



Lax–Wendroff consistency of finite volume schemes for systems of non linear conservation laws: extension to staggered schemes

T. Gallouët¹ · R. Herbin¹ · J.-C. Latché²

Received: 12 March 2021 / Accepted: 30 June 2021 / Published online: 7 August 2021
© The Author(s), under exclusive licence to Sociedad Española de Matemática Aplicada 2021

Abstract

We prove in this paper the Lax–Wendroff consistency of a general finite volume convection operator acting on discrete functions which are possibly not piecewise-constant over the cells of the mesh and over the time steps. It yields an extension of the Lax–Wendroff theorem for general collocated or non-collocated schemes. This result is obtained for general polygonal or polyhedral meshes, under assumptions which, for usual practical cases, essentially boil down to a flux-consistency constraint; this latter is, up to our knowledge, novel and compares the discrete flux at a face to the mean value over the adjacent cell of the continuous flux function applied to the discrete unknown function. We first briefly show how this result copes with multipoint collocated schemes on general meshes. We then apply it to prove the consistency of a finite volume discretisation of a convection operator featuring a (convected) scalar variable and a (convecting) velocity field, with a staggered approximation, i.e. with a cell-centred approximation of the scalar variable and a face-centred approximation of the velocity.

Keywords Finite-volume schemes · Convection · Consistency

Mathematics Subject Classification Primary 65M08 · 76N15 ; Secondary 65M12 · 76N19

1 Introduction

The well-known Lax–Wendroff theorem [11] states that, on uniform 1D grids, if the approximate solutions of a flux-consistent and conservative cell-centred finite volume (FV) scheme

✉ R. Herbin
raphaele.herbin@univ-amu.fr

T. Gallouët
thierry.gallonet@univ-amu.fr

J.-C. Latché
jean-claude.latche@irsn.fr

¹ I2M UMR 7373, Aix-Marseille Université, CNRS, Ecole Centrale de Marseille, 39 rue Joliot Curie, 13453 Marseille, France

² IRSN, BP 13115, St-Paul-lez-Durance Cedex, France

for a system of conservation laws converge a.e. and boundedly as the mesh and time steps tend to zero, then the limit is a weak solution of the conservation law; we call this property “*Lax–Wendroff consistency*” (or LW-consistency for short); it is also stated in a different form [12, Section 12.10], with a BV bound assumption on the scheme. The Lax–Wendroff theorem is an “if-theorem”, which fails to solve the convergence issue for FV schemes since compactness is lacking. Nevertheless, it introduces two crucial tools for the analysis of FV schemes, namely the conservativity and consistency of the numerical fluxes; note that this analysis cannot be handled by the famous Lax–Richtmyer theorem, even in the linear case and even if the exact solution is assumed to be regular, as soon as the mesh is non uniform, see e.g. [4] for more on this subject. Moreover, the Lax–Wendroff theorem remains a useful tool to check whether a particular scheme gives a reasonable approximation when no estimates on the approximate solutions are available to yield some compactness, such as in the case of general hyperbolic systems. The Lax–Wendroff theorem was generalised to non uniform 1D or Cartesian meshes in [3, Theorem 21.2]. In a recent work [1], the Lax–Wendroff theorem is extended to obtain some error estimates for higher order schemes on uniform 1D meshes. The case of general (and, in particular, unstructured) discretisations has been also been tackled over the past decades: [10], [6, Section 4.2.2] [2], [5]. In [2], a quasi-uniformity assumption is required on the mesh, but the flux is only required to be continuous, while in [5], there is no uniformity assumption on the mesh but the flux is supposed to be locally Lipschitz continuous or at least locally “Lipschitz-diagonal” continuous, see Sect. 3 below and [5, Remark 5.2]. In all the above cited works, the scheme is supposed to be collocated, in the sense that the discrete unknowns are associated to the cells of the mesh, so these results may not be used directly to cope with staggered approximations, for instance.

The aim of this paper is to address more general approximations, including those of co-located or staggered type; indeed, we prove the LW-consistency of a generic finite volume convection operator acting on discrete functions that are possibly not piecewise-constant over the cells of the mesh and over the time steps; the result is obtained under sufficient conditions which, in usual cases, turn to essentially boil down to a new flux consistency requirement; this LW-consistency result is stated in Theorem 2.1 below. The flux consistency constraint, formulated by Assertion (11), demands a control on the difference between the discrete flux at a face (or edge) and the mean value over the adjacent cell of the continuous flux function applied to the discrete unknown function. Theorem 2.1 is valid for general polygonal or polyhedral meshes without any supplementary assumptions on the mesh; as a by product of this work, we thus also obtain a consistency result for collocated schemes (i.e. schemes using only piecewise-constant per cell unknowns) with possibly relaxed assumptions for the mesh compared to [5]. However, let us note that the proof that the assumption (11) is satisfied is usually based on the control of the difference between the numerical solution and its space or time translates, see [5, Section 4] and that these latter results may require some regularity assumptions on the mesh, see also Remark 2.2.

This paper is organized as follows. We state and prove the general consistency result in Sect. 2. We then apply it in Sect. 3 to the collocated case and then, in Sect. 4, to a staggered discretisation; precisely speaking, we show the consistency of a finite volume discretisation of a nonlinear convection operator for a scalar variable ρ of the form $\partial_t \beta(\rho) + \operatorname{div}(g(\rho)\mathbf{u})$, where β and g are regular functions and \mathbf{u} is a velocity field, and where we use a cell-centred approximation for ρ and a face-centred approximation of \mathbf{u} .

2 The general LW-consistency result

The aim is to prove the LW-consistency of finite volume approximations of nonlinear convective terms which appear in most models of fluid flow. The general context is the following. Given a numerical scheme which yields some approximate solutions to the system of conservative partial differential equations, we assume that these approximate solutions converge to some functions strongly in L^1 , and we wish to show that the limit is indeed a solution to the system, at least in a weak sense. In order to do so, the usual idea is to multiply the numerical scheme by an interpolate of a smooth function, sum over the cells of the mesh and over the time steps and show that passing to the limit, we get a weak formulation of the system of partial differential equations. The theorem that we prove below is a mean to prove that one may indeed pass to the limit in the terms that involve nonlinear convection operators. Let us begin with an example. Consider the barotropic Euler equations, which read:

$$\partial_t \bar{\rho} + \operatorname{div}(\bar{\rho} \bar{\mathbf{u}}) = 0, \tag{1a}$$

$$\partial_t(\bar{\rho} \bar{\mathbf{u}}) + \operatorname{div}(\bar{\rho} \bar{\mathbf{u}} \otimes \bar{\mathbf{u}}) + \nabla \bar{p} = 0, \tag{1b}$$

where $\bar{\rho}$ is the density, $\bar{\mathbf{u}}$ the velocity and \bar{p} the pressure, which, for barotropic flows, is a function of $\bar{\rho}$ only: $\bar{p} = \mathfrak{p}(\bar{\rho})$. Here and in the remainder of the paper, we use overlined letters when referring to the solution of the continuous problem, while non overlined letters will be used for discrete unknowns. This system of equations is supplemented by an initial condition and suitable boundary conditions if Ω is bounded. An entropy weak solution of the system satisfies the Eq. (1) and also satisfies (in a weak sense, which includes the initial condition) the following entropy condition:

$$\partial_t \bar{E} + \operatorname{div}((\bar{E} + \bar{p}) \bar{\mathbf{u}}) \leq 0, \text{ with } \bar{E} = \frac{1}{2} \bar{\rho} |\bar{\mathbf{u}}|^2 + \mathcal{H}(\bar{\rho}) \text{ and } \mathcal{H}(s) = s \int \frac{\mathfrak{p}(s)}{s^2} ds. \tag{2}$$

The weak consistency of staggered finite volume schemes for this system of equations discretised on multi-dimensional Cartesian or unstructured meshes has been the object of several recent papers, see e.g. [8,9]. The system (1) and (2) may be written as

$$\bar{\mathcal{C}}_1(\bar{\rho}, \bar{\mathbf{u}}) = 0, \tag{3a}$$

$$\bar{\mathcal{C}}_2(\bar{\rho}, \bar{\mathbf{u}}) + \nabla \bar{p} = 0, \tag{3b}$$

$$\bar{\mathcal{C}}_3(\bar{E}, \bar{\mathbf{u}}) + \operatorname{div}(\bar{p} \bar{\mathbf{u}}) \leq 0, \tag{3c}$$

with $\bar{\mathcal{C}}_1(\bar{\rho}, \bar{\mathbf{u}}) = \partial_t \bar{\rho} + \operatorname{div}(\bar{\rho} \bar{\mathbf{u}})$, $\bar{\mathcal{C}}_2(\bar{\rho}, \bar{\mathbf{u}}) = \partial_t(\bar{\rho} \bar{\mathbf{u}}) + \operatorname{div}(\bar{\rho} \bar{\mathbf{u}} \otimes \bar{\mathbf{u}})$, and $\bar{\mathcal{C}}_3(\bar{E}, \bar{\mathbf{u}}) = \partial_t \bar{E} + \operatorname{div}(\bar{E} \bar{\mathbf{u}})$. In the above cited works, the system is discretised with an explicit or implicit in time scheme, and the convection operators \mathcal{C}_1 and \mathcal{C}_2 by a first or second order finite volume scheme. In fact, the system of the barotropic equations can be discretised by different schemes: explicit or implicit, colocated meshes or staggered meshes, using a Riemann solver or using an equation-by-equation procedure. In all cases, the consistency study will have to deal with each of the discrete non linear convection operator \mathcal{C}_i associated to $\bar{\mathcal{C}}_i$. The present work aims at simplifying the proofs of consistency by giving a general result for any nonlinear convection term, discretised on a colocated or staggered mesh, thereby extending our previous result of [5] to staggered meshes. Theorem 2.1 below is an efficient tool to this purpose; it may be used for any of the terms in (3), and is specifically useful to tackle the terms featuring discrete variables with different space approximations, as the operators \mathcal{C}_i of these equations in case of staggered discretizations. We emphasize that both implicit or explicit schemes may be addressed, since the proof deals separately with the discrete time operator and the discrete space divergence operator.

Let us then turn to the general setting; we suppose that:

$$\Omega \subset \mathbb{R}^d, \quad d = 1, 2, 3, \quad T \in (0, +\infty), \quad p \in \mathbb{N}^*, \quad \beta \in C^0(\mathbb{R}^p, \mathbb{R}), \quad \mathbf{f} \in C^0(\mathbb{R}^p, \mathbb{R}^d). \quad (4)$$

We consider the conservative convection operator $\bar{\mathcal{C}}(\bar{U})$ acting on a vector $\bar{U} \in \mathbb{R}^p$ of functions, real-valued, and defined (in the distributional sense), for $\bar{U} \in L^\infty(\Omega \times (0, T), \mathbb{R}^p)$, by:

$$\begin{aligned} \bar{\mathcal{C}}(\bar{U}) : \quad & \Omega \times (0, T) \rightarrow \mathbb{R}, \\ (\mathbf{x}, t) \mapsto & \partial_t(\beta(\bar{U}))(\mathbf{x}, t) + \operatorname{div}(\mathbf{f}(\bar{U}))(\mathbf{x}, t). \end{aligned} \quad (5)$$

Note that, here and throughout the paper, we use $\beta(\bar{U})$ (resp. $\mathbf{f}(\bar{U})$) to denote the function $\beta \circ \bar{U}$ obtained by composition of β and \bar{U} (resp. \mathbf{f} and \bar{U}), so, for instance, $\beta(\bar{U})(\mathbf{x}, t)$ stands for $\beta(\bar{U}(\mathbf{x}, t))$. In the above example of the barotropic Euler equations (1), we have, for $i = 1, 2$, $\bar{\mathcal{C}}_i(\bar{U}) = \partial_t(\beta_i(\bar{U})) + \operatorname{div}(\mathbf{f}_i(\bar{U}))$, with $\bar{U} = (\bar{\rho}, \bar{\mathbf{u}})$, $\beta_1(\bar{U}) = \bar{\rho}$, $\mathbf{f}_1(\bar{U}) = \bar{\rho}\bar{\mathbf{u}}$, $\beta_2(\bar{U}) = \bar{\rho}\bar{\mathbf{u}}$, and $\mathbf{f}_2(\bar{U}) = \bar{\rho}\bar{\mathbf{u}} \otimes \bar{\mathbf{u}}$ (in fact, to match precisely the formalism of Eq. (4), these last two functions have to be considered as d functions, one for each velocity component, associated to d convection operators, which has no consequence for the matter at hand).

Let us denote by \mathcal{P} a mesh of the domain Ω , consisting of a set of disjoint open polyhedral or polygonal subsets of Ω , called cells, whose union of closures is $\bar{\Omega}$ (Fig. 1). To avoid cumbersome notations, we assume that any pair of adjacent cells shares a whole face (in 3D) or edge (in 2D), and not only a part of it; however this assumption is not necessary for the result of Theorem 2.1 to hold. Throughout the paper and when the space dimension is not specified, we use "face" to define the interface between two cells; for a face ζ , $|\zeta|$ stands for its $(d - 1)$ -dimensional measure in 2D and 3D, and we set $|\zeta| = 1$ by convention in one space dimension. The notation $|P|$ stands for the d -dimensional measure of a cell P . We denote by $\delta(\mathcal{P})$ the space step, defined by

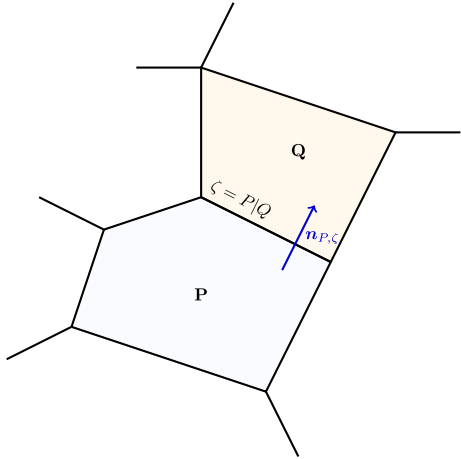
$$\delta(\mathcal{P}) = \max_{P \in \mathcal{P}} \operatorname{diam}(P).$$

Let \mathfrak{F} denote the set of faces of the mesh, and $\mathfrak{F}_{\text{int}}$ denote the set of faces that are not located on the boundary $\partial\Omega$; for a given cell $P \in \mathcal{P}$, let $\mathfrak{F}(P)$ be the set of faces of P . Let $t_0 = 0 < t_1 < \dots < t_N = T$ be a partition of $(0, T)$, denoted by \mathcal{T} ; for such a partition \mathcal{T} , we define the time step by $\delta t = \max\{t_{n+1} - t_n, n \in \llbracket 0, N - 1 \rrbracket\}$, where $\llbracket 0, N - 1 \rrbracket$ denotes the set of integers n such that $0 \leq n \leq N - 1$.

The unknown is supposed to be represented by a function $U \in L^\infty(\Omega \times (0, T), \mathbb{R}^p)$. For a collocated FV scheme, it is the piecewise constant function defined by $U(x, t) = U_K^n$ for $x \in K$ and $t \in]t_n, t_{n+1}[$. For a FV staggered scheme for a system of equations, each component of U is piecewise constant on each associated mesh. But U could also be a piecewise affine function, for instance if say, a DG scheme is used. We emphasize that for non collocated schemes, some unknowns are not piecewise-constant over the cells of the mesh and over the time steps. For instance, when using staggered discretisations in fluid flow simulations, if P is a primal cell, the velocity is possibly discontinuous along surfaces or lines included in P (see the example developed in Sect. 4). The discrete convection operator that we consider here takes the following form:

$$\begin{aligned} \mathcal{C}(U) : \quad & \Omega \times (0, T) \rightarrow \mathbb{R}, \\ (\mathbf{x}, t) \mapsto & \mathcal{C}(U)_p^n, \quad \text{for } \mathbf{x} \in P, \quad P \in \mathcal{P}, \quad \text{and } t \in (t_n, t_{n+1}), \quad n \in \llbracket 0, N - 1 \rrbracket, \end{aligned}$$

Fig. 1 An example of a two-dimensional mesh and associated notations: P and Q are two generic cells, $P, Q \in \mathcal{P}$, with \mathcal{P} the set of cells, $\zeta = P|Q$ is the face separating P and Q , $\zeta \in \mathfrak{F}$, with \mathfrak{F} the set of faces and $\mathbf{n}_{P,\zeta}$ is the normal to ζ pointing outward P



with

$$\mathcal{C}(U)_P^n = (\partial_t \beta)_P^n + \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \mathbf{F}_\zeta^n \cdot \mathbf{n}_{P,\zeta},$$

where $\{\beta_P^n, P \in \mathcal{P}, n \in \llbracket 0, N \rrbracket\}$ is a family of real numbers,

$$(\partial_t \beta)_P^n = \frac{\beta_P^{n+1} - \beta_P^n}{t_{n+1} - t_n}, \quad n \in \llbracket 0, N - 1 \rrbracket, \tag{6}$$

$\{\mathbf{F}_\zeta^n, \zeta \in \mathfrak{F}, n \in \llbracket 0, N - 1 \rrbracket\}$ is a family of real vectors of \mathbb{R}^d and $\mathbf{n}_{P,\zeta}$ stands for the normal vector to ζ pointing outward P . Note that this form of the flux implies that the scheme is conservative. Of course, the real numbers $\{\beta_P^n, P \in \mathcal{P}, n \in \llbracket 0, N \rrbracket\}$ and $\{\mathbf{F}_\zeta^n, \zeta \in \mathfrak{F}, n \in \llbracket 0, N - 1 \rrbracket\}$ are related to the unknown U ; it is the object of Theorem 2.1 below to state precisely the assumptions that must be satisfied by these quantities to ensure the consistency of the discrete convection operator.

Theorem 2.1 (LW-consistency for a multi-dimensional conservative convection operator)

Under the assumptions (4), let $(\mathcal{P}^{(m)}, \mathcal{T}^{(m)})_{m \in \mathbb{N}}$ be a sequence of possibly non uniform space-time discretisations, with $\delta(\mathcal{P}^{(m)})$ and $\delta t^{(m)}$ tending to zero as $m \rightarrow +\infty$, and let $(U^{(m)})_{m \in \mathbb{N}}$ be the associated sequence of discrete functions. We suppose that the sequence $(U^{(m)})_{m \in \mathbb{N}}$ is bounded and converges to a limit:

$$\exists C^u \in \mathbb{R}_+^* \text{ s.t. } \|U^{(m)}\|_\infty \leq C^u, \quad \forall m \in \mathbb{N}, \tag{7}$$

$$\exists \bar{U} \in L^\infty(\Omega \times (0, T), \mathbb{R}^P) \text{ s.t. } \|U^{(m)} - \bar{U}\|_{L^1(\Omega \times (0, T), \mathbb{R}^P)} \rightarrow 0 \text{ as } m \rightarrow +\infty. \tag{8}$$

We also assume that the family $\{(\beta^{(m)})_P^n, P \in \mathcal{P}^{(m)}, n \in \llbracket 0, N^{(m)} \rrbracket, m \in \mathbb{N}\}$ is bounded. In addition, let $U_0 \in L^\infty(\Omega, \mathbb{R}^P)$ and let us suppose that, as $m \rightarrow +\infty$,

$$\sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_P ((\beta^{(m)})_P^0 - \beta(U_0)(\mathbf{x})) \varphi(\mathbf{x}) \, d\mathbf{x} \rightarrow 0, \text{ for any } \varphi \in C_c^\infty(\Omega), \tag{9}$$

$$\sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P ((\beta^{(m)})_P^n - \beta(U^{(m)})(\mathbf{x}, t)) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \rightarrow 0,$$

$$\text{for any } \varphi \in C_c^\infty(\Omega \times [0, T]), \tag{10}$$

$$\sum_{n=0}^{N^{(m)}-1} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \frac{\text{diam}(P)}{|P|} \sum_{\zeta \in \mathcal{F}(P)} |\zeta| \int_{t_n}^{t_{n+1}} \int_P \left| \left((\mathbf{F}^{(m)})_{\zeta}^n - \mathbf{f}(U^m)(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| \, d\mathbf{x} \, dt \rightarrow 0, \tag{11}$$

where $\mathcal{P}_{\text{int}}^{(m)}$ denotes the set of cells of $\mathcal{P}^{(m)}$ that have no face on the boundary $\partial\Omega$. Then, for any $\varphi \in C_c^\infty(\Omega \times [0, T])$,

$$\begin{aligned} \int_0^T \int_{\Omega} \mathcal{C}^{(m)}(U^{(m)}) \mathcal{J}^{(m)}(\varphi)(\mathbf{x}, t) \, d\mathbf{x} \, dt &\rightarrow - \int_{\Omega} \beta(U_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} \\ &- \int_0^T \int_{\Omega} \left(\beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right) \, d\mathbf{x} \, dt \quad \text{as } m \rightarrow +\infty \end{aligned} \tag{12}$$

where $\mathcal{J}^{(m)}(\varphi)$ is an interpolate of φ defined a.e. by

$$\begin{aligned} \mathcal{J}^{(m)}(\varphi)(\mathbf{x}, t) &= \varphi_P^n \text{ for } \mathbf{x} \in P \text{ and } t \in (t_n, t_{n+1}), \\ \text{with } \varphi_P^n &= \frac{1}{|P|} \int_P \varphi(\mathbf{x}, t_n) \, d\mathbf{x}, \quad \text{for } P \in \mathcal{P} \text{ and } n \in \llbracket 0, N-1 \rrbracket. \end{aligned} \tag{13}$$

Before we give the proof of Theorem 2.1, let us first briefly comment on its assumptions.

Remark 2.2 (Flux consistency) The required flux consistency is stated by Assertion (11), which requires the flux $(\mathbf{F}^{(m)})_{\zeta}^n$ through a face ζ of a cell P to be close to the mean value over P of the actual flux function \mathbf{f} applied to the unknown. For a scheme involving only cell unknowns, for instance, the quantity $(\mathbf{F}^{(m)})_{\zeta}^n$ is generally a function of the unknowns in the cell P and in the neighbouring cells, and checking the assumption (11) amounts to bound the difference between the unknowns and their translates. Note that, while Theorem 2.1 holds for very general meshes, as we have already mentioned in the introduction, some regularity assumptions on the sequence of meshes may be required at this step.

To clarify this point, let us consider a simple one-dimensional problem for the scalar unknown u , with $\beta(u) = f(u) = u$, leading to the linear convection operator $\mathcal{C}(u) = \partial_t u + \partial_x u$, which we discretise with the first-order explicit-in-time upwind scheme. Let us suppose that the discrete functions are defined by $u(x, t) = u_P^n$ for $x \in P$ and $t \in (t_n, t_{n+1})$. Then, for $\mathbf{x} \in P$ and $t \in (t_n, t_{n+1})$, $\left| \left((\mathbf{F}^{(m)})_{\zeta}^n - \mathbf{f}(U^m)(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| = \left| (u^{(m)})_{P^-}^n - (u^{(m)})_P^n \right|$ where P^- is the left cell to P when ζ is its left face, and $\left| \left((\mathbf{F}^{(m)})_{\zeta}^n - \mathbf{f}(U^m)(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| = 0$ otherwise (disregarding the boundary cells, according to the formulation of the theorem). Checking Assumption (11) thus consists in proving that the term $R^{(m)}$ defined by

$$R^{(m)} = \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{P \in \mathcal{P}^{(m)}} \text{diam}(P) |u_P^n - u_{P^-}^n|$$

tends to zero as m tends to $+\infty$. This is implied by the convergence in $L^1(\Omega \times (0, T))$ of the sequence of discrete solutions provided that the ratio $|P|/|P^-|$ is bounded independently of m for the sequence of meshes under consideration [5, Section 4]. More elaborate example of application, using collocated then staggered meshes, are provided below, in Sects. 3 and 4 respectively.

Remark 2.3 (Disregarding boundary cells in Assumption (11)) Since the support of the test function φ is compact in $\Omega \times [0, T]$, for $\delta(\mathcal{P}^{(m)})$ small enough, φ vanishes in the boundary cells. Consequently, it is clear from the proof of the theorem below (see the expression (15)

of the term $X_2^{(m)}$) that boundary cells may be excluded in the sum in Assertion (11). This is the reason why only the cells in $\mathcal{P}_{\text{int}}^{(m)}$ are considered in Assumption (11). For numerical fluxes involving wider stencils (for instance in the case of higher order schemes), one could in fact reduce the set of involved cells furthermore.

Remark 2.4 (Regularity of β and f) The proof of Theorem 2.1 holds if β and f are only continuous functions, which is the assumption made in the present section; however, to prove Assertions (10) and (11), a locally Lipschitz-diagonal continuity (see Definition 3.1 below) is often required, as in Sects. 3 and 4.

Remark 2.5 (Stronger convergence assumptions on $\{(\beta^{(m)})_{m \in \mathbb{N}}\}$) In most situations, stronger convergence properties hold for $(\beta^{(m)})_{m \in \mathbb{N}}$, in the sense that the LW-convergence assumptions (9) and (10) are implied by the following strong convergence assumptions:

$$\sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_P |(\beta^{(m)})_P^0 - \beta(U_0(\mathbf{x}))| \, d\mathbf{x} \rightarrow 0 \text{ as } m \rightarrow +\infty,$$

$$\sum_{n=0}^{N_m-1} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_n}^{t_{n+1}} \int_P |(\beta^{(m)})_P^n - \beta(U^{(m)}(\mathbf{x}, t))| \, d\mathbf{x} \, dt \rightarrow 0 \text{ as } m \rightarrow +\infty.$$

This is the case, for instance, for the convection operators considered in Sects. 3 and 4 below. However, there are cases where the convergence of $\beta^{(m)}$ is only weak, see for instance the reconstructed kinetic energy for the full compressible Euler equations in [8].

Remark 2.6 (On the interpolate of the test function) Note that in the definition (13) of $\mathcal{J}^{(m)}(\varphi)$ in (12), the quantities φ_P^n , $n \in \llbracket 0, N \rrbracket$, may be also defined as

$$\varphi_P^n = \frac{1}{|P|} \int_P \varphi(\mathbf{x}, t_{n+1}) \, d\mathbf{x},$$

with minor changes in the arguments of the present section, essentially a slightly different assumption (10), which reads:

$$\sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P \left((\beta^{(m)})_P^{n-1} - \beta(U^{(m)})(\mathbf{x}, t) \right) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \rightarrow 0,$$

$$\text{for any } \varphi \in C_c^\infty(\Omega \times [0, T]).$$

For instance, for a scalar problem, if the discrete function is defined as $u(\mathbf{x}, t) = u_P^{n-1}$ for $\mathbf{x} \in P$ and $t \in [t_{n-1}, t_n)$ (choice often used in explicit schemes) and β_P^{n-1} is defined in the scheme as $\beta(u_P^{n-1})$, this assumption is trivially satisfied, since $(\beta^{(m)})_P^{n-1} = \beta(U^{(m)})(\mathbf{x}, t)$ in $P \times (t_{n-1}, t_n)$, while checking the original assumption (10) needs to bound the time translates of the discrete solution. This is however an easy task, under a very mild regularity assumption for the time discretisation (see Sect. 4 below). The opposite situation occurs (i.e. this is Assumption (10) which is now trivially satisfied) if the discrete function is defined as $u(\mathbf{x}, t) = u_P^n$ for $\mathbf{x} \in P$ and $t \in [t_{n-1}, t_n)$, which is often done for implicit schemes.

Proof of Theorem 2.1 Theorem 2.1 is the consequence of the two following lemmas, which prove respectively the convergence of the time derivative part and the space derivative part. Let us decompose

$$\int_0^T \int_{\Omega} \mathcal{C}^{(m)}(U^{(m)}) \mathcal{J}^{(m)}(\varphi)(\mathbf{x}, t) \, d\mathbf{x} \, dt = X_1^{(m)} + X_2^{(m)}, \text{ with}$$

$$X_1^{(m)} = \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{P \in \mathcal{P}^{(m)}} |P| (\bar{\partial}_t \beta^{(m)})_P^n \varphi_P^n, \tag{14}$$

$$X_2^{(m)} = \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{P \in \mathcal{P}^{(m)}} \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| (\mathbf{F}^{(m)})_{\zeta}^n \cdot \mathbf{n}_{P,\zeta} \varphi_P^n. \tag{15}$$

Then, by Lemma 2.7 below,

$$X_1^{(m)} \rightarrow - \int_{\Omega} \beta(U_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \int_0^T \int_{\Omega} \beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty,$$

and by Lemma 2.8 below,

$$X_2^{(m)} \rightarrow - \int_0^T \int_{\Omega} \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty,$$

which concludes the proof. □

Lemma 2.7 (LW-consistency, time derivative) *Let the sequence $(X_1^{(m)})_{m \in \mathbb{N}}$ be defined by (14). Then, under the assumptions and notations of Theorem 2.1,*

$$X_1^{(m)} \rightarrow - \int_{\Omega} \beta(U_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \int_0^T \int_{\Omega} \beta(\bar{U})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty.$$

Proof By the definition (6) of $\bar{\partial}_t^n \beta^{(m)}(\mathbf{x}, t)$ and thanks to a discrete integration by parts, we get that

$$X_1^{(m)} = - \sum_{P \in \mathcal{P}^{(m)}} |P| (\beta^{(m)})_P^0 \varphi_P^0 - \sum_{n=1}^{N^{(m)}} (t_n - t_{n-1}) \sum_{P \in \mathcal{P}^{(m)}} |P| (\beta^{(m)})_P^n \frac{\varphi_P^n - \varphi_P^{n-1}}{t_n - t_{n-1}}.$$

Let us write the first term of the right-hand side as

$$\begin{aligned} - \sum_{P \in \mathcal{P}^{(m)}} |P| (\beta^{(m)})_P^0 \varphi_P^0 &= - \sum_{P \in \mathcal{P}^{(m)}} \int_P (\beta^{(m)})_P^0 (\varphi_P^0 - \varphi(\mathbf{x}, 0)) \, d\mathbf{x} \\ &\quad - \sum_{P \in \mathcal{P}^{(m)}} \int_P (\beta^{(m)})_P^0 \varphi(\mathbf{x}, 0) \, d\mathbf{x}. \end{aligned}$$

On the one hand, the piecewise-constant function equal to φ_P^0 on each cell $P \in \mathcal{P}^{(m)}$ converges to $\varphi(\mathbf{x}, 0)$ in $L^\infty(\Omega)$ as m tends to $+\infty$, and $(\beta^{(m)})^0$ is supposed to be bounded; the first integral at the right-hand side thus tends to zero. Hence, invoking Assumption (9) for the second integral,

$$- \sum_{P \in \mathcal{P}^{(m)}} |P| (\beta^{(m)})_P^0 \varphi_P^0 \rightarrow - \int_{\Omega} \beta(U_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} \text{ as } m \rightarrow +\infty.$$

Let the piecewise constant function $\bar{\partial}_t^{(m)} \varphi : \Omega \times (0, T) \rightarrow \mathbb{R}^d$ be defined by

$$\bar{\partial}_t^{(m)} \varphi(\mathbf{x}, t) = \frac{\varphi_P^{n+1} - \varphi_P^n}{t_{n+1} - t_n} \text{ for } (\mathbf{x}, t) \in P \times (t_n, t_{n+1}).$$

The function $\tilde{\partial}_t^{(m)}\varphi$ converges uniformly to $\partial_t\varphi$ in $L^\infty(\Omega \times (0, T))$. The second term of $X_1^{(m)}$ may be decomposed as

$$-\sum_{n=1}^{N^{(m)}}(t_n - t_{n-1}) \sum_{P \in \mathcal{P}^{(m)}} |P| (\beta^{(m)})_P^n \frac{\varphi_P^n - \varphi_P^{n-1}}{t_n - t_{n-1}} = Y_1^{(m)} + Y_2^{(m)}$$

with

$$Y_1^{(m)} = -\sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P (\beta^{(m)})_P^n (\tilde{\partial}_t^{(m)}\varphi(\mathbf{x}, t) - \partial_t\varphi(\mathbf{x}, t)) \, \mathbf{x} \, dt,$$

$$Y_2^{(m)} = -\sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P (\beta^{(m)})_P^n \partial_t\varphi(\mathbf{x}, t) \, \mathbf{x} \, dt.$$

Since the family $\{(\beta^{(m)})_P^n, P \in \mathcal{P}^{(m)}, n \in \llbracket 0, N^{(m)} \rrbracket, m \in \mathbb{N}\}$ is assumed to be bounded, the uniform convergence of $\tilde{\partial}_t^{(m)}\varphi$ to $\partial_t\varphi$ yields that the sequence $(Y_1^{(m)})_{m \in \mathbb{N}}$ tends to zero. Invoking the assumption (10), the continuity of β and the convergence of $(U^{(m)})_{m \in \mathbb{N}}$ to \bar{U} , we get that

$$\lim_{m \rightarrow +\infty} X_1^{(m)} = \lim_{m \rightarrow +\infty} Y_2^{(m)} = -\int_0^T \int_\Omega \beta(\bar{U})(\mathbf{x}, t) \partial_t\varphi(\mathbf{x}, t) \, \mathbf{x} \, dt.$$

□

Lemma 2.8 (LW-consistency, space derivative) *Let the sequence $(X_2^{(m)})_{m \in \mathbb{N}}$ be defined by (15). Then, under the assumptions and notations of Theorem 2.1,*

$$X_2^{(m)} \rightarrow -\int_0^T \int_\Omega \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla\varphi(\mathbf{x}, t) \, \mathbf{x} \, dt \quad \text{as } m \rightarrow +\infty.$$

Proof Since φ is compactly supported and since $\delta(\mathcal{P}^{(m)}) \rightarrow 0$ as $m \rightarrow 0$, there exists $M \in \mathbb{N}$ such that for $m \geq M$, $\varphi_P^n = 0$ for all $\mathbf{x} \in \mathcal{P}^{(m)} \setminus \mathcal{P}_{\text{int}}^{(m)}$. Moreover, since for a face ζ separating P and P' , one has $\mathbf{n}_{P,\zeta} = -\mathbf{n}_{P',\zeta}$, we get that

$$\begin{aligned} X_2^{(m)} &= \sum_{n=0}^{N^{(m)}-1} (t_n - t_{n-1}) \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| (\mathbf{F}^{(m)})_\zeta^n \cdot \mathbf{n}_{P,\zeta} \varphi_P^n \\ &= \sum_{n=0}^{N^{(m)}-1} (t_n - t_{n-1}) \sum_{P_{\text{int}} \in \mathcal{P}^{(m)}} A_P^n \end{aligned}$$

with

$$A_P^n = \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| (\mathbf{F}^{(m)})_\zeta^n \cdot \mathbf{n}_{P,\zeta} (\varphi_P^n - \varphi_\zeta^n),$$

where φ_ζ^n denotes the mean value of $\varphi(\mathbf{x}, t_n)$ over ζ . For any $\mathbf{x} \in P, t \in [t_n, t_{n+1})$, we decompose A_P^n as

$$\begin{aligned} A_P^n &= B_P^n(\mathbf{x}, t) + R_P^n(\mathbf{x}, t) \text{ with} \\ B_P^n(\mathbf{x}, t) &= \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \mathbf{f}(U^{(m)})(\mathbf{x}, t) \cdot \mathbf{n}_{P,\zeta} (\varphi_P^n - \varphi_\zeta^n), \\ R_P^n(\mathbf{x}, t) &= \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \left((\mathbf{F}^{(m)})_\zeta^n - \mathbf{f}(U^{(m)})(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} (\varphi_P^n - \varphi_\zeta^n). \end{aligned} \tag{16}$$

Since $\sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \mathbf{n}_{P,\zeta} = 0$, we have

$$\begin{aligned} B_P^n(\mathbf{x}, t) &= - \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \mathbf{f}(U^{(m)})(\mathbf{x}, t) \cdot \mathbf{n}_{P,\zeta} \varphi_\zeta^n = -|P| \mathbf{f}(U^{(m)})(\mathbf{x}, t) \cdot (\nabla\varphi)_P^n, \\ \text{with } (\nabla\varphi)_P^n &= \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \varphi_\zeta^n \mathbf{n}_{P,\zeta} = \frac{1}{|P|} \int_P \nabla\varphi(\mathbf{x}, t_n) \, d\mathbf{x}. \end{aligned} \tag{17}$$

Note that the piecewise constant function $\nabla^{(m)}\varphi : \Omega \times (0, T) \rightarrow \mathbb{R}^d$ defined by

$$\nabla^{(m)}\varphi(\mathbf{x}, t) = (\nabla\varphi)_P^n \text{ for } (\mathbf{x}, t) \in P \times (t_n, t_{n+1})$$

converges uniformly to $\nabla\varphi$ in $L^\infty(\Omega \times (0, T))^d$.

Owing to (16), we have

$$A_P^n = \frac{1}{(t_{n+1} - t_n) |P|} \left(\int_{t_n}^{t_{n+1}} \int_P B_P^n(\mathbf{x}, t) \, d\mathbf{x} \, dt + \int_{t_n}^{t_{n+1}} \int_P R_P^n(\mathbf{x}, t) \, d\mathbf{x} \, dt \right),$$

and, thanks to (17),

$$\begin{aligned} X_2^{(m)} &= \sum_{n=0}^{N^{(m)}-1} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \frac{1}{|P|} \left(\int_{t_n}^{t_{n+1}} \int_P B_P^n(\mathbf{x}, t) \, d\mathbf{x} \, dt + \int_{t_n}^{t_{n+1}} \int_P R_P^n(\mathbf{x}, t) \, d\mathbf{x} \, dt \right) \\ &= - \int_0^T \int_\Omega \mathbf{f}(U^{(m)})(\mathbf{x}, t) \cdot \nabla^{(m)}\varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \\ &\quad + \sum_{n=0}^{N^{(m)}-1} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \frac{1}{|P|} \int_{t_n}^{t_{n+1}} \int_P R_P^n(\mathbf{x}, t) \, d\mathbf{x} \, dt. \end{aligned} \tag{18}$$

Then, thanks to the boundedness and convergence assumptions on $U^{(m)}$ and to the uniform convergence of $\nabla^{(m)}\varphi$ to $\nabla\varphi$, the first term tends to $-\int_0^T \int_\Omega \mathbf{f}(\bar{U})(\mathbf{x}, t) \cdot \nabla\varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt$ as $m \rightarrow +\infty$. Since $|\varphi_\zeta^n - \varphi_P^n| \leq C_\varphi \text{diam}(P)$, with C_φ depending only on φ , we get, for any $\mathbf{x} \in P$ and $t \in (t_n, t_{n+1})$,

$$|R_P^n(\mathbf{x}, t)| \leq C_\varphi \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \left| \left((\mathbf{F}^{(m)})_\zeta^n - \mathbf{f}(U^{(m)})(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| \text{diam}(P).$$

The second term of the right-hand side of Relation (18) thus tends to 0 as $m \rightarrow +\infty$ thanks to the assumption (11), which concludes the proof. \square

3 Application to a collocated scheme

In this section, we apply Theorem 2.1 to a first specific case, namely the case of a convection operator for a single scalar unknown \bar{u} , with a piecewise-constant discretisation of the unknown. We are going to prove a general consistency result for multipoint schemes, assuming minimal regularity of the mesh.

The considered convection operator reads:

$$\begin{aligned} \bar{\mathcal{C}}(\bar{u}) : \Omega \times (0, T) &\rightarrow \mathbb{R}, \\ (\mathbf{x}, t) &\mapsto \partial_t(\beta(\bar{u}))(\mathbf{x}, t) + \operatorname{div}(\mathbf{f}(\bar{u}))(\mathbf{x}, t), \end{aligned} \tag{19}$$

where $\Omega \subset \mathbb{R}^d, d = 1, 2, 3, T \in (0, +\infty), \beta \in C^0(\mathbb{R}, \mathbb{R}), \mathbf{f} \in C^0(\mathbb{R}, \mathbb{R}^d)$. The functions β and \mathbf{f} are supposed to be locally-Lipschitz. The discrete unknowns are $(u_P^n)_{P \in \mathcal{P}, n \in \llbracket 0, N-1 \rrbracket}$ and the discrete convection operator reads:

$$\begin{aligned} \mathcal{C}(u) : \Omega \times (0, T) &\rightarrow \mathbb{R}, \\ (\mathbf{x}, t) &\mapsto \frac{\beta(u_P^{n+1}) - \beta(u_P^n)}{t_{n+1} - t_n} + \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \mathbf{F}_\zeta^n \cdot \mathbf{n}_{P,\zeta}, \end{aligned}$$

for $\mathbf{x} \in P, P \in \mathcal{P}$, and $t \in (t_n, t_{n+1}), n \in \llbracket 0, N - 1 \rrbracket$. (20)

For a face ζ of the mesh, we denote by \mathcal{S}_ζ a set of cells, and suppose that the flux \mathbf{F}_ζ^n reads

$$\mathbf{F}_\zeta^n = \mathbf{F}_{\zeta,n}((u_L^n)_{L \in \mathcal{S}_\zeta}).$$

The set \mathcal{S}_ζ is often referred to as the stencil of the scheme. We denote by $(a)_{L \in \mathcal{S}_\zeta}$ the family of real numbers whose cardinal is the same as \mathcal{S}_ζ and whose elements are all equal to a . The flux is supposed to satisfy the usual consistency assumption:

$$\text{for } a \in \mathbb{R}, \forall \zeta \in \mathfrak{F}, \quad \mathbf{F}_{\zeta,n}((a)_{L \in \mathcal{S}_\zeta}) = \mathbf{f}(a). \tag{21}$$

In addition, we suppose the following local ‘‘Lip-diag’’ (for Lipschitz-diagonal) property of the flux:

Definition 3.1 (Local Lipschitz-diagonal continuity) The numerical flux function $\mathbf{F}_{\zeta,n}$ is said to be locally Lipschitz-diagonal continuous, or Lip-diag, if for any bounded interval $I \subset \mathbb{R}$, there exists $C_I \in \mathbb{R}_+$ depending only on I such that, for any $n \in \llbracket 0, N - 1 \rrbracket$, for any face $\zeta \in \mathfrak{F}$, for any family $(u_L)_{L \in \mathcal{S}_\zeta} \subset I$, and for any P adjacent cell to ζ ,

$$|\mathbf{F}_{\zeta,n}((u_L^n)_{L \in \mathcal{S}_\zeta}) - \mathbf{f}(u_P^n)| \leq C_I \sum_{L \in \mathcal{S}_\zeta} |u_L^n - u_P^n|. \tag{22}$$

Note that this condition is weaker than the local Lipschitz-continuity of the numerical flux function $\mathbf{F}_{\zeta,n}$. For $P \in \mathcal{P}$, we denote by $\mathcal{N}_1(P)$ the set of the neighbours of P , i.e. the cells of \mathcal{P} sharing a face with P ; for $\ell > 1$, we define $\mathcal{N}_\ell(P)$ as the set of the cells sharing a face with a cell of $\mathcal{N}_{\ell-1}(P)$. We assume that there exists $\bar{\ell} \in \mathbb{N}$ such that

$$\forall \zeta \in \mathfrak{F}, \text{ for any cell } P \text{ adjacent to } \zeta, \mathcal{S}_\zeta \subset \mathcal{N}_{\bar{\ell}}(P). \tag{23}$$

The integer $\bar{\ell}$ characterizes the compactness of the stencil of the scheme. The initial value for the scalar unknown u is defined by

$$u_P^0 = \frac{1}{|P|} \int_P u_0(\mathbf{x}) \, d\mathbf{x}, \quad \forall P \in \mathcal{P}, \tag{24}$$

where u_0 is a given function of $L^1(\Omega)$. Finally, we define the discrete function associated to the unknowns as:

$$u(\mathbf{x}, t) = u_P^n \quad \text{for } \mathbf{x} \in P, P \in \mathcal{P}, t \in [t_n, t_{n+1}), n \in \llbracket 0, N - 1 \rrbracket.$$

The consistency result for the discrete convection operator is given in the next lemma; it uses the following regularity parameters of the mesh:

$$\theta_1(\mathcal{P}) = \max_{P \in \mathcal{P}} \frac{\text{diam}(P)^d}{|P|}, \quad \theta_2(\mathcal{P}) = \max \left\{ \frac{|P|}{|Q|}, P \text{ and } Q \text{ adjacent cells of } \mathcal{P} \right\}.$$

We also measure the regularity of the time discretisation by the parameter $\theta_3(\mathcal{T})$ defined by

$$\theta_3(\mathcal{T}) = \max_{1 \leq n \leq N-1} \left\{ \frac{t_{n+1} - t_n}{t_n - t_{n-1}}, \frac{t_n - t_{n-1}}{t_{n+1} - t_n} \right\}.$$

Lemma 3.2 [Consistency, colocated scheme] *Consider a sequence of space and time discretisations $(\mathcal{P}^{(m)})_{m \in \mathbb{N}}$ and $(\mathcal{T}^{(m)})_{m \in \mathbb{N}}$, with $\delta(\mathcal{P}^{(m)})$ and $\delta t^{(m)}$ tending to zero; let $(u^{(m)})_{m \in \mathbb{N}}$ be an associated sequence of discrete functions, and let $\mathcal{C}^{(m)}(u^{(m)})$ be the associated sequence of discrete convection operators defined by (20). We assume that for each $m \in \mathbb{N}$, (21)–(24) hold with, in (23), $\bar{\ell}$ independent of m . We also suppose that*

$$\exists \theta \in \mathbb{R} \text{ such that } \max\{\theta_1(\mathcal{P}^{(m)}), \theta_2(\mathcal{P}^{(m)}), \theta_3(\mathcal{T}^{(m)}), m \in \mathbb{N}\} \leq \theta, \tag{25}$$

and that the number of faces of each cell of the meshes $\mathcal{P}^{(m)}$ is bounded independently of m . Finally, we suppose that the sequence $(u^{(m)})_{m \in \mathbb{N}}$ is bounded in $L^\infty(\Omega \times (0, T))$, and that, when m tends to $+\infty$, it converges in $L^1(\Omega \times (0, T))$ to $\bar{u} \in L^\infty(\Omega \times (0, T))$. Then, for any function $\varphi \in C_c^\infty(\Omega \times [0, T))$,

$$\begin{aligned} \int_0^T \int_\Omega \mathcal{C}^{(m)}(U^{(m)})(\mathbf{x}, t) \mathcal{J}^{(m)}(\varphi) \, d\mathbf{x} \, dt &\rightarrow - \int_\Omega \beta(u_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} \\ &- \int_0^T \int_\Omega \left(\beta(\bar{u})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + (\mathbf{f}(\bar{u}))(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right) \, d\mathbf{x} \, dt \quad \text{as } m \rightarrow +\infty. \end{aligned}$$

Proof Let $m_u \in \mathbb{R}$ and $M_u \in \mathbb{R}$ be two real numbers such that

$$m_u \leq (u^{(m)})_P^n \leq M_u, \quad \forall P \in \mathcal{P}^{(m)}, n \in \llbracket 0, N^{(m)} \rrbracket, \forall m \in \mathbb{N},$$

and let C_β be the Lipschitz modulus of β over the interval $[m_u, M_u]$. We check the assumptions of Theorem 2.1. The consistency of the initialization with the initial condition (Assumption (9)) follows from its definition (24); indeed, for any $\varphi \in C_c^\infty(\Omega)$,

$$\left| \sum_{P \in \mathcal{P}^{(m)}} \int_P \left((\beta^{(m)})_P^0 - \beta(u_0)(\mathbf{x}) \right) \varphi(\mathbf{x}) \, d\mathbf{x} \right| \leq C_\beta \|\varphi\|_{L^\infty(\Omega)} \sum_{P \in \mathcal{P}^{(m)}} \int_P |u_0(\mathbf{x}) - u_P^0| \, d\mathbf{x},$$

and thus tends to zero for any function $u_0 \in L^1(\Omega)$. Since $(\beta^{(m)})_P^n = \beta((u^{(m)})_P^n)$, the left-hand side of Assertion (10) reads, with $\varphi \in C_c^\infty(\Omega \times [0, T))$:

$$\begin{aligned} R_t^{(m)} &= \sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P \left((\beta^{(m)})_P^n - \beta(u^{(m)})(\mathbf{x}, t) \right) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt \\ &= \sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P \left(\beta((u^{(m)})_P^n) - \beta((u^{(m)})_P^{n-1}) \right) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt. \end{aligned}$$

We thus have

$$|R_t^{(m)}| \leq C_\beta \|\varphi\|_{L^\infty(\Omega \times [0, T])} \sum_{n=1}^{N^{(m)}} (t_n - t_{n-1}) \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} |P| |(u^{(m)})_P^n - (u^{(m)})_P^{n-1}|,$$

and thus $R_t^{(m)}$ tends to zero thanks to the assumed regularity of the sequence of time discretisations, invoking the bound of the time-translates of a converging sequence of functions of $L^1(\Omega \times (0, T))$ stated by Lemma A.1 in Appendix. We now check Assumption (11). For $n \in \llbracket 0, N^{(m)} \rrbracket$, $P \in \mathcal{P}_{\text{int}}^{(m)}$ and $\zeta \in \mathfrak{F}(P)$, let

$$R_{P,\zeta}^n = \frac{1}{|P|} \int_{t_n}^{t_{n+1}} \int_P \left| \left((F^{(m)})_\zeta^n - f(u^{(m)})(x, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| dx dt$$

and let

$$R^{(m)} = \sum_{n=0}^{N^{(m)}-1} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \text{diam}(P) \sum_{\zeta \in \mathfrak{F}} |\zeta| R_{P,\zeta}^n.$$

By definition of the discrete flux and the discrete functions, we get:

$$R_{P,\zeta}^n = \frac{1}{|P|} \int_{t_n}^{t_{n+1}} \int_P \left| \left(F_{\zeta,n} \left((u^{(m)})_L^n \right)_{L \in \mathcal{S}_\zeta} - f((u^{(m)})_P^n) \right) \cdot \mathbf{n}_{P,\zeta} \right| dx dt.$$

Thanks to Assumption (22), we thus have

$$R_{P,\zeta}^n \leq (t_{n+1} - t_n) C_{[m_u, M_u]} \sum_{L \in \mathcal{S}_\zeta} |u_L^n - u_P^n|.$$

The remainder term $R^{(m)}$ is thus a collection of differences between the values taken by the unknown in two different cells. In order to apply Lemma A.3 of the appendix, we need to evaluate, for $P \in \mathcal{P}^{(m)}$, the sum ω_P of the weights multiplying the jumps where $(u^{(m)})_P^n$ appears. We first notice that, thanks to the assumed regularity of the mesh, for $Q \in \mathcal{P}^{(m)}$ and ζ a face of Q , $\text{diam}(Q) |\zeta| \leq \text{diam}(Q)^d \leq \theta |Q|$. For P to appear in a difference associated to a face of a cell Q , we need, by assumption, that $P \in \mathcal{N}_{\bar{\zeta}}(Q)$; this in turn requires that $Q \in \mathcal{N}_{\bar{\zeta}}(P)$. The sum ω_P thus satisfies:

$$\omega_P \leq \theta \sum_{L \in \mathcal{N}_{\bar{\zeta}}(P)} |L|,$$

and, invoking once again the regularity constraints on the mesh:

$$\omega_P \leq \theta^{\bar{\ell}+1} \text{card}(\mathcal{N}_{\bar{\zeta}}(P)) |P|.$$

The proof is thus complete thanks to Lemma A.3, since we have supposed that the number of faces of the cells is uniformly bounded, and then so is $\text{card}(\mathcal{N}_{\bar{\zeta}}(P))$. \square

4 An example of application for staggered discretisations

The interest of Theorem 2.1 lies in the fact that it may deal with terms combining several variables, associated to different meshes and time discretisations. A typical example of a such a case is the balance equation for the entropy in barotropic compressible flows (2), where

the entropy E is a nonlinear function of the density ρ and the velocity \mathbf{u} which, in staggered discretisation, are approximated on different meshes, and may also be evaluated at different time levels. Hence, Theorem 2.1 is a suitable tool to prove the consistency of this equation. In this section, we focus on a similar but simpler problem, namely a staggered discretisation of a convection operator combining the time derivative of the function of a single scalar variable and a space divergence term, with a flux obtained as the product of another function of this scalar variable with the velocity.

We suppose that Ω is an open bounded polygonal set of \mathbb{R}^2 , and consider the following convection operator:

$$\begin{aligned} \mathcal{C}(\bar{U}) : \Omega \times (0, T) &\rightarrow \mathbb{R}, \\ (\mathbf{x}, t) &\mapsto \partial_t(\beta(\bar{q}))(\mathbf{x}, t) + \operatorname{div}(g(\bar{q}) \bar{\mathbf{v}})(\mathbf{x}, t), \end{aligned} \tag{26}$$

with $\bar{U} = (\bar{q}, \bar{\mathbf{v}}) : \Omega \times (0, T) \rightarrow \mathbb{R} \times \mathbb{R}^2$, $\mathbf{f}(\bar{U}) = \mathbf{f}(\bar{q}, \bar{\mathbf{v}}) = g(\bar{q}) \bar{\mathbf{v}}$, where $\beta : \mathbb{R} \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ are locally Lipschitz-continuous real functions. Note that, for instance, the convection term of Eq. (1a) may be written as (26) with $\bar{U} = (\bar{\rho}, \bar{\mathbf{u}})$, $\beta(s) = s$ and $g(s) = s$.

In order to discretise this convection operator, we consider two types of staggered arrangements. In both arrangements, the scalar unknowns are located at the centre of the cells. However, these arrangements differ in the use of the vector unknowns. The first discretisation uses the whole velocity vector unknown on each edge of the mesh; this corresponds to the Rannacher–Turek (RT) discrete unknowns in the finite element setting [13]. The second discretisation uses only the normal component of the velocity on each edge; this latter arrangement of the discrete unknowns is very often referred to as the Marker-and-Cell (MAC) scheme [7]. Hence we will refer to the first arrangement as the RT case, and the second as the MAC case. Such discretisations are called staggered and are widely used in computational fluid dynamics; an example of the implementation of a staggered discretisation for the solution of the barotropic and full Euler equations may be found e.g. in [8,9].

We suppose that the mesh is composed either of general quadrangles (RT case), or of rectangles (MAC case). We recall that \mathfrak{F} stands for the set of edges of the mesh, and the internal edge separating the cells P and Q is denoted by $\zeta = P|Q$. This mesh will be referred to in the following as the primal mesh.

We also introduce now one or two dual meshes, depending on the case.

- *RT case* In this case, the (unique) dual mesh consists in a new partition of Ω indexed by the elements of \mathfrak{F} , i.e. $\Omega = \cup_{\zeta \in \mathfrak{F}} D_\zeta$. For an internal edge $\zeta = P|Q$, the set D_ζ is supposed to be a subset of $P \cup Q$ and we define $D_{P,\zeta} = D_\zeta \cap P$, so that $D_\zeta = D_{P,\zeta} \cup D_{Q,\zeta}$ (see Fig. 2); for an external edge ζ of a cell P , D_ζ is a subset of P , and $D_\zeta = D_{P,\zeta}$. The cells $(D_\zeta)_{\zeta \in \mathfrak{F}}$ are referred to as the dual or diamond cells, and $D_{P,\zeta}$ as half dual cells or half diamond cells. For a rectangular cell P , we define $D_{P,\zeta}$ as the simplex having the mass centre of P as vertex and the edge ζ as basis; this definition is extended to general primal meshes by supposing that $|D_{P,\zeta}|$ is still equal to $|P|/4$ and that the sub-cells connectivities (i.e. the way the half-dual cells share a common edge) is left unchanged. Note that the actual geometry of the dual cells does not need to be specified (and a dual cell may not be a polytope, a dual edge being possibly curved).
- *MAC case* In this case, two dual meshes are considered, each of them consisting in a partition of Ω indexed by the vertical and horizontal elements of \mathfrak{F} , i.e. $\Omega = \cup_{\zeta \in \mathfrak{F}^{(i)}} D_\zeta$, $i = 1, 2$, where $\mathfrak{F}^{(1)}$ (resp. $\mathfrak{F}^{(2)}$) denotes the set of vertical (resp. horizontal) edges. The cells $(D_\zeta)_{\zeta \in \mathfrak{F}}$ are still referred to as the dual cells. They are no longer diamond shaped; indeed, a half dual cell $D_{P,\zeta}$ is now half of the rectangle P with side ζ (see Fig. 3). As in the former case, for an internal edge $\zeta = P|Q$, the dual cell D_ζ is the subset of $P \cup Q$

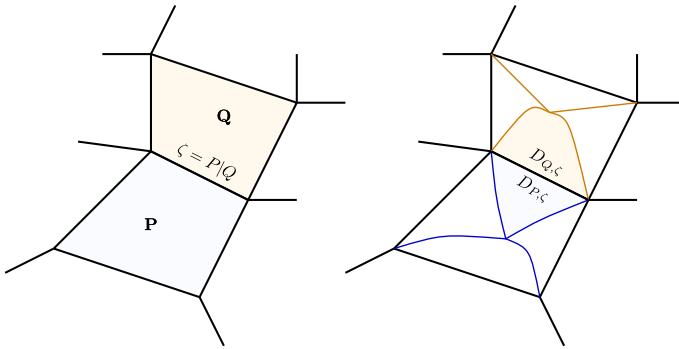


Fig. 2 Primal and dual meshes and associated notations for the quadrangular mesh and Rannacher–Turek like unknowns

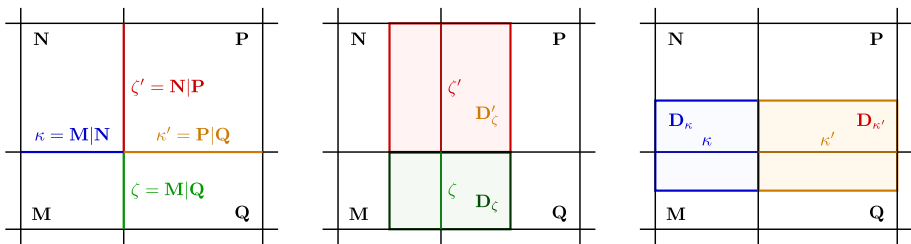


Fig. 3 Primal and dual meshes and associated notations for the MAC case. Left: the primal cells; the edges ζ and ζ' belong to $\mathfrak{F}^{(1)}$ and the edges κ and κ' to $\mathfrak{F}^{(2)}$ Center: the dual cells associated to $\mathfrak{F}^{(1)}$ Right: the dual cells associated to $\mathfrak{F}^{(2)}$

defined as $D_\zeta = D_{P,\zeta} \cup D_{Q,\zeta}$; for an external edge ζ of a cell P , D_ζ is a subset of P , and $D_\zeta = D_{P,\zeta}$.

The scalar unknown q is associated to the primal cells:

$$q(\mathbf{x}, t) = q_P^n \text{ for } \mathbf{x} \in P, P \in \mathcal{P}, t \in [t_n, t_{n+1}), n \in \llbracket 0, N - 1 \rrbracket.$$

The unknowns corresponding to the vector-valued unknown \mathbf{v} are located at the centre of the edges in the RT case; in the MAC case, the unknowns associated to the i th component of \mathbf{v} are located at the centre of the edges of the i th dual mesh. Hence the associated approximate vector function reads:

- RT case—the whole vector unknown is associated to each dual cell :

$$\mathbf{v}(\mathbf{x}, t) = \mathbf{v}_\zeta^n \text{ for } \mathbf{x} \in D_\zeta, \zeta \in \mathfrak{F}, t \in [t_n, t_{n+1}), n \in \llbracket 0, N - 1 \rrbracket.$$

- MAC case—the i th component of the vector unknown is associated to the cells of the i th dual mesh, so that $\mathbf{v}(\mathbf{x}, t) = (v_1(\mathbf{x}, t), v_2(\mathbf{x}, t))^t$ where, for $i = 1, 2$,

$$v_i(\mathbf{x}, t) = v_\zeta^n, \text{ for } \mathbf{x} \in D_\zeta, \zeta \in \mathfrak{F}^{(i)} \text{ and } t \in [t_n, t_{n+1}), n \in \llbracket 0, N - 1 \rrbracket.$$

Let $\mathbf{e}^{(i)}$ denote the i th unit vector of the canonical basis of \mathbb{R}^2 ; with the notations of the previous section, the considered discrete convection operator reads:

$$\mathbb{C}_{\mathcal{P}}(q, \mathbf{v})_P^n = (\partial_i \beta)_P^n + \frac{1}{|P|} \sum_{\zeta \in \mathfrak{F}(P)} |\zeta| \mathbf{F}_\zeta^n \cdot \mathbf{n}_{P,\zeta}, \text{ with } \beta_P^n = \beta(q_P^n) \text{ and } \mathbf{F}_\zeta^n$$

$$= f(q_\zeta^n, \mathbf{v}_\zeta^n) = g(q_\zeta^n) \mathbf{v}_\zeta^n$$

where \mathbf{v}_ζ^n is $\begin{cases} \text{the vector of discrete unknowns} & \text{in the RT case,} \\ \text{defined as } v_\zeta^n \mathbf{e}^{(i)} \text{ for } \zeta \in \mathfrak{F}^{(i)}, i = 1 \text{ or } 2, & \text{in the MAC case,} \end{cases}$

and, for $\zeta = P|Q$, q_ζ^n stands for a convex combination of q_P^n and q_Q^n . For instance the upwind choice would be $q_\zeta^n = q_P^n$ if $v_\zeta^n \geq 0$ and $q_\zeta^n = q_Q^n$ otherwise. Note that for the LW-consistency result, any convex combination works, but this is not so for the stability of the scheme, for which some unpwinding is required.

The initial value for the scalar unknown q is defined by

$$q_P^0 = \frac{1}{|P|} \int_P q_0(\mathbf{x}) \, d\mathbf{x}. \tag{27}$$

The consistency result for the discrete convection operator is given in the next lemma; it uses the same regularity parameters of the mesh as in the collocated case, which we recall:

$$\theta_1(\mathcal{P}) = \max_{P \in \mathcal{P}} \frac{\text{diam}(P)^2}{|P|}, \quad \theta_2(\mathcal{P}) = \max \left\{ \frac{|P|}{|Q|}, P \text{ and } Q \text{ adjacent cells of } \mathcal{P} \right\}.$$

Note that in the MAC case (in fact, for a Cartesian grid), the regularity parameter $\theta_1(\mathcal{P})$ controls the ratio between the two dimensions (i.e. the height and the width) of a cell. For a rectangular computational domain, we thus observe that the ratio $|\zeta|/|\zeta'|$, for $(\zeta, \zeta') \in (\mathfrak{F}^{(i)})^2, i = 1, 2$, is bounded by $\theta_1(\mathcal{P})^2$, which is a quasi-uniformity property of the mesh. This also implies that $\theta_2(\mathcal{P}) \leq \theta_1(\mathcal{P})^2$, and so the second regularity parameter is useless. It may easily be checked that similar relations holds for a general MAC scheme, i.e. a union of matching Cartesian grids, with powers of $\theta_1(\mathcal{P})$ possibly higher than 2. Hence, the regularity of a MAC mesh (or of a Cartesian grid) may be equivalently measured by

$$\theta(\mathcal{P}) = \max \left\{ \frac{\bar{h}^{(1)}}{\underline{h}^{(2)}}, \frac{\bar{h}^{(2)}}{\underline{h}^{(1)}} \right\},$$

with, for $i = 1, 2, \bar{h}^{(i)} = \max\{|\zeta|, \zeta \in \mathfrak{F}^{(i)}\}$ and $\underline{h}^{(i)} = \min\{|\zeta|, \zeta \in \mathfrak{F}^{(i)}\}$.

We also measure the regularity of the time discretisation by the parameter $\theta_3(\mathcal{T})$ defined by

$$\theta_3(\mathcal{T}) = \max_{1 \leq n \leq N-1} \left\{ \frac{t_{n+1} - t_n}{t_n - t_{n-1}}, \frac{t_n - t_{n-1}}{t_{n+1} - t_n} \right\}.$$

Lemma 4.1 (Consistency, staggered scheme) *Let a sequence of discretisations $(\mathcal{P}^{(m)})_{m \in \mathbb{N}}$ and $(\mathcal{T}^{(m)})_{m \in \mathbb{N}}$ be given, with $\delta(\mathcal{P}^{(m)})$ and $\delta t^{(m)}$ tending to zero, and let $(q^{(m)}, \mathbf{v}^{(m)})_{m \in \mathbb{N}}$ be the associated sequence of discrete functions. We suppose that*

$$\exists \theta \in \mathbb{R} \text{ such that } \max\{\theta_1(\mathcal{P}^{(m)}), \theta_2(\mathcal{P}^{(m)}), \theta_3(\mathcal{T}^{(m)}), m \in \mathbb{N}\} \leq \theta. \tag{28}$$

We suppose that the sequences $(q^{(m)})_{m \in \mathbb{N}}$ and $(\mathbf{v}^{(m)})_{m \in \mathbb{N}}$ are bounded in $L^\infty(\Omega \times (0, T))$ and $L^\infty(\Omega \times (0, T))^2$ respectively, and that, when m tends to $+\infty$, they converge in $L^p(\Omega \times (0, T))$ and $L^p(\Omega \times (0, T))^2, 1 \leq p < +\infty$, to $\bar{q} \in L^\infty(\Omega \times (0, T))$ and $\bar{\mathbf{v}} \in L^\infty(\Omega \times (0, T))^2$ respectively. Then, for any function $\varphi \in C_c^\infty(\Omega \times [0, T))$,

$$\begin{aligned} \int_0^T \int_\Omega \mathfrak{C}^{(m)}(U^{(m)})(\mathbf{x}, t) \mathcal{J}^{(m)}(\varphi) \, d\mathbf{x} \, dt &\rightarrow - \int_\Omega \beta(q_0)(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} \\ &- \int_0^T \int_\Omega \left(\beta(\bar{q})(\mathbf{x}, t) \partial_t \varphi(\mathbf{x}, t) + (g(\bar{q}) \bar{\mathbf{v}})(\mathbf{x}, t) \cdot \nabla \varphi(\mathbf{x}, t) \right) \, d\mathbf{x} \, dt \text{ as } m \rightarrow +\infty. \end{aligned}$$

Proof In this proof, we denote by C_β and C_g the Lipschitz modulus of β and g respectively on the interval $[m_q, M_q]$, where $m_q \in \mathbb{R}$ and $M_q \in \mathbb{R}$ are such that

$$m_q \leq (q^{(m)})_P^n \leq M_q, \forall P \in \mathcal{P}^{(m)}, n \in \llbracket 0, N^{(m)} \rrbracket, \forall m \in \mathbb{N}.$$

The proof of this lemma relies on Theorem 2.1. The consistency of the initialization with the initial condition (Assumption (9)) follows from its definition (27); indeed, for any $\varphi \in C_c^\infty(\Omega)$,

$$\left| \sum_{P \in \mathcal{P}^{(m)}} \int_P \left((\beta^{(m)})_P^0 - \beta(q_0)(x) \right) \varphi(x) \, dx \right| \leq C_\beta \|\varphi\|_{L^\infty(\Omega)} \sum_{P \in \mathcal{P}^{(m)}} \int_P |q_0(x) - q_P^0|,$$

and thus tends to zero for any function $q_0 \in L^1(\Omega)$. Since $(\beta^{(m)})_P^n = \beta((q^{(m)})_P^n)$, the left-hand side of Assertion (10) reads, with $\varphi \in C_c^\infty(\Omega \times [0, T])$:

$$\begin{aligned} R_t^{(m)} &= \sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P \left((\beta^{(m)})_P^n - \beta(U^{(m)})(x, t) \right) \varphi(x, t) \, dx \, dt \\ &= \sum_{n=1}^{N^{(m)}} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \int_{t_{n-1}}^{t_n} \int_P \left(\beta((q^{(m)})_P^n) - \beta((q^{(m)})_P^{n-1}) \right) \varphi(x, t) \, dx \, dt. \end{aligned}$$

We thus have

$$|R_t^{(m)}| \leq C_\beta \|\varphi\|_{L^\infty(\Omega \times [0, T])} \sum_{n=1}^{N^{(m)}} (t_n - t_{n-1}) \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} |(q^{(m)})_P^n - (q^{(m)})_P^{n-1}|,$$

and thus $R_t^{(m)}$ tends to zero thanks to the assumed regularity of the sequence of time discretisations, invoking the bound of the time-translates of a converging sequence of functions of $L^1(\Omega \times (0, T))$ stated by Lemma A.1 in Appendix.

We now check Assumption (11). For $n \in \llbracket 0, N^{(m)} \rrbracket$, $P \in \mathcal{P}_{\text{int}}^{(m)}$ and $\zeta \in \mathfrak{F}(P)$, let

$$R_{P,\zeta}^n = \frac{1}{|P|} \int_{t_n}^{t_{n+1}} \int_P \left| \left((F^{(m)})_\zeta^n - f(q^{(m)}, \mathbf{v}^{(m)})(x, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| \, dx \, dt$$

and let

$$R^{(m)} = \sum_{n=0}^{N^{(m)}-1} \sum_{P \in \mathcal{P}_{\text{int}}^{(m)}} \text{diam}(P) \sum_{\zeta \in \mathfrak{F}} |\zeta| R_{P,\zeta}^n.$$

We now express $R_{P,\zeta}^n$, for the RT and MAC discretisations successively.

- RT case—in the case of general quadrangular meshes with the whole vector unknowns located on the edges, we have

$$(F^{(m)})_\zeta^n = g(q_\zeta^n) \mathbf{v}_\zeta^n \quad \text{and} \quad f(q^{(m)}, \mathbf{v}^{(m)})(x, t) = g(q_P^n) \mathbf{v}_{\zeta'}^n \quad \text{for } x \in D_{P,\zeta'}, \zeta' \in \mathfrak{F}(P).$$

We thus get, denoting by $|\mathbf{a}|$ the Euclidean norm of any vector $\mathbf{a} \in \mathbb{R}^2$,

$$\left| \left((F^{(m)})_\zeta^n - f(U^{(m)})(x, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| \leq |g(q_\zeta^n) \mathbf{v}_\zeta^n - g(q_P^n) \mathbf{v}_{\zeta'}^n| \quad \text{for } x \in D_{P,\zeta'}, \zeta' \in \mathfrak{F}(P).$$

Let Q be the primal cell such that $\zeta = P|Q$. Since q_ζ^n is a convex combination of q_P^n and q_Q^n , we thus get, for $\mathbf{x} \in P$, and $t \in [t_n t_{n+1})$,

$$\left| \left((\mathbf{F}^{(m)})_\zeta^n - \mathbf{f}(U^{(m)})(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| \leq C \left(|q_P^n - q_Q^n| + \sum_{\zeta' \in \mathfrak{F}(P)} |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| \right),$$

where C only depends on $\|q^{(m)}\|_{L^\infty(\Omega \times (0, T))}$, $\|\mathbf{v}^{(m)}\|_{L^\infty(\Omega \times (0, T))^2}$ and C_g . Integrating over $P \times (t_n, t_{n+1})$, we obtain

$$R_{P,\zeta}^n \leq C (t_{n+1} - t_n) \left(|q_P^n - q_Q^n| + \sum_{\zeta' \in \mathfrak{F}(P)} |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| \right).$$

- MAC case—in this case, the velocity components are piecewise constant on different grids. Let i be the index such that $\zeta \in \mathfrak{F}^{(i)}$, and let ζ' be the other edge of P normal to $\mathbf{e}^{(i)}$, i.e. the opposite of ζ in P . We have

$$(\mathbf{F}^{(m)})_\zeta^n \cdot \mathbf{n}_{P,\zeta} = g(q_\zeta^n) v_\zeta^n \delta_\zeta \quad \text{and} \quad \mathbf{f}(q^{(m)}, \mathbf{v}^{(m)})(\mathbf{x}, t) = \begin{cases} g(q_P^n) v_\zeta^n \delta_\zeta & \text{if } \mathbf{x} \in D_{P,\zeta}, \\ g(q_P^n) v_{\zeta'}^n \delta_{\zeta'} & \text{if } \mathbf{x} \in D_{P,\zeta'}, \end{cases}$$

with $\delta_\zeta = \mathbf{n}_{P,\zeta} \cdot \mathbf{e}^{(i)}$. We thus get

$$\left| \left((\mathbf{F}^{(m)})_\zeta^n - \mathbf{f}(U^{(m)})(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| = \begin{cases} |g(q_\zeta^n) v_\zeta^n - g(q_P^n) v_\zeta^n| & \text{if } \mathbf{x} \in D_{P,\zeta}, \\ |g(q_\zeta^n) v_\zeta^n - g(q_P^n) v_{\zeta'}^n| & \text{if } \mathbf{x} \in D_{P,\zeta'}, \end{cases}$$

and hence, for $\mathbf{x} \in P$, and $t \in [t_n t_{n+1})$, denoting by Q the primal cell such that $\zeta = P|Q$,

$$\left| \left((\mathbf{F}^{(m)})_\zeta^n - \mathbf{f}(U^{(m)})(\mathbf{x}, t) \right) \cdot \mathbf{n}_{P,\zeta} \right| \leq C \left(|q_P^n - q_Q^n| + |v_\zeta^n - v_{\zeta'}^n| \right),$$

where C only depends on $\|q^{(m)}\|_{L^\infty(\Omega \times (0, T))}$, $\|\mathbf{v}^{(m)}\|_{L^\infty(\Omega \times (0, T))^2}$ and C_g . Therefore, integrating over $P \times (t_n, t_{n+1})$, we finally get

$$R_{P,\zeta}^n \leq C (t_{n+1} - t_n) \left(|q_P^n - q_Q^n| + |v_\zeta^n - v_{\zeta'}^n| \right).$$

Note that, in these computations, we have not addressed the case where ζ is an external edge, taking benefit of the fact that, in the expression of $R^{(m)}$, the sum is restricted to the internal cells.

From the definition of $R^{(m)}$, we thus get that, for both cases, it satisfies the following inequality:

$$R^{(m)} \leq C (R_1^{(m)} + R_2^{(m)}),$$

with

$$R_1^{(m)} = \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{P \in \mathcal{P}^{(m)}} \text{diam}(P) \sum_{\substack{\zeta \in \mathfrak{F}(P), \\ \zeta = P|Q}} |\zeta| |q_P^n - q_Q^n|,$$

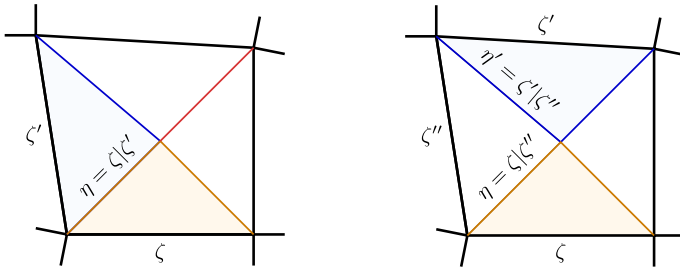


Fig. 4 Left: the primal edges ζ and ζ' are adjacent. Right: the primal edges ζ and ζ' are opposite

and

$$R_2^{(m)} = \begin{cases} \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{P \in \mathcal{P}^{(m)}} \text{diam}(P) \sum_{\{\zeta, \zeta'\} \subset \mathfrak{F}(P)^2} (|\zeta| + |\zeta'|) |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| & \text{in the RT case,} \\ \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{P \in \mathcal{P}^{(m)}} \text{diam}(P) \sum_{\substack{i=1,2, \\ \{\zeta, \zeta'\} \subset \mathfrak{F}^{(i)}(P)^2}} (|\zeta| + |\zeta'|) |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| & \text{in the MAC case.} \end{cases}$$

There only remains to prove that $R_1^{(m)}$ and $R_2^{(m)}$ tend to zero as m tends to $+\infty$. Reordering the summation in $R_1^{(m)}$, we get that

$$R_1^{(m)} = \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{\substack{\zeta \in \mathfrak{F}_{\text{int}}^{(m)}, \\ \zeta = P|Q}} \omega_\zeta |q_P^n - q_Q^n|, \quad \text{with } \omega_\zeta = \left(\text{diam}(P) + \text{diam}(Q) \right) |\zeta|.$$

Lemma A.1 states that $R_1^{(m)}$ tends to zero if the weight ω_ζ is controlled by both $|P|$ and $|Q|$; since we have $\omega_\zeta \leq 2(\max(\text{diam}(P), \text{diam}(Q)))^2$, this is easily obtained using Assumption (28).

As to the term $R_2^{(m)}$, let us start with the RT case. We have:

$$\sum_{\{\zeta, \zeta'\} \subset \mathfrak{F}(P)^2} (|\zeta| + |\zeta'|) |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| \leq 2 \text{diam}(P) \sum_{\{\zeta, \zeta'\} \subset \mathfrak{F}(P)^2} |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n|,$$

We distinguish two cases for the subsets $\{\zeta, \zeta'\} \subset \mathfrak{F}(P)^2$ that appear in the summation: either the dual cells D_ζ and $D_{\zeta'}$ share a common (dual) edge $\eta = \zeta|\zeta' \in \mathfrak{F}^*$, where \mathfrak{F}^* denotes the set of edges of the dual mesh, or they are opposite edges in the quadrilateral cell P ; in this latter case, we may write that

$$|\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| \leq |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta''}^n| + |\mathbf{v}_{\zeta''}^n - \mathbf{v}_{\zeta'}^n|,$$

where $\zeta'' \in \mathfrak{F}(P)$ is such that the dual cell $D_{\zeta''}$ shares a common (dual) edge η (resp. η') $\in \mathfrak{F}^*$ with D_ζ (resp. $D_{\zeta'}$) as shown in Fig. 4. There is one jump between two adjacent faces that appears directly in the summation over $\{\zeta, \zeta'\} \subset \mathfrak{F}(P)^2$, and at most two coming from the decompositions of the jumps needed for pairs of opposite edges, so that altogether,

$$\sum_{(\zeta, \zeta') \in \mathfrak{F}(P)^2} (|\zeta| + |\zeta'|) |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n| \leq 6 \text{diam}(P) \sum_{\eta = \zeta|\zeta' \in \mathfrak{F}^*(P)} |\mathbf{v}_\zeta^n - \mathbf{v}_{\zeta'}^n|,$$

with $\mathfrak{F}^*(P)$ the edges of the dual mesh included in P . We thus get

$$R_2^{(m)} \leq 6 \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{\eta=\zeta|\zeta' \in \mathfrak{F}^*} \text{diam}(P_\eta)^2 |v_\zeta^n - v_{\zeta'}^n|,$$

where P_η stands for the primal cell in which η is included. The right-hand side of this inequality is thus a collection of jumps across the dual edges, with, for an edge η , a weight given by

$$\omega_\eta = 6 \text{diam}(P_\eta)^2.$$

Thanks to Lemma A.1, $R_2^{(m)}$ tends to zero when m tends to $+\infty$ if ω_η is controlled by both $|D_\zeta|$ and $|D_{\zeta'}|$; this is indeed the case thanks to Assumption (28), since $|D_\zeta| \geq |P_\eta|/4$ and $|D_{\zeta'}| \geq |P_\eta|/4$.

Let us now turn to the MAC case, which is in fact simpler; indeed, the differences of velocities appearing in the expression of $R_2^{(m)}$ are all jumps across dual edges, and we may thus recast $R_2^{(m)}$ as

$$R_2^{(m)} = \sum_{n=0}^{N^{(m)}-1} (t_{n+1} - t_n) \sum_{i=1}^2 \sum_{\eta=\zeta|\zeta' \in \mathfrak{F}^{(i,*)}} \text{diam}(P_\eta) (|\zeta| + |\zeta'|) |v_\zeta^n - v_{\zeta'}^n|,$$

where, once again, P_η is the primal cell in which lies η (note that this sum only involves a subset of the dual edges, which corresponds of the dual edges included in primal cell), and $\mathfrak{F}^{(1,*)}$ (resp. $\mathfrak{F}^{(2,*)}$) denotes the set of vertical (respectively horizontal) dual edges. We thus again have a collection of jumps across the dual edges, with, for an edge η included in a primal cell P_η and separating the dual cells D_ζ and $D_{\zeta'}$, a weight given by

$$\omega_\eta = \text{diam}(P_\eta) (|\zeta| + |\zeta'|).$$

Thus, again thanks to Lemma A.1, $R_2^{(m)}$ tends to zero when m tends to $+\infty$ since, remarking that $|D_\zeta| \geq |P_\eta|/2$, $|D_{\zeta'}| \geq |P_\eta|/2$ and $\omega_\eta \leq 2\text{diam}(P_\eta)^2$, so the weight ω_η is controlled by both $|D_\zeta|$ and $|D_{\zeta'}|$ thanks to Assumption (28). \square

Remark 4.2 (On the required regularity of the time discretisation) The assumption $\theta_3(\mathcal{T}^{(m)}) \leq \theta$, for $m \in \mathbb{N}$, may be avoided thanks to a different choice of the interpolation of the test function (see Remark 2.6). However, this assumption is very mild (in fact, we do not have in mind any scheme where the ratio between two consecutive time-steps is likely to blow up when refining the discretisation).

Appendix A. Convergence of discrete functions in L^1

We recall a result proven in [5, Lemma 4.3]. To facilitate its use in the proof of Lemma 4.1, it is rephrased here under more general forms than in [5] (see Remark A.2 below for the differences).

Let \mathcal{M} be a conforming mesh of the domain Ω of \mathbb{R}^d , $d = 1, 2, 3$, in polygonal or polyhedral subsets, and $\mathcal{T} = (t_i)_{i \in \llbracket 0, N \rrbracket}$ be a time discretisation of the interval $(0, T)$, i.e. a sequence of real numbers such that $0 = t_0 < \dots < t_n < \dots < t_N = T$. We denote by $\delta t_{\mathcal{T}}$ the time step, defined by $\delta t_{\mathcal{T}} = \max\{t_{n+1} - t_n, n \in \llbracket 0, N - 1 \rrbracket\}$. For $u \in L^1(\Omega \times (0, T))$, $K \in \mathcal{M}$ and n such that $n \in \llbracket 0, N - 1 \rrbracket$, let u_K^n be the mean value of u over $K \times (t_n, t_{n+1})$.

We denote by \mathcal{E}_{int} the set of internal faces of the mesh and the face $\sigma \in \mathcal{E}_{\text{int}}$ separating the cells K and L is denoted by $\sigma = K|L$. We define the following quantity:

$$T_{\mathcal{M}, \mathcal{T}} u = \sum_{n=0}^{N-1} (t_{n+1} - t_n) \sum_{\substack{\sigma \in \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \omega_{\sigma} |u_K^n - u_L^n| + \sum_{n=0}^{N-2} \delta_{n+1/2} \sum_{K \in \mathcal{M}} |K| |u_K^{n+1} - u_K^n|, \tag{29}$$

where $(\omega_{\sigma})_{\sigma \in \mathcal{E}_{\text{int}}}$ and $(\delta_{n+1/2})_{n \in \llbracket 0, N-2 \rrbracket}$ are two sets of non-negative weights. We introduce the two following parameters:

$$\theta_{\mathcal{M}} = \max_{K \in \mathcal{M}} \max_{\sigma \in \mathcal{E}_{\text{int}}(K)} \frac{\omega_{\sigma}}{|K|}, \quad \theta_{\mathcal{T}} = \max_{n \in \llbracket 0, N-2 \rrbracket} \left\{ \frac{\delta_{n+1/2}}{t_{n+1} - t_n}, \frac{\delta_{n+1/2}}{t_{n+2} - t_{n+1}} \right\}, \tag{30}$$

with $\mathcal{E}_{\text{int}}(K)$ the set of internal faces of K . We denote by $\delta(\mathcal{M})$ the space step characterizing \mathcal{M} , i.e. $\delta(\mathcal{M}) = \max_{K \in \mathcal{M}} \text{diam}(K)$. Then the following convergence result holds.

Lemma A.1 *Let $\theta > 0$ and $(\mathcal{M}^{(m)})_{m \in \mathbb{N}}$ be a sequence of meshes and for each $m \in \mathbb{N}$, $\theta_{\mathcal{M}^{(m)}}$ be defined by (30). We assume that $\theta_{\mathcal{M}^{(m)}} \leq \theta$ for all $m \in \mathbb{N}$ and $\lim_{m \rightarrow +\infty} \delta(\mathcal{M}^{(m)}) = 0$. We suppose that the number of faces of a cell $K \in \mathcal{M}^{(m)}$ is bounded independently from $m \in \mathbb{N}$. For $m \in \mathbb{N}$, we suppose given a time discretisation $\mathcal{T}^{(m)}$, and suppose that $\delta t_{\mathcal{T}^{(m)}}$ also tends to zero when m tends to $+\infty$, and that $\theta_{\mathcal{T}^{(m)}} \leq \theta$ for all $m \in \mathbb{N}$. Let $u \in L^1(\Omega \times (0, T))$ and $(u_p)_{p \in \mathbb{N}}$ be a sequence of functions of $L^1(\Omega \times (0, T))$ such that $u_p \rightarrow u$ in $L^1(\Omega \times (0, T))$ as $p \rightarrow +\infty$.*

Then $T_{\mathcal{M}^{(m)}, \mathcal{T}^{(m)}} u_p$ defined by (29) tends to zero when m tends to $+\infty$ uniformly with respect to $p \in \mathbb{N}$.

Remark A.2 The difference between Lemma A.1 and the formulation of the same convergence result in [5] lies in the definition of the weight of the jumps, which is more general in Lemma A.1. Indeed, the weight of the jumps through the faces featured in the definition of $T_{\mathcal{M}, \mathcal{T}} u$ are defined in [5, Lemma 4.3] as a function of the volume of some dual cells associated to the faces, but a careful examination of the proof itself shows that the introduction of a dual mesh is in fact useless. Therefore, the proof of Lemma [5, Lemma 4.3] readily extends to prove Lemma A.1.

This generalization is in most cases sufficient. However, we may go one step further, still with minor modifications of the proof of [5], as follows. Let \mathcal{S}_x be a set of cardinal 2 - subsets of \mathcal{M} , and \mathcal{S}_t be a set of cardinal 2 - subsets of $\llbracket 0, N - 1 \rrbracket$. Let $\tilde{T}_{\mathcal{M}, \mathcal{T}} u$ be defined by

$$\tilde{T}_{\mathcal{M}, \mathcal{T}} u = \sum_{n=0}^{N-1} (t_{n+1} - t_n) \sum_{\{K, L\} \in \mathcal{S}_x} \omega_{K, L} |u_L^n - u_K^n| + \sum_{\{p, q\} \in \mathcal{S}_t} \delta_{p, q} \sum_{K \in \mathcal{M}} |K| |u_K^p - u_K^q|, \tag{31}$$

where $(\omega_{K, L})_{\{K, L\} \in \mathcal{S}_x}$ and $(\delta_{p, q})_{\{p, q\} \in \mathcal{S}_t}$ are two sets of non-negative weights. We introduce the two following parameters:

$$\theta_{\mathcal{M}} = \max_{K \in \mathcal{M}} \frac{1}{|K|} \sum_{\substack{L \in \mathcal{M} \\ \{K, L\} \in \mathcal{S}_x}} \omega_{K, L}, \quad \theta_{\mathcal{T}} = \max_{n \in \llbracket 0, N-1 \rrbracket} \frac{1}{t_{n+1} - t_n} \sum_{\substack{p \in \llbracket 0, N-1 \rrbracket \\ \{n, p\} \in \mathcal{S}_t}} \delta_{n, p}. \tag{32}$$

For $\{K, L\} \in \mathcal{S}_x$ and $\{p, q\} \in \mathcal{S}_t$, let

$$\vartheta(\{K, L\}) = \max_{(x, y) \in K \times L} |y - x|, \quad \vartheta(\{p, q\}) = \begin{cases} t_{q+1} - t_p & \text{if } q > p, \\ t_{p+1} - t_q & \text{otherwise} \end{cases}$$

and let

$$\mathfrak{d}(\mathcal{M}) = \max_{\{K, L\} \in \mathcal{S}_x} \mathfrak{d}(\{K, L\}), \quad \mathfrak{d}(\mathcal{T}) = \max_{\{p, q\} \in \mathcal{S}_t} \mathfrak{d}(\{p, q\}).$$

Then the following convergence result holds.

Lemma A.3 *Let $(\mathcal{M}^{(m)})_{m \in \mathbb{N}}$ and $(\mathcal{T}^{(m)})_{m \in \mathbb{N}}$ be a given sequence of meshes and time discretisations. Let us suppose there exists $\theta > 0$ such that $\theta_{\mathcal{M}^{(m)}} \leq \theta$ and $\theta_{\mathcal{T}^{(m)}} \leq \theta$ for all $m \in \mathbb{N}$, with $\theta_{\mathcal{M}^{(m)}}$ and $\theta_{\mathcal{T}^{(m)}}$ given by Eq. (32). Let us assume that $\mathfrak{d}(\mathcal{M}^{(m)})$ and $\mathfrak{d}(\mathcal{T}^{(m)})$ tend to zero when m tends to $+\infty$. Let $u \in L^1(\Omega \times (0, T))$ and $(u_p)_{p \in \mathbb{N}}$ be a sequence of functions of $L^1(\Omega \times (0, T))$ such that $u_p \rightarrow u$ in $L^1(\Omega \times (0, T))$ as $p \rightarrow +\infty$.*

Then $\tilde{T}_{\mathcal{M}^{(m)}, \mathcal{T}^{(m)}} u_p$ defined by (31) tends to zero when m tends to $+\infty$ uniformly with respect to $p \in \mathbb{N}$.

References

1. Ben-Artzi, M., Li, J.: Consistency and convergence of finite volume approximations to nonlinear hyperbolic balance laws. *Math. Comput.* **90**, 141–169 (2021)
2. Elling, V.: A Lax–Wendroff type theorem for unstructured quasi-uniform grids. *Math. Comput.* **76**, 251–272 (2007)
3. Eymard, R., Gallouët, T., Herbin, R.: Finite volume methods. In: Ciarlet, P., Lions, J. (eds) *Handbook of Numerical Analysis*, vol. VII, pp. 713–1020. North Holland (2000). <https://hal.archives-ouvertes.fr/hal-02100732v2/>
4. Eymard, R., Gallouët, T., Herbin, R., Latché, J.C.: Finite volume schemes and Lax–Wendroff consistency. Submitted <https://arxiv.org/abs/2106.06380>
5. Gallouët, T., Herbin, R.: Latché: on the weak consistency of finite volume schemes for conservation laws on general meshes. *SeMA J.* **76**, 581–594 (2019)
6. Godlewski, E., Raviart, P.-A.: Numerical approximation of hyperbolic systems of conservation laws, *Applied Mathematical Sciences*, pp. 118. Springer, New York (1996)
7. Harlow, F., Welsh, J.: Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface. *Phys. Fluids* **8**, 2182–2189 (1965)
8. Herbin, R., Latché, J.-C., Minjeaud, S., Therme, N.: Conservativity and weak consistency of a class of staggered finite volume methods for the Euler equations. *Math. Comput.* **90**, 1155–1177 (2021)
9. Herbin, R., Latché, J.C., Nguyen, T.: Consistent segregated staggered schemes with explicit steps for the isentropic and full Euler equations. *ESAIM: Mathematical Modelling and Numerical Analysis* **52**, 893–944 (2018)
10. Kroner, D., Rokyta, M., Wierse, M.: A Lax–Wendroff type theorem for upwind finite volume schemes in 2-D. *East West J. Numer. Math.* **4**, 279–292 (1996)
11. Lax, P., Wendroff, B.: Systems of conservation laws. *Commun. Pure Appl. Math.* **13**, 217–237 (1960)
12. Leveque, R.: *Finite Volume Methods for Hyperbolic Problems*. Cambridge Texts in Applied Mathematics. Cambridge University Press, Cambridge (2002)
13. Rannacher, R., Turek, S.: Simple nonconforming quadrilateral Stokes element. *Numer. Methods Partial Differ. Equ.* **8**, 97–111 (1992)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.