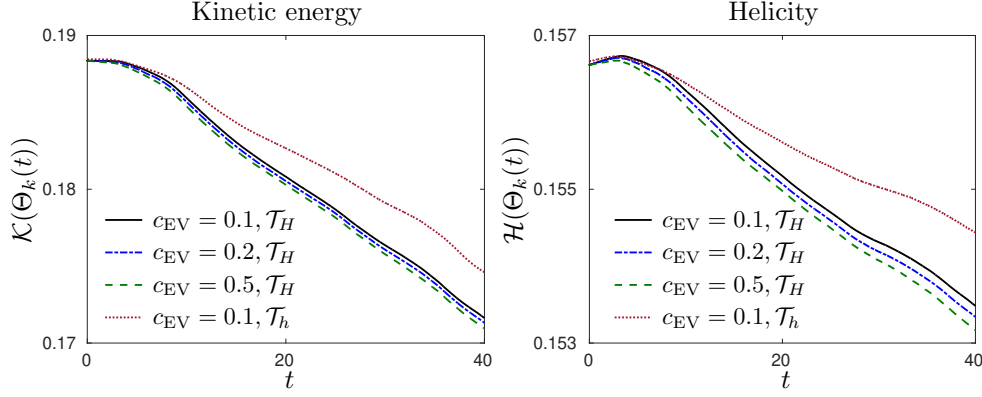


Fig. 2: Single vortex rotation: evolution of the kinetic energy and helicity for different nonlinear viscosity parameters c_{EV} and two different triangulations \mathcal{T}_H (351×351 vertices) and \mathcal{T}_h (512×512 vertices - indicated “fine” in the caption).

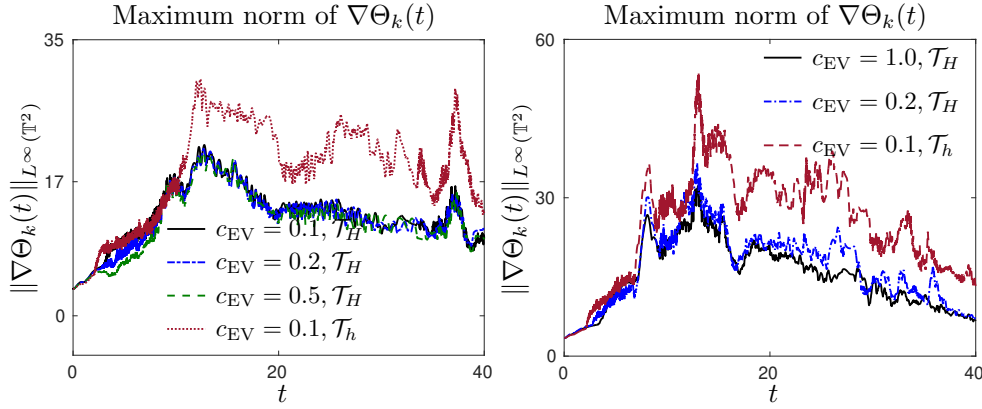


Fig. 3: Evolution of $\|\nabla\Theta_k(t)\|_{L^\infty(\mathbb{T}^2)}$ for the (left) single and (right) double vortex rotation for the different nonlinear viscosity parameters c_{EV} and two different triangulations \mathcal{T}_H (351×351 vertices) and \mathcal{T}_h (512×512 vertices).

We now investigate the interaction between two rotating vortices and consider the following initial temperature

$$\theta_0(x_1, x_2) = e^{-16(x_1 - \pi - \frac{1}{2})^2 - (x_2 - \pi)^2} + e^{-16(x_1 - \pi + \frac{1}{2})^2 - (x_2 - \pi)^2}.$$

Again, we stop the simulation at $T = 40$. As for the single vortex simulation, we select the time step with $CFL = 0.4$ on the coarse mesh \mathcal{T}_h . The entropy viscosity parameter c_{EV} is set to 0.1.

Snapshots of the buoyancy Θ_k^n are provided in Figure 4. We observe the development of a sharp layer between the two vortices around $t = 8$. Then this sharp layer reduces over time but the tip of two vortices do not merge. We also note the presence

of other small scale vortices that develop in time confirming the ability of the method to produce fine details without over-resolving. Notice that the above mentioned sharp layer does not appear to be a singularity in view of the evolution of $\|\nabla\Theta_k(t)\|_{L^\infty(\mathbb{T}^2)}$ provided in Figure 3 (right).

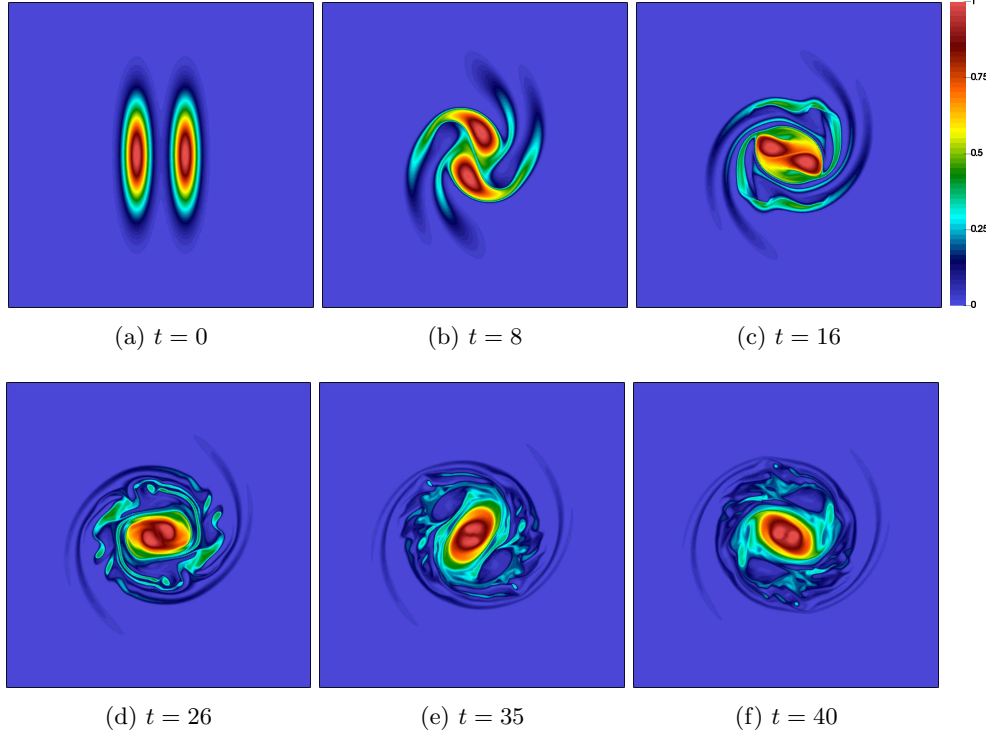


Fig. 4: Snapshots of the double vortex rotation on a 351×351 space resolution and with $\text{CFL} = 0.4$. A sharp layer develop between the two vortices around $t = 8$. Although the intensity of the layer separating the two vortices is reducing over time, the vortices do not merge.

4.4. Viscous SQG with Sharp Transitions. In this section, we consider the case $\varkappa > 0$ and solve a benchmark problem that was previously investigated in [20, 48] and [19, 56]. We set \varkappa are interested in approximating the solution to the viscous SQG system with $\varkappa = 0.001$. Recall that the Ekman pumping number \varkappa is typically small in our physical setting. We have already mentioned that singularities can develop only when a saddle structure is present in the initial buoyancy as for

$$\theta_0(x_1, x_2) = \sin x_1 \sin x_2 + \cos x_2.$$

The space discretization consists of 351×351 vertices and we chose a time step so that $\text{CFL} = 0.25$. In addition, the viscosity coefficient c_{EV} is taken to be 1.

The evolution of the potential temperature is depicted in Figure 5. The initial data contains two smooth waves that evolve in time and at $t \approx 7$ a sharp transition appears between them. A similar sharp transition develops further in other parts

of the computational domain. In addition, we also report in Figure 6 the evolution of $\|\nabla\Theta_k(t)\|_{L^\infty(\mathbb{T}^2)}$. As in the above mentioned previous works, we observe that the latter grows when sharp transitions appear and then oscillates.

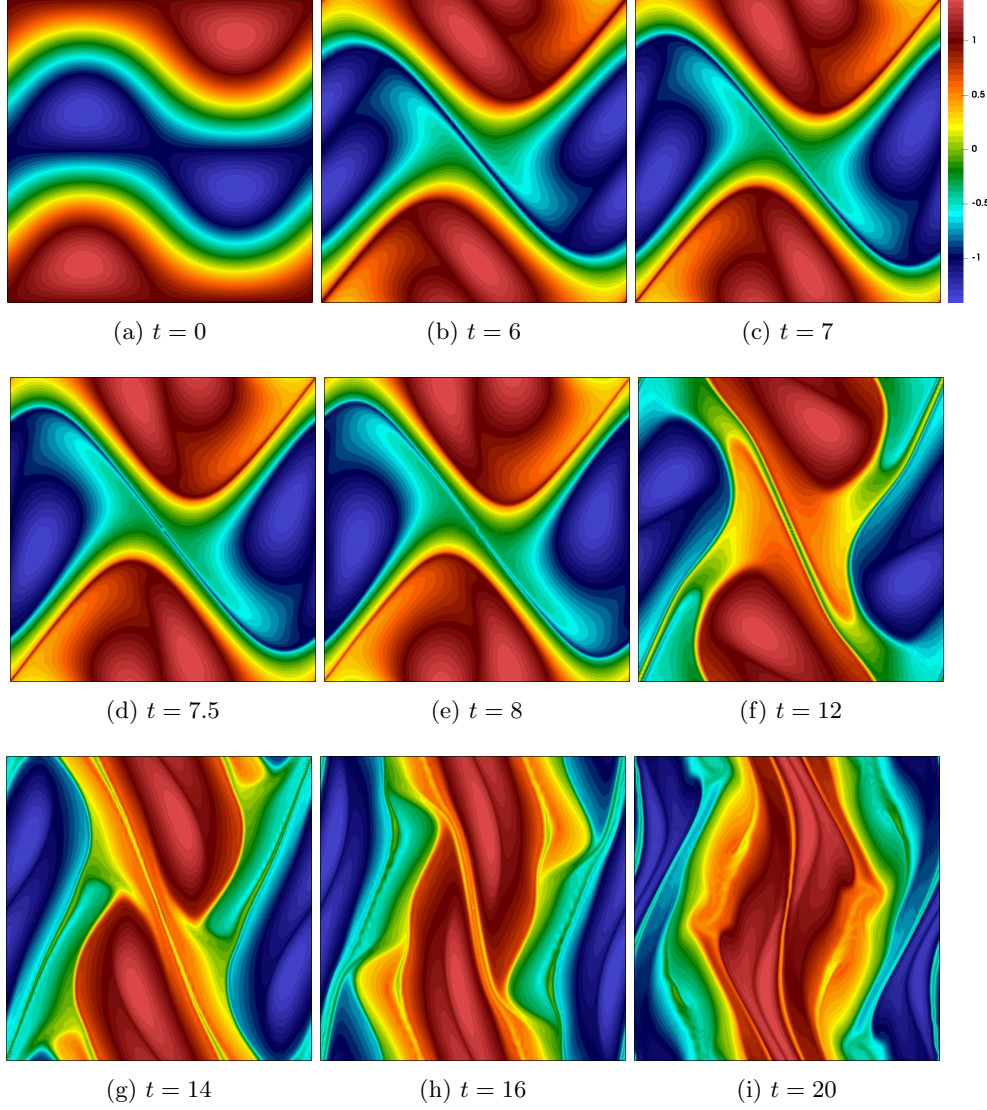


Fig. 5: Snapshots of the potential temperature for the viscous SQG. The mesh consists of 351×351 vertices, $c_{EV} = 1$ and $CFL = 0.25$. Sharp layers are developing out of the saddle configuration present in the initial data.

4.5. Freely Decaying Turbulence and Kolmogorov Energy Cascade. Experimentally [58] and numerically [43, 44] it is observed that starting from a noisy initial data to represent incoherent vortices, vortices appear and merge with other vortices with the same rotation direction to form bigger vortices. This process con-

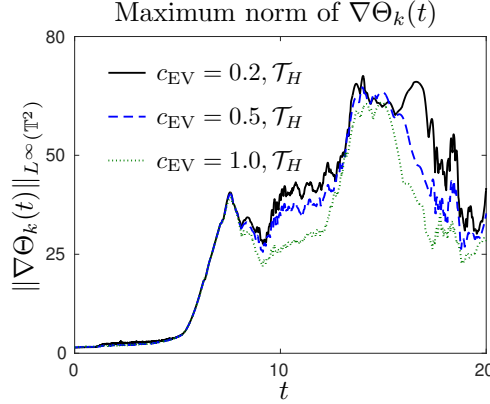


Fig. 6: Viscous SQG with sharp layer: Evolution of $\|\nabla\Theta_k(t)\|_{L^\infty(\mathbb{T}^2)}$.

tinuous until only two vertices with opposing rotating velocities are left and decay diffusively.

We perform a simulation up to time $T = 80$ of freely decaying turbulence for the SQG system using a triangulation of the domain \mathbb{T}^2 with vertices of coordinates $(2\pi n/512, 2\pi m/512)$, $n, m = 0, \dots, 512$. The time step is chosen so that the CFL number is 0.4 and $c_{EV} = 0.1$. Incoherent vortices are represented here by an initial buoyancy whose value at each vertex coordinates is randomly chosen from a uniform partition over $[-10, 10]$, see Figure 7 (a). As expected, vortices emerge and merge with other vortices with the same rotation direction to form bigger vortices as observed in Figure 7. In order to make all small scale vortices visible, we plot the solution in Figure 7 in Schlieren gray-scale diagram:

$$\sigma = \exp \left(-10 \frac{|\nabla\Theta_k^n| - \min_{\mathbb{T}^2} |\nabla\Theta_k^n|}{\max_{\mathbb{T}^2} |\nabla\Theta_k^n| - \min_{\mathbb{T}^2} |\nabla\Theta_k^n|} \right).$$

Kolmogorov energy cascade describes the energy transfer from larger scale vortices to the smaller ones. For isotropic flows like in this setting, it suffices to consider the kinetic energy $\mathcal{K}(t)$ in (2.5) and determine the energy distribution $\hat{R}(k, t)$ for $k = 1, 2, \dots$ so that

$$\mathcal{K}(t) = \sum_{k=0}^{\infty} \hat{R}(k, t).$$

We briefly describe the process and refer to [51] for additional details and the Kolmogorov assumptions. We denote by $R(y, t)$, $y \in \mathbb{R}$, $t \geq 0$, the two point correlation function in the first variable

$$R(y, t) := \frac{1}{2} \int_{\mathbb{T}^2} \theta((x_1, x_2), t) \theta((x_1 + y, x_2), t) dx_1 dx_2.$$

Note that because the fluid is assumed to be isotropic, it suffices to compute the correlation with respect to one variable (the first here) and that

$$\mathcal{K}(t) = R(0, t).$$

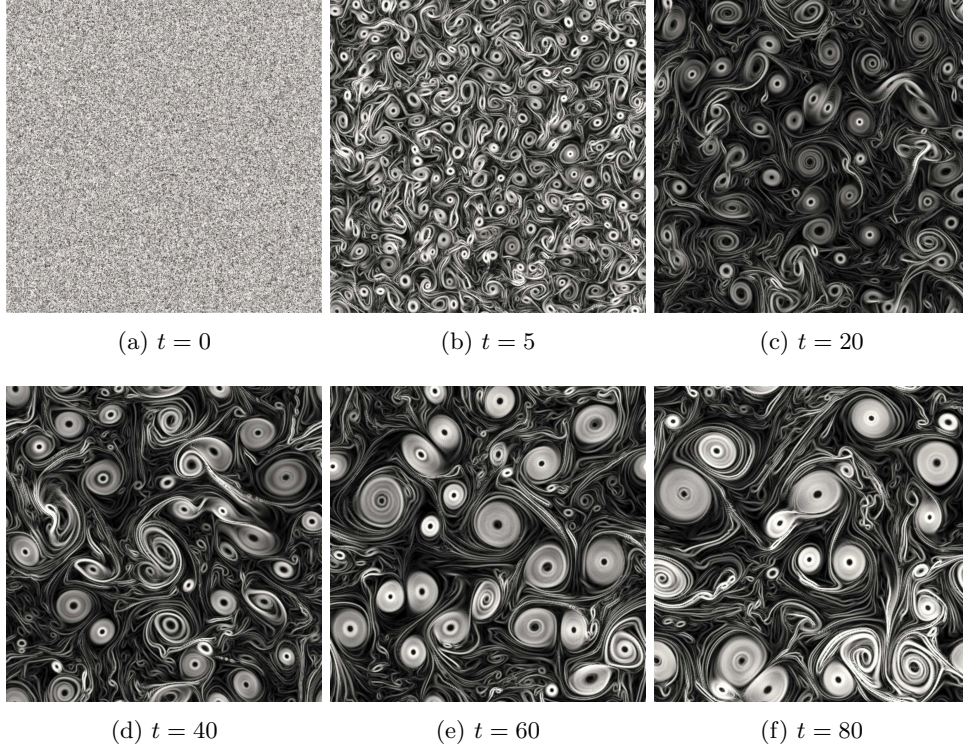


Fig. 7: Freely decaying turbulence: Schlieren diagram of the potential temperature starting from a white noise initial data.

The wavelength decomposition readily follows from the Fourier series of $R(y, t)$:

$$(4.2) \quad R(y, t) = \sum_{n=-\infty}^{\infty} \hat{R}_n(t) e^{-iny}$$

or, regrouping the terms

$$R(y, t) = \sum_{k=0}^{\infty} \sum_{|n|=k} \hat{R}_n(t) e^{-iny}.$$

Whence, we obtain the desired expression

$$\mathcal{K}(t) = \sum_{k=0}^{\infty} \hat{R}(k, t) \quad \text{with} \quad \hat{R}(k, t) := \sum_{|n|=k} \hat{R}_n(t).$$

In practice, we use discrete Fourier Transform with $N = 512$ terms to approximate the Fourier expansion (4.2)

$$\tilde{R}_m(t) = \frac{1}{N} \sum_{n=0}^{N-1} R(2\pi m/N, t) e^{-\frac{2\pi i}{N} m},$$

with $m = 0, \dots, N - 1$ so that

$$\mathcal{K}(t) \approx \frac{1}{N} \sum_{m=0}^{N-1} |\tilde{R}_m(t)|.$$

We present the energy cascade $m \mapsto |\tilde{R}_m(t)|$ for several times in the fully inviscid SQG case in Figure 8. At large scales, the theoretical prediction of the energy decay $-\frac{5}{3}$ as in full three-dimensional turbulence was obtained by [31, 50]. However inverse cascade of energy typical of two dimensional flows are predicted at small scales thereby leading to an energy decay of -3 [59, 37, 52, 12]. The left panel of Figure 8 (left) depicts the energy decay for the SQG simulation. For smaller wave numbers we observe the theoretical $-\frac{5}{3}$ rate, however, for bigger wavenumbers, the slope becomes steeper. We note that similar decays are observed in the viscous case. The energy decay when $\varkappa = 0.001$ and $s = \frac{1}{2}$ is reported in Figure 8 (middle).

For comparison, we mention that in the quasi-geostrophic (QG) system, the stream function is computed via the relation

$$(4.3) \quad (-\Delta)\psi = \theta$$

instead of (2.4). This system is widely used to mostly study 2D turbulence, see for instance [43, 44, 13] and references therein. Our numerical experiments reproduce the expected decay of approximately -5 predicted in [43, 44], see Figure 8 (right).

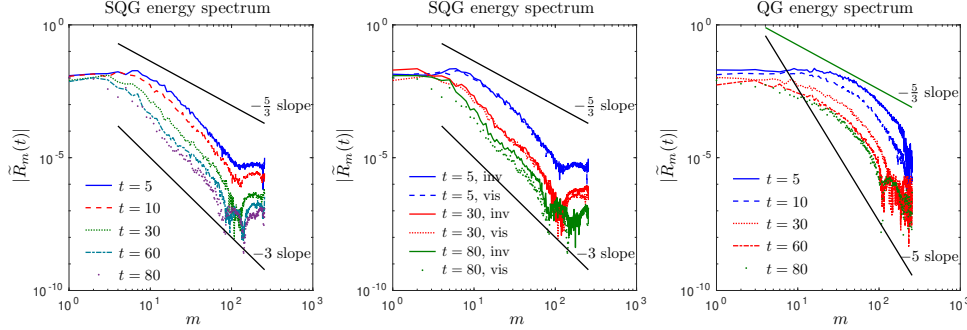


Fig. 8: 2D freely decaying turbulence: the energy spectrum $|\tilde{R}_m(t)|$ vs m , $m = 0, \dots, 256$ in the inviscid case for (left) the inviscid SQG system, (middle) the viscous SQG system and (right) the QG system. The observed decay rate $-5/3$ for small wavenumbers and -3 for large wavenumbers for the SQG systems are in accordance with the theoretical predictions. Compare with the QG system which exhibits an energy decay rate of -5 instead.

REFERENCES

- [1] A. BONITO, J. P. BORTHAGARAY, R. H. NOCHETTO, E. OTÁROLA, AND A. J. SALGADO, *Numerical methods for fractional diffusion*, Computing and Visualization in Science, 19 (2018), pp. 19–46.
- [2] A. BONITO AND J.-L. GUERMOND, *Maximum principle preserving approximation of scalar conservation equations with fractional laplacian*, in preparation.
- [3] A. BONITO, J.-L. GUERMOND, AND B. POPOV, *Stability analysis of explicit entropy viscosity methods for non-linear scalar conservation equations*, Math. Comp., 83 (2014),

- pp. 1039–1062, <https://doi.org/10.1090/S0025-5718-2013-02771-8>, <https://doi.org/10.1090/S0025-5718-2013-02771-8>.
- [4] A. BONITO, W. LEI, AND J. E. PASCIAK, *Numerical approximation of the integral fractional laplacian*, Numerische Mathematik, 142 (2019), pp. 235–278.
 - [5] A. BONITO, W. LEI, AND J. E. PASCIAK, *On sinc quadrature approximations of fractional powers of regularly accretive operators*, Journal of Numerical Mathematics, 27 (2019), pp. 57–68.
 - [6] A. BONITO AND J. PASCIAK, *Numerical approximation of fractional powers of elliptic operators*, Mathematics of Computation, 84 (2015), pp. 2083–2110.
 - [7] A. BONITO AND J. E. PASCIAK, *Numerical approximation of fractional powers of regularly accretive operators*, IMA Journal of Numerical Analysis, 37 (2017), pp. 1245–1273.
 - [8] A. BONITO AND P. WEI, *Electroconvection of thin liquid crystals: Model reduction and numerical simulations*, Journal of Computational Physics, 405 (2020), p. 109140.
 - [9] J. P. BORIS AND D. L. BOOK, *Flux-corrected transport. I. SHASTA, a fluid transport algorithm that works* [J. Comput. Phys. **11** (1973), no. 1, 38–69], J. Comput. Phys., 135 (1997), pp. 170–186. With an introduction by Steven T. Zalesak, Commemoration of the 30th anniversary of J. Comput. Phys.
 - [10] T. BUCKMASTER, S. SHKOLLER, AND V. VICOL, *Nonuniqueness of weak solutions to the sqg equation*, Communications on Pure and Applied Mathematics, 72 (2019), pp. 1809–1874.
 - [11] L. A. CAFFARELLI AND A. VASSEUR, *Drift diffusion equations with fractional diffusion and the quasi-geostrophic equation*, Annals of Mathematics, (2010), pp. 1903–1930.
 - [12] X. CAPET, P. KLEIN, B. L. HUA, G. LAPEYRE, AND J. C. MCWILLIAMS, *Surface kinetic energy transfer in surface quasi-geostrophic flows*, Journal of Fluid Mechanics, 604 (2008), pp. 165–174.
 - [13] G. F. CARNEVALE, J. C. MCWILLIAMS, Y. POMEAU, J. B. WEISS, AND W. R. YOUNG, *Evolution of vortex statistics in two-dimensional turbulence*, Phys. Rev. Lett., 66 (1991), pp. 2735–2737, <https://doi.org/10.1103/PhysRevLett.66.2735>, <https://link.aps.org/doi/10.1103/PhysRevLett.66.2735>.
 - [14] J. G. CHARNEY, *Geostrophic turbulence*, Journal of the Atmospheric Sciences, 28 (1971), pp. 1087–1095.
 - [15] J. G. CHARNEY, *On the scale of atmospheric motions*, in The Atmosphere—A Challenge, Springer, 1990, pp. 251–265.
 - [16] P. G. CIARLET, *The finite element method for elliptic problems*, vol. 40, Siam, 2002.
 - [17] P. CLÉMENT, *Approximation by finite element functions using local regularization*, Rev. Française Automat. Informat. Recherche Opérationnelle Sér., 9 (1975), pp. 77–84.
 - [18] P. CONSTANTIN, D. CORDOBA, AND J. WU, *On the critical dissipative quasi-geostrophic equation*, Indiana University mathematics journal, (2001), pp. 97–107.
 - [19] P. CONSTANTIN, M.-C. LAI, R. SHARMA, Y.-H. TSENG, AND J. WU, *New numerical results for the surface quasi-geostrophic equation*, Journal of Scientific Computing, 50 (2012), pp. 1–28.
 - [20] P. CONSTANTIN, A. J. MAJDA, AND E. TABAK, *Formation of strong fronts in the 2-d quasi-geostrophic thermal active scalar*, Nonlinearity, 7 (1994), p. 1495.
 - [21] P. CONSTANTIN, Q. NIE, AND N. SCHÖRGHOFER, *Nonsingular surface quasi-geostrophic flow*, Physics Letters A, 241 (1998), pp. 168–172.
 - [22] A. ERN AND J.-L. GUERMOND, *Theory and practice of finite elements*, vol. 159 of Applied Mathematical Sciences, Springer-Verlag, New York, 2004, <https://doi.org/10.1007/978-1-4757-4355-5>, <https://doi.org/10.1007/978-1-4757-4355-5>.
 - [23] A. E. GILL, *Atmosphere—ocean dynamics*, Elsevier, 2016.
 - [24] J.-L. GUERMOND AND M. NAZAROV, *A maximum-principle preserving C^0 finite element method for scalar conservation equations*, Comput. Methods Appl. Mech. Engrg., 272 (2014), pp. 198–213, <https://doi.org/10.1016/j.cma.2013.12.015>, <http://dx.doi.org/10.1016/j.cma.2013.12.015>.
 - [25] 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, <https://doi.org/10.1137/17M1149961>, <https://doi.org/10.1137/17M1149961>.
 - [26] J.-L. GUERMOND, M. NAZAROV, B. POPOV, AND Y. YANG, *A second-order maximum principle preserving Lagrange finite element technique for nonlinear scalar conservation equations*, SIAM J. Numer. Anal., 52 (2014), pp. 2163–2182.
 - [27] J.-L. GUERMOND, R. PASQUETI, AND B. POPOV, *Entropy viscosity method for nonlinear conservation laws*, J. Comput. Phys., 230 (2011), pp. 4248–4267.
 - [28] J.-L. GUERMOND AND B. POPOV, *Error estimates of a first-order Lagrange finite element*

- technique for nonlinear scalar conservation equations, SIAM J. Numer. Anal., 54 (2016), pp. 57–85, <https://doi.org/10.1137/140990863>, <https://doi.org/10.1137/140990863>.
- [29] 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, <https://doi.org/10.1137/16M1074291>, <https://doi.org/10.1137/16M1074291>.
- [30] 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, <https://doi.org/10.1137/16M1106560>, <https://doi.org/10.1137/16M1106560>.
- [31] I. M. HELD, R. T. PIERREHUMBERT, S. T. GARNER, AND K. L. SWANSON, *Surface quasi-geostrophic dynamics*, Journal of Fluid Mechanics, 282 (1995), pp. 1–20.
- [32] B. J. HOSKINS, *The geostrophic momentum approximation and the semi-geostrophic equations*, Journal of the Atmospheric Sciences, 32 (1976), pp. 233–242.
- [33] A. KISELEV, F. NAZAROV, AND A. VOLBERG, *Global well-posedness for the critical 2d dissipative quasi-geostrophic equation*, Inventiones mathematicae, 167 (2007), pp. 445–453.
- [34] D. KUZMIN, R. LÖHNER, AND S. TUREK, *Flux-Corrected Transport*, Scientific Computation, Springer, 2005. 3-540-23730-5.
- [35] D. KUZMIN AND S. TUREK, *Flux correction tools for finite elements*, Journal of Computational Physics, 175 (2002), pp. 525–558.
- [36] J. H. LACASCE AND A. MAHADEVAN, *Estimating subsurface horizontal and vertical velocities from sea-surface temperature*, Journal of Marine Research, 64 (2006), pp. 695–721.
- [37] G. LAPEYRE, *Surface quasi-geostrophy*, Fluids, 2 (2017), p. 7.
- [38] G. LAPEYRE AND P. KLEIN, *Dynamics of the upper oceanic layers in terms of surface quasi-geostrophy theory*, Journal of physical oceanography, 36 (2006), pp. 165–176.
- [39] C. LEITH, *Nonlinear normal mode initialization and quasi-geostrophic theory*, Journal of the Atmospheric Sciences, 37 (1980), pp. 958–968.
- [40] A. LISCHKE, G. PANG, M. GULIAN, F. SONG, C. GLUSA, X. ZHENG, Z. MAO, W. CAI, M. M. MEERSCHAERT, M. AINSWORTH, ET AL., *What is the fractional laplacian?*, arXiv preprint arXiv:1801.09767, (2018).
- [41] L. LU, M. NAZAROV, AND P. FISCHER, *Nonlinear artificial viscosity for spectral element methods*, C. R. Math. Acad. Sci. Paris, 357 (2019), pp. 646–654, <https://doi.org/10.1016/j.crma.2019.07.006>, <https://doi.org/10.1016/j.crma.2019.07.006>.
- [42] A. MAJDA, *Introduction to PDEs and Waves for the Atmosphere and Ocean*, vol. 9, American Mathematical Soc., 2003.
- [43] J. C. MCWILLIAMS, *The emergence of isolated coherent vortices in turbulent flow*, Journal of Fluid Mechanics, 146 (1984), p. 21–43, <https://doi.org/10.1017/S0022112084001750>.
- [44] J. C. MCWILLIAMS, *The vortices of two-dimensional turbulence*, Journal of Fluid Mechanics, 219 (1990), p. 361–385, <https://doi.org/10.1017/S0022112090002981>.
- [45] M. NAZAROV, *Convergence of a residual based artificial viscosity finite element method*, Computers & Mathematics with Applications, 65 (2013), pp. 616 – 626, <https://doi.org/10.1016/j.camwa.2012.11.003>, <http://www.sciencedirect.com/science/article/pii/S0898122112006499>.
- [46] M. NAZAROV AND J. HOFFMAN, *On the stability of the dual problem for high Reynolds number flow past a circular cylinder in two dimensions*, SIAM J. Sci. Comput., 34 (2012), pp. A1905–A1924, <https://doi.org/10.1137/110836213>, <https://doi.org/10.1137/110836213>.
- [47] M. NAZAROV AND A. LARCHER, *Numerical investigation of a viscous regularization of the Euler equations by entropy viscosity*, Comput. Methods Appl. Mech. Engrg., 317 (2017), pp. 128–152, <https://doi.org/10.1016/j.cma.2016.12.010>, <https://doi.org/10.1016/j.cma.2016.12.010>.
- [48] K. OHKITANI AND M. YAMADA, *Inviscid and inviscid-limit behavior of a surface quasi-geostrophic flow*, Physics of Fluids, 9 (1997), pp. 876–882.
- [49] J. PEDLOSKY, *Geophysical fluid dynamics*, Springer Science & Business Media, 2013.
- [50] R. T. PIERREHUMBERT, I. M. HELD, AND K. L. SWANSON, *Spectra of local and nonlocal two-dimensional turbulence*, Chaos, Solitons and Fractals, 4 (1994), pp. 1111–1116.
- [51] S. B. POPE, *Turbulent flows*, Cambridge University Press, Cambridge, 2000, <https://doi.org/10.1017/CBO9780511840531>, <https://doi.org/10.1017/CBO9780511840531>.
- [52] F. RAGONE AND G. BADIN, *A study of surface semi-geostrophic turbulence: freely decaying dynamics*, J. Fluid Mech., 792 (2016), pp. 740–774, <https://doi.org/10.1017/jfm.2016.116>, <https://doi.org/10.1017/jfm.2016.116>.
- [53] S. G. RESNICK, *Dynamical problems in non-linear advective partial differential equations.*, PhD thesis, University of Chicago, 1996.

- [54] L. R. SCOTT AND S. ZHANG, *Finite element interpolation of nonsmooth functions satisfying boundary conditions*, Math. Comp., 54 (1990), pp. 483–493, <https://doi.org/10.2307/2008497>, <https://doi.org/10.2307/2008497>.
- [55] C.-W. SHU AND S. OSHER, *Efficient implementation of essentially non-oscillatory shock-capturing schemes*, J. Comput. Phys., 77 (1988), pp. 439 – 471.
- [56] F. SONG AND G. E. KARNIADAKIS, *Fractional spectral vanishing viscosity method: Application to the quasi-geostrophic equation*, Chaos, Solitons & Fractals, 102 (2017), pp. 327–332.
- [57] P. R. STINGA AND J. L. TORREA, *Extension problem and harnack’s inequality for some fractional operators*, Communications in Partial Differential Equations, 35 (2010), pp. 2092–2122.
- [58] P. TABELING, S. BURKHART, O. CARDOSO, AND H. WILLAIME, *Experimental study of freely decaying two-dimensional turbulence*, Phys. Rev. Lett., 67 (1991), pp. 3772–3775, <https://doi.org/10.1103/PhysRevLett.67.3772>, <https://link.aps.org/doi/10.1103/PhysRevLett.67.3772>.
- [59] R. TULLOCH AND K. S. SMITH, *A theory for the atmospheric energy spectrum: Depth-limited temperature anomalies at the tropopause*, Proceedings of the National Academy of Sciences, 103 (2006), pp. 14690–14694, <https://doi.org/10.1073/pnas.0605494103>, <https://www.pnas.org/content/103/40/14690>, <https://arxiv.org/abs/https://www.pnas.org/content/103/40/14690.full.pdf>.
- [60] G. K. VALLIS, *Atmospheric and oceanic fluid dynamics*, Cambridge University Press, 2017.
- [61] P. WEI, *Numerical approximation of time dependent fractional diffusion with drift: applications to surface quasi-geostrophic dynamics and electroconvection*, PhD thesis, Texas A&M University, 2019.
- [62] K. YOSIDA, *Functional analysis*, Classics in Mathematics, Springer-Verlag, Berlin, 1995, <https://doi.org/10.1007/978-3-642-61859-8>, <https://doi.org/10.1007/978-3-642-61859-8>. Reprint of the sixth (1980) edition.
- [63] S. T. ZALESAK, *Fully multidimensional flux-corrected transport algorithms for fluids*, J. Comput. Phys., 31 (1979), pp. 335–362.