

An Analysis of Finite Volume, Finite Element, and Finite Difference Methods Using Some Concepts from Algebraic Topology

Claudio Mattiussi

CLAMPCO Sistemi s.r.l.—NIRLAB, AREA Science Park, Padriciano 99, 34012 Trieste, Italy

Received May 23, 1996; revised January 17, 1997

In this paper we apply the ideas of algebraic topology to the analysis of the finite volume and finite element methods, illuminating the similarity between the discretization strategies adopted by the two methods, in the light of a geometric interpretation proposed for the role played by the weighting functions in finite elements. We discuss the intrinsic discrete nature of some of the factors appearing in the field equations, underlining the exception represented by the constitutive term, the discretization of which is maintained as the key issue for numerical methods devoted to field problems. We propose a systematic technique to perform this task, present a rationale for the adoption of two dual discretization grids and point out some optimization opportunities in the combined selection of interpolation functions and cell geometry for the finite volume method. Finally, we suggest an explanation for the intrinsic limitations of the classical finite difference method in the construction of accurate high order formulas for field problems. © 1997 Academic Press

1. INTRODUCTION

To solve a field problem by means of the finite element method (FE), we start by partitioning the domain of the problem into *elements* and assigning a certain number of *nodes* to each element [1]. The starting point of the finite volume method (FV) is very similar, but for the introduction of an additional staggered grid of *cells*, usually defining one cell around each node [2]. Despite the grid differences, a system of equations fully equivalent to the FV one can be obtained with FE using as weighting functions the characteristic functions of FV cells, i.e., functions equal to unity inside the cell and zero outside [3]. Once ascertained that particular FE weighting functions can be used to define FV cells, we can be tempted to ask if a similar role can be ascribed to generic weighting functions. This paper will show that we can answer positively to this question, provided we recognize that the cells defined by generic weighting functions are not necessarily crisp, but can be *spread* (Fig. 1). The discussion to follow, while showing the relevance of algebraic topologic concepts in this field of investigation, throws some light on the nature of both FE and FV methods and underlines some optimization opportunities for the latter. In particular, this paper ex-

plains why it is expedient to use two distinct and dual discretization grids, it shows how they must be staggered to achieve optimal performance and proposes a technique for the construction of high order algorithms which comply with the physics of the problem on regular and irregular grids. A final section devoted to an analysis of the discretization strategies adopted by the finite difference methods (FD) underlines the absence in the classical version of FD [4] of the distinct geometric flavor of FE and FV, suggesting how this reflects in the performance of formulas obtained with it. New approaches to FD [13, 14] are also briefly analyzed and commented

For concreteness, in the course of the exposition we will almost always refer to a steady heat conduction model problem. The choice of thermostatics was made in order to present the results in a context which, for its simplicity, should be familiar to the widest possible audience. Nonetheless, the discussion applies to every field theory whose field equation admits a factorization in the spirit of the one presented in Section 2. Reference [5] shows that a great number of physical equations admit this factorization. Moreover, even if our model problem translates in a boundary value problem, the analysis performed in this paper applies also to problems which translate in initial-boundary value problems, the extension merely requiring the introduction of the necessary time-like geometric objects.

2. A DIGRESSION ON EQUATIONS

2.1. The Factorization of Field Equations

Let us start by defining the terminology adopted in this paper, considering the case of thermostatics. Apart from boundary conditions, the field equation

$$-\operatorname{div}(\lambda \operatorname{grad} T) = \sigma \quad (1)$$

(λ is the thermal conductivity of the medium) constitutes the link between the unknown field (the temperature T) and the given source field (the rate of heat generation σ).

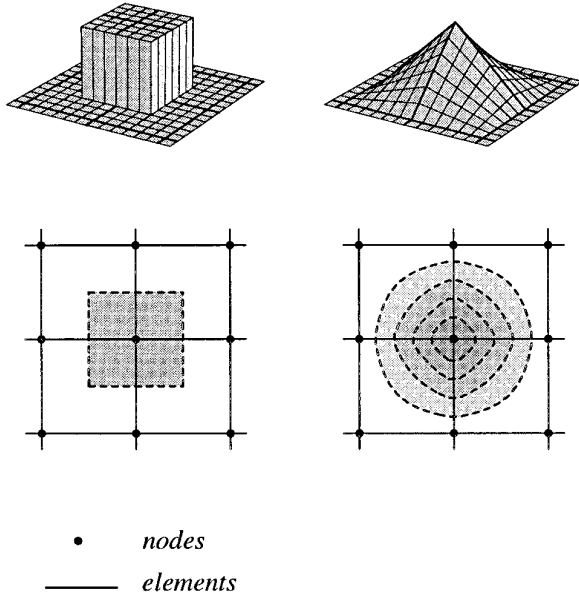


FIG. 1. FE weighting functions and corresponding crisp and spread cells.

It is expedient to factorize the field equation decomposing it into three parts: balance equation, constitutive equation, and kinematic equation (this last name is not standard and has only a mnemonic purpose, inspired by elasticity [1]). An expressive representation for the factorization is achieved with the diagram of Fig. 2 where, to allow the reader to think in terms of a field problem of his choice, conventional names were assigned to the physical quantities.

This paper will show that the structure of the field equation exhibited in Fig. 2 finds a direct correspondence in the strategies adopted by FV and FE to replace Eq. (1) with a system of algebraic equations. In the next three sections each factor of the field equation is examined in this discretization perspective.

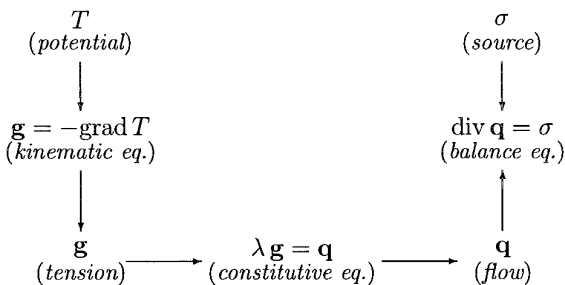


FIG. 2. The factorization diagram for the field equation of thermostatics, showing the terminology adopted in this paper for the fields and the equations.

2.2. The Balance Equation

Given a balance equation in differential terms, its discretization amounts to writing it for domains D of finite extension. For thermostatics the correspondence is

$$\sigma = \text{div } \mathbf{q} \quad \leftrightarrow \quad Q_{\text{source}}(D) = Q_{\text{flow}}(\partial D), \quad (2)$$

(differential or *local*) (discrete or *global*)

where ∂ stands for “the boundary of.” We will call *local* the equations and physical quantities of the differential case and *global* those of the discrete case. Note that the transition from local to global takes place without any approximation error. This happens because the balance equation does not require for its validity uniform fields, homogeneous materials or any other condition holding in general only in domains of infinitesimal extension. Let us express this concept by calling it a *topological equation* to suggest an idea of invariance under arbitrary homeomorphic transformations. For a topological equation, the discrete version appears therefore as the fundamental one, with the differential statement proceeding from it if additional hypotheses are fulfilled.

2.3. The Constitutive Equation

Contrary to the case of the balance equation it is generally impossible to put in discrete form a local constitutive equation without limiting its applicability to particular field configurations. This happens because a general form for a constitutive equation, valid over regions of finite extension, can be given only if the field is supposed uniform and the material homogeneous, which is generally not the case. The exact rendering of constitutive equations requires the use of *metrical* concepts like length, area, volume, and angle, along with terms describing the properties of the medium.

2.4. The Kinematic Equation

The operator appearing in the differential kinematic equation of thermostatics is the gradient, and we know that its discrete counterpart involves a simple difference. It is therefore possible to state the kinematic equation in discrete form:

$$\mathbf{g} = -\text{grad } T \quad \leftrightarrow \quad G(D) = T(\partial D), \quad (3)$$

(differential or *local*) (discrete or *global*)

Here G is the global thermal tension, the domain is a line, and its boundary consists of two points. $T(\partial D)$ is the difference of two temperatures (this will be considered in more detail later). Summing up, we can say that the kinematic equation is a topological equation, since there is no

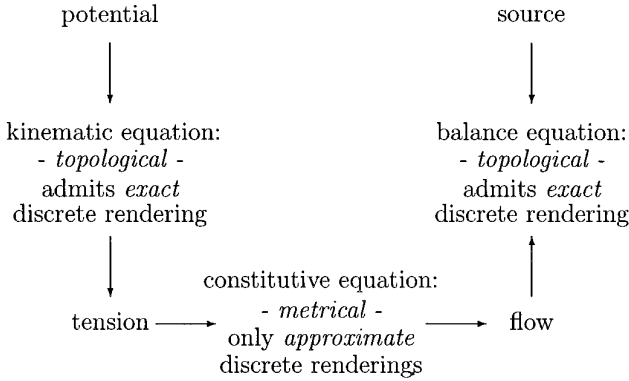


FIG. 3. The discretizability of the factors of the field equation.

approximation involved in its discrete rendering. Figure 3 recapitulates our analysis of the discretizability of each term of the factorized field equation.

3. THE REPRESENTATION OF DISCRETIZED DOMAINS

3.1. Chains

To give a formal enunciation to the discrete version of topological equations, we introduce a tool aimed at the representation of discretized domains: the concept of *chain*. In FV a balance equation is written for each cell of the grid. An implicit orientation of the cells is assumed, usually such that the heat *generated* is to be taken with positive sign (in contrast with the heat *absorbed*). We will say that the cell is oriented as a *source* (opposed to a cell oriented as a *sink*). If we label a particular oriented cell as \mathbf{c} , it is natural to denote the same cell with opposite orientation by $-\mathbf{c}$. In this way we can represent an arbitrary ensemble of oriented cells belonging to the grid, by simply attributing them the coefficient 0 (cell not in the ensemble), 1 (cell in the ensemble with its orientation unchanged) or -1 (cell in the ensemble with its orientation reversed). We can enlarge further the capabilities of this representation by admitting arbitrary integer multiplicities n for cells; this will permit, for example, the representation of a multiple loop (Fig. 4). We obtain a collection of cells with integer multiplicity, an object that in algebraic topology is called a *chain* with integer coefficients [7].

The concept of chain plays a fundamental role in the establishment of a point of view comprising both FV and FE, since we can interpret chains as the discrete counterpart of oriented domains with a weighting function defined over them (abridged below as *weighted domains*). In this way the FE’s weighting functions will acquire a geometric interpretation.

3.2. A Formal Notation for Chains

To proceed in our treatment of chains, we need a reasonably compact notation for them. Let us start by labeling each oriented cell of the grid with a superscript univocally identifying it. The multiplicity with which the cell labeled i appears in a chain will be denoted by n_i . With these choices we can represent a chain as a *formal sum*:

$$\mathbf{C} = \sum_i n_i \mathbf{c}^i. \tag{4}$$

This notation has some link with the intuitive idea of composing a domain by “adding” its parts. A chain is in fact an element of a free module which has the cells as generators and chains generated by a given ensemble of cells can be added, subtracted, and multiplied by integers, allowing the algebraic manipulation of domains. In a 3D space the formal sum (4) can be used to represent ensembles of oriented and weighted volumes, surfaces, lines, and points. To prevent the use of many names for a unique concept, provided those geometric objects satisfy some additional condition (e.g., of being simply connected [7]), topologists speak in all cases of *cells*, prefixing the name with the appropriate dimension number. So volumes become *3-dimensional cells*—in short, *3-cells*—surfaces become *2-cells*, lines *1-cells*, and points *0-cells*. Chains formed with p -cells are called *p-dimensional chains* or *p-chains*. This convention will be adopted and the notation adapted accordingly, writing $\mathbf{c}_{(p)}$ for a p -cell and $\mathbf{C}_{(p)}$ for a p -chain (5):

$$\mathbf{C}_{(p)} = \sum_i n_i \mathbf{c}_{(p)}^i. \tag{5}$$

In conclusion, *chains can be used to represent the discretized geometry of a problem.*

3.3. Grids and Cell-Complexes

Consider a 3D domain discretized by partitioning it into 3-cells. If the discretization is “properly performed” (in a topological sense, that can be easily formalized [7]), two

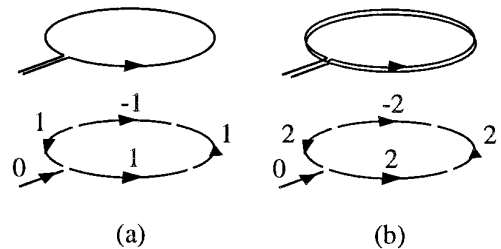


FIG. 4. (a) A single loop represented as an ensemble of oriented lines with multiplicities 0, 1, and -1 . (b) A multiple loop requires generic integer multiplicities.

3-cells intersect on a 2-cell or have an empty intersection, two 2-cells intersect on a 1-cell or have an empty intersection and finally two 1-cells intersect on a 0-cell, i.e., a point, or have no points in common. All these cells of various orders constitute what in algebraic topology is called a tridimensional *cell-complex*, and—when the cells of all orders are oriented—a 3D *oriented cell complex*. The presence of an oriented cell-complex is a prerequisite to the very definition of chains and was assumed implicitly in the former discussion. Note that to be a cell-complex, a discretization grid must satisfy certain conditions (which exclude, for example, overlapping cells). The term *grid* is therefore more general than *cell-complex* and will be used to refer to discretized domains in a broader sense.

3.4. The Boundary of a Chain

The boundary of a domain is a fundamental notion in the enunciation of physical laws and therefore it is advisable to define this concept for the chains.

The boundary $\partial \mathbf{c}_{(p)}$ of an oriented p -cell $\mathbf{c}_{(p)}$ is defined as the $(p - 1)$ -chain composed by the $(p - 1)$ -cells of the cell-complex having nonempty intersection with $\mathbf{c}_{(p)}$, endowed with the orientation *induced* on them by $\mathbf{c}_{(p)}$ [7]. Building on this definition, the linear extension (6) represents the procedure to calculate the boundary of a chain as a combination of its cells' boundaries.

$$\partial \left(\sum_i n_i \mathbf{c}_{(p)}^i \right) = \sum_i n_i (\partial \mathbf{c}_{(p)}^i) \quad (6)$$

This defines the *boundary operator* ∂ , which transforms p -chains in $(p - 1)$ -chains and is compatible with the additive and the (external) multiplicative structure of chains; in other words, it is a *linear transformation* of the module of p -chains into the module of $(p - 1)$ -chains over the same cell-complex:

$$\{\mathbf{C}_{(p)}\} \xrightarrow{\partial} \{\mathbf{C}_{(p-1)}\}. \quad (7)$$

3.5. Boundaries with Internal Vestiges

In this section we will show that the boundary of a chain has certain peculiarities with respect to the traditional geometric notion of boundary of a domain and that only for particular chains do the two concepts coincide.

Consider first two adjacent 2-cells (i.e., having a 1-cell of their boundary in common) with compatible orientation (i.e., inducing on the common 2-cell opposite orientations) and both with multiplicity 1 (Fig. 5a). If we apply the boundary operator to this 2-chain, we find that the common 1-cell does not appear in the result. This telescoping property is a consequence of the opposite orientations induced on the common 1-cell by the two 2-cells. Now consider a

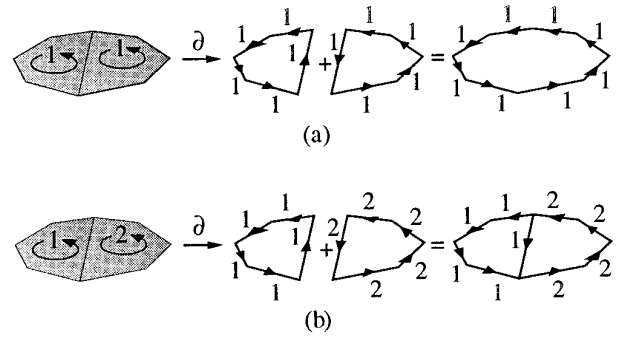


FIG. 5. (a) The boundary of a 2-chain. (b) Appearance of internal vestiges in the boundary of a 2-chain.

chain composed again by two adjacent 2-cells with compatible orientation but this time with different multiplicity (Fig. 5b). When we apply to this chain the boundary operator, the common 1-cell receives from its two adjacent cells different multiplicities, the sum of which does not vanish. To the boundary of our 2-chain, belongs in this case a 1-cell that we are used to considering internal to the subdomain composed by the two 2-cells.

In general, only in the case of p -chains composed by adjacent p -cells with *compatible* orientation and *uniform* multiplicity, the telescoping property works to cancel all the internal p -cells, and the boundary operator generates a $(p - 1)$ -chain which lies on what we are used to considering the boundary of the domain individualized by the ensemble of p -cells appearing with nonnull multiplicity in the p -chain.

The reason for this long digression on a seemingly minor point lies in our desire to interpret the weighted domains as continuous counterparts of chains. To build a complete correspondence, it is mandatory to define (at least implicitly) the boundary of a weighted domain. The present result anticipates that with a weighting function that is not constant (on its support), the boundary of the corresponding weighted domain will appear to be spread over the whole domain.

4. FIELDS AND DISCRETIZED DOMAINS

4.1. The Representation of Fields

FE and FV discretize the domain of the problem; in this perspective we have reviewed some tools allowing a formal description of discretized domains. We need a similar set of tools to describe the *fields* over such domains. The approach adopted can be better understood considering that, from an operative point of view, the continuous representation of a field in terms of a *field function* is but an abstraction. This appears obvious as soon as we realize that from a measurement we always obtain the value of a quantity

associated with a finite region of space, i.e., a *global* quantity (for example, the magnetic flux associated with a surface of finite extension and not the magnetic induction in a point). In this light, a field function defined in a region of space should be considered an abstraction representing the measurements of a given global quantity that can be performed over *all* the suitable (which means “with the proper dimension and kind of orientation”) subdomains of the region. In this perspective, when we discretize a region of space, we are implicitly deciding to consider only a given subset of all its possible subdomains. Consequently we no longer need the full representation of the field—the *local* representation—but we can content ourselves with a representation containing only the *global quantities* we are possibly going to deal with, i.e., those associated with the cells of our discretization. We emphasize the fact that such a representation does *not* constitute or involve any *approximation* of the field.

4.2. Cochains

Given a field problem defined in a discretized region, we have to deal with various fields which, as pointed out in the preceding section, reveal themselves as global quantities associated with suitable p -cells. For example, in 3D thermostatics the source field manifests itself as a rate of heat generation (or absorption) within oriented 3-cells, whereas the flow field manifests itself as a rate of heat flow through oriented 2-cells. This means that to each 2-cell of the complex, the flow field associates a well-defined value of heat (8) and the same happens for each 3-cell as a consequence of the presence of the source field (9):

$$\mathbf{c}_{(2)}^i \xrightarrow{\mathbf{q}} Q_{\text{flow}}^i \quad (8)$$

$$\mathbf{c}_{(3)}^j \xrightarrow{\sigma} Q_{\text{source}}^j \quad (9)$$

We can therefore represent a field on a cell-complex as a function associating global quantities with all the p -cells of the complex having a given value of p and a given kind of orientation (as will be explained in Section 5.2) both characteristic of the field. The global quantities can be scalars (as in thermostatics, where they are values of energy or temperature, and in electromagnetics, where they are values of charge or ratios of action and charge), vectors (as in fluid-dynamics) or other mathematical entities.

Let us examine the properties of these functions. Consider first two adjacent 3-cells with the same orientation. Think of them as a 3-chain over a suitable cell-complex, with multiplicities 1 for these two 3-cells and 0 for all the other 3-cells of the complex. The heat generated within the two cells is the sum of the heat generated within each one (10):

$$\mathbf{c}_{(3)}^j + \mathbf{c}_{(3)}^k \xrightarrow{\sigma} Q_{\text{source}}^j + Q_{\text{source}}^k \quad (10)$$

Inverting the orientation of a cell, the sign of the associated quantity changes (11):

$$-\mathbf{c}_{(3)}^j \xrightarrow{\sigma} -Q_{\text{source}}^j \quad (11)$$

Finally, let us consider the behavior with respect to cells multiplicity. We will present only the following heuristic argument. Given the magnetic field associated with a loop, if the loop is doubled, the field associated with it doubles. In the same spirit it is reasonable to assume that the quantity associated with a cell with multiplicity n , is n times the quantity associated with the same cell with multiplicity 1 (12):

$$n_j \mathbf{c}_{(3)}^j \xrightarrow{\sigma} n_j Q_{\text{source}}^j \quad (12)$$

Putting it all together we obtain

$$\mathbf{C}_{(3)} = \sum_j n_j \mathbf{c}_{(3)}^j \xrightarrow{\sigma} \sum_j n_j Q_{\text{source}}^j = Q_{\text{source}}^{\mathbf{C}} \quad (13)$$

In short, the field associates a global quantity with each cell of the complex; a chain is a weighted sum of cells and therefore the field associates a global quantity with each chain; moreover, this association is linear over the module of all the chains constructed over the same cell-complex. For example, the heat source field σ manifests itself as a linear transformation of the module of the 3-chains into the field of the real numbers:

$$\{\mathbf{C}_{(3)}\} \xrightarrow{\sigma} \mathcal{R}. \quad (14)$$

In algebraic topology, such a transformation is called a *real-valued 3-dimensional cochain* or, in short, a *3-cochain*. To emphasize the joint role of the domain and of the field in the generation of the global quantity, the representation (15) is often used,

$$\langle \mathbf{C}_{(3)}, \mathbf{Q}_{\text{source}}^{(3)} \rangle = Q_{\text{source}}^{\mathbf{C}}, \quad (15)$$

where $\mathbf{C}_{(3)}$ is a 3-chain and $\mathbf{Q}_{\text{source}}^{(3)}$ is the heat source 3-cochain, which is the representation of the source field over the cell-complex. In the same spirit, the heat flow field \mathbf{q} manifests itself as a real-valued 2-cochain $\mathbf{Q}_{\text{flow}}^{(2)}$,

$$\langle \mathbf{C}_{(2)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle = Q_{\text{flow}}^{\mathbf{C}}, \quad (16)$$

and, in general, a field manifests itself on a cell-complex as a (not necessarily real-valued) p -cochain. We can re-

phrase all this, saying that *cochains constitute a representation for fields over discretized domains*. As field functions can be added and multiplied by a scalar, so can cochains defined over a same cell-complex. Collectively all these cochains constitute a module.

5. THE REPRESENTATION OF TOPOLOGICAL EQUATIONS

5.1. General Remarks

Consider the two topological equations of thermostatics (Fig. 2). In their discrete form, both equations assert the existence of a relation between a global quantity associated with an oriented domain and another global quantity associated with the boundary of that domain (Eqs. (2) and (3)). On a discretized domain, if we represent a volume as a chain $\mathbf{V}_{(3)}$, the source field as a cochain $\mathbf{Q}_{\text{source}}^{(3)}$, and the flow field as a cochain $\mathbf{Q}_{\text{flow}}^{(2)}$, we can write the balance equation of thermostatics as

$$\langle \mathbf{V}_{(3)}, \mathbf{Q}_{\text{source}}^{(3)} \rangle = \langle \partial \mathbf{V}_{(3)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle. \quad (17)$$

Similarly, the kinematic equation asserts the equality of the tension associated with an oriented line and of a combination of the potentials associated with the two oriented points which constitute its boundary. Since the concept of an *oriented point* may sound unfamiliar to the reader, we shall consider to some greater extent the concept of orientation.

5.2. Internal and External Orientation

Deriving the boundary of a cell and writing a topological equation requires the concept of orientation *induced* by an oriented domain on its boundary. This concept implies the possibility of comparing the orientation of the domain with the orientation of the boundary. For example, in a 3D ambient space the source/sink orientation of a 3-cell (which is, in fact, an orientation with dichotomic symbols meaning “in” and “out”) can be compared with the “through” direction which constitutes the orientation of the 2-cells lying on its boundary. To calculate the boundary of a 2-cell oriented with a “through” direction, its boundary 1-cells must be oriented by means of a *sense of rotation* around them, and this kind of orientation of a 1-cell can be compared with that of 0-cells endowed with a tridimensional *screw-sense*, or *vortex* (Fig. 6a). Similarly, to associate a thermal tension with a 1-cell, the cell must be oriented by means of a running direction along it. The boundary of such a 1-cell must be 0-cells with source/sink orientation. This kind of oriented 1-cell can be the boundary of a 2-cell oriented by means of a sense of rotation on it and, in turn, this 2-cell can be the boundary of a 3-cell oriented by a tridimensional vortex (Fig. 6b).

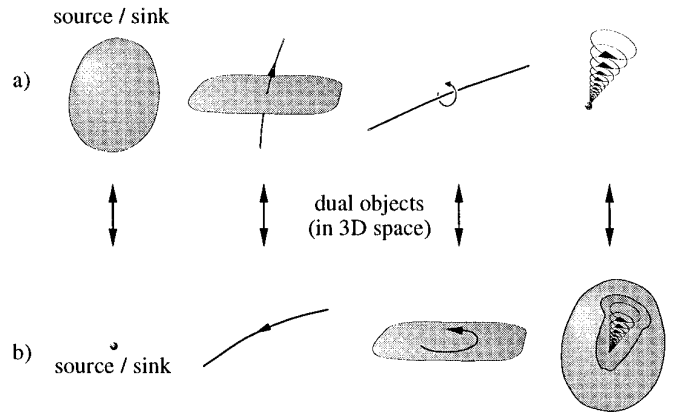


FIG. 6. (a) External orientation for 3-, 2-, 1-, and 0-cells. (b) Internal orientation for 0-, 1-, 2-, and 3-cells.

In conclusion, we have *two kinds of orientations* for geometric objects: those of Fig. 6b are called *internal* orientations; those of Fig. 6a are called *external* orientations. The same distinction holds in every “well-behaved” n -dimensional ambient space [9]. Note that the symbol which gives internal orientation to p -cells gives by definition external orientation to $(n - p)$ -cells. This fact permits the erection in an n -dimensional ambient space of two dual cell complexes, with each p -cell with internal orientation of the first matched by a $(n - p)$ -cell with external orientation of the second.

From the existence of two kinds of orientation follows the need of two discretization grids, one with internal orientation and the other with external orientation (the importance of this distinction of orientations and the benefits deriving from the adoption of two dual cell-complexes as grids, are usually underestimated [17]). To refer compactly and unambiguously to them, the former will be called *primary grid* or—when the term “cell-complex” applies—*primary cell-complex*; the latter will be called *secondary grid* or *secondary cell-complex*. Correspondingly we will speak of primary and secondary cells, chains, and cochains.

Let us consider how this distinction of orientation applies to some familiar case. In 3D thermostatics, the temperature is associated with internally oriented points, the thermal tension with internally oriented lines, the heat flow with externally oriented surfaces, and the heat generation with externally oriented volumes. In 3D magnetostatics the vector potential is associated with internally oriented lines, the magnetic induction with internally oriented surfaces, the magnetic field with externally oriented lines, and the charge current with externally oriented surfaces. Note that a topological equation always involves quantities associated with domains that, being one the boundary of the other, have the same kind of orientation. Later we will observe that constitutive equations link quantities associated with domains having different kinds of orientation.

We can now resume the discussion concerning the kinematic equation of thermostatics. If we represent an internally oriented line as a chain $\mathbf{L}_{(1)}$, the thermal field as a cochain $\mathbf{T}^{(0)}$, and the thermal tension field as a cochain $\mathbf{G}^{(1)}$, we can write the discrete kinematic equation (3) as

$$\langle \mathbf{L}_{(1)}, \mathbf{G}^{(1)} \rangle = \langle \partial \mathbf{L}_{(1)}, \mathbf{T}^{(0)} \rangle. \quad (18)$$

Incidentally, note that since the line induces a source-orientation on its starting point and a sink-orientation on its endpoint, if originally the points are sink-oriented (as is implicit in *calculus* and, therefore, in the definition of the gradient operator), we have

$$\langle \partial \mathbf{L}_{(1)}, \mathbf{T}^{(0)} \rangle = T_{\text{endpoint}} - T_{\text{startingpoint}}. \quad (19)$$

Due to the fact that in *heat theory*, the points to which we associate temperatures are source-oriented, we have instead

$$\langle \partial \mathbf{L}_{(1)}, \mathbf{T}^{(0)} \rangle = -(T_{\text{endpoint}} - T_{\text{startingpoint}}). \quad (20)$$

The possibility of this difference in the default orientation assigned to points in mathematics and in physics, is the reason for the presence of the minus sign in kinematic equations like $\mathbf{g} = -\text{grad } T$ in thermostatics, $\mathbf{E} = -\text{grad } \varphi$ in electrostatics, and of many other “minus” signs appearing in textbooks of physics.

5.3. The Coboundary of a Cochain

A topological equation asserts the equality of two global quantities, one associated with a geometric object and the other with its boundary. For example, on a discretized domain we can write the heat balance equation (2) as

$$\langle \mathbf{c}_{(3)}, \mathbf{Q}_{\text{source}}^{(3)} \rangle = \langle \partial \mathbf{c}_{(3)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle \quad \forall \mathbf{c}_{(3)} \in \text{Cell complex}. \quad (21)$$

We can write this topological equation in a more compact way by *defining* a 3-cochain $\delta \mathbf{Q}_{\text{flow}}^{(2)}$ which satisfies the relation

$$\langle \mathbf{c}_{(3)}, \delta \mathbf{Q}_{\text{flow}}^{(2)} \rangle \stackrel{\text{def}}{=} \langle \partial \mathbf{c}_{(3)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle \quad \forall \mathbf{c}_{(3)} \in \text{Cell complex}. \quad (22)$$

In other words, the 3-cochain $\delta \mathbf{Q}_{\text{flow}}^{(2)}$ associates with each 3-cell the rate of heat flow that the 2-cochain $\mathbf{Q}_{\text{flow}}^{(2)}$ associates with the boundary of that 3-cell. This definition allows the restatement of (21) in the simpler form (23), where we no longer need to quote explicitly the p -cells nor to assert “for all 3-cells,” since both are implicit in the cochain concept:

$$\mathbf{Q}_{\text{source}}^{(3)} = \delta \mathbf{Q}_{\text{flow}}^{(2)}. \quad (23)$$

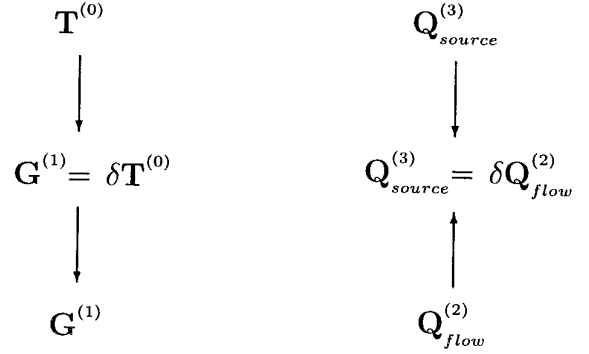


FIG. 7. The discrete representation of topological equations employing the concept of cochain and the definition of the coboundary operator. The left and right columns refer to quantities associated with objects endowed with internal and external orientation, respectively.

We can extend the definition of δ from single cells to generic chains of any order

$$\langle \mathbf{C}_{(p+1)}, \delta \mathbf{C}^{(p)} \rangle \stackrel{\text{def}}{=} \langle \partial \mathbf{C}_{(p+1)}, \mathbf{C}^{(p)} \rangle. \quad (24)$$

Equation (24) defines an *operator* δ , which is called the *coboundary operator* and transforms p -cochains in $(p + 1)$ -cochains (25). It can be shown [7] that it is a linear transformation of the module of p -cochains into the module of $(p + 1)$ -cochains over the same cell-complex:

$$\{\mathbf{C}^{(p)}\} \xrightarrow{\delta} \{\mathbf{C}^{(p+1)}\}. \quad (25)$$

Making use of the definitions of cochain and coboundary, we can redraw the diagram of the factorized field equation (Fig. 2) with an explicit discrete representation for the fields and the topological equations (Fig. 7).

Note that we have defined the coboundary operator without reference to any differential operator. The connection will be established in the next section.

5.4. Coboundary and Differential Operators

In vector analysis there are familiar differential operators that act as the coboundary does. For example, the operator *divergence* transforms a (vector) field that can be integrated over oriented surfaces in a (pseudoscalar) field that can be integrated over oriented volumes, allowing the substitution of

$$\int_V \sigma \, dv = \int_{\partial V} \mathbf{q} \cdot \mathbf{ds} \quad \forall V \subset \text{Domain} \quad (26)$$

with (compare with (21) and (23), respectively)

$$\sigma = \text{div } \mathbf{q}. \quad (27)$$

The operators *curl* and *gradient* act in an analogous way for the transition from oriented lines to oriented surfaces and from oriented points to oriented lines, respectively. Therefore the coboundary operator can be considered as the discrete counterpart of the three differential operators *grad*, *curl*, *div*. It indeed satisfies the property $\delta \delta \equiv 0$ (corresponding to $\text{curl grad} \equiv 0$ and $\text{div curl} \equiv 0$) and its converse which, on a simply p -connected complex, leads to the construction of a “potential” [7]. On a pair of n -dimensional *dual cell-complexes* the coboundary operator acting between p and $(p + 1)$ cochains of the primary complex is the *adjoint*—relative to a natural duality between cochain spaces (see the Appendix)—of the coboundary acting between $(n - p - 1)$ and $(n - p)$ cochains of the secondary complex (this corresponds, for example, to the adjointness of $-\text{grad}$ and div relative to the duality of spaces established by $\int_D T \sigma \, dv$ and $\int_D \mathbf{g} \cdot \mathbf{q} \, dv$). Finally, the very definition of the coboundary guarantees that *conservation* of physical quantities possibly expressed by the topological equation is preserved in the discrete equation. This means that the use of the coboundary operator to render in discrete form the topological equations preserves in the discrete operator properties that other approaches—for example, the Support-Operators method [13, 14]—are obliged to enforce explicitly.

To complete the parallelism between the coboundary and the differential operators, we should substitute a *weighted volume* to V in Eq. (26). We will show in Sections 5.5 and 5.7 that integration by parts as is used in FE exactly fulfills this goal.

5.5. Balance Equations in FV and FE

For FV, the enforcement of the 3D heat balance equation

$$\int_{\partial V} \mathbf{q} \cdot d\mathbf{s} = \int_V \sigma \, dv \quad \forall V \subset \mathcal{D} \quad (28)$$

simply amounts to its application to each 3-cell of the grid. With the discrete symbolism this amounts to

$$\langle \partial \mathbf{c}_{(3)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle = \langle \mathbf{c}_{(3)}, \mathbf{Q}_{\text{source}}^{(3)} \rangle \quad \forall \mathbf{c}_{(3)}. \quad (29)$$

Note that only *global quantities* are involved in these equations; *there is no need to work at this point with average values of fields over cells, nor to involve in the topological equations the extensions of cells*, as is often done in FV practice [10].

FE takes a small detour to discretize the balance equation. In the case of weighted residues, the method starts from the local version of the balance equation

$$\text{div } \mathbf{q} = \sigma \quad (30)$$

and enforces the corresponding weighted residual equation for each node n of the grid, endowed with a given weighting function w^n with support $\mathbf{d}_{(3)}^n$:

$$\int_{\mathbf{d}_{(3)}^n} w^n \text{div } \mathbf{q} \, dv = \int_{\mathbf{d}_{(3)}^n} w^n \sigma \, dv \quad \forall n. \quad (31)$$

The next step is the integration by parts of the left-hand side of (31):

$$\int_{\partial \mathbf{d}_{(3)}^n} w^n \mathbf{q} \cdot d\mathbf{s} - \int_{\mathbf{d}_{(3)}^n} \text{grad } w^n \cdot \mathbf{q} \, dv = \int_{\mathbf{d}_{(3)}^n} w^n \sigma \, dv. \quad (32)$$

Let us give a geometric interpretation to Eq. (31) which parallels the obvious one of Eq. (28). Equation (31) presents integrations over a domain $\mathbf{d}_{(3)}^n$ endowed with a weighting function; think of this weighted domain as a chain composed by infinitesimal cells with real multiplicity. To underline this interpretation we can write (31) in the form

$$\int_{w^n \mathbf{d}_{(3)}^n} \text{div } \mathbf{q} \, dv = \int_{w^n \mathbf{d}_{(3)}^n} \sigma \, dv \quad \forall n. \quad (33)$$

In the spirit of the divergence theorem we can write (33) in the form

$$\int_{\partial(w^n \mathbf{d}_{(3)}^n)} \mathbf{q} \cdot d\mathbf{s} = \int_{w^n \mathbf{d}_{(3)}^n} \sigma \, dv \quad \forall n. \quad (34)$$

Equation (34) *constitutes the balance equations as written by FE*. The undefined left-hand side term—an integral over the boundary of a weighted domain—can be defined by comparing Eqs. (34) and (32):

$$\int_{\partial(w^n \mathbf{d}_{(3)}^n)} \mathbf{q} \cdot d\mathbf{s} \stackrel{\text{def}}{=} \int_{\partial \mathbf{d}_{(3)}^n} w^n \mathbf{q} \cdot d\mathbf{s} - \int_{\mathbf{d}_{(3)}^n} \text{grad } w^n \cdot \mathbf{q} \, dv. \quad (35)$$

Although the definition is an implicit one, it is apparent that, as anticipated in Section 3.5, the boundary of a weighted domain is spread over the whole domain $\mathbf{d}_{(3)}^n$, except when the weighting function is constant on $\mathbf{d}_{(3)}^n$, in which case the term containing $\text{grad } w^n$ on the right-hand side of Eq. (35) vanishes and no “internal vestiges” remain.

Comparison of (34) with (29) shows that FE writes balance equations that are very similar in spirit to the ones written by FV; while FV writes them on crisp cells, FE does it on spread cells (this parallelism can be a starting point for enquires on the issue of quantity conservation in FE). Note that FV writes the balance equation in terms of global quantities, whereas FE, making use of spread cells with spread boundary, is forced to use field functions defined over the whole domain. Despite this difference, neither FV nor FE discretize the operator representing the

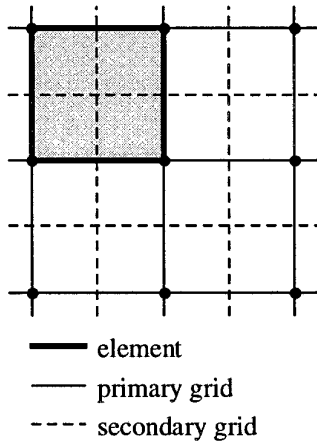


FIG. 8. The two grids and the elements mesh used in FV (the orientation of cells is not represented).

local version of the balance equation. Instead they resort both to the *global* version of it, and with reason, since a topological equation applies directly to regions with finite extension.

5.6. The Two Discretization Grids and the Elements Mesh

FV makes use of two staggered discretization grids (Fig. 8) while, apparently, FE does not make use of two grids. In fact, FE achieves a similar goal by means of the weighting functions which define a spread cell around each node (Fig. 1). Note that spread cells relative to different nodes of the same element, overlap, while the definition of cell-complexes does not admit cell overlapping. In any case, even if the characterization of a secondary grid is for them far from complete, the charge brought against FE of “reducing all to nodes” must be reconsidered. In many cases, quantities apparently referred to nodes are, in fact, associated with the spread cells surrounding the nodes or to their boundary. Weighting functions should not be considered merely as analytical tools necessary to calculate residues, since they appear endowed with a significant geometric meaning.

In addition to the discretization grids, an additional geometric structure emerges from the distinction that must be made between *cells* and *elements*. Cells are expedient to discrete field representation, for we associate them with global quantities. Elements, as will become clear in Section 6.1, constitute the approximation regions used to perform the transition from global quantities to local field representations, required by the FE *and also by the FV* discretization technique of constitutive equations. In an n -dimensional ambient space, elements and primary n -cells often coincide, but it may happen—especially if the elements have internal nodes—that each element be composed by the union of more than one n -cell (Fig. 11). So we have

actually a superposition of two distinct structures, a cell grid (usually the primary) and the *elements mesh*.

5.7. The Role of Integration by Parts in FE

In Section 5.5 an implicit definition for the boundary of a tridimensional weighted domain was given. By means of the identities employed to perform integration by parts, similar formulas for the other cases can be easily obtained. For 3D problems the boundary of weighted 2D and 1D domains are implicitly defined by (36) and (37):

$$\int_{\partial(w^n \mathbf{a}_{(2)}^n)} \mathbf{A} \cdot d\mathbf{l} \stackrel{\text{def}}{=} \int_{\partial \mathbf{a}_{(2)}^n} w^n \mathbf{A} \cdot d\mathbf{l} - \int_{\mathbf{a}_{(2)}^n} \text{grad } w^n \times \mathbf{A} \cdot d\mathbf{s} \quad (36)$$

$$\int_{\partial(w^n \mathbf{a}_{(1)}^n)} T d\mathbf{p} \stackrel{\text{def}}{=} \int_{\partial \mathbf{a}_{(1)}^n} w^n T d\mathbf{p} - \int_{\mathbf{a}_{(1)}^n} T \text{grad } w^n \cdot d\mathbf{l}. \quad (37)$$

In both cases the ‘internal’ part vanishes only if the weighting function is constant on \mathbf{d}^n . The deduction of the corresponding formulas for 2D and 1D ambient spaces is straightforward.

The fundamental role played in FE by the technique of integration by parts appears, therefore, as a manifestation of the necessity to operate with the boundary of spread cells in order to write topological equations. Under this light Eqs. (38)–(40) below play in FE the role played in FV by the generalized Stokes theorem (41):

$$\int_{w^n \mathbf{a}_{(3)}^n} \text{div } \mathbf{q} d\mathbf{v} = \int_{\partial(w^n \mathbf{a}_{(3)}^n)} \mathbf{q} \cdot d\mathbf{s} \quad (38)$$

$$\int_{w^n \mathbf{a}_{(2)}^n} \text{curl } \mathbf{A} \cdot d\mathbf{s} = \int_{\partial(w^n \mathbf{a}_{(2)}^n)} \mathbf{A} \cdot d\mathbf{l} \quad (39)$$

$$\int_{w^n \mathbf{a}_{(1)}^n} \text{grad } T \cdot d\mathbf{l} = \int_{\partial(w^n \mathbf{a}_{(1)}^n)} T d\mathbf{p} \quad (40)$$

$$\langle \mathbf{C}_{(p+1)}, \delta \mathbf{C}^{(p)} \rangle = \langle \partial \mathbf{C}_{(p+1)}, \mathbf{C}^{(p)} \rangle. \quad (41)$$

As a final remark, observe that integrals defined over weighted domains should be considered *Stieltjes’ integrals*. Lebesgue made the point in his celebrated lectures [12] that, owing to its profound geometric meaning, Stieltjes’ integral—and the implied concept of quantities associated with geometric objects—should be the tool of choice for mathematical modeling in physics.

6. THE DISCRETIZATION OF CONSTITUTIVE EQUATIONS

6.1. General Remarks

Constitutive equations connect the left and right columns of the factorization diagrams and represent therefore a bridge between field variables associated with primary cells and field variables associated with secondary cells.

Topological equations and constitutive equations together compose the complete field equation and we know that from a discretized formulation of a field problem we usually obtain only an approximation to the solution obtainable (at least potentially) from the local formulation. Therefore, having stated that topological equations admit exact discrete representation, we can conclude by means of a *reductio ad absurdum* that the discretization of constitutive equations implies necessarily some approximation in the transition from the local to the discrete model. The reasons behind this inevitable error can be understood by trying a direct discretization of the local constitutive equation of thermostatics for an isotropic material:

$$\lambda \mathbf{g} = \mathbf{q}. \quad (42)$$

Calling Q_{flow} the rate of heat flow through a secondary 2-cell $\mathbf{c}_{(2)}^s$ of extension S and G the thermal tension across a primary 1-cell $\mathbf{c}_{(1)}^p$ of extension L , which is the dual of $\mathbf{c}_{(2)}^s$ and shares its orientation, the simpler discrete relation which mimics (42) is

$$\lambda \frac{G}{L} = \frac{Q_{\text{flow}}}{S}. \quad (43)$$

This relation holds in general only if the material is homogeneous, the fields are uniform, $\mathbf{c}_{(2)}^s$ is planar, $\mathbf{c}_{(1)}^p$ is straight, and $\mathbf{c}_{(1)}^p$ is orthogonal to $\mathbf{c}_{(2)}^s$. Note that in infinitesimal regions and away from abrupt material discontinuities, all these conditions are automatically satisfied (except the orthogonality of cells, which is separately expressed by the parallelism of \mathbf{g} and \mathbf{q} implicit in (42)); this accounts for the success of local representations in physics.

Equation (43) and its validity conditions remind us that to write constitutive equations we have to take care—along with the properties of the medium—of extension, curvature, relative position, and other metrical characteristics of cells. As noted in Section 2.3, constitutive equations possess a *metrical* nature, not merely a topological one.

6.2. Local Representations

A local representation for fields which reflects the natural association of field quantities with geometric objects endowed with internal and external orientation is given by ordinary and twisted *differential forms*, with the operators appearing in topological equations represented by the exterior differential d [8]. For example, in thermostatics the potential can be represented by ordinary 0-forms, the tension by ordinary 1-forms, the flow by twisted 2-forms, and the source by twisted 3-forms. For historical reasons the representation most widely adopted for fields with scalar global quantities uses instead scalars, (contravariant) vectors, and the three differential operators *grad*, *curl*, *div*.

The transition from the representation in terms of forms to the usual one requires the introduction of two reference entities: a *metric tensor* (from which derives also a *unit volume*) and a *screw-sense* [9].

Due to the metrical nature of constitutive equations we know that a metric tensor must appear in their local representation. A real algebraic relation between the fields—a constitutive tensor composed by the metric tensor and some material parameters—is often considered a general enough representation, but for example, it is not adequate when the medium involves a nonlocal relation between the fields. The presence of the metric tensor in both the transition between the representations of fields and the constitutive equations, along with the erroneous assumption that constitutive equations are always representable as algebraic relations, induce some authors to represent both with the same operator \star . If this gives an elegant mathematical appearance to the formalism obtained, it is obviously a nonsense that confuses a mathematical manipulation void of physical content, with the transformation representing the medium, depriving the constitutive equations of their physical content. The use of the star operator to represent constitutive equations and of the metric conjugate δ of d to write equations of physics are both examples of this kind of confusion [17].

6.3. The Structure of Constitutive Equations in FV

Consider the constitutive equation of thermostatics. In the local representation it is a transformation Λ which links the field functions \mathbf{g} and \mathbf{q} and translates the properties of the medium. As pointed out in the previous section, Λ is a generic transformation, not necessarily representable with a tensor:

$$\mathbf{g} \xrightarrow{\Lambda} \mathbf{q}. \quad (44)$$

Correspondingly, in the FV formulation the objects put in relation by the constitutive equation are the two cochains $\mathbf{G}^{(1)}$ and $\mathbf{Q}_{\text{flow}}^{(2)}$,

$$\mathbf{G}^{(1)} \xrightarrow{\Lambda} \mathbf{Q}_{\text{flow}}^{(2)}. \quad (45)$$

A natural representation of a cochain $\mathbf{C}^{(p)}$ is the vector of its values $\langle \mathbf{c}_{(p)}^i, \mathbf{C}^{(p)} \rangle$, i.e., the global quantities associated with the p -cells of the complex over which $\mathbf{C}^{(p)}$ is defined. With this representation, if n_1 is the number of primary 1-cells and n_2 that of secondary 2-cells, (45) becomes

$$\begin{pmatrix} G_1 \\ \vdots \\ G_{n_1} \end{pmatrix} \xrightarrow{\Lambda} \begin{pmatrix} Q_1 \\ \vdots \\ Q_{n_2} \end{pmatrix} \quad (46)$$

which corresponds to n_2 functions Λ_i expressing the heat flow through each secondary 2-cell as a function of the thermal tensions across primary 1-cells. The simplest yet nontrivial form that (46) can assume is that of a linear transformation:

$$\sum_{j=1}^{n_1} \Lambda_{ij} G_j = Q_i, \quad i = 1, \dots, n_2. \quad (47)$$

The generic coefficient Λ_{ij} of (47) expresses the influence of the thermal tension G_j across the j th primary 1-cell $\mathbf{c}_{(1)}^j$ on the rate of heat flow Q_i through the i th secondary 2-cell $\mathbf{c}_{(2)}^i$. There are actually good reasons to put to zero the greater part of the coefficients Λ_{ij} ; the greater the number of zero coefficients in (47), the greater the sparseness of the matrix of the global system of equations. On the other hand, the greater the number of terms involved, the better we can expect to approximate the local constitutive equation with the discrete link Λ .

6.4. The Structure of Constitutive Equations in FE

To write FE balance equations we need a local representation of the flow. As a consequence, in the model problem, the objects put in relation by the constitutive equation will no longer be the two cochains $\mathbf{G}^{(1)}$ and $\mathbf{Q}_{\text{flow}}^{(2)}$, but the cochain $\mathbf{G}^{(1)}$ and the field function \mathbf{q} ,

$$\mathbf{G}^{(1)} \xrightarrow{\Lambda} \mathbf{q}. \quad (48)$$

The first term of the link is a discrete representation, whereas the second is a local one; the link appears as a kind of approximation of a field function. For this reason we avoid writing the FE link in general form, waiting for the definition of cochain approximation to do this.

6.5. Strategies for Constitutive Equation Discretization

In principle, any conceivable method capable of supplying a set of coefficients Λ_{ij} for Eq. (47) is acceptable as discretization strategy. For example, one could run a genetic algorithm having the values of the coefficients Λ_{ij} as parameters and the fitness measure of each individual linked to the errors resulting from the solution of a given ensemble of ‘‘adaptation’’ problems. Note that with the approach presented in this paper, once the domain is discretized the rendering of topological equations is univocally determined (in terms of coboundary). Therefore the only degree of freedom left is the choice of the discretization technique for the constitutive equations. Still, this freedom permits the construction of many algorithms, with different matrix structures and computational properties. In the following pages we will show how to approach in general terms the discretization of constitutive equations,

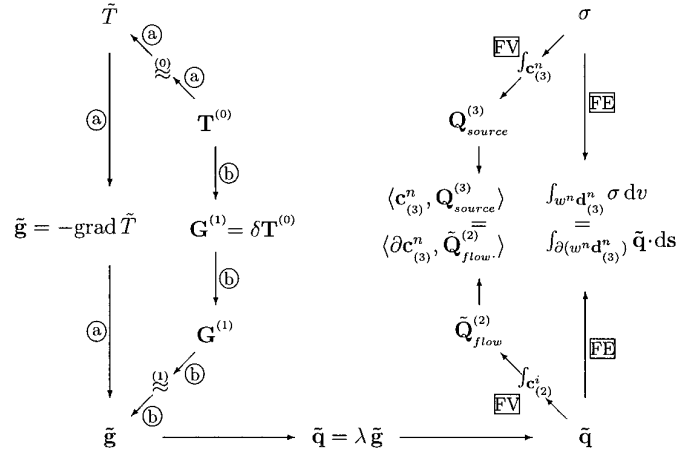


FIG. 9. Alternative FV/FE strategies for the discretization of the constitutive equation of thermostatics: (a) 0-cochain approximation followed by differential kinematic equation; (b) discrete kinematic equation followed by 1-cochain approximation.

complementing the exposition with a couple of elementary examples. We will not undertake the analysis of the computational properties of the algorithms obtained in these examples.

To derive the discretized constitutive equation, FV and FE resort explicitly to the local constitutive equation. For a 3D thermostatics problem, the two methods proceed as follows: an approximation $\tilde{\mathbf{g}}$ of the field function \mathbf{g} is obtained directly from the cochain $\mathbf{G}^{(1)}$ or indirectly (via the differential kinematic equation) from the cochain $\mathbf{T}^{(0)}$. The local constitutive equation $\lambda \mathbf{g} = \mathbf{q}$ is applied to $\tilde{\mathbf{g}}$ to obtain an approximate flow density $\tilde{\mathbf{q}}$. Note that from the previous step we obtain an actual expression (with the coefficients of $\mathbf{T}^{(0)}$ as parameters) for the field function $\tilde{\mathbf{g}}$ and we can apply to this expression a generic transformation, for example a tensor or an integral relation. Given $\tilde{\mathbf{q}}$, FV integrates it over secondary 2-cells, to obtain the cochain $\tilde{\mathbf{Q}}_{\text{flow}}^{(2)}$ which is necessary for the enforcement of the balance equation $\delta \tilde{\mathbf{Q}}_{\text{flow}}^{(2)} = \mathbf{Q}_{\text{source}}^{(3)}$; whereas FE writes the balance equation in terms of $\tilde{\mathbf{q}}$, of the source density σ , and of the weighting functions w which define the shape of the spread cells (Fig. 9).

The final result is the establishment of a direct link expressing the flow cochain $\tilde{\mathbf{Q}}_{\text{flow}}^{(2)}$ (for FV) or the flow density $\tilde{\mathbf{q}}$ (for FE) as a function of the unknown cochain $\mathbf{T}^{(0)}$. The fact that neither the intermediate fields $\tilde{\mathbf{g}}$ and $\tilde{\mathbf{q}}$, nor the cochains $\mathbf{G}^{(1)}$ and $\tilde{\mathbf{Q}}_{\text{flow}}^{(2)}$ appear explicitly in the final expression of the balance equations (with the exception of their role in the enforcement of boundary conditions, which will be considered later) explains why the detailed nature of the discretization process and of its main obstacle—the constitutive equation—tend to remain hidden. Indeed, even in FV, an explicit expression for the discrete constitutive equation linking $\mathbf{G}^{(1)}$ to $\tilde{\mathbf{Q}}_{\text{flow}}^{(2)}$, is seldom given.

The reader should be aware of the fact that, depending on the theory involved and the phenomena considered, constitutive equations can link quantities other than tension and flow (using the conventional names of Fig. 2) and that more than one constitutive equation can appear simultaneously in the field equation. Moreover, the quantity occupying the position of the tension in the diagram might be no longer associated with 1-cells (this happens, for example, in 3D magnetostatics, where the magnetic flux—associated with primary 2-cells—takes the place of our conventional tension). In each of these cases the discretization procedure sketched above, which starts with the reconstruction from cochains of an approximated field function, still applies with minor changes. The fundamental step is always the approximation of a continuous representation of a field associated with p -dimensional domains, based on the information constituted by a p -cochain: an operation indicated in Fig. 9 with the symbol $\overset{(p)}{\approx}$, and called from now on *p -cochain-based field function approximation* or, concisely, *p -cochain approximation*.

6.6. 0-Cochain Approximation and Equation Assembly

For both FV and FE the approximation based on a 0-cochain is the usual point interpolation. Remember that the kind of mathematical object associated with each p -cell (a scalar, a vector, etc.) depends on the kind of cochains we are dealing with (scalar-valued, vector-valued, etc.). For example, in electromagnetics the global quantities are scalars and therefore only scalar-valued p -cochain approximation should be performed, avoiding instead the much used interpolations based on nodal values of local vector quantities like \mathbf{E} and \mathbf{H} . As done until now, only scalar-valued cochains will be considered in the following.

In the case of thermostatics, the 0-cochain we start with is $T^{(0)}$, i.e., the temperature values on primary 0-cells (which correspond to the nodes of the FE terminology). Within each element we construct an approximation of the field function T by means of *interpolation functions* τ_e , where the subscript e indicates an approximation holding only within a particular element. As anticipated in Section 3.3, *elements* play the role of *approximation regions* for the construction of a continuous representation \tilde{T} for the field function. The application to \tilde{T} of the differential kinematic equation first, and of the local constitutive equation next, is performed element by element. Note that a separate expression for \tilde{T} , $\tilde{\mathbf{g}}$, and $\tilde{\mathbf{q}}$ is available within each element, but that to write the balance equations we must consider instead *secondary n -cells* (n is the dimension of the problem), which can overlap more than one element. In the case of FV, it is therefore advisable to distinguish between secondary n -cells which are completely contained in a single element (let us call them *interior cells*), and secondary n -cells which lie across neighboring elements (Fig. 12). For

FE, this corresponds to the distinction between nodes, or better, spread secondary n -cells, which are *interior* to the element, and spread n -cells which lie on its boundary and are therefore in common with other elements (with the exception of nodes lying on the boundary of the problem domain). To write the balance equation for an interior secondary n -cell (crisp for FV, spread for FE) we need only the approximate flow density calculated within the element which contains the cell. Conversely, for secondary n -cells shared between elements, we have to *assemble* the contributions derived from the flow densities calculated within all the elements involved. As a consequence, to interior cells correspond balance equations involving a minimum number of terms, namely only those corresponding to the temperatures of the nodes of a single element, whereas a very “shared” cell is characterized by a large number of nodal temperatures appearing in the corresponding balance equation.

6.7. p -Cochain Approximation

The idea behind p -cochain approximation is the natural extension to global quantities associated with p -cells, of the idea upon which 0-cochain approximation is based. For example, taking the 1-cochain $\mathbf{G}^{(1)}$ as starting point, we have the global quantity “thermal tension” on primary 1-cells and from this knowledge we want to build an approximation to the field function \mathbf{g} . To this purpose, within each element we assign an *approximation function* Φ_e , taking care to select it with the same vectorial nature of \mathbf{g} (note once again that it would be more appropriate to represent \mathbf{g} as a differential 1-form). The number of degrees of freedom of Φ_e must be equal to the number of primary 1-cells within the e th element. If the approximation function is properly selected in relation to the element’s shape, the values of $\mathbf{G}^{(1)}$ on the primary 1-cells belonging to the element determine univocally the approximating function $\tilde{\mathbf{g}}$ within it. In other words, while in 0-cochain approximation we ask the approximating functions τ_e to take the values of $\mathbf{T}^{(0)}$ on the primary 0-cells belonging to the element e , in 1-cochain approximation we ask the integral of Φ_e on the primary 1-cells belonging to the element e to take the values of $\mathbf{G}^{(1)}$ on them. The extension of this approximation technique to generic p -cochains is straightforward. The procedure will be clarified by the examples in Section 8.

It is worth stressing that with the point of view adopted in this paper, the reconstruction of field functions from p -cochains is only expedient to the discretization of the constitutive equations and is not made with the aim of obtaining an expression for field functions over the whole domain of the problem. In this light, the issue of interelement continuity for approximation functions becomes the automatically satisfied requirement of consistency in the

association of global quantities with p -cells, when this association is made by pairs of adjacent elements.

More important than the mere technique of p -cochain approximation, is the discussion about the pros and cons of its adoption (with $p > 0$) in place of 0-cochain approximation. When—as happens in thermostatics (Fig. 9)—both strategies are available, the choice is a matter of habit and taste, but there are field problems where there are no 0-cochains to base the approximation on. For example, in 3D magnetostatics the quantities which correspond to the potential and the tension in the left column of Fig. 3 are the line integral of the vector potential \mathbf{A} (let us call it Π) and the magnetic flux Φ , respectively. This means that we have a 1-cochain $\Pi^{(1)}$ and a 2-cochain $\Phi^{(2)}$, but no 0-cochains, to start the interpolation. From a mathematical point of view we can always resort to point-interpolation based on nodal values of vector quantities like \mathbf{A} and \mathbf{B} . However, by so doing we will neglect the correct association of quantities with oriented geometric objects, losing its rich geometric content and strong adherence to the physical nature of the problem. Therefore, the widespread practice of nodal interpolation of *local* vector quantities should be avoided and it comes as little surprise the fact that the idea behind p -cochain approximation is at the heart of some of the techniques adopted by the FE community (for example, the use of edge elements [15]), in order to prevent the appearance of nonphysical terms in the solution. Note also that the p -cochain approximation approach avoids the traditional dilemma concerning the *location* of local vector field quantities, since it works with global quantities that we know are associated with p -cells.

6.8. Examples

In this section, the theory expounded will be substantiated with some example. For simplicity only 2D domains partitioned in rectangular primary 2-cells which coincide with elements (Fig. 8) will be considered. Later, the case of rectangular elements containing more than one primary 2-cell will be examined. We strongly emphasize the fact that *the theory presented in the previous sections applies also to nonorthogonal and to irregular grids*, these cases require only an increased bookkeeping effort to take care—and only in the discretization of the *constitutive* equations—of the extension and relative position of cells. In all cases, particular care must be applied in the placement of the discretization grids with respect to discontinuities in the properties of the medium and singularities in the source term—for example, requiring their placement along the boundaries between elements, so that the interpolation functions do not impose too demanding continuity conditions.

With the choice of Fig. 8, element edges coincide with primary 1-cells, while the four primary 0-cells of each pri-

mary 2-cell are the nodes of the element (remember that all these geometric objects are oriented, even if the orientation is not represented in the figure). If we take the route of 0-cochain approximation, we must define within each element an interpolating function with the degrees of freedom appropriate to the four temperature values associated with the 0-cells. For example, the bilinear polynomial (49) with its four coefficients and its geometric isotropy will do:

$$p_{\text{bil}}(x, y) = a + bx + cy + dxy. \quad (49)$$

Obviously, if there are reasons to believe that a different family of functions (for example, with exponential or rational terms) is particularly suited to a given problem, the choice of the approximating functions should comply with this additional information. To perform the interpolation, we need to know the extent of the elements. Let us call Δx and Δy the discretization steps and set a local coordinate system having its origin in a primary 0-cell of the element. Writing $T_{i,j}$ for $T(i \Delta x, j \Delta y)$, the $\mathbf{T}^{(0)}$ -based 0-cochain approximation performed within the element e , corresponds to finding a bilinear polynomial $\tilde{T}_{e,\text{bil}}(x, y)$ which satisfies the conditions

$$\begin{aligned} \tilde{T}_{e,\text{bil}}(0, 0) &= T_{0,0} \\ \tilde{T}_{e,\text{bil}}(\Delta x, 0) &= T_{1,0} \\ \tilde{T}_{e,\text{bil}}(\Delta x, \Delta y) &= T_{1,1} \\ \tilde{T}_{e,\text{bil}}(0, \Delta y) &= T_{0,1}. \end{aligned} \quad (50)$$

This gives

$$\begin{aligned} \tilde{T}_{e,\text{bil}}(x, y) &= T_{0,0} + \frac{(-T_{0,0} + T_{1,0})x}{\Delta x} + \frac{(-T_{0,0} + T_{0,1})y}{\Delta y} \\ &\quad + \frac{(T_{0,0} - T_{1,0} + T_{1,1} - T_{0,1})xy}{\Delta x \Delta y} \end{aligned} \quad (51)$$

Applying to $\tilde{T}_{e,\text{bil}}(x, y)$ the differential kinematic equation $\mathbf{g} = -\text{grad } T$ and the constitutive equation $\mathbf{q} = \lambda \mathbf{g}$ we arrive at

$$\begin{aligned} \tilde{\mathbf{q}}_e(x, y) &= \lambda(x, y) \left(\frac{(T_{0,0} - T_{1,0})}{\Delta x} \right. \\ &\quad \left. + \frac{(-T_{0,0} + T_{1,0} - T_{1,1} + T_{0,1})y}{\Delta x \Delta y} \right) \hat{\mathbf{i}} \\ &\quad + \lambda(x, y) \left(\frac{(T_{0,0} - T_{0,1})}{\Delta y} \right. \\ &\quad \left. + \frac{(-T_{0,0} + T_{1,0} - T_{1,1} + T_{0,1})x}{\Delta x \Delta y} \right) \hat{\mathbf{j}}. \end{aligned} \quad (52)$$

where $\hat{\mathbf{i}}$ and $\hat{\mathbf{j}}$ are the unit base vectors. A relation of this kind holds within each element; collectively these relations constitute the FE link between the primary 0-cochain $\mathbf{T}^{(0)}$ and the field function $\tilde{\mathbf{q}}$ anticipated in Section 6.4.

In the case of FV, we have to reconstitute the secondary 2-cochain $\tilde{\mathbf{Q}}_{\text{flow}}^{(2)}$ from the field function $\tilde{\mathbf{q}}$. To this end we must perform an integration of $\tilde{\mathbf{q}}$ on secondary 1-cells, and this requires the definition of their position. This is the subject of cell optimization that will be discussed in Section 8. For the time being, let us just assume that the secondary 2-cells are rectangular, symmetrically staggered with respect to primary 2-cells, so that each secondary 1-cell intersects orthogonally in the middle a primary 1-cell (Fig. 8). With this choice, a secondary 1-cell lies across two elements and to calculate the flow $\tilde{Q}_{(i,j) \rightarrow (h,k)}$ associated with it, we need to integrate the approximated flow density $\tilde{\mathbf{q}}$ calculated over two adjacent elements. Let us call e_1 and e_2 the two elements involved in the determination of $\tilde{Q}_{(0,0) \rightarrow (1,0)}$ and $\tilde{\mathbf{q}}_{e_1}$ and $\tilde{\mathbf{q}}_{e_2}$ the corresponding approximate flow densities. Of course, the result of the integration depends on the function $\lambda(x, y)$ which defines the material. To simplify the calculations we consider a homogeneous material and $\Delta x = \Delta y = \Delta$, obtaining the relation

$$\begin{aligned} \tilde{Q}_{(0,0) \rightarrow (1,0)}^{\text{flow}} &= \int_0^{\Delta/2} \tilde{\mathbf{q}}_{e_1}(x, y) \cdot \hat{\mathbf{i}} dy + \int_{-\Delta/2}^0 \tilde{\mathbf{q}}_{e_2}(x, y) \cdot \hat{\mathbf{i}} dy \\ &= \lambda \left(+\frac{1}{8}T_{0,1} - \frac{1}{8}T_{1,1} + \frac{3}{4}T_{0,0} - \frac{3}{4}T_{1,0} + \frac{1}{8}T_{0,-1} - \frac{1}{8}T_{1,-1} \right). \end{aligned} \quad (53)$$

Once a similar calculation has been performed for each 1-cell of the secondary grid, the balance equations can be written. For FV this amounts to equating the heat source $\mathbf{Q}_{\mathbf{c}(2)}^{\text{source}} = \langle \mathbf{c}(2), \mathbf{Q}_{\text{source}}^{(2)} \rangle$ within each secondary 2-cell, to the sum of the four heat flows associated with the 1-cells which constitute its (oriented) boundary. In local coordinates centered within the cell this means

$$\begin{aligned} \mathbf{Q}_{0,0}^{\text{source}} &= \tilde{Q}_{(0,0) \rightarrow (1,0)}^{\text{flow}} + \tilde{Q}_{(0,0) \uparrow (0,1)}^{\text{flow}} \\ &\quad - \tilde{Q}_{(-1,0) \rightarrow (0,0)}^{\text{flow}} - \tilde{Q}_{(0,-1) \uparrow (0,0)}^{\text{flow}}. \end{aligned} \quad (54)$$

We obtain for each secondary 2-cell, the equation

$$\begin{aligned} \mathbf{Q}_{\mathbf{c}(2)}^{\text{source}} &= \lambda \left(-\frac{1}{4}T_{-1,1} - \frac{1}{2}T_{0,1} - \frac{1}{4}T_{1,1} - \frac{1}{2}T_{-1,0} \right. \\ &\quad \left. + 3T_{0,0} - \frac{1}{2}T_{1,0} - \frac{1}{4}T_{-1,-1} - \frac{1}{2}T_{0,-1} - \frac{1}{4}T_{1,-1} \right). \end{aligned} \quad (55)$$

which involves the nine temperatures associated with the 0-cells lying in the four elements which “share” the secondary 2-cell. The collection of these equations, constitute the FV linear system.

Writing the FE balance equation requires an integration of the flow density over all the elements belonging to the

support of the weighting functions. This happens because the weighting functions define a spread secondary 2-cell, which in general has a spread boundary. For example, consider Galerkin FE, i.e., FE with weighting functions obtained assembling functions of the same kind of the interpolating functions. With bilinear interpolating functions, the secondary 2-cells are spread over four elements (Fig. 1)—which we call collectively $\mathbf{d}_{(2)}$ —and we have

$$\mathbf{Q}_{w_{\text{bil}}\mathbf{d}_{(2)}}^{\text{source}} = \int_{w_{\text{bil}}\mathbf{d}_{(2)}} \sigma(x, y) dx dy \quad (56)$$

$$\tilde{Q}_{\partial(w_{\text{bil}}\mathbf{d}_{(2)})}^{\text{flow}} = \int_{\partial(w_{\text{bil}}\mathbf{d}_{(2)})} \tilde{\mathbf{q}}(x, y) \cdot d\mathbf{l} \quad (57)$$

with the following defining relation for the r.h.s. of Eq. (57):

$$\begin{aligned} \int_{\partial(w_{\text{bil}}\mathbf{d}_{(2)})} \tilde{\mathbf{q}}(x, y) \cdot d\mathbf{l} &\stackrel{\text{def}}{=} \int_{\partial\mathbf{d}_{(2)}} w_{\text{bil}}(x, y) \tilde{\mathbf{q}}(x, y) \cdot d\mathbf{l} \\ &\quad - \int_{\mathbf{d}_{(2)}} \text{grad } w_{\text{bil}}(x, y) \cdot \tilde{\mathbf{q}}(x, y) dx dy. \end{aligned} \quad (58)$$

The balance equation is

$$\mathbf{Q}_{w_{\text{bil}}\mathbf{d}_{(2)}}^{\text{source}} = \tilde{Q}_{\partial(w_{\text{bil}}\mathbf{d}_{(2)})}^{\text{flow}} \quad (59)$$

and the final result (with $\Delta x = \Delta y = \Delta$) once again involves the nine temperatures defined within the four elements which constitute the support of the spread secondary 2-cell

$$\begin{aligned} \int_{w_{\text{bil}}\mathbf{d}_{(2)}} \sigma(x, y) dx dy &= \lambda \left(-\frac{1}{3}T_{-1,1} - \frac{1}{3}T_{0,1} - \frac{1}{3}T_{1,1} \right. \\ &\quad \left. - \frac{1}{3}T_{-1,0} + \frac{8}{3}T_{0,0} - \frac{1}{3}T_{1,0} \right. \\ &\quad \left. - \frac{1}{3}T_{-1,-1} - \frac{1}{3}T_{0,-1} - \frac{1}{3}T_{1,-1} \right). \end{aligned} \quad (60)$$

Obviously in Eq. (55) the source term and the coefficients of the r.h.s. depend on the choice of the shape of the secondary 2-cell, whereas in Eq. (60) they depend on the choice of the weighting functions, which define the “shape” of the spread 2-cell.

Deciding to take the route of 1-cochain approximation, the first step requires the application of the discrete kinematic equation. With all the primary 0-cells oriented as sources, the equation is

$$G_{(i,j) \rightarrow (h,k)} = T_{i,j} - T_{h,k}. \quad (61)$$

In comparison to 0-cochain interpolation, note that by applying first the topological equation we defer the appearance of metrical notions like length and angle. The next step is the estimation of the field function \mathbf{g} based on the

cochain $\mathbf{G}^{(1)}$. Within each element there are four primary 1-cells with their thermal tensions

$$\begin{array}{ccc} & G_{(0,1) \rightarrow (1,1)} & \\ G_{(0,0) \uparrow (0,1)} & & G_{(1,0) \uparrow (1,1)} \\ & G_{(0,0) \rightarrow (1,0)} & \end{array} \quad (62)$$

therefore we can choose as the approximating function a vector-valued function with four coefficients and suitable symmetry, for example,

$$\mathbf{p}_{\text{lin}}(x, y) = (a + by)\hat{\mathbf{i}} + (c + dx)\hat{\mathbf{j}}. \quad (63)$$

The 1-cochain approximation amounts to requiring that the line integral of the approximating function $\tilde{\mathbf{g}}_{c,\text{lin}}(x, y)$ performed on the four primary 1-cells be equal to the corresponding thermal tension:

$$\begin{aligned} \int_0^{\Delta x} \tilde{\mathbf{g}}_{c,\text{lin}}(x, 0) \cdot \hat{\mathbf{i}} \, dx &= G_{(0,0) \rightarrow (1,0)} \\ \int_0^{\Delta x} \tilde{\mathbf{g}}_{c,\text{lin}}(x, \Delta y) \cdot \hat{\mathbf{i}} \, dx &= G_{(0,1) \rightarrow (1,1)} \\ \int_0^{\Delta y} \tilde{\mathbf{g}}_{c,\text{lin}}(0, y) \cdot \hat{\mathbf{j}} \, dy &= G_{(0,0) \uparrow (0,1)} \\ \int_0^{\Delta y} \tilde{\mathbf{g}}_{c,\text{lin}}(\Delta x, y) \cdot \hat{\mathbf{j}} \, dy &= G_{(1,0) \uparrow (1,1)}. \end{aligned} \quad (64)$$

The result is

$$\begin{aligned} \tilde{\mathbf{g}}_{c,\text{lin}}(x, y) &= \left(\frac{G_{(0,0) \rightarrow (1,0)}}{\Delta x} + \frac{(G_{(0,1) \rightarrow (1,1)} - G_{(0,0) \rightarrow (1,0)})y}{\Delta x \Delta y} \right) \hat{\mathbf{i}} \\ &+ \left(\frac{G_{(0,0) \uparrow (0,1)}}{\Delta y} + \frac{(G_{(1,0) \uparrow (1,1)} - G_{(0,0) \uparrow (0,1)})x}{\Delta x \Delta y} \right) \hat{\mathbf{j}}. \end{aligned} \quad (65)$$

From here we proceed as for 0-cochain interpolation, obtaining the same final equation. It is worth noting that in the case of FV, we can write explicitly the discrete constitutive equation. The prototype of the link is the expression

$$\begin{aligned} Q_{(0,0) \rightarrow (1,0)}^{\text{flow}} &= \frac{1}{8}G_{(0,1) \rightarrow (1,1)} + \\ &+ 0G_{(0,0) \uparrow (0,1)} + 0G_{(1,0) \uparrow (1,1)} + \\ &+ \frac{3}{4}G_{(0,0) \rightarrow (1,0)} + \\ &+ 0G_{(0,-1) \uparrow (0,0)} + 0G_{(1,-1) \uparrow (1,0)} + \\ &+ \frac{1}{8}G_{(0,-1) \rightarrow (1,-1)} \end{aligned} \quad (66)$$

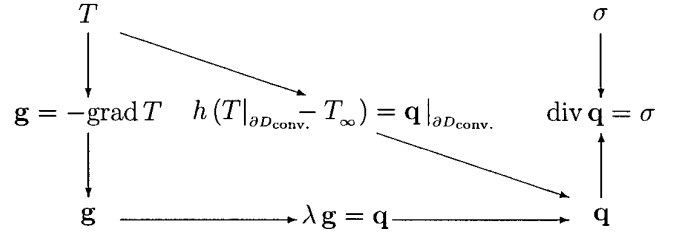


FIG. 10. Mixed boundary conditions originate additional constitutive equations. In the thermostatics example considered here, they are due to the presence of a convective heat exchange across the boundary.

7. BOUNDARY CONDITIONS

We consider first the case of boundary conditions which assign some of the unknown (possibly discretized) fields along the boundary of the domain. For FV the enforcement of this kind of boundary conditions requires simply the positioning along the boundary of the proper kind of primary or secondary p -cells. The global quantities associated with these cells appear as known terms in the final equations. In FE the cells and their boundaries can be spread; to include the boundary terms in the equations we must assure that the weighting functions defining the cells do not vanish along the parts of the boundary where the fields are assigned. This is usually obtained with the placement of nodes along the boundary, but our analysis shows that this is not mandatory. Note also that the boundary terms are naturally partitioned among bordering *cells* and not—as is usually assumed in FE—among bordering *nodes*.

To explain how mixed boundary conditions fit in our discretization scheme, we must consider the physical phenomena which originates them. In thermostatics they can appear as a consequence of convective heat exchange across the boundary. We can write the rate of heat exchange as

$$h(T|_{\partial D_{\text{conv}}} - T_{\infty}) = \mathbf{q}|_{\partial D_{\text{conv}}}, \quad (67)$$

where ∂D_{conv} are the parts of the boundary where the exchange takes place, h is the coefficient of convective heat exchange, and T_{∞} is the ambient temperature. To enforce this kind of boundary condition we have only to treat (67) as a constitutive equation valid on ∂D_{conv} (Fig. 10). The discretization of the new term proceeds for both FV and FE as described above for constitutive equations.

8. CELL OPTIMIZATION

In the examples of Section 6.8, we saw that once the primary grid, the elements mesh, and the approximation functions have been selected, in order to write the balance

equations it is necessary to set the shape of the secondary 2-cells. In FE, with the interpretation proposed in this paper, this corresponds to the choice of the weighting functions. The options available in the joint selection of shape and weighting functions are well documented in FE literature, along with the existence of *optimal* choices for particular problems.

In FV literature, on the other hand, there is seldom mention of optimality criteria in the choice of cell shape. Instead, the cells are usually constructed on the basis of symmetry considerations and the very adoption of two staggered grid is often considered an oddity, justified *a posteriori* by the superior performances obtained [11, 16]. The purpose of this section is to show that FV, like FE, can benefit greatly of a properly performed joint selection of approximation function and cell shape. The examples to follow are not intended to exhaust the topic of cell optimization in FV, but only to attract interest to the problem. Therefore we will consider only 2D examples where, mimicking FE, a primary grid, the elements mesh, and the approximation functions have been arbitrarily set and it remains only to define the shape of secondary 2-cells. Moreover, as in the examples of Section 6.8, all primary 2-cells and elements are rectangular.

The basic idea behind our optimization strategy is borrowed from polynomial approximation theory and can be explained with the following 1D example. Given a real function $f(x)$ defined on an interval $[x_0, x_1]$, a simple approximation of it is constituted by a linear interpolation polynomial $\tilde{f}_{\text{lin}}(x)$, taking the values $f(x_0)$ and $f(x_1)$ at the endpoints of the interval:

$$\tilde{f}_{\text{lin}}(x) = f(x_0) + \frac{(f(x_1) - f(x_0))(x - x_0)}{x_1 - x_0}. \quad (68)$$

The derivative of $f(x)$ can be approximated by the derivative of the interpolation polynomial, but since this last derivative is a constant function, we expect this approximation of $f'(x)$ achieved by $\tilde{f}'(x)$ to be of a lower order, compared to that of $f(x)$ achieved by $\tilde{f}(x)$. It happens, however, that *all* the infinitely many quadratic polynomials $\tilde{f}_{\text{quad}}^\alpha(x)$, taking the values $f(x_0)$ and $f(x_1)$ in x_0, x_1 ,

$$\begin{aligned} \tilde{f}_{\text{quad}}^\alpha(x) = f(x_0) + \frac{(f(x_1) - f(x_0))(x - x_0)}{x_1 - x_0} \\ + \alpha(x - x_0)(x - x_1), \end{aligned} \quad (69)$$

have a derivative which in the central point $x_c = (x_0 + x_1)/2$ of the interval takes the *same* value of the derivative of the linear interpolation polynomial.

Let us apply this principle to the example of Section 6.8, having a primary grid with rectangular 2-cells which coincide with elements (Fig. 8). In this case to approximate

the temperature distribution within the elements we adopt the bilinear polynomial (49); we know (Eq. (51)) the result $\tilde{T}_{e,\text{bil}}(x, y)$ of 0-cochain approximation performed within the element. For reasons of symmetry, we can aim in our search for higher order approximations in the calculation of the flow, at quadratic polynomials

$$p_{\text{quad}}(x, y) = a' + b'x + c'y + d'xy + \alpha x^2 + \beta y^2. \quad (70)$$

We can write the infinitely many polynomials taking the temperature values of an elements' four 0-cells as

$$\begin{aligned} \tilde{T}_{e,\text{quad}}^{\alpha,\beta}(x, y) = T_{0,0} + \frac{(-T_{0,0} + T_{1,0})x}{\Delta x} + \frac{(-T_{0,0} + T_{0,1})y}{\Delta y} \\ + \frac{(T_{0,0} - T_{1,0} + T_{1,1} - T_{0,1})xy}{\Delta y \Delta y} \\ + \alpha x(x - \Delta x) + \beta y(y - \Delta y), \end{aligned} \quad (71)$$

where the two parameters α and β can take arbitrary values. We must establish if there exists, within the element, a locus which can be taken as a piece of the boundary of a secondary 2-cell and such that the flow through it, calculated from *all* these quadratic functions, equals the one calculated with the bilinear function.

As a first approach, we can look for a curve $\gamma(s)$ ($s_0 \leq s \leq s_1$) lying in the element and satisfying in each point the condition

$$\text{grad } \tilde{T}_{e,\text{bil}}(\gamma(s)) \cdot \hat{\mathbf{n}}_{\gamma(s)} = \text{grad } \tilde{T}_{e,\text{quad}}^{\alpha,\beta}(\gamma(s)) \cdot \hat{\mathbf{n}}_{\gamma(s)} \quad \forall \alpha, \beta, \quad (72)$$

where $\hat{\mathbf{n}}_{\gamma(s)}$ is the normal to the curve γ in its point of parameter s . In addition, we ask the curve to constitute the boundary of a 2-cell, or to contribute to the formation of such a boundary when joined with curves calculated in adjacent elements. An alternative is to impose directly, instead of the equality of the orthogonal component of the two gradients in each point of the curve, the equality through the curves of the values of the *flow* calculated with the two approximating functions: in other words, in place of (72) we can impose

$$\begin{aligned} \int_{\gamma(s)} \tilde{\mathbf{q}}_{\text{bil}}(\gamma(s)) \cdot \hat{\mathbf{n}}_{\gamma(s)} \, ds \\ = \int_{\gamma(s)} \tilde{\mathbf{q}}_{e,\text{quad}}^{\alpha,\beta}(\gamma(s)) \cdot \hat{\mathbf{n}}_{\gamma(s)} \, ds \quad \forall \alpha, \beta, \end{aligned} \quad (73)$$

with

$$\tilde{\mathbf{q}}_{e,\text{bil}}(x, y) = -\lambda(x, y) \text{grad } \tilde{T}_{e,\text{bil}}(x, y) \quad (74)$$

$$\tilde{\mathbf{q}}_{\text{quad}}^{\alpha,\beta}(x, y) = -\lambda(x, y) \text{grad } \tilde{T}_{e,\text{quad}}^{\alpha,\beta}(x, y). \quad (75)$$

This last formulation is more adherent to the spirit of FV discretization and imposes a lesser constraint on the variables of the problem but it appears more complex to apply. To keep things simple let us enforce Eq. (72). Giving to the curve γ which constitutes the as-yet-unknown cell boundary, the parametric representation

$$\gamma \equiv \begin{cases} x \\ y(x), \end{cases} \quad (76)$$

we have the following cartesian expression for $\hat{\mathbf{n}}_\gamma(s)$:

$$\hat{\mathbf{n}}_\gamma(x) = \frac{-y'(x)\hat{\mathbf{i}} + \hat{\mathbf{j}}}{\sqrt{1 + y'(x)^2}}. \quad (77)$$

Substituting (51), (71), and (77) in (72) we obtain an equation which can be reduced to

$$\alpha(2x - \Delta x)y'(x) - \beta(2y(x) - \Delta y) = 0 \quad \forall \alpha, \beta. \quad (78)$$

This equation is satisfied for arbitrary α, β with

$$y(x) = \frac{\Delta y}{2}. \quad (79)$$

For reasons of symmetry, assigning to γ the parametric representation

$$\gamma \equiv \begin{cases} x(y) \\ y, \end{cases} \quad (80)$$

we obtain the solution

$$x(y) = \frac{\Delta x}{2}. \quad (81)$$

This means that the axes of symmetry of the primary 2-cells are optimal loci for the evaluation of the normal flow and, therefore, for the placement of secondary 1-cells. The traditional symmetrically staggered secondary grid with 0-cells placed in the barycentre of primary 2-cells, adopted tentatively in Section 6.8 (Fig. 8) in this case, is indeed optimal, in the sense that the results obtained with these cells using bilinear interpolation polynomials, have the same degree of accuracy obtainable interpolating with complete second-order polynomials. In other words, the adoption of the optimal secondary 2-cells in place of generic ones, increases by one the rate of convergence of the method.

As second example we consider a domain with the same regular primary grid as the first example but with a different elements mesh; now each element no longer coincides with

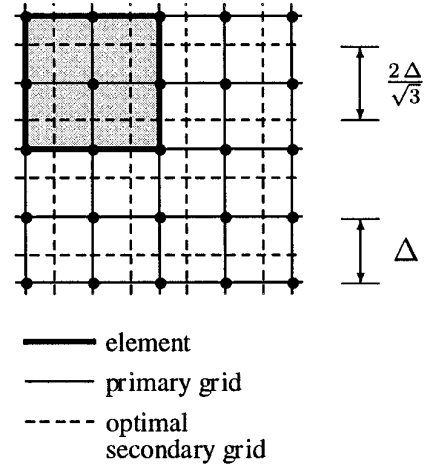


FIG. 11. Grids and elements mesh for biquadratic 0-cochain approximation (the orientation of cells is not represented).

a primary 2-cell but is composed by four of them (Fig. 11). Instead of the four primary 0-cells per element of the preceding example, we have now nine primary 0-cells per element and for this reason the approximation functions for temperature within the elements can be biquadratic polynomials:

$$p_{\text{biq}}(x, y) = a + bx + cy + dxy + ex^2 + fy^2 + gx^2y^2 + rx^2y + sxy^2. \quad (82)$$

We can aim in our search for a better approximation, at polynomials containing the two cubical terms missing in (82)

$$p_{\text{cub}}^{\alpha, \beta}(x, y) = a' + b'x + c'y + d'xy + e'x^2 + f'y^2 + g'x^2y^2 + r'x^2y + s'xy^2 + \alpha x^3 + \beta y^3. \quad (83)$$

Applying 0-cochain approximation based on polynomials (82) and (83), we obtain the field functions $\tilde{T}_{e, \text{biq}}(x, y)$ and $\tilde{T}_{e, \text{cub}}^{\alpha, \beta}(x, y)$. The optimization equation is

$$\text{grad } \tilde{T}_{e, \text{biq}}(\gamma(s)) \cdot \hat{\mathbf{n}}_{\gamma(s)} = \text{grad } \tilde{T}_{e, \text{cub}}^{\alpha, \beta}(\gamma(s)) \cdot \hat{\mathbf{n}}_{\gamma(s)} \quad \forall \alpha, \beta \quad (84)$$

and gives rise to the equation (referred to a local coordinate system with origin in the central primary 0-cell)

$$\alpha(3x^2 - \Delta x^2)y'(x) - \beta(3y(x)^2 - \Delta y^2) = 0 \quad \forall \alpha, \beta \quad (85)$$

which is satisfied for arbitrary α, β with

$$y(x) = \pm \frac{\Delta y}{\sqrt{3}}. \quad (86)$$

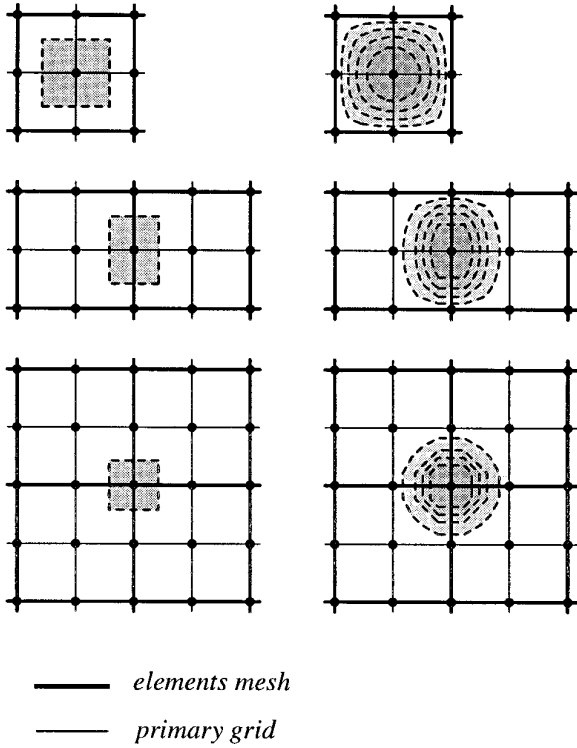


FIG. 12. Optimal FV secondary 2-cells (left column) and Galerkin FE spread 2-cells (right column) for the same primary grid and elements mesh and with the same biquadratic interpolation functions.

As before we also have the solution

$$x(y) = \pm \frac{\Delta x}{\sqrt{3}}. \quad (87)$$

Equations (86) and (87) mean that adopting biquadratic interpolation on nine-point square elements (i.e., with $\Delta x = \Delta y$), we have three kinds of optimal secondary 2-cells (Fig. 12): large square cells interior to the element, to which corresponds a 9-point formula; medium-sized rectangular cells astride the element edges (and therefore shared among two adjacent elements) to which correspond 15-point formulas and, finally, small square cells centered in the element vertices (i.e., shared among four elements), to which correspond 25-point formulas. The coefficients of the formulas can be calculated following the procedure described in Section 6.8. Proceeding with Galerkin FE with the same primary grid and elements mesh and adopting the same interpolation functions, we obtain again three kinds of secondary 2-cells, with corresponding 9-, 15-, and 25-point formulas. Figure 12 shows the optimal FV cells, along with the corresponding Galerkin FE spread cells, and the “influence region” for each kind of cell. Note that in this case, to place, as is usually done, the FV secondary

1-cells on the axes of symmetry of the primary 2-cells (obtaining a regular secondary grid with cells of uniform extension) is *not* an optimal choice, i.e., will not give, using biquadratic interpolation polynomials, the degree of accuracy obtainable by interpolating with complete third-order polynomials. The ensuing absence of improvement in the rate of convergence, passing (unawarely) from optimal bilinear to nonoptimal biquadratic interpolation, may account for the widespread—but wrong—feeling that FV, while performing well with low-order interpolation, becomes less attractive for approximations of higher order.

8.1. Taylor Expansion Applied to the Determination of FV Optimal Cells

In the spirit of Eq. (72) and (84), we try to determine a locus for the secondary 2-cell boundary such that the approximation of the component of $\text{grad } T(x, y)$ orthogonal to it is optimal. Suppose that the curve γ which constitutes the boundary of the cell, is given the functional representation (76). Calling $\hat{\mathbf{n}}_\gamma(x)$ the normal to γ in its generic point, we can write the approximation we are trying to achieve as

$$\text{grad } T(x, y) \cdot \hat{\mathbf{n}}_\gamma(x) \approx \sum_{i,j=-1}^1 a_{i,j}(x, y(x)) T(i \Delta x, j \Delta y). \quad (88)$$

Expanding $T(x, y)$ in Taylor series around the generic point of γ we can write

$$\begin{aligned} T(i \Delta x, j \Delta y) &= \sum_{h,k} \frac{1}{h!k!} (i \Delta x - x)^h (j \Delta y - y(x))^k T^{(h,k)}(x, y(x)). \end{aligned} \quad (89)$$

Combining (89) with (88), we obtain

$$\text{grad } T(x, y) \cdot \hat{\mathbf{n}}_\gamma(x) \approx \sum_{h,k} b_{i,j}(x, y(x)) T^{(h,k)}(x, y(x)), \quad (90)$$

where the coefficients $b_{i,j}$ are functions of the $a_{i,j}$ s. Substituting the cartesian expression (77) for $\hat{\mathbf{n}}_\gamma(x)$, Eq. (90) becomes

$$\begin{aligned} & - \frac{T^{(1,0)}(x, y(x)) y'(x)}{\sqrt{1 + y'(x)^2}} + \frac{T^{(0,1)}(x, y(x))}{\sqrt{1 + y'(x)^2}} \\ & \approx \sum_{h,k} b_{i,j}(x, y(x)) T^{(h,k)}(x, y(x)). \end{aligned} \quad (91)$$

To obtain a perfect approximation we should achieve equality in (91). In fact, the best we can do is to impose the equality for as many terms of low order as possible. This gives rise to the system of equations

$$\begin{aligned}
b_{0,0} &= 0 & \text{and} \\
b_{1,0} &= -\frac{y'(x)}{\sqrt{1+y'(x)^2}} & k_y &= -\frac{\Delta y}{\sqrt{3}} \\
b_{0,1} &= \frac{1}{\sqrt{1+y'(x)^2}} & a_{-1,-1}(x) &= \dots \\
b_{2,0} &= 0 & & \vdots \\
&\vdots & & \\
b_{2,2} &= 0 & &
\end{aligned} \tag{92}$$

which is a mixed system of algebraic and differential equations. Equations (92) must be solved for the coefficients $a_{i,j}(x)$ of Eq. (88) and for the function $y(x)$. To simplify matters, let us look for solutions only among curves parallel to the x axis, that is, with $y(x) = k_y$, where k_y is a constant we must determine. With these choices, (92) becomes

$$\begin{aligned}
b_{0,0} &= 0 \\
b_{1,0} &= 0 \\
b_{0,1} &= 1 \\
b_{2,0} &= 0 \\
&\vdots \\
b_{2,2} &= 0.
\end{aligned} \tag{93}$$

Writing (93) explicitly in terms of the $a_{i,j}$ s we obtain

$$\begin{aligned}
a_{-1,-1} + \dots + a_{1,1} &= 0 \\
(-\Delta x - x)a_{-1,-1} + \dots + (\Delta x - x)a_{1,1} &= 0 \\
(-\Delta y - k_y)a_{-1,-1} + \dots + (\Delta y - k_y)a_{1,1} &= 0 \\
&\vdots & \vdots & \vdots
\end{aligned} \tag{94}$$

which, solved for the $a_{i,j}$ s and k_y gives two solutions:

$$\begin{aligned}
k_y &= \frac{\Delta y}{\sqrt{3}} \\
a_{-1,-1}(x) &= \frac{(3 - 2\sqrt{3})(\Delta x - x)x}{12 \Delta x^2 \Delta y} \\
&\vdots \\
a_{1,1}(x) &= \frac{(3 + 2\sqrt{3})(\Delta x + x)x}{12 \Delta x^2 \Delta y}
\end{aligned} \tag{95}$$

For reasons of symmetry, to a curve γ with a parametric representation (80) there correspond the solutions $x(y) = k_x = \pm \Delta x/\sqrt{3}$, with their set of coefficients $a_{i,j}(x)$. With these curves we can construct the boundaries $\partial \mathbf{c}_{(2)}^i$ of the optimal secondary 2-cells. Once the optimal cells are found, substituting the calculated coefficients $a_{i,j}$ in (88) we obtain the expression of the optimal approximation of $\text{grad } T \cdot \hat{\mathbf{n}}_{\partial \mathbf{c}_{(2)}^i}$ along the boundary of the optimal cells and then we calculate the flux to enforce, finally, the discrete balance equation.

8.2. FV Optimal Cells and FE Superconvergent Points

The introduction of optimal cells should sound familiar to FE practitioners, reminding them of FE *optimal flux-sampling* or *superconvergent* points. In fact both concepts have a common root, but while in FE we look for isolated points where the value of the flow density is of higher precision, in FV we look for a higher precision global flow through secondary $(n - 1)$ -cells forming the boundary of secondary n -cells. This makes the FV optimization task locally easier, because we look for a higher precision only for the component of the flow density (or for its integral) orthogonal to the $(n - 1)$ -cells, but globally much harder, since we look for optimal bounding $(n - 1)$ -dimensional *loci* instead of optimal isolated points. Therefore, while FE superconvergent points can be easily found for an ample category of cell shapes, it may well happen that with some choice of the primary grid, of the elements mesh, and of the approximation functions, there are no optimal loci capable of constituting the boundary of secondary cells. This is often the case for irregular primary grids, with the optimization strategy sketched above. In these cases it may be advisable to try a joint optimization of primary and secondary grids, or to content oneself with an approximation of lower degree for the constitutive equations. Note that while FE superconvergent points are used to obtain more accurate numerical data once the calculations have been performed, FV optimal cells permit setting an improved system of equations prior to the actual execution of numerical calculations.

9. THE FINITE DIFFERENCE METHODS

FD approaches the discretization problem from a very different perspective than FE and FV. The classical FD

discretization process aims to determine a *local* discrete approximation for the *operator* constituting the complete *field equation*. For example, in the case of thermostatics (for simplicity, with homogeneous and isotropic material) the field equation is

$$\lambda \nabla^2 T = \sigma \quad (97)$$

and it is the differential operator ∇^2 that FD tries to write in discrete terms. Performing on Eq. (97) the discretization steps for a nine-point formula on a regular grid with square cells of width $\Delta x = \Delta y = \Delta$, we obtain

$$\begin{aligned} \sigma(0,0)\Delta^2 = & \lambda(0T_{-1,1} - 1T_{0,1} + 0T_{1,1} - 1T_{-1,0} + 4T_{0,0} \\ & - 1T_{1,0} + 0T_{-1,-1} - 1T_{0,-1} + 0T_{1,-1}). \end{aligned} \quad (98)$$

Equation (98), which is the optimal nine-point FD formula, appears in fact to be the traditional FD five-point formula for the laplacian. In other words, with the FD approach we are suggested to altogether ignore the four ‘‘corner’’ temperatures. Conversely, the optimal FV nine-point formula (55) (and also the FE nine-point formula (60)) for the same problem and the same discretization grid makes use of this information, assigning nonzero weights to the information carried by these 0-cells. This happens because the two formulas are aimed at different goals: the FV one optimizes the approximation of the *constitutive equation*, and therefore of the flow, *through the boundary* of the secondary cells, while the FD formula constitutes an optimal approximation for the *laplacian operator in the central point* of the patch. This local versus global approach is also apparent in the l.h.s. source terms of the two equations.

Differently rephrased, the classical FD approach gives an optimal solution to a problem which is *not* the one the discretization methods for field problems are expected to face. This is the fundamental weakness of this approach, and it explains why it is difficult to obtain well-performing high-order formulas for field problems with it. Note, on the other hand, that in 1D problems the balance equations are written for secondary 1-cells, the boundaries of which are formed by 0-cells. Therefore, in this case the FD strategy which optimizes on points can coincide with the FV strategy which optimizes on 1-cell boundaries. In this light the historical primacy of the FD approach for the discretization of 2D and 3D field problems might be ascribed to this confusion between two different discretization paradigms and to an unfortunate choice of the wrong one in the generalization from 1D to higher dimensions.

Other kinds of FD schemata—like the support-operator FD methods [13, 14]—discretize separately in a nonlocal way the first-order differential operators. This approach would be similar to the one proposed in this paper, but for the failure to acknowledge the intrinsic discrete nature

of topological equations. This is reflected in the recourse to *field functions* and *metrical tools* in the relative discretization step, where global quantities and topological tools are more appropriate. Consequently the discretization efforts are partially diverted from the constitutive equations, where they should be focused, in order to select a discretization for the differential operators and to enforce properties on it that the (univocally determined) coboundary operator satisfies automatically (see Section 5.4).

10. CONCLUSIONS

We have shown that, introducing the concept of chain and that of spread cell, a strict parallelism can be built between the finite element method and the finite volume method. It appears that both methods can be better analyzed by distinguishing the topological and constitutive parts of the field equation and recognizing the intrinsic discrete nature of the *topological equations*. A consequence of this is the realization that the discretization of *constitutive equations* is the central issue in the construction of effective discretization schemata and the only place where the recourse to local representations is fully justified. For topological equations a discrete representation for geometry, fields, and operators is preferable and, for operators, it is uniquely determined in terms of coboundaries. This underlines the importance of the correct attribution of physical quantities to geometric objects with appropriate orientation and dimension. To the lack of such attribution the poor performance of high order FD formulas for field problems and some troubles afflicting FE methods are ascribed.

The paper also shows how the recognition of the opportunity of employing both *internal and external orientation* constitutes a first argument for the utilization of two distinct discretization grids and it points out that the use of two *dual cell-complexes* as grids assures the preservation of the fundamental algebraic properties of topological operators. A distinction is made between the discretization grids and the elements mesh. It is shown that the latter is not a tool devoted to the representation in discrete terms of geometry and fields, but it is instead instrumental to the field function approximation, which is a step in the discretization of the constitutive equations. A further argument supporting not only the distinction of primary and secondary grids, but also their staggered placement, lies in the minimization thus obtainable of the error implicit in the discretization of constitutive equations.

With all this in mind, the choice of the grids and of the elements mesh, the enforcement of boundary conditions and the whole discretization process can be conducted in a more systematic and efficient way and with close adherence to the physical nature of the problems, resulting in greater insight and better performance.

APPENDIX

We will show in this appendix that on n -dimensional dual cell-complexes (primary and secondary complexes, with internal and external orientation, respectively) the coboundary operator acting on primary $(p - 1)$ -cochains is the adjoint, under a natural duality of cochain spaces, of the coboundary operator acting on secondary $(n - p)$ -cochains.

For concreteness we will sketch the proof for the coboundary operators appearing in the discretization of 3D thermostatics (Fig. 7). We can assume without loss of generality that all cells be oriented and univocally labeled in such a way that each pair of dual cells have the same default orientation and label. We represent a p -cochain on a complex with the vector of the global quantities associated with the p -cells of the complex:

$$\begin{aligned} \mathbf{T}^{(0)} &= [T_1 \dots T_N], & \mathbf{Q}_{\text{source}}^{(3)} &= [Q_1^s \dots Q_N^s] \\ \mathbf{G}^{(1)} &= [G_1 \dots G_M], & \mathbf{Q}_{\text{flow}}^{(2)} &= [Q_1^f \dots Q_M^f]. \end{aligned} \quad (99)$$

We put in duality each space of primary p -cochains with the space of secondary $(n - p)$ -cochains, by means of the bilinear forms (100) and (101)

$$\langle \mathbf{T}^{(0)}, \mathbf{Q}_{\text{source}}^{(3)} \rangle = \sum_{i=1}^N T_i Q_i^s \quad (100)$$

$$\langle \mathbf{G}^{(1)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle = \sum_{i=1}^M G_i Q_i^f \quad (101)$$

that are a discrete counterpart of

$$\langle T, \sigma \rangle = \int_D T \sigma \, dv \quad (102)$$

$$\langle \mathbf{g}, \mathbf{q} \rangle = \int_D \mathbf{g} \cdot \mathbf{q} \, dv. \quad (103)$$

We want to prove that

$$\langle \delta \mathbf{T}^{(0)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle = \langle \mathbf{T}^{(0)}, \bar{\delta} \mathbf{Q}_{\text{flow}}^{(2)} \rangle, \quad (104)$$

where the overbar distinguishes the coboundary acting on the secondary cell-complex from the one acting on the primary one.

To show that (104) holds, we start by representing the coboundary operators δ and $\bar{\delta}$ by means of the *incidence matrices* I and \bar{I} [7].

$$\langle \delta \mathbf{T}^{(0)}, \mathbf{Q}_{\text{flow}}^{(2)} \rangle = \sum_{i=1}^M \left(\sum_{j=1}^N I_{i,j} T_j \right) Q_i^f \quad (105)$$

$$\langle \mathbf{T}^{(0)}, \bar{\delta} \mathbf{Q}_{\text{flow}}^{(2)} \rangle = \sum_{i=1}^N T_i \left(\sum_{j=1}^M \bar{I}_{i,j} Q_j^f \right). \quad (106)$$

Then we observe that with our choice of *dual* cell-complexes and our pairing of orientation and labels of dual cells, we have

$$\bar{I} = I^T. \quad (107)$$

Substituting (107) in (106), exchanging running indices, and reordering, the identity of (105) and (106) is easily verified. This proves that $\bar{\delta}$ is the adjoint of δ under the duality of cochain spaces given by (100) and (101).

ACKNOWLEDGMENT

Many thanks to Gian Guido Folena for his help.

REFERENCES

1. D. S. Burnett, *Finite Element Analysis* (Addison-Wesley, Reading, MA, 1987).
2. S. Ramadhyani and S. V. Patankar, Solution of the Poisson equation: Comparison of the Galerkin and control-volume methods, *Int. J. Methods Eng.* **15**, 1395 (1980).
3. E. Oñate and S. R. Idelsohn, A comparison between finite element and finite volume methods in CFD, *Comput. Fluid Dyn.* **1**, 93 (1992).
4. G. D. Smith, *Numerical Solution of Partial Differential Equations* (Oxford Univ. Press, London, 1965).
5. E. Tonti, *On the Formal Structure of Physical Theories* (Consiglio Nazionale delle Ricerche, Milano, 1975).
6. P. Penfield Jr. and H. A. Haus, *Electrodynamics of Moving Media* (The MIT Press, Cambridge, MA, 1966).
7. W. Franz, *Algebraic Topology* (Ungar, New York, 1968).
8. W. L. Burke, *Applied Differential Geometry* (Cambridge Univ. Press, Cambridge, 1985).
9. J. A. Schouten, *Tensor Analysis for Physicists* (Dover, New York, 1989).
10. J. M. Hyman, R. J. Knapp, and J. C. Scovel, High order finite volume approximations of differential operators on nonuniform grids, *Physica D* **60**, 112 (1992).
11. M. Vinokur, An analysis of finite-difference and finite-volume formulations of conservation laws, *J. Comput. Phys.* **81**, 1 (1989).
12. H. Lebesgue, *Leçons sur l'Intégration* (Chelsea, New York, 1973), p. 292.
13. A. A. Samarskii, V. F. Tishkin, A. P. Favorskii and M. Yu. Shashkov, Operational finite-difference schemes, *Differential Equations* **17**, 854 (1981).
14. M. Shashkov and S. Steinberg, Support-operator finite-difference algorithms for general elliptic problems, *J. Comput. Phys.* **118**, 131 (1995).
15. D. Sun, J. Manges, X. Yuan, and Z. Cendes, Spurious modes in finite-element methods, *IEEE Antennas Propag. Mag.* **37**(5), 12 (1995).
16. N. K. Madsen, Divergence preserving discrete surface integral methods for Maxwell curl equations using non-orthogonal unstructured grids, *J. Comput. Phys.* **119**, 34 (1995).
17. A. A. Dezin, *Multidimensional Analysis and Discrete Models* (CRC Press, Boca Raton, 1995).