

FFT-based Homogenization Methods

Professor:

Francisco Manuel Andrade Pires

Student:

José Luís Passos Vila-Chã

Report presented under the scope of the
Doctoral Program in Mechanical Engineering

Porto, September 2021

Page intentionally left blank.

Contents

List of Figures	vii
List of Tables	ix
1 Continuum Mechanics and Finite Element Method	1
1.1 Kinematics of Deformation	1
1.1.1 Motion	1
1.1.2 Material and spatial descriptions	2
1.1.3 Deformation gradient	2
1.1.3.1 Isochoric/Volumetric decomposition	3
1.1.3.2 Polar decomposition	3
1.2 Strain tensors	4
1.3 Forces and stress measures	4
1.3.0.1 Cauchy stress tensor	5
1.3.0.2 First Piola-Kirchhoff stress tensor	5
1.3.0.3 Kirchhoff stress tensor	5
1.3.0.4 Deviatoric/Hydrostatic decomposition	6
1.4 Heat	6
1.4.0.1 Heat flux vector	6
1.5 Fundamental conservation principles	6
1.5.1 Principle of mass conservation	6
1.5.2 Principle of linear momentum conservation	6
1.5.3 First principle of thermodynamics	7
1.5.4 Second principle of thermodynamics	8
1.5.5 Clausius-Duhem inequality	8
1.6 Mechanical constitutive initial value problem	9
1.6.1 Thermodynamics with internal variables	9
2 Mechanical problem	13
2.0.1 Mechanical constitutive initial value problem	13
2.0.2 Weak equilibrium. The principle of virtual work	14
2.0.3 Mechanical constitutive initial boundary value problem	15
2.1 Time discretization	16
2.2 Finite Element Method	18
2.2.1 Finite element concept	18
2.2.2 Interpolation functions	18
2.2.3 Interpolation matrix and discrete gradient operators	19
2.2.4 Spatial discretization	19

2.2.5	Numerical integration	22
2.3	Linearisation	22
2.4	Constitutive laws	23
3	Thermo field	25
3.1	Governing equations	25
3.2	Thermal constitutive initial value problem	25
3.3	Weak energy balance equation	28
3.4	The thermal initial boundary value problem	28
3.5	Finite Element Method	29
4	Thermo-mechanical problem	31
4.0.1	Thermo-mechanical constitutive initial value problem	32
4.0.2	Weak equilibrium. The principle of virtual work	34
4.0.3	Mechanical constitutive initial boundary value problem	34
4.1	Time discretization	36
4.2	Finite Element Method	38
4.2.1	Interpolation	38
4.2.2	Spatial discretization	39
4.3	Linearisation	41
5	Solution procedures for coupled fields	43
5.1	Context field elimination	43
5.2	Monolithic	44
5.2.1	Numerical considerations	44
5.2.2	Usage examples	46
5.3	Partitioned	47
5.3.1	Operator splits	48
5.3.2	Loosely vs. Strongly coupled schemes	48
5.3.3	Loosely coupled	49
5.3.4	Strongly coupled	52
5.4	Comparison of solution techniques	53
6	Strongly coupled methods for coupled fields	57
6.1	Equations to be solved	57
6.2	A classification scheme for iterative methods	58
6.2.1	Predictor	59
6.2.2	Global Approaches	60
6.2.3	Convergence criteria	62
6.3	One-point iteration function	63
6.3.1	Fixed-point approaches	63
6.3.1.1	Block Jacobi or Schwarz additive	63
6.3.1.2	Block Gauss-Seidel or Schwarz multiplicative	63
6.3.2	Newton's method	64
6.3.3	Constant Underrelaxation	67
6.4	One-point iteration function with memory	67
6.4.0.1	Aitken relaxation	67
6.4.1	Multi-secant methods	69
6.4.1.1	Generalized Broyden	71
6.4.1.2	Anderson mixing	72

6.4.1.3	Generalized Broyden's family	73
6.4.1.4	Anderson's family	73
6.4.1.5	The Broyden-like class	74
6.4.2	Practical considerations	74
6.5	Multipoint iteration functions	80
6.5.1	Finite-Difference Newton Method	80
6.5.2	Newton-Krylov methods	81
6.5.3	Extrapolation techniques with cycling	83
6.6	Multipoint iteration functions with memory	88
6.7	Conclusions	88
Bibliography		91

List of Figures

5.1	Devices of partitioned analysis time-stepping (Felippa et al., 2001).	49
6.1	Geometric interpretation of the fixed-point iteration method in one dimension. The fixed-point of f is sought, which is equivalent to the root of $x - f(x)$	64
6.2	Geometric interpretation of the Newton method in one dimension for an example function f , whose derivative is denoted by f'	65
6.3	Geometric interpretation of the Aitken relaxation in one dimension for an example function f and corresponding interpretation as the secant method.	69
6.4	Geometrical interpretation of Aitken's Δ^2 method.	84

List of Tables

5.1	Summary of the comparison between the FFT-Galerkin method.	55
6.1	Summary of the comparison between method for the solution methods of non-linear systems of equations. n here denotes the number of unknowns and m denotes depending on the context the number of previous iterates considered, the number of fixed point evaluations or the size of the Krylov subspace.	89
6.2	Summary of the update formulas for the solution methods of non-linear systems of equations.	90

Page intentionally left blank.

Chapter 1

Continuum Mechanics and Finite Element Method

This chapter deals with the concepts needed to describe the behavior of a solid undergoing large deformation, as well as, the conservation principles that ensure its mechanical equilibrium. It also presents a succinct overview of the Finite Element Method as a tool to solve mechanical initial value equilibrium problem. These topics are broadly covered in the literature and here the approach used follows [1],

1.1 Kinematics of Deformation

1.1.1 Motion

Let a deformable body \mathcal{B} occupy an open region Ω_0 of the tridimensional Euclidean space \mathcal{E} with a regular boundary $\partial\Omega_0$ in its reference configuration. Its motion, depicted in Figure 1.1, is defined by a smooth one-to-one function

$$\boldsymbol{\varphi}: \Omega \times \mathcal{R} \rightarrow \mathcal{E}, \quad (1.1)$$

mapping each material particle of coordinates \mathbf{X} in the reference configuration to its position \mathbf{x} in the deformed configuration, for a given instant of time t , as

$$\mathbf{x} = \boldsymbol{\varphi}(\mathbf{X}, t) = \boldsymbol{\varphi}_t(\mathbf{X}). \quad (1.2)$$

Thus, the displacement field is defined as

$$\mathbf{u}(\mathbf{X}, t) = \boldsymbol{\varphi}(\mathbf{X}, t) - \mathbf{X}, \quad (1.3)$$

and, since the function that defines the motion is one-to-one, the reference configuration can be recovered as

$$\mathbf{X} = \boldsymbol{\varphi}^{-1}(\mathbf{x}, t) = \mathbf{x} - \mathbf{u}(\boldsymbol{\varphi}^{-1}(\mathbf{x}, t), t), \quad (1.4)$$

where $\boldsymbol{\varphi}^{-1}$ is the reference mapping function.

1.1.2 Material and spatial descriptions

Dealing with finite deformations, the behavior of the body under analysis can be described with respect to the reference configuration, using the so-called material or Lagrangian description, or to the deformed configuration, using the so-called spatial or Eulerian description.

In the Lagrangian description any field, be it scalar, vectorial or tensorial defined over the body is expressed as a function of the reference configuration, $\mathbf{X} \in \Omega_0$. On the other hand, the Eulerian description of same field is done using the deformed configuration, $\mathbf{x} \in \Omega$.

As such let $\alpha(\mathbf{x}, t)$ be a spatial field and $\beta(\mathbf{X}, t)$ a material field. Their material α_m and spatial β_s descriptions are given by

$$\alpha_m(\mathbf{X}, t) = \alpha(\boldsymbol{\varphi}(\mathbf{X}, t), t), \quad (1.5)$$

$$\beta_s(\mathbf{x}, t) = \beta(\boldsymbol{\varphi}^{-1}(\mathbf{x}, t), t), \quad (1.6)$$

noting that any field associated with a motion of \mathcal{B} can be expressed as a function of time and material or spatial position.

The same distinction between material and spatial descriptions applies to operators such as the divergence and the gradient. The spatial and material gradients, ∇ and ∇_0 , respectively, are defined as

$$\nabla \alpha = \frac{\partial}{\partial \mathbf{x}} \alpha(\mathbf{x}, t), \quad \nabla_0 \beta = \frac{\partial}{\partial \mathbf{X}} \beta(\mathbf{X}, t), \quad (1.7)$$

where the derivatives are taken with respect to the spatial and reference configuration accordingly.

1.1.3 Deformation gradient

The deformation gradient, a second order tensor denoted by \mathbf{F} , is defined as

$$\mathbf{F}(\mathbf{X}, t) \equiv \nabla_0 \boldsymbol{\varphi}(\mathbf{X}, t) = \frac{\partial \mathbf{x}}{\partial \mathbf{X}}, \quad (1.8)$$

or, taking into account that

$$\mathbf{x} = \mathbf{X} + \mathbf{u}(\mathbf{X}, t), \quad (1.9)$$

it can be expressed as

$$\mathbf{F}(\mathbf{X}, t) = \mathbf{I} + \nabla_0 \mathbf{u}. \quad (1.10)$$

The deformation gradient relates the relative position between two neighboring material particles before and after deformation. To see this let \mathbf{X} be the coordinates of some material particle in the reference configuration and $\mathbf{X} + d\mathbf{X}$ the coordinates of some material particle in its neighborhood, their corresponding coordinates in the deformed configuration are given by

$$\mathbf{X} = \mathbf{x} - \mathbf{u}(\mathbf{X}, t), \quad (1.11)$$

$$\mathbf{X} + d\mathbf{X} = \mathbf{x} + d\mathbf{x} - \mathbf{u}(\mathbf{X} + d\mathbf{X}, t). \quad (1.12)$$

Subtracting Equation (1.11) to Equation (1.12), it is found that

$$d\mathbf{X} = d\mathbf{x} + \mathbf{u}(\mathbf{X}, t) - \mathbf{u}(\mathbf{X} + d\mathbf{X}, t) \quad (1.13)$$

$$= (\mathbf{I} + \nabla_0 \mathbf{u}(\mathbf{X}, t)) d\mathbf{x} \quad (1.14)$$

$$= \mathbf{F} d\mathbf{x}. \quad (1.15)$$

Due to this relation, it can be shown that the determinant of the deformation gradient has a physical meaning. It is the local unit volume change, that is,

$$J \equiv \det \mathbf{F} = \frac{dv}{dv_0}, \quad (1.16)$$

where dv_0 is an infinitesimal volume of the body in its reference configuration and dv the infinitesimal volume after deformation.

1.1.3.1 Isochoric/Volumetric decomposition

Any deformation can be locally decomposed in volumetric and isochoric (or distortional) components. From Equation (1.16) it can be gathered that an isochoric deformation is characterized by $J = 1$. As such, the deformation gradient can be decomposed as

$$\mathbf{F} = \mathbf{F}_{\text{iso}} \mathbf{F}_{\text{vol}} = \mathbf{F}_{\text{vol}} \mathbf{F}_{\text{iso}}, \quad (1.17)$$

where the isochoric and volumetric components are defined by

$$\mathbf{F}_{\text{iso}} = (\det \mathbf{F})^{-\frac{1}{3}}, \quad \mathbf{F}_{\text{vol}} = (\det \mathbf{F})^{\frac{1}{3}} \mathbf{I}. \quad (1.18)$$

1.1.3.2 Polar decomposition

The deformation gradient can also be decomposed in rotation and stretch components, the so-called polar decomposition, defined as

$$\mathbf{F} = \mathbf{R}\mathbf{U} = \mathbf{V}\mathbf{R}, \quad (1.19)$$

where \mathbf{R} is the proper orthogonal rotation tensor and \mathbf{U} and \mathbf{V} are the symmetric positive right and left stretch tensors, respectively.

Equation (1.19) has a physical interpretation with the right polar decomposition ($\mathbf{F} = \mathbf{R}\mathbf{U}$) corresponding to a stretch mapping followed by a rotation, and the left polar decomposition ($\mathbf{F} = \mathbf{V}\mathbf{R}$) corresponding to a rotation followed by a stretch mapping. The right \mathbf{U} and left \mathbf{V} stretch tensors are related through the rotation matrix \mathbf{R} as

$$\mathbf{V} = \mathbf{R}\mathbf{U}\mathbf{R}^T, \quad (1.20)$$

and can be obtained from deformation gradient by

$$\mathbf{C} \equiv \mathbf{U}^2 = \mathbf{F}^T \mathbf{F}, \quad \mathbf{B} \equiv \mathbf{V}^2 = \mathbf{F} \mathbf{F}^T, \quad (1.21)$$

where \mathbf{C} and \mathbf{B} are the right and left Cauchy-Green strain tensors.

Since \mathbf{U} and \mathbf{V} are symmetric tensors, they admit the spectral decomposition

$$\mathbf{U} = \sum_{i=1}^3 \lambda_i \mathbf{E}_i^* \otimes \mathbf{E}_i^*, \quad \mathbf{V} = \sum_{i=1}^3 \lambda_i \mathbf{e}_i^* \otimes \mathbf{e}_i^*, \quad (1.22)$$

where λ_i , $i = 1, 2, 3$, are the eigenvalues of both \mathbf{U} and \mathbf{V} and \mathbf{E}_i^* and \mathbf{e}_i^* are the respective eigenvectors.

The eigenvectors of left \mathbf{V} and right \mathbf{U} stretch tensors are related through

$$\mathbf{e}_i^* = \mathbf{R} \mathbf{E}_i^*. \quad (1.23)$$

forming two orthogonal bases. These vectors define the Lagrangian and Eulerian principal directions, respectively, allowing for the expression of the local stretching from a material particle, associated with any deformation, as a superposition of stretches along the three mutual orthogonal directions.,

1.2 Strain tensors

In Continuum Mechanics there are two main families of strain tensors derived from the deformation gradient and used to describe the body deformation. The Lagrange family strain tensors are defined as

$$\mathbf{E}^{(m)} = \begin{cases} \frac{1}{m} (\mathbf{U}^m - \mathbf{I}), & m \neq 0, \\ \ln(\mathbf{U}), & m = 0, \end{cases} \quad (1.24)$$

where m is a real number, and likewise, the Euler family strain tensors are defined as

$$\mathbf{e}^{(m)} = \begin{cases} \frac{1}{m} (\mathbf{V}^m - \mathbf{I}), & m \neq 0, \\ \ln(\mathbf{V}), & m = 0, \end{cases} \quad (1.25)$$

where m is also real number.

In particular, choosing $m = 0$, one obtains the so-called material and spatial logarithmic strain tensors

$$\mathbf{E}^{(0)} \equiv \ln[\mathbf{U}] = \sum_{i=1}^3 \ln \lambda_i \mathbf{E}_i^* \otimes \mathbf{E}_i^*, \quad (1.26)$$

$$\mathbf{e}^{(0)} \equiv \ln[\mathbf{V}] = \sum_{i=1}^3 \ln \lambda_i \mathbf{e}_i^* \otimes \mathbf{e}_i^*. \quad (1.27)$$

1.3 Forces and stress measures

The deformation of a body is intrinsically related to the forces acting on it. These forces can be divided in two classes, from a purely mechanical point of view: volume (or body) forces, proportional to the mass contained in a volume element, as such measured in force per unit volume, and surface forces, acting on the surface of a volume element, measured as force per unit area. Related to the latter is the concept of stress, that can be described mathematically by second order tensors with different definitions.

1.3.0.1 Cauchy stress tensor

According to Cauchy's theorem the relation between the so-called Cauchy stress vector, $\mathbf{t}(\mathbf{x}, \mathbf{n})$, and the unitary outward vector normal to the deformed surface under analysis, \mathbf{n} , is linear and given by

$$\mathbf{t}(\mathbf{x}, \mathbf{n}) \equiv \boldsymbol{\sigma}(\mathbf{x}) \mathbf{n}, \quad (1.28)$$

where $\boldsymbol{\sigma}$ is the second order Cauchy stress tensor.

The Cauchy stress vector is naturally associated with the deformed configured and thus, expressed in a spatial description and measured in force per unit deformed area. It must also be noted that as a consequence of the balance of angular momentum, the Cauchy stress tensor is symmetric.

1.3.0.2 First Piola-Kirchhoff stress tensor

The First Piola-Kirchhoff stress tensor, \mathbf{P} , can be regarded as the material counterpart of the Cauchy stress tensor, as it establishes a linear dependence between the stress vector $\mathbf{t}_0(\mathbf{X}, \mathbf{m})$, measured in force per unit reference area, and the unitary outward vector normal to the undeformed surface under analysis, \mathbf{m} ,

$$\mathbf{t}_0 = \mathbf{P} \mathbf{m}, \quad (1.29)$$

which must related to the Cauchy stress vector by

$$\mathbf{t}_0 = \frac{da}{da_0} \mathbf{t} = \frac{da}{da_0} \boldsymbol{\sigma} \mathbf{n}, \quad (1.30)$$

where da is the infinitesimal deformed area normal to the unitary vector \mathbf{n} and da_0 the corresponding undeformed area normal to \mathbf{m} . It can be shown that the relation between da and da_0 is

$$\frac{da}{da_0} \mathbf{n} = J \mathbf{F}^{-T} \mathbf{m}, \quad (1.31)$$

and substituting on the equation above motivates the following definition

$$\mathbf{P} \equiv J \boldsymbol{\sigma} \mathbf{F}^{-T}, \quad (1.32)$$

where J is the determinant of the deformation gradient \mathbf{F} and $\boldsymbol{\sigma}$ is the Cauchy stress tensor. From Equation (1.32), one gathers that, in general, the First Piola-Kirchhoff stress tensor is not symmetric.

1.3.0.3 Kirchhoff stress tensor

The Kirchhoff stress tensor, $\boldsymbol{\tau}$, is a widely used symmetric tensor, defined as

$$\boldsymbol{\tau} \equiv J \boldsymbol{\sigma}. \quad (1.33)$$

1.3.0.4 Deviatoric/Hydrostatic decomposition

The Cauchy stress tensor, $\boldsymbol{\sigma}$, can split as

$$\boldsymbol{\sigma} = \boldsymbol{\sigma}_d - p\mathbf{I}, \quad (1.34)$$

where p is the hydrostatic pressure defined as

$$p \equiv -\frac{1}{3} \text{tr} [\boldsymbol{\sigma}], \quad (1.35)$$

and $\boldsymbol{\sigma}_d$ is the deviatoric stress defined as

$$\boldsymbol{\sigma}_d \equiv \boldsymbol{\sigma} - p\mathbf{I}. \quad (1.36)$$

1.4 Heat

Changes in temperature of a body are very often related to heat flowing inside, entering or leaving it. In continuum mechanics, heat is measured in power per unit surface.

1.4.0.1 Heat flux vector

According to Cauchy's theorem the relation between the heat flux across a surface, $h(\mathbf{x}, \mathbf{n})$, and the unitary outward normal to the deformed surface under analysis, \mathbf{n} , is linear and given by

$$h(\mathbf{x}, \mathbf{n}) = -\mathbf{q}(\mathbf{x}) \cdot \mathbf{n}. \quad (1.37)$$

where \mathbf{q} is the heat flux vector.

1.5 Fundamental conservation principles

In Continuum Mechanics, there is a set of conservation principles and thermodynamic laws, that irrespective of the quantities used to describe the mechanical behavior of a body undergoing large deformations must always be satisfied.

1.5.1 Principle of mass conservation

The principle of mass conservation can be stated as

$$\dot{\rho} + \rho \text{div } \dot{\mathbf{u}} = 0, \quad (1.38)$$

where ρ is the material density measured in mass per unit deformed volume.

1.5.2 Principle of linear momentum conservation

The principle of linear momentum conservation can be stated in both material and spatial description. In a spatial description it reads

$$\begin{cases} \operatorname{div} \boldsymbol{\sigma} + \mathbf{b} = \rho \ddot{\mathbf{u}}, & \forall \mathbf{x} \in \Omega, \\ \mathbf{t} = \boldsymbol{\sigma} \mathbf{n}, & \forall \mathbf{x} \in \partial\Omega, \end{cases} \quad (1.39)$$

where \mathbf{b} is the body forces field measured as per unit deformed volume.

One can also write the principle of linear momentum conservation in material coordinates, as

$$\begin{cases} \operatorname{div}_0 \mathbf{P} + \mathbf{b}_0 = \rho_0 \ddot{\mathbf{u}}, & \forall \mathbf{X} \in \Omega_0, \\ \mathbf{t}_0 = \mathbf{P} \mathbf{m}, & \forall \mathbf{X} \in \partial\Omega_0, \end{cases} \quad (1.40)$$

where \mathbf{b}_0 is the body forces field, measured in force per unit undeformed volume, and ρ_0 is the material density, measured in mass per unit undeformed volume. Both these quantities can be found from their spatial counterparts as

$$\mathbf{b}_0 = J\mathbf{b}, \quad \rho_0 = J\rho. \quad (1.41)$$

Take notice of the abuse of language regarding functions defined on the reference configuration Ω_0 and on the deformed configuration Ω . The same symbol, f , is used for a function f defined on Ω and the function $f \circ \boldsymbol{\varphi}$ defined on Ω_0 .

Equations (1.39) and (1.40) are the so-called strong, point-wise or local equilibrium equations, as they enforce the mechanical equilibrium at every material particle of the body.

1.5.3 First principle of thermodynamics

Let e be the internal energy per unit mass, r the heat supply per unit mass and \mathbf{q} the heat flux, then the first principle of thermodynamics pertaining to the balance of energy can be written in the spatial description as

$$\begin{cases} \rho \dot{e} = \boldsymbol{\sigma} : \mathbf{D} + \rho r - \operatorname{div} \mathbf{q}, & \forall \mathbf{x} \in \Omega, \\ \mathbf{t} = \boldsymbol{\sigma} \mathbf{n}, & \forall \mathbf{x} \in \partial\Omega, \\ h = \mathbf{q} \cdot \mathbf{n}, & \forall \mathbf{x} \in \partial\Omega. \end{cases} \quad (1.42)$$

The second order tensor \mathbf{D} denotes a strain rate measure, such that the double contraction $\boldsymbol{\sigma} : \mathbf{D}$ represents the stress power per unit volume in the deformed configuration of body. In material coordinates, it reads

$$\begin{cases} \rho_0 \dot{e} = \mathbf{P} : \dot{\mathbf{F}} + \rho_0 r - \operatorname{div}_0 \mathbf{q}_0, & \forall \mathbf{X} \in \Omega_0, \\ \mathbf{t}_0 = \mathbf{P} \mathbf{m}, & \forall \mathbf{X} \in \partial\Omega_0, \\ h_0 = \mathbf{q}_0 \cdot \mathbf{m}, & \forall \mathbf{X} \in \partial\Omega_0, \end{cases} \quad (1.43)$$

where \mathbf{q}_0 is the Piola transformation of \mathbf{q} , i.e.,

$$\mathbf{q}_0 = J\mathbf{F}^{-T} \mathbf{q}, \quad (1.44)$$

and

$$h_0 = Jh. \quad (1.45)$$

1.5.4 Second principle of thermodynamics

The local entropy balance can be written as

$$\rho \dot{s} = -\operatorname{div} \left[\frac{\mathbf{q}}{\theta} \right] + \frac{\rho r}{\theta} + \hat{s}, \quad (1.46)$$

where \hat{s} is the entropy production. The second principle of thermodynamics postulates that the changes in the entropy in the universe can never be negative, which is mathematically expressed by

$$\hat{s} \geq 0, \quad (1.47)$$

yielding

$$\rho \dot{s} + \operatorname{div} \left[\frac{\mathbf{q}}{\theta} \right] - \frac{\rho r}{\theta} \geq 0, \quad (1.48)$$

where θ and s are the temperature and specific entropy fields respectively. In a material description, it reads

$$\rho_0 \dot{s} + \operatorname{div}_0 \left[\frac{\mathbf{q}_0}{\theta} \right] - \frac{\rho_0 r}{\theta} \geq 0. \quad (1.49)$$

1.5.5 Clausius-Duhem inequality

Combining the first and second thermodynamic principles yields

$$\rho \dot{s} + \operatorname{div} \left[\frac{\mathbf{q}}{\theta} \right] - \frac{1}{\theta} (\rho \dot{e} - \boldsymbol{\sigma} : \mathbf{D} + \operatorname{div} \mathbf{q}) \geq 0, \quad (1.50)$$

From the definition of the specific Helmholtz free energy

$$\psi \equiv e - \theta s, \quad (1.51)$$

and defining the temperature field gradient as $\mathbf{g} = \nabla \theta$, it is possible to establish the so-called Clausius-Duhem inequality in the spatial description as

$$\underbrace{\boldsymbol{\sigma} : \mathbf{D} - \rho (\dot{\psi} + s \dot{\theta})}_{\mathcal{D}_{\text{mech}}} - \underbrace{\frac{1}{\theta} \mathbf{q} \cdot \mathbf{g}}_{\mathcal{D}_{\text{cond}}} \geq 0, \quad (1.52)$$

where the identity

$$\operatorname{div} \left[\frac{\mathbf{q}}{\theta} \right] = \frac{1}{\theta} \operatorname{div} \mathbf{q} - \frac{1}{\theta^2} \mathbf{q} \cdot \nabla \theta. \quad (1.53)$$

is used.

From a physical point of view, the Clausius-Duhem inequality states that the energy dissipation per unit deformed volume is always non-negative. Moreover the terms in the inequality can be split between the mechanical internal dissipation $\mathcal{D}_{\text{mech}}$ and the dissipation due to heat conduction $\mathcal{D}_{\text{cond}}$. From

$$\hat{s} = \boldsymbol{\sigma} : \mathbf{D} - \rho (\dot{\psi} + s \dot{\theta}) - \frac{1}{\theta} \mathbf{q} \cdot \mathbf{g}, \quad (1.54)$$

assuming that the process leads to an uniform temperature distribution, yields for the mechanical dissipation $\mathcal{D}_{\text{mech}}$,

$$\mathcal{D}_{\text{mech}} = \hat{s}|_{\theta \text{ uniform}} = \boldsymbol{\sigma} : \mathbf{D} - \rho (\dot{\psi} + s\dot{\theta}), \quad (1.55)$$

since conduction is excluded and only mechanical and temperature transient effects remain. If on the other hand, only conduction effects are retained, i.e., assuming the process to be isothermic, isochoric and isovolumetric, the dissipation due to conduction, $\mathcal{D}_{\text{cond}}$, is obtained as

$$\mathcal{D}_{\text{cond}} = -\frac{1}{\theta} \mathbf{q} \cdot \mathbf{g}. \quad (1.56)$$

Equation (1.52) can also be written as

$$\mathbf{P} : \dot{\mathbf{F}} - \rho_0 (\dot{\psi} + s\dot{\theta}) - \frac{1}{\theta} \mathbf{q}_0 \cdot \mathbf{g}_0 \geq 0, \quad (1.57)$$

where $\mathbf{g}_0 = \nabla_0 \theta$, aplying the Piola transformation, and as

$$\boldsymbol{\tau} : \mathbf{D} - \rho_0 (\dot{\psi} + s\dot{\theta}) - \frac{J}{\theta} \mathbf{q} \cdot \mathbf{g} \geq 0, \quad (1.58)$$

multiplying it by J and attending to the definition of the Kirchhoff stress tensor, where the left hand side represents now the energy dissipation per unit reference volume.

1.6 Mechanical constitutive initial value problem

In Continuum Mechanics, a constitutive model is a set of equations, also called constitutive equations, establishing the stress-strain relation for some material. Before going further, it is important to define a thermokinetic process of a body \mathcal{B} as

$$\text{thermokinetic process: } \{\boldsymbol{\varphi}(\mathbf{X}, t), \theta(\mathbf{X}, t)\}, \quad (1.59)$$

ans a calordynamic process of \mathcal{B} as

$$\text{calorodynamic process: } \{\boldsymbol{\sigma}(\mathbf{X}, t), e(\mathbf{X}, t), s(\mathbf{X}, t), r(\mathbf{X}, t), \mathbf{b}(\mathbf{X}, t), \mathbf{q}(\mathbf{X}, t)\}, \quad (1.60)$$

which satisfies the fundamental conservation principles previously introduced.

It is also important to note that any constitutive model must satisfy a set of constitutive axioms, explained in detail by ? . As these are too general to be used directly in practice, a particular case of the general history functional-based constitutive theory based on the thermodynamics with internal variables approach is used.

1.6.1 Thermodynamics with internal variables

The values of $\boldsymbol{\sigma}$, ψ , s and \mathbf{q} at a material particle define its thermodynamic state, assuming \mathbf{b} follows from the balance of linear momentum and r from the energy

balance equation. In thermodynamics with interval variables approach, that thermodynamic state is assumed to be completely defined by the instantaneous values of a finite number of state variables

$$\{\mathbf{F}, \theta, \mathbf{g}, \boldsymbol{\alpha}\}. \quad (1.61)$$

at a given instant of the calorodynamic process, where

$$\boldsymbol{\alpha} = \{\alpha_k\} \quad (1.62)$$

is a set of internal variables, scalar or tensorial in nature, associated with dissipative mechanisms. As such, the accuracy of the constitutive model depends strongly on the appropriate choice of internal variables, as these contain the relevant information about the thermodynamical history of the material.

Accordingly, the specific Helmholtz free energy is postulated to follow

$$\psi = \psi(\mathbf{F}, \theta, \boldsymbol{\alpha}). \quad (1.63)$$

To find the constitutive equations for the stress tensor and the entropy, one can substitute

$$\dot{\psi} = \frac{\partial \psi}{\partial \mathbf{F}} : \dot{\mathbf{F}} + \frac{\partial \psi}{\partial \theta} \dot{\theta} + \frac{\partial \psi}{\partial \alpha_k} \dot{\alpha}_k, \quad (1.64)$$

found from the chain rule, on the Clausius-Duhem equation, Equation (1.52), obtaining

$$\left(\boldsymbol{\sigma} \mathbf{F}^{-T} - \rho \frac{\partial \psi}{\partial \mathbf{F}} \right) : \dot{\mathbf{F}} - \rho \left(s + \frac{\partial \psi}{\partial \theta} \right) \dot{\theta} - \rho \frac{\partial \psi}{\partial \alpha_k} \dot{\alpha}_k - \frac{1}{\theta} \mathbf{q} \cdot \mathbf{g} \geq 0, \quad (1.65)$$

where the velocity gradient was adopted to set the work conjugacy as

$$\boldsymbol{\sigma} : \mathbf{D} = \boldsymbol{\sigma} : \mathbf{L} = \boldsymbol{\sigma} : \dot{\mathbf{F}} \mathbf{F}^{-1} = \boldsymbol{\sigma} \mathbf{F}^{-T} : \dot{\mathbf{F}}. \quad (1.66)$$

Since the Clausius-Duhem inequality must hold for any thermokinetic process and so remain valid for any set $\{\dot{\mathbf{F}}(t), \dot{\theta}(t)\}$, the Cauchy stress and entropy constitutive equations must be

$$\boldsymbol{\sigma} = \rho \frac{\partial \psi}{\partial \mathbf{F}} \mathbf{F}^T, \quad (1.67)$$

$$s = -\frac{\partial \psi}{\partial \theta}. \quad (1.68)$$

It is also possible to write the constitutive equations for the Kirchhoff stress tensor as

$$\boldsymbol{\tau} = J \rho \frac{\partial \psi}{\partial \mathbf{F}} \mathbf{F}^T, \quad (1.69)$$

multiplying Equation (1.67) by J , and the first Piola-Kirchhoff stress tensor as

$$\mathbf{P} = \rho_0 \frac{\partial \psi}{\partial \mathbf{F}} \quad (1.70)$$

multiplying Equation (1.65) also by J .

For each internal variable α_k of the set α of internal variables, the conjugate thermodynamical forces are defined to be

$$A_k \equiv \rho_0 \frac{\partial \psi}{\partial \alpha_k}, \quad (1.71)$$

so that the Clausius-Duhem equation can be written in a reduced form as

$$- \mathbf{A} * \dot{\boldsymbol{\alpha}} - \frac{J}{\theta} \mathbf{q} \cdot \mathbf{g} \geq 0, \quad (1.72)$$

where \mathbf{A} is the set of conjugate thermodynamical forces and $*$ denotes the appropriate product operation.

To completely define the constitutive model, one still needs to postulate the constitutive equations for the flux variables $\dot{\boldsymbol{\alpha}}$ and $\frac{1}{\theta} \mathbf{q}$. These are given by

$$\dot{\boldsymbol{\alpha}} = f(\mathbf{F}, \theta, \mathbf{g}, \boldsymbol{\alpha}), \quad (1.73)$$

$$\frac{1}{\theta} \mathbf{q} = g(\mathbf{F}, \theta, \mathbf{g}, \boldsymbol{\alpha}). \quad (1.74)$$

A sufficient condition for the previous constitutive functions to satisfy the Clausius-Duhem inequality is the hypothesis of normal dissipativity, whereby one defines the constitutive functions for the flux variables as

$$\dot{\alpha}_k = - \frac{\partial \Xi}{\partial A_k}, \quad \frac{1}{\theta} \mathbf{q} = - \frac{\partial \Xi}{\partial \mathbf{g}}, \quad (1.75)$$

where the dissipation potential is

$$\Xi = \Xi(\mathbf{A}, \mathbf{g}; \mathbf{F}, \theta, \boldsymbol{\alpha}), \quad (1.76)$$

a convex function with respect to each A_k and \mathbf{g} , and zero valued at the origin, $\{\mathbf{A}, \mathbf{g}\} = \{\mathbf{0}, \mathbf{0}\}$. Note that in the previous definition the state variables appear only as parameters.

Chapter 2

Mechanical problem

In the following chapter, the general framework presented in the previous chapter is applied to a purely mechanical analysis, neglecting the thermal terms.

2.0.1 Mechanical constitutive initial value problem

In the purely mechanical case, with all the quantities related to the thermal domain removed, a constitutive model based on internal variables is established by the following set of equations

$$\mathbf{P} = \rho_0 \frac{\partial \psi}{\partial \mathbf{F}}, \quad (2.1)$$

$$\psi = \psi(\mathbf{F}, \boldsymbol{\alpha}), \quad (2.2)$$

$$\dot{\boldsymbol{\alpha}} = f(\mathbf{F}, \boldsymbol{\alpha}). \quad (2.3)$$

Thus, the spatial mechanical constitutive initial value problem can be stated as follows

Problem 2.1 | Spatial mechanical constitutive initial value problem.

Given the initial values of the internal variables, $\boldsymbol{\alpha}(t_0)$, and the history of the deformation gradient

$$\mathbf{F}(t), \quad t \in [t_0, t_{\text{end}}], \quad (2.4)$$

find the functions for $\boldsymbol{\sigma}(t)$ and $\boldsymbol{\alpha}(t)$ such that the constitutive equations

$$\boldsymbol{\sigma} = \rho \frac{\partial \psi}{\partial \mathbf{F}} \mathbf{F}^T, \quad (2.5)$$

$$\psi = \psi(\mathbf{F}, \boldsymbol{\alpha}), \quad (2.6)$$

$$\dot{\boldsymbol{\alpha}} = f(\mathbf{F}, \boldsymbol{\alpha}). \quad (2.7)$$

are satisfied for every $t \in [t_0, t_{\text{end}}]$.

Likewise, in a material description it can be stated as

Problem 2.2 | Material mechanical constitutive initial value problem.

Given the initial values of the internal variables, $\alpha(t_0)$, and the history of the deformation gradient

$$\mathbf{F}(t), \quad t \in [t_0, t_{\text{end}}], \quad (2.8)$$

find the functions for $\mathbf{P}(t)$ and $\alpha(t)$ such that the constitutive equations

$$\mathbf{P} = \rho_0 \frac{\partial \psi}{\partial \mathbf{F}}, \quad (2.9)$$

$$\psi = \psi(\mathbf{F}, \alpha), \quad (2.10)$$

$$\dot{\alpha} = f(\mathbf{F}, \alpha). \quad (2.11)$$

are satisfied for every $t \in [t_0, t_{\text{end}}]$.

2.0.2 Weak equilibrium. The principle of virtual work

The strong equations that enforce the equilibrium of a body can be written using the spatial description as

$$\rho \ddot{\mathbf{u}} = \text{div } \boldsymbol{\sigma} + \mathbf{b} \quad \text{in } \Omega, \quad (2.12)$$

and the material description as

$$\rho_0 \ddot{\mathbf{u}} = \text{div}_0 \mathbf{P} + \mathbf{b}_0 \quad \text{in } \Omega_0. \quad (2.13)$$

From a practical standpoint, finding the exact solution to the strong equilibrium equations in the context of real engineering problems is most often nearly or entirely impossible. Most numerical methods obtain only approximate solutions to the so-called weak equilibrium equations to circumvent this problem. These result from relaxing the strong equilibrium equations so that the solutions need only satisfy the equilibrium equations in an average sense instead of satisfying them pointwise. This is achieved through an integration over the body volume. The weak equilibrium equations can be found making use of several energetic and weighted residual methods, such as the Virtual Work Principle used here.

Problem 2.3 | Principle of virtual work (spatial version).

The Virtual Work Principle states, in a spatial description, that the body is in equilibrium if and only if the Cauchy stress field satisfies

$$\int_{\Omega} [\boldsymbol{\sigma} : \nabla \boldsymbol{\eta} - (\mathbf{b} - \rho \ddot{\mathbf{u}}) \cdot \boldsymbol{\eta}] dv - \int_{\partial\Omega} \mathbf{t} \cdot \boldsymbol{\eta} da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_u, \quad (2.14)$$

where \mathcal{V}_u is the space of virtual displacement of the body, defined by the space of sufficiently regular arbitrary displacements

$$\boldsymbol{\eta} : \Omega \rightarrow \mathcal{U} \quad (2.15)$$

where \mathcal{U} is the n -dimension vector associated with \mathcal{E} .

The principle of virtual work can be expressed in a completely equivalent way using a material description.

Problem 2.4 | Principle of virtual work (material version).

The Virtual Work Principle states, in a material description, that the body is in equilibrium if and only if the First Piola-Kirchhoff stress field satisfies

$$\int_{\Omega_0} [\mathbf{P} : \nabla_0 \boldsymbol{\eta} - (\mathbf{b}_0 - \rho_0 \ddot{\mathbf{u}}) \cdot \boldsymbol{\eta}] dv - \int_{\partial\Omega_0} \mathbf{t}_0 \cdot \boldsymbol{\eta} da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_{u,0}, \quad (2.16)$$

where $\mathcal{V}_{u,0}$ is the space of virtual displacement of the body, defined by the space of sufficiently regular arbitrary displacements

$$\boldsymbol{\eta} : \Omega_0 \rightarrow \mathcal{U}. \quad (2.17)$$

2.0.3 Mechanical constitutive initial boundary value problem

It is now possible to pose the mechanical constitutive initial value problem in its weak form. Assume that a body \mathcal{B} is made from a generic material, characterized by a given constitutive model, whose internal variables are known at the initial time, as presented in Figure ?? . In addition, it is assumed that the interior of the body was subjected to a prescribed history of body forces, $\mathbf{b}(\mathbf{X}, t)$, $t \in [t_0, t_{\text{end}}]$, and to the following boundary conditions:

- **Natural (or Neumann) boundary condition:** The boundary portion $\Omega_{\text{traction},0}$ of \mathcal{B} is subjected to a prescribed history of traction forces, $\mathbf{t}_{\text{presc}}(\mathbf{X}, t)$, $\mathbf{X} \in \partial\Omega_{\text{traction},0}$, $t \in [t_0, t_{\text{end}}]$,
- **Essential (or Dirichlet) boundary condition:** The boundary portion $\Omega_{\text{motion},0}$ of \mathcal{B} is subjected to a prescribed displacement field history, $\mathbf{u}_{\text{presc}}(\mathbf{X}, t)$, such that

$$\boldsymbol{\varphi}(\mathbf{X}, t) = \mathbf{X} + \mathbf{u}_{\text{presc}}(\mathbf{X}, t), \quad \mathbf{X} \in \partial\Omega_{\text{motion},0}, \quad t \in [t_0, t_{\text{end}}].$$

It is also convenient to define the set of kinematically admissible displacements of \mathcal{B} as the set of all sufficiently regular displacement functions that satisfy the essential boundary condition (?),

$$\mathcal{K}_u \equiv \{\mathbf{u} : \Omega_0 \times \mathcal{R} \rightarrow \mathcal{U} \mid \mathbf{u}(\mathbf{X}, t) = \mathbf{u}_{\text{presc}}(\mathbf{X}, t), \quad \mathbf{X} \in \partial\Omega_{\text{motion},0}, \quad t \in [t_0, t_{\text{end}}]\}. \quad (2.18)$$

So the weak form of the quasi-static mechanical constitutive initial boundary value problem can be stated in a spatial description as follows

Problem 2.5 | Spatial mechanical initial BVP.

Find a kinematically admissible displacement function, $\mathbf{u} \in \mathcal{K}_u$, such that for every $t \in [t_0, t_{\text{end}}]$, the body \mathcal{B} is in equilibrium as stated by the Virtual Work Principle

$$\int_{\Omega} [\boldsymbol{\sigma} : \nabla \boldsymbol{\eta} - (\mathbf{b} - \rho \ddot{\mathbf{u}}) \cdot \boldsymbol{\eta}] dv - \int_{\partial\Omega} \mathbf{t} \cdot \boldsymbol{\eta} da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_u, \quad (2.19)$$

where the space of virtual displacements at time t is defined by

$$\mathcal{V}_u \equiv \{\boldsymbol{\eta} : \Omega \rightarrow \mathcal{U} \mid \boldsymbol{\eta} = \mathbf{0} \text{ in } \boldsymbol{\varphi}(\partial\Omega_{\text{motion},0}, t)\}, \quad (2.20)$$

and at each point of \mathcal{B} , the Cauchy stress tensor is the solution of spatial mechanical constitutive initial values problem.

and in the material description as

Problem 2.6 | Material mechanical initial BVP.

Find a kinematically admissible displacement function, $\mathbf{u} \in \mathcal{K}_u$, such that for every $t \in [t_0, t_{\text{end}}]$, the body \mathcal{B} is in equilibrium as stated by the Virtual Work Principle

$$\int_{\Omega_0} [\mathbf{P} : \nabla_0 \boldsymbol{\eta} - (\mathbf{b}_0 - \rho_0 \ddot{\mathbf{u}}) \cdot \boldsymbol{\eta}] d\nu - \int_{\partial\Omega_0} \mathbf{t}_0 \cdot \boldsymbol{\eta} da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_{u,0}, \quad (2.21)$$

where the space of virtual displacements at time t is defined by

$$\mathcal{V}_{u,0} \equiv \{\boldsymbol{\eta} : \Omega_0 \rightarrow \mathcal{U} \mid \boldsymbol{\eta} = \mathbf{0} \text{ in } \partial\Omega_{\text{motion},0}\}, \quad (2.22)$$

and at each point of \mathcal{B} , the First Piola-Kirchhoff stress tensor is the solution of material mechanical constitutive initial values problem.

2.1 Time discretization

Given a generic path-dependent model, i.e., a model in which the stress state does not depend only on the instantaneous deformation state but also on the deformation history, the solution of the constitutive initial value problem for a given set of initial conditions is usually not known for complex strain paths $\mathbf{F}(t)$. Thus, there is a need to use an appropriate numerical algorithm to integrate the rate constitutive equations.

In general, the algorithms for integrating rate constitutive equations are obtained adopting some time (or pseudo-time) discretization and some hypothesis on the deformation path between adjacent time stations. In the present document, an algorithm is adopted based on approximated incremental constitutive functions. Attending to the mechanical constitutive initial boundary value problem and considering the time increment $[t_n, t_{n+1}]$, this approach is comprised by the following two requirements:

- **Cauchy and First Piola-Kirchhoff stress tensors.** Considering a time increment $[t_n, t_{n+1}]$ and given the set $\boldsymbol{\alpha}_n$ of internal variables at t_n , the deformation gradient \mathbf{F}_{n+1} at time t_{n+1} determines the stress $\boldsymbol{\sigma}_{n+1}$ uniquely through

$$\boldsymbol{\sigma}_{n+1} = \hat{\boldsymbol{\sigma}}(\boldsymbol{\alpha}_n, \mathbf{F}_{n+1}), \quad (2.23)$$

where $\hat{\boldsymbol{\sigma}}$ is the incremental constitutive function for the Cauchy stress tensor.

Similarly, the First Piola-Kirchhoff stress tensor \mathbf{P}_{n+1} must be uniquely determined by the prescribed deformation gradient \mathbf{F}_{n+1} prescribed at t_{n+1} as

$$\mathbf{P}_{n+1} = \hat{\mathbf{P}}(\boldsymbol{\alpha}_n, \mathbf{F}_{n+1}), \quad (2.24)$$

where $\hat{\mathbf{P}}$ is the incremental constitutive function for the First Piola-Kirchhoff stress tensor.

- **Set of internal variables.** Assuming that the set of internal variables $\boldsymbol{\alpha}_n$ is known at t_n , the set of internal variables must be uniquely determined by the prescribed deformation gradient \mathbf{F}_{n+1} prescribed at t_{n+1} as

$$\boldsymbol{\alpha}_{n+1} = \hat{\boldsymbol{\alpha}}(\boldsymbol{\alpha}_n, \mathbf{F}_{n+1}), \quad (2.25)$$

where $\hat{\boldsymbol{\alpha}}$ is the incremental constitutive function for the set of internal variables.

Generally, the numerical constitutive laws are nonlinear and path-independent within one increment. In other words, within each increment, $\boldsymbol{\sigma}_{n+1}$ and $\boldsymbol{\alpha}_{n+1}$, they are functions of \mathbf{F}_{n+1} alone with the argument $\boldsymbol{\alpha}_n$ constant within the same time interval.

Making use of the aforementioned time discretization, one can state the weak form of the mechanical constitutive initial boundary value problem in the spatial description as

Problem 2.7 | Spatial incremental mechanical initial BVP

Given the set of internal variables $\boldsymbol{\alpha}_n$ at t_n , the prescribed body and traction force fields \mathbf{b}_{n+1} and \mathbf{t}_{n+1} at t_{n+1} , and the prescribed deforming gradient \mathbf{F}_{n+1} at t_{n+1} , find the kinematically admissible displacement field $\mathbf{u}_{n+1} \in \mathcal{K}_{u,n+1}$ such that the body \mathcal{B} is in equilibrium as stated by the virtual Work Principle

$$\int_{\Omega_{n+1}} [\hat{\boldsymbol{\sigma}}(\mathbf{F}_{n+1}, \boldsymbol{\alpha}_n) : \nabla \boldsymbol{\eta} - (\mathbf{b}_{n+1} - \rho \ddot{\mathbf{u}}_{n+1}) \cdot \boldsymbol{\eta}] dV - \int_{\partial\Omega_{n+1}} \mathbf{t}_{n+1} \cdot \boldsymbol{\eta} dA = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_u, \quad (2.26)$$

where the space of kinematically admissible displacement fields $\mathcal{K}_{u,n+1}$ is defined by

$$\mathcal{K}_{u,n+1} \equiv \{\mathbf{u} : \Omega_0 \times \mathcal{B} \rightarrow \mathcal{U} \mid \mathbf{u}_{n+1}(\mathbf{X}) = \mathbf{u}_{\text{presc},n+1}(\mathbf{X}), \mathbf{X} \in \partial\Omega_{\text{motion},0}\}. \quad (2.27)$$

and in the material description as

Problem 2.8 | Material incremental mechanical initial BVP

Given the set of internal variables $\boldsymbol{\alpha}_n$ at t_n , the prescribed body and traction force fields $\mathbf{b}_{0,n+1}$ and $\mathbf{t}_{0,n+1}$ at t_{n+1} , and the prescribed deforming gradient \mathbf{F}_{n+1} at t_{n+1} , find the kinematically admissible displacement field $\mathbf{u}_{n+1} \in \mathcal{K}_{u,n+1}$ such that the body \mathcal{B} is in equilibrium as stated by the virtual Work Principle

$$\int_{\Omega_0} [\hat{\mathbf{P}}(\mathbf{F}_{n+1}, \boldsymbol{\alpha}_n) : \nabla_0 \boldsymbol{\eta} - (\mathbf{b}_{0,n+1} - \rho_0 \ddot{\mathbf{u}}_{n+1}) \cdot \boldsymbol{\eta}] dV - \int_{\partial\Omega_{0,n+1}} \mathbf{t}_{0,n+1} \cdot \boldsymbol{\eta} dA = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_{u,0}, \quad (2.28)$$

where the space of kinematically admissible displacement fields $\mathcal{K}_{u,n+1}$ is defined by

$$\mathcal{K}_{u,n+1} \equiv \{\mathbf{u} : \Omega_0 \times \mathcal{B} \rightarrow \mathcal{U} \mid \mathbf{u}_{n+1}(\mathbf{X}) = \mathbf{u}_{\text{presc},n+1}(\mathbf{X}), \mathbf{X} \in \partial\Omega_{\text{motion},0}\}. \quad (2.29)$$

2.2 Finite Element Method

With the incremental weak form of the mechanical constitutive initial boundary value problem now established, an approximated solution can be found using the Finite Element Method.

2.2.1 Finite element concept

The first in the Finite Element method is to discretize the continuum domain Ω in a finite set of n_{elem} mutually exclusive subdomains called finite elements $\Omega^{(e)}$. The discretized domain, ${}^h\Omega$, is therefore an approximation to the continuum domain expressed by

$$\Omega \approx {}^h\Omega \equiv \bigcup_{e=1}^{n_{\text{elem}}} \Omega^{(e)}. \quad (2.30)$$

The spaces of virtual displacements \mathcal{V}_u and $\mathcal{V}_{u,0}$ as well as the space of kinematically admissible displacement fields \mathcal{K}_u are also discretized in the same way, with their discretized forms denoted by ${}^h\mathcal{V}_u$, ${}^h\mathcal{V}_{u,0}$ and ${}^h\mathcal{K}_u$.

2.2.2 Interpolation functions

Let e be a generic finite element with n_{nodes} nodes, where each node i of coordinates \mathbf{x}^i is associated with an interpolation function $N_i^{(e)}$. These interpolation functions are often called shape functions and perform the required field interpolations inside the element domain $\Omega^{(e)}$.

Letting $a(\mathbf{x})$ be a generic field defined over $\Omega^{(e)}$, its interpolation at any point \mathbf{x} inside the element is defined by the element shape functions as

$$a(\mathbf{x}) \approx {}^h a(\mathbf{x}) \equiv \sum_{i=1}^{n_{\text{nodes}}} a(\mathbf{x}_i) N_i^{(e)}(\mathbf{x}). \quad (2.31)$$

If instead $a(\mathbf{x})$ is instead a generic field defined over the global domain Ω , the interpolation of $a(\mathbf{x})$ at any point \mathbf{x} is defined by the global shape functions as

$$a(\mathbf{x}) \approx {}^h a(\mathbf{x}) \equiv \sum_{i=1}^{n_{\text{points}}} a(\mathbf{x}_i) N_i^g(\mathbf{x}), \quad (2.32)$$

where n_{points} is the total number of nodes of the finite element mesh. The discretized spaces ${}^h\mathcal{V}_u$ and ${}^h\mathcal{K}_u$ can now be defined as

$${}^h\mathcal{K}_u \equiv \left\{ {}^h\mathbf{u}(\mathbf{x}) = \sum_{i=1}^{n_{\text{points}}} \mathbf{u}(\mathbf{x}_i) N_i^g(\mathbf{x}) \mid \mathbf{u}(\mathbf{x}_i) = \mathbf{u}_{\text{presc}}(\mathbf{x}_i) \text{ if } \mathbf{x}_i \in \partial\Omega_{\text{motion},0} \right\}, \quad (2.33)$$

$${}^h\mathcal{V}_u \equiv \left\{ {}^h\boldsymbol{\eta}(\mathbf{x}) = \sum_{i=1}^{n_{\text{points}}} \boldsymbol{\eta}(\mathbf{x}_i) N_i^g(\mathbf{x}) \mid \boldsymbol{\eta}(\mathbf{x}_i) = \mathbf{0} \text{ if } \mathbf{x}_i \in \partial\Omega_{\text{motion},0} \right\} \quad (2.34)$$

Quantities defined on the reference configuration Ω_0 accepted a treatment entirely similar to the one described above, and thus is omitted.

2.2.3 Interpolation matrix and discrete gradient operators

The global shape functions can be conveniently assembled in the so-called global interpolation matrix as

$$\mathbf{N}^g(\mathbf{x}) \equiv \left[\text{diag}[N_1^g(\mathbf{x})] \text{diag}[N_2^g(\mathbf{x})] \cdots \text{diag}[N_{n_{\text{points}}}^g(\mathbf{x})] \right], \quad (2.35)$$

where $\text{diag}[N_i^g]$ is a diagonal matrix $n_{\text{dim}} \times n_{\text{dim}}$

$$\text{diag}[N_i^g(\mathbf{x})] \equiv \begin{bmatrix} N_i^g & 0 & \cdots & 0 \\ 0 & N_i^g & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & N_i^g \end{bmatrix} \quad (2.36)$$

where n_{dim} is the number of degrees of freedom per node.

Defining the global vector of nodal displacements as

$$\mathbf{u} = \left[u_1^1, \dots, u_{n_{\text{dim}}}^1, \dots, u_1^{n_{\text{points}}}, \dots, u_{n_{\text{dim}}}^{n_{\text{points}}} \right]^T, \quad (2.37)$$

the displacement field $\mathbf{u}(\mathbf{x})$ defined over the global domain Ω , can be found from Equation (2.32) at any point \mathbf{x} as

$${}^h\mathbf{u}(\mathbf{x}) \equiv \mathbf{N}^g(\mathbf{x})\mathbf{u}, \quad {}^h\mathbf{u} \in {}^h\mathcal{K}_u. \quad (2.38)$$

2.2.4 Spatial discretization

Applying the aforementioned finite element discretization to the incremental mechanical constitutive initial boundary value problem, we can then write in the spatial description

$$\int_{h\Omega} \left[\hat{\boldsymbol{\sigma}}^T \mathbf{B}^g \boldsymbol{\eta} - (\mathbf{b}_{n+1} - \rho \ddot{\mathbf{u}}_{n+1}) \cdot \mathbf{N}^g \boldsymbol{\eta} \right] d\nu - \int_{\partial^h \Omega_{\text{traction}}} \mathbf{t}_{n+1} \cdot \mathbf{N}^g \boldsymbol{\eta} da = 0, \quad \forall \boldsymbol{\eta} \in {}^h\mathcal{V}_u, \quad (2.39)$$

where \mathbf{B}^g is the discrete symmetric global gradient operator, defined for a 2D problem in cartesian coordinates as

$$\mathbf{B}^g \equiv \begin{bmatrix} \frac{\partial N_1^g}{\partial x} & 0 & \frac{\partial N_2^g}{\partial x} & 0 & \cdots & \frac{\partial N_{n_{\text{points}}}^g}{\partial x} & 0 \\ 0 & \frac{\partial N_1^g}{\partial y} & 0 & \frac{\partial N_2^g}{\partial y} & \cdots & 0 & \frac{\partial N_{n_{\text{points}}}^g}{\partial y} \\ \frac{\partial N_1^g}{\partial y} & \frac{\partial N_1^g}{\partial x} & \frac{\partial N_2^g}{\partial y} & \frac{\partial N_2^g}{\partial x} & \cdots & \frac{\partial N_{n_{\text{points}}}^g}{\partial y} & \frac{\partial N_{n_{\text{points}}}^g}{\partial x} \end{bmatrix}. \quad (2.40)$$

Equation (2.39) can be rewritten as

$$\left\{ \int_{h\Omega} \left[\mathbf{B}^{gT} \hat{\boldsymbol{\sigma}}(\boldsymbol{\alpha}_n, \mathbf{F}_{n+1}) - \mathbf{N}^{gT} \mathbf{b}_{n+1} + \mathbf{N}^{gT} \rho \ddot{\mathbf{u}}_{n+1} \right] d\nu - \int_{\partial^h \Omega_{\text{traction}}} \mathbf{N}^{gT} \mathbf{t}_{n+1} da \right\}^T \boldsymbol{\eta} = 0, \quad \forall \boldsymbol{\eta} \in {}^h\mathcal{V}_u, \quad (2.41)$$

and, since it must be satisfied for any $\boldsymbol{\eta} \in {}^h\mathcal{V}_u$, the incremental quasi-static discretized mechanical constitutive initial boundary value problem can thus be stated in the spatial description as

Problem 2.9 | Spatial incremental discretized mechanical initial BVP.

Given the set of internal variables $\boldsymbol{\alpha}_n$ at t_n , the prescribed body and traction force fields \mathbf{b}_{n+1} and \mathbf{t}_{n+1} , and the prescribed deformation gradient \mathbf{F}_{n+1} at t_{n+1} , find the kinematically admissible nodal displacement field $\mathbf{u}_{n+1} \in {}^h\mathcal{K}_{u,n+1}$ such that the body \mathcal{B} is in equilibrium as stated by the Virtual Work Principle

$$\mathbf{M}\ddot{\mathbf{u}}_{n+1} + \mathbf{f}^{\text{int}}(\mathbf{u}_{n+1}) - \mathbf{f}_{n+1}^{\text{ext}} = \mathbf{0}, \quad (2.42)$$

where \mathbf{f}^{int} e $\mathbf{f}_{n+1}^{\text{ext}}$ are the global vectors of internal and external forces defined as

$$\mathbf{f}^{\text{int}} \equiv \int_{h\Omega} \mathbf{B}^g{}^T \hat{\boldsymbol{\sigma}}(\mathbf{F}_{n+1}, \boldsymbol{\alpha}_n) d\nu, \quad (2.43)$$

$$\mathbf{f}_{n+1}^{\text{ext}} \equiv \int_{h\Omega} \mathbf{N}^g{}^T \mathbf{b}_{n+1} d\nu + \int_{\partial^h\Omega_{\text{traction}}} \mathbf{N}^g{}^T \mathbf{t}_{n+1} da, \quad (2.44)$$

and \mathbf{M} is the mass matrix defined as

$$\mathbf{M} = \int_{h\Omega} \rho \mathbf{N}^g{}^T \mathbf{N}^g d\nu. \quad (2.45)$$

In a material description, Equation (2.41) is written as

$$\left\{ \int_{h\Omega_0} \left[\mathbf{G}^g{}^T \hat{\mathbf{P}}(\boldsymbol{\alpha}_n, \mathbf{F}_{n+1}) - \mathbf{N}^g{}^T \mathbf{b}_{0,n+1} + \mathbf{N}^g{}^T \rho_0 \ddot{\mathbf{u}}_{n+1} \right] d\nu - \int_{\partial^h\Omega_{\text{traction},0}} \mathbf{N}^g{}^T \mathbf{t}_{0,n+1} da \right\}^T \boldsymbol{\eta} = 0, \quad \forall \boldsymbol{\eta} \in {}^h\mathcal{V}_{u,0}, \quad (2.46)$$

where \mathbf{G}^g is the discrete global gradient operator, defined for a 2D problem in cartesian coordinates as

$$\mathbf{G}^g \equiv \begin{bmatrix} \frac{\partial N_1^g}{\partial x} & 0 & \frac{\partial N_2^g}{\partial x} & 0 & \dots & \frac{\partial N_{n_{\text{points}}}^g}{\partial x} & 0 \\ 0 & \frac{\partial N_1^g}{\partial x} & 0 & \frac{\partial N_2^g}{\partial x} & 0 & \dots & \frac{\partial N_{n_{\text{points}}}^g}{\partial x} \\ \frac{\partial N_1^g}{\partial y} & 0 & \frac{\partial N_2^g}{\partial y} & \dots & 0 & \frac{\partial N_{n_{\text{points}}}^g}{\partial y} & 0 \\ 0 & \frac{\partial N_1^g}{\partial y} & 0 & \frac{\partial N_2^g}{\partial y} & \dots & 0 & \frac{\partial N_{n_{\text{points}}}^g}{\partial y} \end{bmatrix}, \quad (2.47)$$

and, as for the spatial description, it must be satisfied for any $\boldsymbol{\eta} \in {}^h\mathcal{V}_{u,0}$, the incremental quasi-static discretized mechanical constitutive initial boundary value problem can thus be stated in the material description as

Problem 2.10 | Material incremental discretized mechanical initial BVP.

Given the set of internal variables α_n at t_n , the prescribed body and traction force fields $\mathbf{b}_{0,n+1}$ and $\mathbf{t}_{0,n+1}$, and the prescribed deformation gradient \mathbf{F}_{n+1} at t_{n+1} , find the kinematically admissible nodal displacement field $\mathbf{u}_{n+1} \in {}^h\mathcal{K}_{u,n+1}$ such that the body \mathcal{B} is in equilibrium as stated by the Virtual Work Principle

$$\mathbf{M}\ddot{\mathbf{u}}_{n+1} + \mathbf{f}^{\text{int}}(\mathbf{u}_{n+1}) - \mathbf{f}_{n+1}^{\text{ext}} = \mathbf{0}, \quad (2.48)$$

where \mathbf{f}^{int} e $\mathbf{f}_{n+1}^{\text{ext}}$ are the global vectors of internal and external forces defined as

$$\mathbf{f}^{\text{int}} \equiv \int_{h\Omega_0} \mathbf{G}^g T \hat{\mathbf{P}}(\mathbf{F}_{n+1}, \alpha_n) d\nu, \quad (2.49)$$

$$\mathbf{f}_{n+1}^{\text{ext}} \equiv \int_{h\Omega_0} \mathbf{N}^g T \mathbf{b}_{0,n+1} d\nu + \int_{\partial^h\Omega_{\text{traction},0}} \mathbf{N}^g T \mathbf{t}_{0,n+1} da. \quad (2.50)$$

and \mathbf{M} is the mass matrix defined as

$$\mathbf{M} = \int_{h\Omega_0} \rho_0 \mathbf{N}^g T \mathbf{N}^g d\nu. \quad (2.51)$$

The global vectors for the internal and external forces are usually obtained by assemblage of their elemental counterparts as

$$\mathbf{f}^{\text{int}} = \bigvee_{e=1}^{n_{\text{elem}}} \left(\mathbf{f}^{\text{int}} \right)^{(e)}, \quad (2.52)$$

$$\mathbf{f}^{\text{ext}} = \bigvee_{e=1}^{n_{\text{elem}}} \left(\mathbf{f}^{\text{ext}} \right)^{(e)}, \quad (2.53)$$

where the elemental vectors in the spatial description are defined as

$$\left(\mathbf{f}^{\text{int}} \right)^{(e)} \equiv \int_{h\Omega^{(e)}} \mathbf{B}^T \hat{\boldsymbol{\sigma}}(\mathbf{F}_{n+1}, \alpha_n) d\nu, \quad (2.54)$$

$$\left(\mathbf{f}_{n+1}^{\text{ext}} \right)^{(e)} \equiv \int_{h\Omega^{(e)}} \mathbf{N}^T \mathbf{b}_{n+1} d\nu + \int_{\partial^h\Omega_{\text{traction}}^{(e)}} \mathbf{N}^T \mathbf{t}_{n+1} da, \quad (2.55)$$

and in material description as

$$\left(\mathbf{f}^{\text{int}} \right)^{(e)} \equiv \int_{h\Omega_0^{(e)}} \mathbf{G}^T \hat{\mathbf{P}}(\mathbf{F}_{n+1}, \alpha_n) d\nu, \quad (2.56)$$

$$\left(\mathbf{f}_{n+1}^{\text{ext}} \right)^{(e)} \equiv \int_{h\Omega_0^{(e)}} \mathbf{N}^T \mathbf{b}_{0,n+1} d\nu + \int_{\partial^h\Omega_{0,\text{traction}}^{(e)}} \mathbf{N}^T \mathbf{t}_{0,n+1} da, \quad (2.57)$$

The matrices \mathbf{N} , \mathbf{B} , and \mathbf{G} are the elemental interpolation matrix, the symmetric elemental gradient operator, and the discrete elemental gradient operator.

In a similar manner, the global mass matrix is also usually obtained by assemblage of their elemental counterparts as

$$\mathbf{M} \equiv \bigvee_{e=1}^{n_{\text{elem}}} \mathbf{M}^{(e)}, \quad (2.58)$$

where the elemental mass matrices in the spatial description are defined as

$$\mathbf{M}^{(e)} = \int_{h\Omega^{(e)}} \rho \mathbf{N}^T \mathbf{N} d\nu, \quad (2.59)$$

and the material description

$$\mathbf{M}^{(e)} = \int_{h\Omega_0^{(e)}} \rho_0 \mathbf{N}^T \mathbf{N} d\nu. \quad (2.60)$$

2.2.5 Numerical integration

In the Finite Element Method, the integrations over the element domain are generally performed numerically using the Gaussian Quadrature Method. Stating it's application succinctly, let $a(\mathbf{x})$ be a generic field, if there is a coordinate transformation from a local (or natural) normalized domain Υ to the element domain $\Omega^{(e)}$, $\mathbf{x}: \Upsilon \rightarrow \Omega^{(e)}$, the integral of $a(\mathbf{x})$ over the domain $\Omega^{(e)}$ can be numerically determined as

$$\int_{\Omega^{(e)}} a(\mathbf{x}) d\mathbf{x} = \int_{\Upsilon} a(\mathbf{x}(\boldsymbol{\zeta})) j(\boldsymbol{\zeta}) d\boldsymbol{\zeta} \approx \sum_{i=1}^{n_{GP}} w_i a(\mathbf{x}(\boldsymbol{\zeta}_i)) j(\boldsymbol{\zeta}_i), \quad (2.61)$$

where $\boldsymbol{\zeta}_i$ and w_i , $i = 1, \dots, n_{GP}$ are the positions and weights of the Gauss sampling points in the domain Υ and $j(\boldsymbol{\zeta})$ is the determinant of the coordinate transformation's Jacobian defined as

$$j(\boldsymbol{\zeta}) = \det \left(\frac{\partial \mathbf{x}}{\partial \boldsymbol{\zeta}} \right). \quad (2.62)$$

2.3 Linearisation

The equilibrium equation, Equation (2.42) in a spatial description and Equation (3.25) in a material description, is generally nonlinear due to geometrical and/or material nonlinearities. The Newton-Raphson Method is an efficient and robust iterative scheme with a quadratic convergence rate often used to solve the equilibrium equation at each time increment, t_n . The residual of the fully discretized balance of linear momentum is defined for an iteration step i of the Newton-Raphson method as

$$\mathbf{r}(\mathbf{u}_{n+1}^i) = \mathbf{M} \ddot{\mathbf{u}}_{n+1}^i + \mathbf{f}^{\text{int}}(\mathbf{u}_{n+1}^i) - \mathbf{f}_{n+1}^{\text{ext}}. \quad (2.63)$$

A Taylor expansion about the current solution \mathbf{u}_{n+1}^i is performed, discarding all terms of higher order than one, yielding the linearised form

$$\text{Lin } \mathbf{r}(\mathbf{u}_{n+1}^i) = \mathbf{r}(\mathbf{u}_{n+1}^i) + \underbrace{\frac{\partial \mathbf{r}(\mathbf{u}_{n+1}^i)}{\partial \mathbf{u}_{n+1}}}_{\mathbf{K}(\mathbf{u}_{n+1}^i)} \delta \mathbf{u}. \quad (2.64)$$

with the dynamic effective tangential stiffness matrix $\mathbf{K}(\mathbf{u}_{n+1}^i)$. The linearisation of the internal forces included in \mathbf{K} is known as the tangential stiffness matrix \mathbf{K}_T , which is defined as

$$\mathbf{K}_T^i = \frac{\partial \mathbf{f}^{\text{int}}}{\partial \mathbf{u}_{n+1}} \bigg|_i. \quad (2.65)$$

Equilibrium is achieved if

$$\text{Lin} \mathbf{r}(\mathbf{u}_{n+1}^i) = \mathbf{0}, \quad (2.66)$$

so that a linear system of equation is given by

$$\mathbf{K}(\mathbf{u}_{n+1}^i) \delta \mathbf{u} = -\mathbf{r}(\mathbf{u}_{n+1}^i). \quad (2.67)$$

Thus, a new solution of the displacement increment $\delta \mathbf{u}$ for current iteration step $i + 1$ is determined, and the final displacement solution of time step $n + 1$ is obtained via updating

$$\mathbf{u}_{n+1}^{i+1} = \mathbf{u}_{n+1}^i + \delta \mathbf{u}. \quad (2.68)$$

A solution of t_{n+1} is found, i.e. an equilibrium state is reached and $\mathbf{u}_{n+1} = \mathbf{u}_{n+1}^{i+1}$, if prescribed, user-defined convergence criteria are fulfilled.

2.4 Constitutive laws

Chapter 3

Thermo field

For the development of the thermomechanical models, the temperature field needs to be considered. This section provides an overview of the governing equations required to describe a temperature field with the finite element method (FEM). The procedure to establish a fully discrete system of equations for the thermal field is comparable to the one for the structural field in chapter ???. Furthermore, the basics of nonlinear continuum thermodynamics have already been featured in chapter ??. Consequently, the detailed derivation are skipped in this chapter.

3.1 Governing equations

Based on the general model presented in Section ??, the balance equations for the temperature field are obtained as a special case by neglecting all mechanical terms. Hence, the energy balance equation (Equation (1.42)), now in material description, reduces to

$$\rho_0 \dot{e} = -\operatorname{div}_0 \mathbf{q}_0 + \rho_0 r \quad \text{in } \Omega_0, \quad (3.1)$$

where all mechanical terms are neglected. The target application of the present work are coupled generally nonlinear thermomechanical interaction problems, where the initial and the current domains are not equal, i.e. $\Omega_0 \neq \Omega$. Thus, for the sake of simplicity and in view of the later coupled problem, all following relations are expressed in material quantities. A purely thermal analysis is independent of the deformation, so that reference and current configuration are identical and the domain remains constant, i.e. $\Omega_0 \equiv \Omega$.

3.2 Thermal constitutive initial value problem

From Section ??, discarding all variables related to the mechanical problem, the general thermal constitutive initial value problem is

Problem 3.1 | General thermal constitutive initial value problem.

Given the initial value of the internal variables $\boldsymbol{\alpha}(t_0)$ and the history of the

temperature distribution

$$\theta(t), \quad t \in [t_0, t_{\text{end}}],$$

find the function for $\mathbf{q}_0(t)$, $s(t)$ and $\boldsymbol{\alpha}(t)$ such that the constitutive equations

$$s = -\frac{\partial \psi}{\partial \theta}, \quad (3.2)$$

$$\psi = \psi(\theta), \quad (3.3)$$

$$\dot{\boldsymbol{\alpha}} = f(\theta, \mathbf{g}_0, \boldsymbol{\alpha}), \quad (3.4)$$

$$\frac{1}{\theta} \mathbf{q}_0 = g(\theta, \mathbf{g}_0, \boldsymbol{\alpha}). \quad (3.5)$$

are satisfied for every $t \in [t_0, t_{\text{end}}]$.

No distinction between spatial and material configurations applies as $\Omega = \Omega_0$. Next, a standard set of assumptions are introduced.

Helmholtz free energy As a first step, the specific heat C_V is established and defined according to the thermodynamical principles to be the amount of heat required to change a unit mass of a substance by one degree in temperature, i.e.

$$C_V = \frac{\partial e}{\partial \theta}. \quad (3.6)$$

The index $(\cdot)_V$ denotes that C_V is measured at constant volume. Its dimensions are energy over temperature, i.e., $[E/\Theta]$, and using the International System of Units (SI), C_V is expressed in joule per kelvin. Using Equation 1.51, the specific heat at constant volume can be written as

$$C_V = -\frac{\partial^2 \psi}{\partial \theta^2} \theta = \frac{\partial s}{\partial \theta} \theta. \quad (3.7)$$

In general, the heat capacity depends on the deformation and on the temperature. A substance whose specific volume (or density) is constant is called an incompressible substance. This incompressibility or constant-volume assumption should be taken to imply that the energy associated with the volume change is negligible compared with other forms of energy. For the application to elastomers, see for instance Netz [96], the heat capacity C_V can be assumed to depend only on the temperature, and the partial derivatives in Equation (3.7) turn into exact derivatives, yielding

$$C_V = \frac{ds}{d\theta} \theta. \quad (3.8)$$

Furthermore, for the application to metals, a constant specific heat capacity (i.e. $C_V = \text{const.}$) is a valid assumption, utilised e.g. in Adam and Ponthot [1], Ghadiani [48], Ibrahimbegovic and Chorfi [61], and Simo and Miehe [122]. Accordingly, the heat capacity is also assumed to be constant (i.e. $C_V = \text{const.}$), since focus in this work is on the application to metals. Thus, the entropy can be written as

$$s(\theta) = C_V \ln \left(\frac{\theta}{\theta_0} \right), \quad (3.9)$$

after integration, where θ_0 and C_V denote the constant initial temperature and the constant specific heat, respectively.

Given the constitutive relation for the entropy (Equation (3.2)), the Helmholtz free energy per unit reference volume is found to be

$$\psi(\theta) = -C_V \left[(\theta - \theta_0) - \theta \ln \left(\frac{\theta}{\theta_0} \right) \right], \quad (3.10)$$

Subsequently, the time derivative of the entropy is

$$\dot{s}(\theta) = \frac{\partial s}{\partial \theta} \dot{\theta} = -\frac{\partial^2 \psi}{\partial \theta^2} \dot{\theta} = C_V \frac{1}{\theta} \dot{\theta}. \quad (3.11)$$

Law for the heat flux As previously mentioned, in a purely thermal analysis the deformation is neglected, consequently the material and spatial heat flux coincide, that is $\mathbf{q}_0 \equiv \mathbf{q}$, which is also valid for the material and spatial gradient, hence $\nabla_0 \theta = \nabla \theta$. To satisfy the dissipation inequality due to conduction (Equation (3.1)), a constitutive law for the heat flux has to be chosen associating the heat flux \mathbf{q}_0 with its dual variable \mathbf{g}_0 and the temperature θ . Accordingly, so-called Fourier's law, which is linear and isotropic is utilised, which is defined as

$$\mathbf{q}_0 = -k \mathbf{g}_0. \quad (3.12)$$

Herein, the thermal conductivity k is assumed constant and positive that is $k \geq 0$. Thus, heat is conducted in the direction of decreasing temperatures. Apart from Fourier's law, different constitutive laws for the heat flux are available in the literature, as e.g. Duhamel's law of heat conduction (see e.g. Holzapfel [58]) which uses a positive semi-definite second-order tensor \mathbf{k} instead of the constant conductivity k . If Duhamel's law is restricted to thermally isotropic behaviour (i.e. no preferred direction), the conductivity tensor reduces to $\mathbf{k} = k \mathbf{I}$. If a constant heat conductivity $k = \text{const.}$ is assumed, Fourier's law is recovered as a special form of Duhamel's law. Moreover, e.g. in Holzapfel and Simo [59] and Sherief and Abd El-Latief [117], a variable conductivity ($k \neq \text{const.}$ is assumed in the context of elastomers). In Bargmann and Steinmann [13] and Bargmann et al. [14], three different constitutive laws for the heat flux \mathbf{q} are proposed based on the Green-Naghdi's non-classical theory. Nevertheless, for the present work Fourier's law yields physical results and hence is exclusively considered in this work.

"Standard" thermal constitutive description No extra internal variables $\boldsymbol{\alpha}$ are considered in the present description of the thermal problem. Thus, the thermal constitutive initial value problem given the standard assumptions laid out above accepts a closed form solution, i.e., the functions for \mathbf{q} and s are known from the outset.

Problem 3.2 | "Standard" thermal constitutive description

Given the history of the temperature distribution

$$\theta(t), \quad t \in [t_0, t_{\text{end}}],$$

compute the functions for $\mathbf{q}_0(t)$ and $s(t)$ at every $t \in [t_0, t_{\text{end}}]$ using the constitutive equations

$$s = C_V \ln \left(\frac{\theta}{\theta_0} \right), \quad (3.13)$$

$$\mathbf{q}_0 = -k \mathbf{g}_0. \quad (3.14)$$

3.3 Weak energy balance equation

The solution of the thermal problem using the FEM requires the use of the weak form of the energy balance equation. Applying to the governing equation (Equation (3.1)) in the strong form, a variational approach, multiplying it by the virtual temperatures ξ followed by integration by parts, one can find the energy balance equation in its weak form.

Problem 3.3 | Weak energy balance equation

There is energy balance in the body if and only if the temperature distribution satisfies

$$\int_{\Omega_0} \left[(\dot{\epsilon} - \rho_0 r) \xi - \mathbf{q}_0 \cdot \nabla_0 \xi \right] dv - \int_{\partial\Omega_0} h_0 \xi da = 0, \quad \forall \xi \in \mathcal{V}_{\theta,0}, \quad (3.15)$$

where $\mathcal{V}_{\theta,0}$ is the space of virtual temperature distributions on the body, defined by the space of sufficiently regular arbitrary temperature distributions.

3.4 The thermal initial boundary value problem

Following the same approach as in Section ??, it is now possible to introduce the thermal initial boundary value problem. Assume that the internal variables governing the body \mathcal{B} are known at the initial time t_0 . In addition, assume that the heat generated in the interior of the body is prescribed, $r(\mathbf{X}, t)$, $t \in [t_0, t_{\text{end}}]$, as well as,

- **Natural (or Neumann) boundary condition.** The boundary portion $\partial\Omega_{\text{heat},0}$ of \mathcal{B} is subject to a prescribed history of heat flux, $h_{\text{presc},0}(\mathbf{X}, t) = \mathbf{q}_{\text{presc},0}(\mathbf{X}, t) \cdot \mathbf{m}(\mathbf{X})$, $\mathbf{X} \in \partial\Omega_{\text{heat},0}$, $t \in [t_0, t_{\text{end}}]$.
- **Essential (or Dirichlet) boundary condition.** The boundary portion $\partial\Omega_{\text{temperature},0}$ of \mathcal{B} is subject to a prescribed temperature history, $\theta_{\text{presc}}(\mathbf{X}, t)$, $\mathbf{X} \in \partial\Omega_{\text{temperature},0}$, $t \in [t_0, t_{\text{end}}]$.

As before the admissible temperature distributions for the body \mathcal{B} are all sufficiently regular temperature fields that satisfy the essential boundary condition,

$$\mathcal{K}_\theta = \{ \theta : \Omega_0 \times \mathbb{R} \rightarrow \mathbb{R} \mid \theta(\mathbf{X}, t) = \theta_{\text{presc}}(\mathbf{X}, t), \quad \mathbf{X} \in \partial\Omega_{\text{temperature},0}, \quad t \in [t_0, t_{\text{end}}] \}. \quad (3.16)$$

Combining the weak energy balance equations with the "standard" thermal constitutive description, the weak form of the "standard" thermal constitutive initial boundary value problem can be stated as follows

Problem 3.4 | "Standard" thermal initial BVP.

Find an admissible temperature distribution, $\theta \in \mathcal{K}_\theta$, such that for every $t \in [t_0, t_{\text{end}}]$, the body \mathcal{B} is in energetic equilibrium

$$\int_{\Omega_0} \left[(C_V \dot{\theta} - \rho_0 r) \xi + k \mathbf{g}_0 \cdot \nabla_0 \xi \right] dv - \int_{\partial\Omega_0} h_0 \xi da = 0, \quad \forall \xi \in \mathcal{V}_{\theta,0}, \quad (3.17)$$

where the space of virtual temperature distributions at time t is defined by

$$\mathcal{V}_{\theta,0} \equiv \left\{ \xi : \Omega_0 \rightarrow \mathbb{R} \mid \xi = 0 \text{ in } \partial\Omega_{\text{temperature},0} \right\}. \quad (3.18)$$

3.5 Finite Element Method

Following a procedure entirely similar to the one described in Section ??, the global shape functions can be conveniently assembled in the so-called global interpolation matrix as

$$\mathbf{N}^g(\mathbf{X}) \equiv \left[N_1^g(\mathbf{X}), N_2^g(\mathbf{X}), \dots, N_{n_{\text{points}}}^g(\mathbf{X}) \right]. \quad (3.19)$$

The vector containing the nodal values of the temperature is denoted by $\boldsymbol{\theta}$ and defined as

$$\boldsymbol{\theta}(t) = \left[\theta^1(t), \dots, \theta^{n_{\text{points}}}(t) \right]^T, \quad (3.20)$$

such that the value of the temperature inside the discretized domain ${}^h\Omega_0$ can be found from

$${}^h\theta(\mathbf{X}, t) \equiv \mathbf{N}^g(\mathbf{X}) \boldsymbol{\theta}(t), \quad {}^h\theta \in {}^h\mathcal{K}_\theta. \quad (3.21)$$

It is also convenient to define the discrete global gradient operator \mathbf{H}^g . For instance, in a 2D problem, where cartesian coordinates are employed, this discrete operator is defined as

$$\mathbf{H}^g \equiv \begin{bmatrix} \frac{\partial N_1^g}{\partial X} & \frac{\partial N_2^g}{\partial X} & \cdots & \frac{\partial N_{n_{\text{points}}}^g}{\partial X} \\ \frac{\partial N_1^g}{\partial Y} & \frac{\partial N_2^g}{\partial Y} & \cdots & \frac{\partial N_{n_{\text{points}}}^g}{\partial Y} \end{bmatrix}. \quad (3.22)$$

Applying the aforementioned finite element discretization to the "standard" thermal initial BVP yields

$$\int_{{}^h\Omega_0} \left[(C_V \dot{\theta} - \rho_0 r) \mathbf{N}^g \xi + k \mathbf{g}_0 \cdot \mathbf{H}^g \xi \right] dv - \int_{{}^h\partial\Omega_0} h_0 \mathbf{N}^g \xi da = 0, \quad \forall \xi \in {}^h\mathcal{V}_{\theta,0}, \quad (3.23)$$

which can be rewritten

$$\left\{ \int_{h\Omega_0} \left[(\mathbf{N}^g)^T (C_V \dot{\theta} - \rho_0 r) + k(\mathbf{H}^g)^T \mathbf{H}^g \theta \right] d\nu - \int_{h\partial\Omega_0} (\mathbf{N}^g)^T h_0 da \right\}^T \xi = 0, \quad \forall \xi \in {}^h\mathcal{V}_{\theta,0}, \quad (3.24)$$

where the relation $\mathbf{g}_0 = \mathbf{H}^g \theta$ is employed. Since Equation (3.24) must be satisfied for any $\xi \in {}^h\mathcal{V}_{\theta,0}$, the discretized "standard" thermal initial boundary value problem can be stated as

Problem 3.5 | Discretized "standard" thermal initial BVP.

Given the prescribed heat sources and heat fluxes $r(\mathbf{X}, t)$ and $h_0(\mathbf{X}, t)$ find the admissible nodal temperatures $\theta(t) \in {}^h\mathcal{V}_\theta$ such that the body \mathcal{B} is in energetic equilibrium

$$\mathbf{C} \dot{\theta}(t) + \mathbf{K} \theta(t) - \mathbf{f}^{\text{ext}}(t) = \mathbf{0}, \quad (3.25)$$

where \mathbf{C} and \mathbf{K} are the temperature damping and stiffness matrix defined as

$$\mathbf{C} = \int_{h\Omega_0} C_V \mathbf{N}^g{}^T \mathbf{N}^g d\nu, \quad (3.26)$$

$$\mathbf{K} = \int_{h\Omega_0} k \mathbf{H}^g{}^T \mathbf{H}^g d\nu. \quad (3.27)$$

and $\mathbf{f}^{\text{ext}}(t)$ is the global vector of external forces defined as

$$\mathbf{f}^{\text{ext}}(t) \equiv \int_{h\Omega_0} \rho \mathbf{N}^g{}^T r(\mathbf{X}, t) d\nu + \int_{\partial^h\Omega_{\text{heat},0}} \mathbf{N}^g{}^T h_0(\mathbf{X}, t) da. \quad (3.28)$$

Chapter 4

Thermo-mechanical problem

In the following chapter, the general framework presented in chapter ?? is applied to the complete thermo-mechanical analysis. Chapters ?? and ?? concerning the standalone mechanical and thermal problems provide a guide for the developments that follow. This coupling is a so-called volume coupling, as in each point of the domain, the two fields are coupled. In contrast, surface-coupled problems, e.g., fluid-structure interactions problems, include problems where there is coupling only at the interface between the fluid and the structural domain.

The standard approach to finite deformation thermomechanical analysis is the thermomechanical potential . It is already included in the description of continuum thermodynamics in Chapter ??, where the Helmholtz free-energy is allowed to depend on the temperature θ . Thus, the temperature is directly included in the constitutive law. Following the setup of finite strain plasticity in the general non-isothermal case, in addition to the plastic intermediate configuration, a thermal intermediate configuration has to be considered , i.e.,

$$\mathbf{F} = \mathbf{F}^t \mathbf{F}^e \mathbf{F}^p. \quad (4.1)$$

In contrast to thermal deformations, which are solely volumetric, the volume remains constant during plastic deformations for most metals. Hence, plastic deformations are assumed isochoric, $J^p = \det \mathbf{F}^p = 1$, yielding the following split for the jacobian of the deformation

$$J = J^t J^e. \quad (4.2)$$

Following [Danowski \(2014\)](#), the thermal expansion is assumed to be

$$J^t = \frac{dv}{dv_0} = \exp(3\alpha_T \Delta\theta). \quad (4.3)$$

Combining the last two equations, the elastic volumetric deformation J^e , can be expressed as

$$J^e = J^e(\theta) = J/J^t. \quad (4.4)$$

Consequently, the additional thermal intermediate configuration can be omitted and volumetric deformations are described only J^e . Thus, if thermal stresses arise due to a temperature change, elastic strains balance the body which implicitly correspond to the thermal strains according to Equation (4.4).

4.0.1 Thermo-mechanical constitutive initial value problem

In the full thermo-mechanical case, considering all quantities related to both the mechanical and the thermal domain, a constitutive model based on internal variables is established by the following set of equations

$$\mathbf{P} = \rho_0 \frac{\partial \psi}{\partial \mathbf{F}}, \quad (4.5)$$

$$s = -\frac{\partial \psi}{\partial \theta}, \quad (4.6)$$

$$\psi = \psi(\mathbf{F}, \theta, \mathbf{g}_0, \boldsymbol{\alpha}), \quad (4.7)$$

$$\dot{\boldsymbol{\alpha}} = \mathbf{f}(\mathbf{F}, \theta, \mathbf{g}_0, \boldsymbol{\alpha}), \quad (4.8)$$

$$\frac{1}{\theta} \mathbf{q}_0 = \mathbf{g}(\mathbf{F}, \theta, \boldsymbol{\alpha}). \quad (4.9)$$

The Helmholtz free energy ψ in Equation 4.7 is expressed with respect to the reference volume, so that ψ is reformulated using potential functions, to emphasize this additive decomposition, as follows

$$\rho_0 \psi(\mathbf{F}, \theta, \boldsymbol{\alpha}) := \hat{\mathbb{U}}(J^e, \theta) + \hat{\mathbb{W}}(\mathbf{F}_{\text{iso}}, \theta) + \hat{\mathbb{M}}(J^e, \theta) + \hat{\mathbb{T}}(\theta) + \hat{\mathbb{K}}(\boldsymbol{\alpha}, \theta), \quad (4.10)$$

where in contrast to the deformation gradient \mathbf{F} , the Jacobi-determinant J^e and the isochoric deformation gradient \mathbf{F}_{iso} are applied. $\hat{\mathbb{U}}$ and $\hat{\mathbb{W}}$ can be identified with the standard hyperelastic materials potentials, with the caveat that now the material parameters can depend on the temperature, whereas $\hat{\mathbb{M}}$ describes the thermomechanical coupling potential. The potential $\hat{\mathbb{T}}$ represents the purely thermal potential and is assumed identical to Equation (??). Finally, $\hat{\mathbb{K}}$ is the convex plastic potential.

Based on the potential functions, the coupling of the two fields, mechanical and thermal, can be explained:

- the temperature enters the structural field via additional thermal stresses and possibly via temperature-dependent material parameters. Herein, $\hat{\mathbb{M}}$ characterizes the thermomechanical coupling potential, leading to thermal stresses and thermal expansion and dilatation, whereas $\hat{\mathbb{K}}$ being temperature-dependent and therefore enables exemplarily von Mises plasticity combined with temperature-dependent isotropic hardening and thermal softening.
- the structure enters the thermal field via coupling terms, arising from $\hat{\mathbb{M}}$ and $\hat{\mathbb{K}}$, in addition to the purely thermal energy $\hat{\mathbb{T}}$, and thus, coupling terms as the internal or mechanical dissipation $\mathcal{D}_{\text{mech}}$ may emerge in the thermal balance equation. Furthermore, for finite deformation thermomechanical problems, where the initial domain Ω_0 deforms to Ω , so that $\Omega \neq \Omega_0$, and a Lagrangian formulation is used, the deformation enters the thermal field additionally due to the mapping of all quantities in the balance equations to the reference configuration.

The constitutive relations found for the first Piola-Kirchhoff stress tensor (Equation (??)) and the entropy (Equation (??)) are the same. The expression for the entropy in full can be written as

$$s = C_V \ln \left(\frac{\theta}{\theta_0} \right) - \frac{1}{\rho_0} \left(\frac{\partial \hat{\mathcal{U}}(J^e, \theta)}{\partial \theta} + \frac{\partial \hat{\mathcal{W}}(\mathbf{F}_{\text{iso}}, \theta)}{\partial \theta} + \frac{\partial \hat{\mathcal{M}}(J^e, \theta)}{\partial \theta} + \frac{\partial \hat{\mathcal{K}}(\boldsymbol{\alpha}, \theta)}{\partial \theta} \right). \quad (4.11)$$

As in Chapter ??, the heat conduction law is chosen to be the Fourier heat conduction law, repeated here for the sake of clarity in the spatial description

$$\mathbf{q} = -k \mathbf{g}. \quad (4.12)$$

Applying the Piola transformation yields for the material description of the heat flux vector,

$$\mathbf{q}_0 = -k_0 \mathbf{C}^{-1} \mathbf{g}_0. \quad (4.13)$$

Thus, the material thermo-mechanical constitutive initial value problem can be stated as

Problem 4.1 | Material thermomechanical constitutive initial value problem.

Given the initial values of the internal variables, $\boldsymbol{\alpha}(t_0)$, the history of the deformation gradient

$$\mathbf{F}(t), \quad t \in [t_0, t_{\text{end}}], \quad (4.14)$$

and the history of the temperature distribution

$$\theta(t), \quad t \in [t_0, t_{\text{end}}], \quad (4.15)$$

find the functions for $\mathbf{P}(t)$, $s(t)$, $\mathbf{q}_0(t)$ and $\boldsymbol{\alpha}(t)$ such that the constitutive equations

$$\mathbf{P} = \rho_0 \frac{\partial \psi}{\partial \mathbf{F}}, \quad (4.16)$$

$$s = C_V \ln \left(\frac{\theta}{\theta_0} \right) - \frac{1}{\rho_0} \left(\frac{\partial \hat{\mathcal{U}}(J^e, \theta)}{\partial \theta} + \frac{\partial \hat{\mathcal{W}}(\mathbf{F}_{\text{iso}}, \theta)}{\partial \theta} + \frac{\partial \hat{\mathcal{M}}(J^e, \theta)}{\partial \theta} + \frac{\partial \hat{\mathcal{K}}(\boldsymbol{\alpha}, \theta)}{\partial \theta} \right), \quad (4.17)$$

$$\psi = \frac{1}{\rho_0} \left(\hat{\mathcal{U}}(J^e) + \hat{\mathcal{W}}(\mathbf{F}_{\text{iso}}) + \hat{\mathcal{M}}(J^e, \theta) + \hat{\mathcal{T}}(\theta) + \hat{\mathcal{K}}(\boldsymbol{\alpha}, \theta) \right), \quad (4.18)$$

$$\mathbf{q}_0 = -k_0 \mathbf{C}^{-1} \mathbf{g}_0, \quad (4.19)$$

$$\dot{\boldsymbol{\alpha}} = f(\mathbf{F}, \theta, \mathbf{g}_0, \boldsymbol{\alpha}), \quad (4.20)$$

are satisfied for every $t \in [t_0, t_{\text{end}}]$, where

$$\hat{\mathcal{T}}(\theta) = -C_V \left[(\theta - \theta_0) - \theta \ln \left(\frac{\theta}{\theta_0} \right) \right]. \quad (4.21)$$

4.0.2 Weak equilibrium. The principle of virtual work

The strong equations that enforce the equilibrium of a body can be written using the material description as

$$\rho_0 \ddot{\mathbf{u}} = \operatorname{div}_0 \mathbf{P} + \mathbf{b}_0, \quad \text{in } \Omega_0, \quad (4.22)$$

$$\rho_0 \dot{e} = \mathbf{P} : \dot{\mathbf{F}} + \rho_0 r - \operatorname{div}_0 \mathbf{q}_0, \quad \text{in } \Omega_0. \quad (4.23)$$

Following an approach entirely similar to the ones presented in Chapters ?? and ??, the weak form of the linear momentum and energy balance equations can be found to be

Problem 4.2 | Weak form of the linear momentum and energy balance equations (material version).

In a material description, the body is in mechanical and energetic equilibrium if and only if the First Piola-Kirchhoff stress field, $\mathbf{P}(t)$, the heat flux vector $\mathbf{q}_0(t)$, satisfy

$$\int_{\Omega_0} [\mathbf{P}(t) : \nabla_0 \boldsymbol{\eta} - (\mathbf{b}_0(t) - \rho_0 \ddot{\mathbf{u}}(t)) \cdot \boldsymbol{\eta}] \, d\nu - \int_{\partial\Omega_0} \mathbf{t}_0(t) \cdot \boldsymbol{\eta} \, da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_{u,0}, \quad (4.24)$$

$$\int_{\Omega_0} \left[(\rho_0 \dot{e}(t) - \mathbf{P}(t) : \dot{\mathbf{F}}(t) - \rho_0 r(t)) \xi - \mathbf{q}_0(t) \cdot \nabla_0 \xi \right] \, d\nu - \int_{\partial\Omega_0} h_0(t) \xi \, da = 0, \quad \forall \xi \in \mathcal{V}_{\theta,0}, \quad (4.25)$$

where $\mathcal{V}_{u,0}$ is the space of virtual displacement of the body, defined by the space of sufficiently regular arbitrary displacements

$$\boldsymbol{\eta} : \Omega_0 \rightarrow \mathcal{U}. \quad (4.26)$$

and $\mathcal{V}_{\theta,0}$ is the space of virtual temperature distributions of the body, defined by the space of sufficiently regular arbitrary temperature distributions

$$\xi : \Omega_0 \rightarrow \mathcal{R}. \quad (4.27)$$

4.0.3 Mechanical constitutive initial boundary value problem

It is now possible to pose the thermo-mechanical constitutive initial value problem in its weak form. Assume that a body \mathcal{B} is made from a generic material, characterized by a given constitutive model, whose internal variables are known at the initial time, as presented in Figure ??. In addition, it is assumed that the interior of the body was subjected to a prescribed history of body forces, $\mathbf{b}(\mathbf{X}, t)$, and a prescribed history of heat sources, $r(\mathbf{X}, t)$, $t \in [t_0, t_{\text{end}}]$, and to the following boundary conditions:

- **Natural (or Neumann) boundary condition:**

- **Mechanical field:** The boundary portion $\Omega_{\text{traction},0}$ of \mathcal{B} is subjected to a prescribed history of traction forces, $\mathbf{t}_{\text{presc},0}(\mathbf{X}, t)$, $\mathbf{X} \in \partial\Omega_{\text{traction},0}$, $t \in [t_0, t_{\text{end}}]$,
- **Thermal field:** The boundary portion $\partial\Omega_{\text{heat},0}$ of \mathcal{B} is subject to a prescribed history of heat flux, $h_{\text{presc},0}(\mathbf{X}, t) = \mathbf{q}_{\text{presc},0}(\mathbf{X}, t) \cdot \mathbf{m}(\mathbf{X})$, $\mathbf{X} \in \partial\Omega_{\text{heat},0}$, $t \in [t_0, t_{\text{end}}]$.

• **Essential (or Dirichlet) boundary condition:**

- **Mechanical field:** The boundary portion $\Omega_{\text{motion},0}$ of \mathcal{B} is subjected to a prescribed displacement field history, $\mathbf{u}_{\text{presc}}(\mathbf{X}, t)$, such that

$$\boldsymbol{\varphi}(\mathbf{X}, t) = \mathbf{X} + \mathbf{u}_{\text{presc}}(\mathbf{X}, t), \quad \mathbf{X} \in \partial\Omega_{\text{motion},0}, \quad t \in [t_0, t_{\text{end}}].$$

- **Thermal field:** The boundary portion $\partial\Omega_{\text{temperature},0}$ of \mathcal{B} is subject to a prescribed temperature history, $\theta_{\text{presc}}(\mathbf{X}, t)$, $\mathbf{X} \in \partial\Omega_{\text{temperature},0}$, $t \in [t_0, t_{\text{end}}]$.

It is also convenient to define the set of kinematically admissible displacements and admissible temperature distributions of \mathcal{B} as the set of all sufficiently regular displacement and temperature functions that satisfy the correspondin essential boundary conditions (??),

$$\begin{aligned} \mathcal{K}_u \equiv \{ \mathbf{u} : \Omega_0 \times \mathcal{R} \rightarrow \mathcal{U} \mid \mathbf{u}(\mathbf{X}, t) = \mathbf{u}_{\text{presc}}(\mathbf{X}, t), \\ \mathbf{X} \in \partial\Omega_{\text{motion},0}, \quad t \in [t_0, t_{\text{end}}] \}. \end{aligned} \quad (4.28)$$

$$\begin{aligned} \mathcal{K}_\theta \equiv \{ \theta : \Omega_0 \times \mathcal{R} \rightarrow \mathcal{R} \mid \theta(\mathbf{X}, t) = \theta_{\text{presc}}(\mathbf{X}, t), \\ \mathbf{X} \in \partial\Omega_{\text{temperature},0}, \quad t \in [t_0, t_{\text{end}}] \}. \end{aligned} \quad (4.29)$$

To obtain the thermo-mechanical initial boundary value problem, express the internal specific energy of the system using its specific Helmholtz free energy, yielding

$$\dot{e} \equiv \dot{\psi} + s\dot{\theta} + \dot{s}\theta. \quad (4.30)$$

Considering Equation (??) for the entropy, its time derivative is

$$\dot{s} = C_V \frac{1}{\theta} \dot{\theta} - \frac{1}{\rho_0 \theta} \mathcal{H}^{\text{ep}}, \quad (4.31)$$

where the definition of $\hat{\mathbb{T}}$ (Equation (??)) and the Gough-Joule effect, \mathcal{H}^{ep} is defined to be

$$\begin{aligned} \mathcal{H}^{\text{ep}} = \theta \left(\frac{\partial^2 \hat{\mathbb{U}}(J^e, \theta)}{\partial \theta \partial J^e} j^e + \frac{\partial^2 \hat{\mathbb{U}}(J^e, \theta)}{\partial \theta^2} \dot{\theta} + \frac{\partial^2 \hat{\mathbb{W}}(\mathbf{F}_{\text{iso}}, \theta)}{\partial \theta^2} \dot{\theta} + \frac{\partial^2 \hat{\mathbb{W}}(\mathbf{F}_{\text{iso}}, \theta)}{\partial \theta \partial \mathbf{F}_{\text{iso}}} : \dot{\mathbf{F}}_{\text{iso}} \right. \\ \left. + \frac{\partial^2 \hat{\mathbb{M}}(J^e, \theta)}{\partial \theta \partial J^e} j^e + \frac{\partial^2 \hat{\mathbb{M}}(J^e, \theta)}{\partial \theta^2} \dot{\theta} + \frac{\partial^2 \hat{\mathbb{K}}(\boldsymbol{\alpha}_k, \theta)}{\partial \theta^2} \dot{\theta} + \frac{\partial^2 \hat{\mathbb{K}}(\boldsymbol{\alpha}_k, \theta)}{\partial \theta \partial \boldsymbol{\alpha}_k} * \dot{\boldsymbol{\alpha}}_k \right). \end{aligned} \quad (4.32)$$

Applying the chain-rule, one finds for $\dot{\psi}$

$$\dot{\psi} = \frac{\partial \psi}{\partial \mathbf{F}} : \dot{\mathbf{F}} + \frac{\partial \psi}{\partial \theta} \dot{\theta} + \frac{\partial \psi}{\partial \boldsymbol{\alpha}} * \dot{\boldsymbol{\alpha}}, \quad (4.33)$$

which can be rewritten considering the constitutive relations for the first Piola-Kirchhoff stress tensor (Equation (??)) and the entropy (Equation (??)), and the definition of the thermodynamical forces (Equation (??)) as

$$\dot{\psi} = \frac{1}{\rho_0} \mathbf{P} : \dot{\mathbf{F}} - s\dot{\theta} - \mathbf{A} * \dot{\boldsymbol{\alpha}}. \quad (4.34)$$

Combining Equation (4.34) and the expression found for the mechanical dissipation $\mathcal{D}_{\text{mech}}$ in Equation (??) yields

$$\mathcal{D}_{\text{mech}} = \rho_0 \mathbf{A} * \boldsymbol{\alpha}, \quad (4.35)$$

and thus $\dot{\psi}$ can be expressed as

$$\dot{\psi} = \frac{1}{\rho_0} \mathbf{P} : \dot{\mathbf{F}} - s \dot{\theta} - \frac{1}{\rho_0} \mathcal{D}_{\text{mech}}. \quad (4.36)$$

Combining Equations (??), Equation (4.31) and Equation (4.36) yields

$$\dot{e} = \frac{1}{\rho_0} \mathbf{P} : \dot{\mathbf{F}} - \frac{1}{\rho_0} \mathcal{D}_{\text{mech}} - C_V \dot{\theta} - \frac{1}{\rho_0} \mathcal{H}^{\text{ep}}. \quad (4.37)$$

So the weak form of the thermo-mechanical constitutive initial boundary value problem can be stated in a spatial description as follows and in the material description as

Problem 4.3 | Material thermo-mechanical initial BVP.

Find a kinematically admissible displacement function, $\mathbf{u} \in \mathcal{K}_u$, and an admissible temperature distribution, $\theta \in \mathcal{K}_\theta$, such that for every $t \in [t_0, t_{\text{end}}]$, the body \mathcal{B} is in mechanical and energetic equilibrium

$$\int_{\Omega_0} [\mathbf{P}(t) : \nabla_0 \boldsymbol{\eta} - (\mathbf{b}_0(t) - \rho_0 \ddot{\mathbf{u}}(t)) \cdot \boldsymbol{\eta}] \, dv - \int_{\partial\Omega_0} \mathbf{t}_0(t) \cdot \boldsymbol{\eta} \, da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_{u,0}, \quad (4.38)$$

$$\begin{aligned} \int_{\Omega_0} \left[\left(\rho_0 C_V \dot{\theta}(t) - \mathcal{D}_{\text{mech}}(t) - \mathcal{H}^{\text{ep}}(t) - \rho_0 r(t) \right) \xi - \mathbf{q}_0(t) \cdot \nabla_0 \xi \right] \, dv \\ - \int_{\partial\Omega_0} h_0(t) \xi \, da = 0, \quad \forall \xi \in \mathcal{V}_{\theta,0}, \end{aligned} \quad (4.39)$$

where the space of virtual displacements at time t is defined by

$$\mathcal{V}_{u,0} \equiv \{ \boldsymbol{\eta} : \Omega_0 \rightarrow \mathcal{U} \mid \boldsymbol{\eta} = \mathbf{0} \text{ in } \partial\Omega_{\text{motion},0} \}, \quad (4.40)$$

the space of virtual temperatures at time t is defined by

$$\mathcal{V}_{\theta,0} \equiv \{ \xi : \Omega_0 \rightarrow \mathcal{R} \mid \xi = 0 \text{ in } \partial\Omega_{\text{temperature},0} \}, \quad (4.41)$$

and at each point of \mathcal{B} , the First Piola-Kirchhoff stress tensor is the solution of material mechanical constitutive initial value problem.

4.1 Time discretization

In the context of thermo-mechanical problems, a generic path-dependent model is a model in which the thermodynamical state does not depend only on the instantaneous deformation and temperature states but also on their history. Under these conditions, the solution of the constitutive initial value problem for a given set of initial conditions

is usually not known for complex strain, $\mathbf{F}(t)$, or temperature paths, $\theta(t)$, paths. Thus, as in the domain of a strictly mechanical problem, there is a need to use an appropriate numerical algorithm to integrate the rate constitutive equations.

Attending to the thermo-mechanical constitutive initial boundary value problem and considering the time increment $[t_n, t_{n+1}]$, this approach is comprised by the following requirements:

- **First Piola-Kirchhoff stress tensors.** Considering a time increment $[t_n, t_{n+1}]$ and given the set $\boldsymbol{\alpha}_n$ of internal variables at t_n , the deformation gradient \mathbf{F}_{n+1} and the temperature distribution θ_{n+1} at time t_{n+1} determines the First Piola-Kirchhoff stress tensor \mathbf{P}_{n+1} uniquely as

$$\mathbf{P}_{n+1} = \hat{\mathbf{P}}(\mathbf{F}_{n+1}, \theta_{n+1}, \boldsymbol{\alpha}_n), \quad (4.42)$$

where $\hat{\mathbf{P}}$ is the incremental constitutive function for the First Piola-Kirchhoff stress tensor.

- **Set of internal variables.** Assuming that the set of internal variables $\boldsymbol{\alpha}_n$ is known at t_n , the set of internal variables must be uniquely determined by the prescribed deformation gradient \mathbf{F}_{n+1} and temperature distribution θ_{n+1} prescribed at t_{n+1} as

$$\boldsymbol{\alpha}_{n+1} = \hat{\boldsymbol{\alpha}}(\mathbf{F}_{n+1}, \theta_{n+1}, \boldsymbol{\alpha}_n), \quad (4.43)$$

where $\hat{\boldsymbol{\alpha}}$ is the incremental constitutive function for the set of internal variables.

- **Mechanical dissipation.** Considering a time increment $[t_n, t_{n+1}]$ and given the set $\boldsymbol{\alpha}_n$ of internal variables at t_n , the deformation gradient \mathbf{F}_{n+1} and the temperature distribution θ_{n+1} at time t_{n+1} determines the mechanical dissipation $\mathcal{D}_{\text{mech}}$ as

$$\mathcal{D}_{\text{mech},n+1} = \hat{\mathcal{D}}_{\text{mech}}(\mathbf{F}_{n+1}, \theta_{n+1}, \boldsymbol{\alpha}_n). \quad (4.44)$$

- **Gough-Joule effect.** Considering a time increment $[t_n, t_{n+1}]$ and given the set $\boldsymbol{\alpha}_n$ of internal variables at t_n , the deformation gradient \mathbf{F}_{n+1} and the temperature distribution θ_{n+1} at time t_{n+1} determines the Gough-Joule effect \mathcal{H}^{ep} as

$$\mathcal{H}_{n+1}^{\text{ep}} = \hat{\mathcal{H}}^{\text{ep}}(\mathbf{F}_{n+1}, \theta_{n+1}, \boldsymbol{\alpha}_n). \quad (4.45)$$

Generally, the numerical constitutive laws are nonlinear and path-independent within one increment.

Making use of the aforementioned time discretization, one can state the weak form of the mechanical constitutive initial boundary value problem in the material description as

Problem 4.4 | Material incremental thermo-mechanical initial BVP.

Given the set of internal variables α_n at t_n , the prescribed body and traction force fields $\mathbf{b}_{0,n+1}$ and $\mathbf{t}_{0,n+1}$, as well as, the prescribed heat sources, \mathbf{r}_{n+1} and heat fluxes, $h_{0,n+1}$, at t_{n+1} , and the prescribed deformation gradient \mathbf{F}_{n+1} and temperature distribution θ_{n+1} at t_{n+1} , find the kinematically admissible displacement field $\mathbf{u}_{n+1} \in \mathcal{K}_{u,n+1}$ and the admissible temperature distribution $\theta_{n+1} \in \mathcal{K}_{\theta,n+1}$ such that the body \mathcal{B} is in mechanical and energetic equilibrium

$$\int_{\Omega_0} [\hat{\mathbf{P}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1}, \alpha_n) : \nabla_0 \boldsymbol{\eta} - (\mathbf{b}_{0,n+1} - \rho_0 \ddot{\mathbf{u}}_{n+1}) \cdot \boldsymbol{\eta}] dv - \int_{\partial\Omega_0} \mathbf{t}_{0,n+1} \cdot \boldsymbol{\eta} da = 0, \quad \forall \boldsymbol{\eta} \in \mathcal{V}_{u,n+1}, \quad (4.46)$$

$$\int_{\Omega_0} \left[\left(\rho_0 C_V \dot{\theta}_{n+1} - \hat{\mathcal{D}}_{\text{mech}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1}, \alpha_n) - \hat{\mathcal{H}}^{\text{ep}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1}, \alpha_n) - \rho_0 r_{n+1} \right) \xi - \mathbf{q}_{0,n+1} \cdot \nabla_0 \xi \right] dv - \int_{\partial\Omega_0} \hat{h}_{0,n+1} \xi da = 0, \quad \forall \xi \in \mathcal{V}_{\theta,n+1}, \quad (4.47)$$

where

$$\mathcal{K}_{u,n+1} \equiv \{\mathbf{u} : \Omega \times \mathcal{R} \rightarrow \mathcal{U} \mid \mathbf{u}_{n+1}(\mathbf{X}) = \mathbf{u}_{\text{presc},n+1}(\mathbf{X}), \mathbf{X} \in \partial\Omega_{\text{motion},0}\}, \quad (4.48)$$

$$\mathcal{K}_{\theta,n+1} \equiv \{\theta : \Omega \times \mathcal{R} \rightarrow \mathcal{R} \mid \theta_{n+1}(\mathbf{X}) = \theta_{\text{presc},n+1}(\mathbf{X}), \mathbf{X} \in \partial\Omega_{\text{temperature},0}\}. \quad (4.49)$$

4.2 Finite Element Method

It is now possible to apply the finite element method to discretize spatial the incremental thermo-mechanical initial boundary value problem. The approach is entirely similar to the one present in the context of the strictly mechanical (see Section ??) and strictly thermal (see Section ??) problems. As such, some details are omitted.

4.2.1 Interpolation

Defining the global vector of nodal displacements \mathbf{u} and the global vector of nodal temperatures θ , like before, as

$$\mathbf{u}(t) = \left[u_1^1(t), \dots, u_{n_{\text{dim}}}^1(t), \dots, u_1^{n_{\text{points}}}(t), \dots, u_{n_{\text{dim}}}^{n_{\text{points}}}(t) \right]^T, \quad (4.50)$$

$$\boldsymbol{\theta}(t) = \left[\theta^1(t), \theta^2(t), \dots, \theta^{n_{\text{points}}}(t) \right]^T, \quad (4.51)$$

the displacement $\mathbf{u}(\mathbf{X}, t)$ and temperature $\theta(\mathbf{X}, t)$ fields defined over the global domain Ω_0 , can be found at any point \mathbf{X} as

$$^h \mathbf{u}(\mathbf{X}, t) \equiv \mathbf{N}^{g,u}(\mathbf{X}) \mathbf{u}(t), \quad ^h \mathbf{u} \in ^h \mathcal{K}_u, \quad (4.52)$$

$$^h \theta(\mathbf{X}, t) \equiv \mathbf{N}^{g,\theta}(\mathbf{X}) \boldsymbol{\theta}(t), \quad ^h \boldsymbol{\theta} \in ^h \mathcal{K}_\theta, \quad (4.53)$$

where $\mathbf{N}^{g,u}$ and $\mathbf{N}^{g,\theta}$ are the global interpolation matrices with appropriate dimensions given the dimensions of \mathbf{u} and θ . The global shape functions used to interpolate between the nodal values of the displacement and the temperature can even be different if need be.

4.2.2 Spatial discretization

The application of finite element discretization to the mechanical part of the incremental thermo-mechanical constitutive initial boundary value problem is exactly the same as the one presented in Section ?? for the strictly mechanical problem, as the equation to be discretized is the same. On the other hand, applying the discretization to the thermal part of the incremental thermo-mechanical constitutive initial boundary value problem, yields

$$\int_{h\Omega_0} \left[\left(\rho_0 C_V \mathbf{N}^{g,\theta} \dot{\theta}_{n+1} - \hat{\mathcal{D}}_{\text{mech}} - \hat{\mathcal{H}}^{\text{ep}} - \rho_0 r_{n+1} \right) \cdot \mathbf{N}^{g,\theta} \boldsymbol{\xi} - \mathbf{q}_{0,n+1} \cdot \mathbf{H}^{g,\theta} \boldsymbol{\xi} \right] d\nu - \int_{h\partial\Omega_0} h_{0,n+1} \mathbf{N}^{g,\theta} \boldsymbol{\xi} da = 0, \quad \forall \boldsymbol{\xi} \in {}^h\mathcal{V}_\theta, \quad (4.54)$$

where $\mathbf{H}^{g,\theta}$ is the discrete global gradient operator for scalars, defined in Equation (??). Equation (4.55) can be rewritten as

$$\left\{ \int_{h\Omega_0} \left[\mathbf{N}^{g,\theta T} \left(\rho_0 C_V \mathbf{N}^{g,\theta} \dot{\theta}_{n+1} - \hat{\mathcal{D}}_{\text{mech}} - \hat{\mathcal{H}}^{\text{ep}} - \rho_0 r_{n+1} \right) - \mathbf{H}^{g,\theta T} \mathbf{q}_{0,n+1} \right] d\nu - \int_{h\partial\Omega_0} \mathbf{N}^{g,\theta T} h_{0,n+1} da \right\}^T \boldsymbol{\xi} = 0, \quad \forall \boldsymbol{\xi} \in {}^h\mathcal{V}_\theta, \quad (4.55)$$

and, since it must be satisfied for any $\boldsymbol{\xi} \in {}^h\mathcal{V}_\theta$, the incremental discretized thermo-mechanical constitutive initial boundary value problem can thus be stated in the material description as

Problem 4.5 | Material incremental discretized thermo-mechanical initial BVP.

Given the set of internal variables $\boldsymbol{\alpha}_n$ at t_n , the prescribed body and traction force fields $\mathbf{b}_{0,n+1}$ and $\mathbf{t}_{0,n+1}$, as well as, the prescribed heat sources and heat flux fields $\mathbf{r}_{0,n+1}$ and $h_{0,n+1}$, and both the prescribed deformation gradient \mathbf{F}_{n+1} and the prescribed temperature θ_{n+1} at t_{n+1} , find the kinematically admissible nodal displacement field $\mathbf{u}_{n+1} \in {}^h\mathcal{K}_{u,n+1}$ and the admissible nodal temperature field $\theta_{n+1} \in {}^h\mathcal{K}_{\theta,n+1}$ such that the body \mathcal{B} is in mechanical and energetic equilibrium

$$\mathbf{M}\ddot{\mathbf{u}}_{n+1} + \mathbf{f}_u^{\text{int}}(\boldsymbol{\theta}_{n+1}, \mathbf{u}_{n+1}) - \mathbf{f}_{u,n+1}^{\text{ext}} = \mathbf{0}, \quad (4.56)$$

$$\mathbf{C}\dot{\boldsymbol{\theta}}_{n+1} + \mathbf{f}_\theta^{\text{int}}(\boldsymbol{\theta}_{n+1}, \mathbf{u}_{n+1}) - \mathbf{f}_{\theta,n+1}^{\text{ext}} = \mathbf{0}, \quad (4.57)$$

where $\mathbf{f}_u^{\text{int}}$ e $\mathbf{f}_{u,n+1}^{\text{ext}}$ are the mechanical global vectors of internal and external forces

defined as

$$\mathbf{f}_u^{\text{int}} \equiv \int_{h\Omega_0} \mathbf{G}^{g,uT} \hat{\mathbf{P}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1} \boldsymbol{\alpha}_n) \, dv, \quad (4.58)$$

$$\mathbf{f}_{u,n+1}^{\text{ext}} \equiv \int_{h\Omega_0} \mathbf{N}^{g,uT} \mathbf{b}_{0,n+1} \, dv + \int_{\partial^h \Omega_{\text{traction},0}} \mathbf{N}^{g,uT} \mathbf{t}_{0,n+1} \, da, \quad (4.59)$$

and $\mathbf{f}_\theta^{\text{int}}$ e $\mathbf{f}_{\theta,n+1}^{\text{ext}}$ the thermal global vectors of internal and external forces defined as

$$\mathbf{f}_\theta^{\text{int}} \equiv \left[\int_{h\Omega_0} \mathbf{N}^{g,\theta T} \hat{\mathcal{D}}_{\text{mech}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1} \boldsymbol{\alpha}_n) + \mathbf{N}^{g,\theta T} \hat{\mathcal{H}}^{\text{ep}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1} \boldsymbol{\alpha}_n) \right. \\ \left. \mathbf{H}^{g,\theta T} \mathbf{q}_0(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1}) \right] \, dv, \quad (4.60)$$

$$\mathbf{f}_{\theta,n+1}^{\text{ext}} \equiv \int_{h\Omega_0} \mathbf{N}^{g,\theta T} \mathbf{r}_{0,n+1} \, dv + \int_{\partial^h \Omega_{\text{heat},0}} \mathbf{N}^{g,\theta T} h_{0,n+1} \, da. \quad (4.61)$$

The mass matrix \mathbf{M} and the thermal capacitance matrix \mathbf{C} are defined by

$$\mathbf{M} = \int_{h\Omega_0} \rho_0 \mathbf{N}^{g,uT} \mathbf{N}^{g,u} \, dv, \quad (4.62)$$

$$\mathbf{C} = \int_{h\Omega_0} \rho_0 C_V \mathbf{N}^{g,\theta T} \mathbf{N}^{g,\theta} \, dv. \quad (4.63)$$

The global vectors for the internal and external forces are usually obtained by assemblage of their elemental counterparts as

$$\mathbf{f}_u^{\text{int}} = \mathbf{A}_{e=1}^{n_{\text{elem}}} \left(\mathbf{f}_u^{\text{int}} \right)^{(e)}, \quad (4.64)$$

$$\mathbf{f}_u^{\text{ext}} = \mathbf{A}_{e=1}^{n_{\text{elem}}} \left(\mathbf{f}_u^{\text{ext}} \right)^{(e)}, \quad (4.65)$$

$$\mathbf{f}_\theta^{\text{int}} = \mathbf{A}_{e=1}^{n_{\text{elem}}} \left(\mathbf{f}_\theta^{\text{int}} \right)^{(e)}, \quad (4.66)$$

$$\mathbf{f}_\theta^{\text{ext}} = \mathbf{A}_{e=1}^{n_{\text{elem}}} \left(\mathbf{f}_\theta^{\text{ext}} \right)^{(e)}, \quad (4.67)$$

where the mechanical elemental vectors in the material description are defined as

$$\mathbf{f}_u^{\text{int}} \equiv \int_{h\Omega_0^{(e)}} \mathbf{G}^{uT} \hat{\mathbf{P}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1} \boldsymbol{\alpha}_n) \, dv, \quad (4.68)$$

$$\mathbf{f}_{u,n+1}^{\text{ext}} \equiv \int_{h\Omega_0^{(e)}} \mathbf{N}^{uT} \mathbf{b}_{0,n+1} \, dv + \int_{\partial^h \Omega_{\text{traction},0}^{(e)}} \mathbf{N}^{uT} \mathbf{t}_{0,n+1} \, da, \quad (4.69)$$

and the thermal vectors, also in the material description, are

$$\mathbf{f}_\theta^{\text{int}} \equiv \left[\int_{h\Omega_0^{(e)}} \mathbf{N}^{\theta T} \hat{\mathcal{G}}_{\text{mech}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1} \boldsymbol{\alpha}_n) + \mathbf{N}^{\theta T} \hat{\mathcal{H}}^{\text{ep}}(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1} \boldsymbol{\alpha}_n) \right. \\ \left. \mathbf{H}^{\theta T} \mathbf{q}_0(\mathbf{F}(\mathbf{u}_{n+1}), \theta_{n+1}) \right] \mathrm{d}v, \quad (4.70)$$

$$\mathbf{f}_{\theta, n+1}^{\text{ext}} \equiv \int_{h\Omega_0^{(e)}} \mathbf{N}^{\theta T} \mathbf{r}_{0, n+1} \mathrm{d}v + \int_{\partial h\Omega_{\text{heat}, 0}^{(e)}} \mathbf{N}^{\theta T} h_{0, n+1} \mathrm{d}a. \quad (4.71)$$

The elemental mechanical interpolation of the matrices \mathbf{N}^u , \mathbf{B}^u , \mathbf{N}^θ and \mathbf{H}^θ are the elemental mechanical interpolation matrix, the mechanical discrete elemental gradient operator, the elemental thermal interpolation matrix, and the thermal discrete elemental gradient operator for scalars.

In a similar manner, the global mass and capacitance matrices are also usually obtained by assemblage of their elemental counterparts as

$$\mathbf{M} \equiv \bigcup_{e=1}^{n_{\text{elem}}} \mathbf{M}^{(e)}, \quad (4.72)$$

$$\mathbf{C} \equiv \bigcup_{e=1}^{n_{\text{elem}}} \mathbf{C}^{(e)}, \quad (4.73)$$

where the elemental mass matrices in the material description are defined as

$$\mathbf{M}^{(e)} = \int_{h\Omega^{(e)}} \rho_0 \mathbf{N}^u T \mathbf{N}^u \mathrm{d}v, \quad (4.74)$$

and the elemental thermal capacitance matrices, also in the material description, are defined as

$$\mathbf{C}^{(e)} = \int_{h\Omega^{(e)}} \rho_0 C_V \mathbf{N}^{\theta T} \mathbf{N}^\theta \mathrm{d}v, \quad (4.75)$$

4.3 Linearisation

The equilibrium equation, Equation (2.42) in a spatial description and Equation (3.25) in a material description, is generally nonlinear due to geometrical and/or material nonlinearities. The Newton-Raphson Method is an efficient and robust iterative scheme with a quadratic convergence rate often used to solve the equilibrium equation at each time increment, t_n . The residual of the fully discretized balance of linear momentum is defined for an iteration step i of the Newton-Raphson method as

$$\mathbf{r}(\mathbf{u}_{n+1}^i) = \mathbf{M} \ddot{\mathbf{u}}_{n+1}^i + \mathbf{f}^{\text{int}}(\mathbf{u}_{n+1}^i) - \mathbf{f}_{n+1}^{\text{ext}}. \quad (4.76)$$

A Taylor expansion about the current solution \mathbf{u}_{n+1}^i is performed, discarding all terms of higher order than one, yielding the linearised form

$$\text{Lin} \mathbf{r}(\mathbf{u}_{n+1}^i) = \mathbf{r}(\mathbf{u}_{n+1}^i) + \underbrace{\frac{\partial \mathbf{r}(\mathbf{u}_{n+1})}{\partial \mathbf{u}_{n+1}}}_{\mathbf{K}(\mathbf{u}_{n+1}^i)} \bigg|_{\mathbf{u}_{n+1}^i}^i \delta \mathbf{u}. \quad (4.77)$$

with the dynamic effective tangential stiffness matrix $\mathbf{K}(\mathbf{u}_{n+1}^i)$. The linearisation of the internal forces included in \mathbf{K} is known as the tangential stiffness matrix \mathbf{K}_T , which is defined as

$$\mathbf{K}_T^i = \left. \frac{\partial \mathbf{f}^{\text{int}}}{\partial \mathbf{u}_{n+1}} \right| ^i. \quad (4.78)$$

Equilibrium is achieved if

$$\text{Lin} \mathbf{r}(\mathbf{u}_{n+1}^i) = \mathbf{0}, \quad (4.79)$$

so that a linear system of equation is given by

$$\mathbf{K}(\mathbf{u}_{n+1}^i) \delta \mathbf{u} = -\mathbf{r}(\mathbf{u}_{n+1}^i). \quad (4.80)$$

Thus, a new solution of the displacement increment $\delta \mathbf{u}$ for current iteration step $i + 1$ is determined, and the final displacement solution of time step $n + 1$ is obtained via updating

$$\mathbf{u}_{n+1}^{i+1} = \mathbf{u}_{n+1}^i + \delta \mathbf{u}. \quad (4.81)$$

A solution of t_{n+1} is found, i.e. an equilibrium state is reached and $\mathbf{u}_{n+1} = \mathbf{u}_{n+1}^{i+1}$, if prescribed, user-defined convergence criteria are fulfilled.

Correção do TPC de DT

Chapter 5

Solution procedures for coupled fields

This chapter presents an overview of solution procedures for coupled fields. It includes techniques applied to various coupled field problems, such as thermomechanical coupling, fluid-structure interaction, and aeroelasticity. Its goal is to support the choice of solution techniques for thermo-plastic problems, which are accurate, stable, efficient, both in terms of memory and computational time, and easy to implement and extend later to further couplings, e.g., electro-thermomechanical problems.

5.1 Context field elimination

Field elimination achieves the solution of a coupled problem by eliminating the variables of the first field and introducing them into the second field. This second field is then solved.

The main advantage of this procedure is the reduction of the number of state variables. Which in turn, leads to smaller systems of equations, which are presumably easier to solve. Furthermore, the analyst can choose the remaining variables such that they are the variables of interest. In this way, the variables eliminated do not need to be recovered. (Felippa and Park, 1980)

On the other hand Felippa and Park (1980) cite as disadvantages

- only possible for problems allowing explicit (and well-conditioned) variable eliminations;
- sparseness and symmetry attributes of matrices associated with the original coupled system can be adversely affected by the eliminations process; and
- available software modules for the isolated fields are not likely to be of much use for processing the reduced system.

The remainder of the chapter disregards these procedures, including in Section ?? where the comparison of the different schemes is discussed.

5.2 Monolithic

Monolithic algorithms solve the coupled nonlinear multi-physics system simultaneously. Predominantly, implicit schemes are applied to achieve good stability properties. In turn, the nonlinear residual equations are often solved using the Newton-Raphson method. A particular challenge for monolithic algorithms is the efficient solution of the large system of equations, including potential nonlinearities or lack of symmetry. Even the units chosen can contribute to the ill-conditioning of the system matrix. One essential aspect for efficient solvers for large-scale problems is a good preconditioning technique.

5.2.1 Numerical considerations

For the solution of a large system of equations, iterative methods are preferable to direct methods, in part, due to memory footprint considerations. The Newton-Krylov methods such as GMRES and the BiCGStab are among the most commonly used in multi-physics problems (Hron and Turek, 2006). However, their use does not suffice for an efficient and robust solution procedure for a multi-physics problem. In addition, the use of preconditioners alleviates the possible large condition numbers of the system matrix. There are several preconditioning techniques for the solution of large systems of equations, e.g., ILU preconditioners, domain decomposition, including multigrid approaches; multilevel recursive Schur complements preconditioners (see Smith et al. (2004) and Chen (2005)).

Heil (2004) is concerned with the fully coupled solution of large-displacement fluid-structure interaction problems by Newton's method. They use block-triangular approximations of the Jacobian matrix, obtained by neglecting selected fluid-structure interaction blocks, and show that they provide suitable preconditioners for the solution of the linear systems with GMRES. A Schur complement approximation for the Navier-Stokes block and multigrid approximations for the solution of the computationally most expensive operations is the basis for the efficient approximate implementation of the preconditioners.

Hron and Turek (2006) propose a method based on a fully implicit, monolithic formulation of the problem in the arbitrary Lagrangian-Eulerian framework to solve the problem of fluid-structure interaction of an incompressible elastic object in laminar incompressible viscous flow. They utilize the standard geometric multigrid approach based on a hierarchy of grids obtained by successive regular refinement of a given coarse mesh. The complete multigrid iteration is performed in the standard defect-correction setup with the V or F-type cycle.

Tezduyar et al. (2006) show how preconditioning techniques more sophisticated than diagonal preconditioning can be used in iterative solutions of the linear equation systems in fluid-structure interaction problems.

In Gee et al. (2011), the authors focus on the strong coupling fluid-structure interaction employing monolithic solution schemes. Therein, a Newton-Krylov method is applied to the monolithic set of nonlinear equations. They propose two preconditioners that apply algebraic multigrid techniques to the entire fluid-structure interaction system of equations. As the first option, the authors employ a standard block Gauss-Seidel approach, where approximate inverses of the individual field blocks are based on an algebraic multigrid hierarchy tailored for the type of the underlying physical problem. A monolithic coarsening scheme for the coupled system that uses prolongation and restriction projections constructed for the individual fields

provides the basis for the second preconditioner. The resulting nonsymmetric monolithic algebraic multigrid method involves coupling the fields on coarse approximations to the problem yielding significantly enhanced performance, claim the authors.

In the context of multi-physics problems, [Lin et al. \(2010\)](#) propose a fully coupled algebraic multilevel preconditioner for Newton-Krylov solution methods. A set of multi-physics partial differential equation (PDE) applications attests its performance: a drift-diffusion approximation for semiconductor devices, a low Mach number formulation for the simulation of coupled flow, transport, and non-equilibrium chemical reactions, a low Mach number formulation for visco-resistive magnetohydrodynamics (MHD) systems. An aggressive-coarsening graph-partitioning of the non-zero block structure of the Jacobian matrix provides the basis for the algebraic multilevel preconditioner. Using a different approach [Badia et al. \(2014\)](#) employ a new family of recursive block LU preconditioners to solve the thermally coupled induction less magnetohydrodynamics problem equations, which model the flow of an electrically charged fluid under the influence of an external electromagnetic field with thermal coupling.

[Netz \(2013\)](#) addresses a thermo-mechanically coupled problem of thermo-viscoelasticity at finite strains using a monolithic approach. The authors solve the system of nonlinear algebraic equations obtained from the spatial (FEM) and temporal (diagonally-implicit Runge-Kutta methods) discretization of the problem monolithically. They employ the Multilevel-Newton algorithm to obtain a high-order result in the space and the time domain. The numerical concept is applied to a constitutive model of finite strain thermo-viscoelasticity. [Rothe et al. \(2015\)](#) also employ in the context of thermo-viscoelasticity, the multilevel Newton algorithm to solve the system of algebraic equations describing the discretized problem.

[Danowski et al. \(2013\)](#) presents a monolithic solution scheme for thermo-structure interaction problem, using right preconditioning and a GMRES. The preconditioner "sub-problem" is solved using a Richardson iteration scheme and a relaxed block Gauss-Seidel method, which uncouples the mechanical and thermal problems. This procedure tackles each problem using an independent algebraic multigrid (AMG) preconditioner. [Verdugo and Wall \(2016\)](#) also considers the procedure just mentioned, as well as a preconditioner based on a semi-implicit method for pressure-linked equations, extended to deal with an arbitrary number of fields. This technique also results in uncoupled problems that can be solved with standard AMG. They also introduce a more sophisticated preconditioner that enforces the coupling at all AMG levels, unlike the other two techniques, which resolve the coupling only at the finest level. These techniques are applied successfully to three different coupled problems: thermo-structure interaction, fluid-structure interaction, and a complex model of the human lung.

[Mayr et al. \(2020\)](#) propose a hybrid interface preconditioner for the monolithic solution of surface-coupled problems. They combine physics-based block preconditioners with an additional additive Schwarz preconditioner, whose subdomains span across the interface on purpose. This approach is motivated by the error assessment of physics-based block preconditioners, revealing an accumulation of the error at the coupling surface, despite their overall efficiency.

5.2.2 Usage examples

Thermo-mechanical coupling In the following paragraph, a small overview of the literature is presented regarding the application of monolithic solvers to the thermo-mechanical coupled problem. [Carter and Booker \(1989\)](#) suggests a monolithic approach to the thermoelastic problem at small strains. The constitutive laws considered do not acknowledge the dependence of the mechanical properties on the temperature and are not deduced from a Helmholtz energy function. [Glaser \(1992\)](#) uses monolithic algorithms for the calculation of thin-walled structures using shell elements and an arc-length method for the TSI solution. While all coupling terms were considered, only a simplified mechanical dissipation was included where the hardening power was neglected (according to [Danowski \(2014\)](#)). [Ibrahimbegovic and Chorfi \(2002\)](#) presents a thermoplasticity covariant formulation within the framework of the principal axis methodology, which the authors claim, leads to a very efficient numerical implementation. The paper contains several numerical simulations dealing with the fully coupled thermomechanical response at large viscoplastic strains, including strain localization and cyclic loading cases, to illustrate the performance of the proposed methodology. The authors consider the von Mises thermoplasticity yield criterion and strain energy depending on logarithmic stretches, a hardening variable, and temperature. A monolithic solver achieves the solution to the coupled problem, but no details about it are given. [Danowski \(2014\)](#) proposes a volume-coupled TSI model based on the finite element method for the structural and thermal field. Various temperature-dependent, isotropic, elastic, and elastoplastic material models for small and finite strains are employed, incorporating the effect of the highly elevated temperatures predominating in rocket nozzles, the practical application focused in the Ph.D. thesis. The author considers both monolithic and partitioned coupling algorithms to solve fully coupled thermomechanical systems. Regarding the former, a novel monolithic Newton-Krylov scheme with problem-specific block Gauss-Seidel preconditioner and algebraic multigrid methods is introduced. Concerning the latter, loosely and strongly coupled partitioned schemes are examined, possibly including acceleration techniques, as, e.g., the Aitken Δ^2 method. [Netz \(2013\)](#) and [Rothe et al. \(2015\)](#) both present monolithic approaches, based on the multilevel Newton method, for the solution of the thermo-mechanical problem. In both contributions, thermo-visco-plastic materials are successfully analyzed. Recently, [Felder et al. \(2021\)](#) have presented a finite strain thermo-mechanically coupled two-surface damage-plasticity theory. The authors obtain the solution for the three coupled fields, displacement, nonlocal damage variable, and temperature, employing an implicit and monolithic solution scheme.

The thermo-mechanical coupling has also been studied in the more specific domain of contact mechanics. [Zavarise et al. \(1992\)](#) present one of the earliest contributions in this direction. They propose a FEM formulation of frictionless contact, accounting for full thermo-elastic coupling. The penalty method is used to enforce the non-penetration conditions. Another contribution, [Hansen \(2011\)](#), advances a standard mortar discretization with Lagrange multipliers to solve the small strain thermo-elasticity problem. The authors consider the heat equation coupled to linear mechanics through a thermal expansion term in their formulation. The solution approach is based on a preconditioned Jacobian-free Newton Krylov solution method, and the use case under analysis is a light water reactor nuclear fuel rod. [Dittmann et al. \(2014\)](#) investigate thermomechanical mortar contact algorithms and their application to NURBS-based Isogeometric Analysis in the context of nonlinear

elasticity. Mortar methods are applied to both the mechanical and thermal fields to model frictional contact, the energy transfer between the surfaces, and frictional heating. A monolithic approach is pursued in solving the nonlinear algebraic equations found after the discretization in time and space. In the Ph.D. thesis by the same first author, [Dittmann \(2017\)](#), this approach is further pursued in multi-field contact problems. More recently, [Seitz et al. \(2018\)](#); [Seitz \(2019\)](#) tackles the numerical treatment of contact problems considering inelastic deformation and thermomechanical coupling. It accounts for plastic spin, visco-plasticity, and thermo-plastic coupling, as well as temperature-dependent material parameters. The authors also opt for a monolithic solver, although no further details are supplied. See also, in the context of contact mechanics, [Oancea and Laursen \(1997\)](#); [Pantuso et al. \(2000\)](#); [Hüeber and Wohlmuth \(2009\)](#); [Hesch and Betsch \(2011\)](#); [Gitterle \(2012\)](#) and [Novascone et al. \(2015\)](#).

Others In the context of fluid-structure interaction, the monolithic approach seems to be more widely used than in thermo-mechanically coupled problems. A few contributions in this domain using a monolith approach are [Blom \(1998\)](#); [Heil \(2004\)](#); [Hübner et al. \(2004\)](#); [Michler et al. \(2004\)](#); [Zhang and Hisada \(2004\)](#); [Dettmer and Perić \(2006\)](#); [Hron and Turek \(2006\)](#); [Tezduyar et al. \(2006\)](#); [Küttler et al. \(2010\)](#); [Gee et al. \(2011\)](#); [Klöppel et al. \(2011\)](#); [Mayr et al. \(2015\)](#) and [Mayr et al. \(2020\)](#). The use of a monolithic approach can also be found in the domain of saturated soils (e.g., [Lewis and Sukirman \(1993\)](#), [Borja et al. \(1998\)](#), [Jha and Juanes \(2007\)](#), [White and Borja \(2008\)](#)). Monolithic solvers are also used in the context of magnetohydrodynamics (e.g., [Shadid et al. \(2010\)](#) and [Badia et al. \(2014\)](#)).

5.3 Partitioned

The following Section presents the partitioned time-stepping algorithms. For a detailed comparison with the monolithic approach and between themselves, see Section ??.

A field partition is a field-by-field decomposition of the space discretization. Partitioning may be algebraic or differential. In algebraic partitioning, the complete coupled system is spatially discretized first and then decomposed. In differential partitioning, the decomposition is done first, and each field is then discretized separately. Differential partitioning often leads to non-matched meshes, as typical of fluid-structure interaction. Algebraic partitioning was initially developed for matched meshes and substructuring ([Felippa et al., 2001](#)).

The earliest contributions regarding the partitioned treatment of coupled systems emerged in the mid 1970s, involving structure-structure interactions and fluid-structure interactions (see e.g. [Belytschko and Mullen \(1976\)](#), [Park et al. \(1977\)](#), [Belytschko and Mullen \(1978\)](#), [Hughes and Liu \(1978\)](#) and [Belytschko et al. \(1979\)](#)).

Given a complex system, there are usually many ways of partitioning it into subsystems or fields. [Felippa and Park \(1980\)](#) provide a very pragmatic and helpful criterion to select the fields to be considered. According to their definition, a field is characterized by computational considerations. It is a segment of the overall problem for which a separable software module is either available or readily prepared if the interaction terms are suppressed. As such, a partitioned approach to the solution of multi-physics problems employs field analyzers specific to each field separately stepped in time. The coupling between the fields is achieved through proper

communication between the individual components using prediction, substitution, and synchronization techniques.

Before moving on, it may be helpful to clear up the difference between partitioned schemes, staggered schemes, operator splits, and fractional-step methods. The first is probably the most general term and includes the others. Its definition has already been given. A staggered scheme is a term most often used for the partitioned schemes where the solution concerning each field is sequential and obtained only once per time step as in the loosely coupled schemes to be introduced. However, it may also include the strongly coupled schemes, as well. An operator split is obtained through the decomposition of the fully coupled problem into subproblems. The structure of the problem is the same, as well as the unknowns considered. The only difference between the subproblems is the physical effects considered. The equation terms concerning each physical effect must be divided exclusively and exhaustively between the subproblems. Finally, according to [Armero and Simo \(1992\)](#) staggered algorithms for coupled problems can be viewed as fractional steps methods, in the sense of [Holt and Yanenko \(2012\)](#), arising from an operator split of the coupled problem of evolution.

5.3.1 Operator splits

The most common operator splits into thermomechanical problems are the isothermic and adiabatic split.

Isothermic The isothermic split is perhaps the most straightforward and natural approach, as noted by [Argyris and Doltsinis \(1981\)](#), one of the earliest contributions on the topic. The scheme achieves the solution of the thermo-mechanical problem, first solving the mechanical problem at a constant temperature, then a purely thermal phase is considered at a fixed configuration.

Adiabatic The adiabatic split is proposed in [Armero and Simo \(1992\)](#). It consists of a first phase where constant entropy is enforced and a second phase of purely thermal conduction with a fixed reference. In terms of implementation complexity, it is comparable to the isothermal split. This is possible because the constant entropy phase can be cast as a mechanical phase, where the stiffness and the external force are adjusted as a function of an intermediate temperature. This temperature is computed considering the strong form of the temperature evolution equation to retain the computational efficiency of the isothermal split, despite the momentum equation being enforced in its weak form. The advantage of this split is that when used in a staggered scheme, it is unconditionally stable (see Section ??).

5.3.2 Loosely vs. Strongly coupled schemes

According to [Felippa et al. \(2001\)](#) there are several basic techniques associated with partitioned schemes (see Figure 5.1). These are

- prediction;
- substitution;
- interfield iteration;

- full step correction;
- lockstep advancing;
- midpoint correction;
- subcycling;
- augmentation.

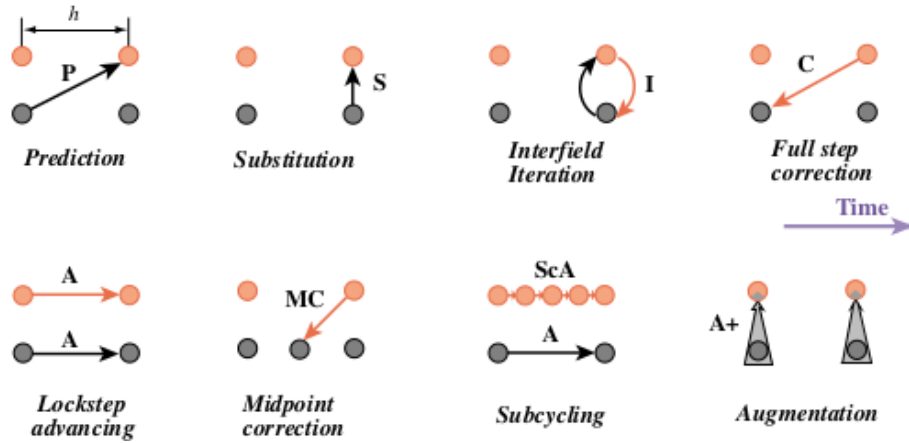


Figure 5.1: Devices of partitioned analysis time-stepping (Felippa et al., 2001).

Inter-field iterations are the primary criterion distinguishing loosely or one-way staggered coupled schemes and strongly or iterative staggered coupled schemes. In the loosely coupled schemes, the integration algorithm proceeds sequentially, solving the problem in each field only once per time step. On the other hand, for strongly coupled schemes, inter-field iterations are present, such that the problems are solved multiple times at the same time instant. This inner loop is repeated until a given tolerance is reached for the unknowns in each field.

The remainder of the techniques listed will be mentioned and explained in the discussion below.

5.3.3 Loosely coupled

The solution for the fully coupled problem is found in loosely coupled schemes by solving each field sequentially. For the thermomechanical problem, the two available schemes are the isothermal split (see e.g. Simo and Armero (1992), Agelet de Saracibar (1998)) and the adiabatic split (see e.g. Armero and Simo (1992) and Armero and Simo (1993)), as mentioned above.

According to Felippa et al. (2001), in linear problems, the first concern with partitioning is the degradation of time-stepping stability. After the analyst has ensured stability, an accuracy analysis of the method should be performed. In strongly nonlinear problems, such as fluid flow, stability and accuracy tend to be intertwined since numerical stability is harder to define. As such, they are usually considered together in method design. The expectation is for a method that operates well at a reasonable timestep.

Felippa and Geers (1988) present a detailed explanation about how to design from scratch a loosely coupled time-stepping algorithm applicable to linear systems of equations. It includes implementation details, such as the choice of the predictor formula, and the design steps, from the formulation of the original field equations and temporal discretization to the stability and accuracy analysis. Other contributions focused mainly on linear systems are Neishlos (1983), Zienkiewicz et al. (1988) and Combescure and Gravouil (2002),

Because the loosely coupled schemes are explicit, they are also often only conditionally stable. The isothermic split is such an example (Armero and Simo, 1992). On the other hand, the adiabatic split proposed in Armero and Simo (1992) is unconditionally stable, despite being explicit. Farhat et al. (1991) also propose a stable staggered scheme, achieved through semi-algebraic augmentation, which, however, is limited to linearized thermoelasticity. In the context of coupled flow and geomechanics, Kim et al. (2011b) show that when the mechanical problem is solved first, the drained split combined with a backward Euler discretization is conditionally stable, and the undrained split is unconditionally stable when combined with midpoint rule. When instead the flow problem is solved first, Kim et al. (2011a) show that the fixed-stress split is conditionally and the fixed-strain split is unconditionally stable for appropriate choices of the generalized midpoint rule.

Moreover, in the domain of fluid-structure interactions, it can be shown that staggered methods are inherently non-conservative. As time progresses in the simulation, these schemes introduce parasitic energy at the boundary, which contributes to their poor numerical stability (Michler et al., 2003). A further problem appears when solving these coupled physical problems, the so-called artificial added-mass effect, which leads to instability. It manifests itself when a slender structure and fluid have similar densities, and the latter is modeled as an incompressible fluid (Causin et al., 2005; Förster, 2007). It can even be shown that for every sequentially staggered scheme and spatial discretization of a problem, a mass ratio between the fluid and structural mass density can be found at which the coupled system becomes unstable (Förster et al., 2007).

Despite this, some contributions detail strategies allowing for the unconditional stability of these schemes. As part of the development loop of commercial tire designs, Gillard (2019) tackles the problem of tire hydroplaning. The author presents a robust explicit coupling scheme that relies on rigorous control of the energy artificially introduced at the interface by the staggering process through a dynamic adaptation of the coupling time step size. Regarding the artificial added-mass effect, Farhat et al. (2010) demonstrates that even for fluid-structure applications with strong added mass effects, a carefully designed staggered and sub-iteration-free time-integrator can achieve numerical stability and robustness concerning the slenderness of the structure, as long as the fluid is justifiably modeled as a compressible medium.

Another technique available to improve the stability of loosely coupled schemes is (algebraic) augmentation. It rests on the injection of one of the coupled equations into the other, after discretization in space, to 'soften' the system, either by reducing the large eigenvalues of the uncoupled stiff equation or by introducing some damping into it. Some examples of this approach include Park et al. (1977) and Park (1983).

Yet another technique to ensure stability in the context of fluid-structure interactions is presented in Fernández et al. (2006). Stability is achieved employing a semi-implicit coupling scheme, splitting the added-mass, viscous effects, and geometrical/convective nonlinearities, through a Chorin-Temam projection scheme

within the fluid.

Regarding accuracy, the loosely coupled schemes do not necessarily inherit the accuracy order of the schemes used in the integration of the separate fields, often being just of the first order in time (Farhat et al., 2006). However, some contributions detail approaches that are second-order time-accurate. In the context of thermo-elasticity, Armero and Simo (1992) show that a double-pass approach using the adiabatic split yields such a second-order accurate time-stepping algorithm. A few approaches yield similar results in the domain of fluid-structure interaction (see Piperno (1997), Farhat et al. (2006) and Farhat et al. (2010)). In any case, whatever the theoretical convergence order of the loosely coupled method, at a given time instant, the fully coupled discretized equations of the problem will never be exactly satisfied by the solutions found. There is a lag between the fields considered, e.g., the mechanical and thermal fields in a thermomechanical problem. In the context of strong coupling, this lag can be conceived as a numerical evaluation error. Solving approximately the exact (i.e., aggregated) equations can be reinterpreted as exactly solving a set of approximate (i.e., segregated) equations. Thus, one can construe loosely-coupled methods as solving a set of segregated equations instead of aggregated equations. Accordingly, the incurred numerical evaluation error can be reinterpreted as a discretization error. Loosely-coupled methods, therefore, satisfy conservation only in an asymptotic sense, i.e., for vanishing mesh width; this is a basic consistency requirement Michler (2005).

Prediction techniques can improve the order of the numerical evaluation error incurred by loosely-coupled partitioned methods. For the sake of explanation, consider the thermo-mechanical problem being solved using the isothermic split. When using predictors, instead of integrating the mechanical equations based on the structure's temperature in the previous time instant, a prediction can be used for the temperature of the structure boundary in the current time instant. Such predictions are generally based on an extrapolation of the solution from the previous time step. Prediction techniques improve the solution accuracy and stability of loosely-coupled methods (Piperno, 1997; Piperno and Farhat, 2001; Michler, 2005; Farhat et al., 2006).

Another technique available to improve the accuracy of the loosely coupled methods is subcycling. It involves solving each field's problems using different time steps since the fields present in a multi-physics problem often have different time scales. In the context of aeroelasticity, Piperno et al. (1995) claims that it can offer substantial computational advantages, including savings in the simulation CPU time because the structural field will be advanced fewer times. Farhat et al. (1997) and Piperno (1997) also argue for this technique along the same lines.

Usage examples The loosely coupled scheme has been used in the context of thermoelasticity (Argyris and Doltsinis, 1981; Armero and Simo, 1992; Johansson and Klarbring, 1993; Miehe, 1995a,b; Holzapfel and Simo, 1996), thermo-plasticity (Armero and Simo, 1992, 1993; Simo and Miehe, 1992; Wriggers et al., 1992; Agelet de Saracibar, 1998; Agelet de Saracibar et al., 1999) and thermo-viscoplasticity (Adam and Ponthot, 2002a,b).

For examples in aeroelasticity see e.g., Piperno et al. (1995); Farhat and Lesoinne (2000) and Farhat et al. (2003), and in fluid-structure interaction more broadly see e.g. Tezduyar et al. (2006) and Miller (2015) Other applications include fluid-soil interaction analysis (Saetta and Vitaliani, 1992; Armero, 1999; Mikelić and Wheeler, 2013).

5.3.4 Strongly coupled

In the strongly coupled scheme, inter-field iterations are performed until a given tolerance for the unknowns of each field is reached. They converge to the solution of the monolithic scheme and are thus able to satisfy discrete versions of the coupled problem exactly (Förster, 2007; Danowski, 2014). In principle, regarding thermomechanics, either the isothermal or the adiabatic slit can be used, but there seems to be no example of the latter. In contrast to the staggered schemes, there is no problem of conditional stability, but the scheme may converge very slowly or not at all. As an example coming from fluid-structure interactions, it has been shown that the number of coupling iterations increases when the time step decreases or when the structure becomes more flexible (Degroote et al., 2008). This can place a severe restriction on the use of these schemes. Several acceleration techniques are available in the literature to speed up convergence.

A straightforward way to improve the convergence behavior of the strongly coupled schemes is using predictors, in contrast to the values found in the last step. Thus, the initial guesses can be improved using well-chosen predictors Michler (2005). Along these lines, Erbts and Düster (2012) employs polynomial prediction methods, and Wendt et al. (2015) uses a line extrapolation method to improve the first guess of the unknown and thus decrease the number of iterations needed to achieve convergence.

Another approach that is well established for series acceleration is the Aitken delta-squared process. It uses previously computed values to obtain more accurate estimates for the unknown. Irons and Tuck (1969) is an early contribution detailing this low-memory convergence acceleration scheme. In the context of thermomechanics, Danowski (2014), Erbts et al. (2015) and Wendt et al. (2015) use this technique, with the last authors also employing a quasi-Newton least squares method. Some examples of contributions in the domain of fluid-structure interactions taking advantage of this approach are Degroote et al. (2008), Küttler and Wall (2008) and Küttler and Wall (2009). The last authors also introduce a vector extrapolation approach that includes more than three previous values of the iteration scheme in the improved estimate.

The strongly coupled approach lends itself to an interpretation as a nonlinear block Jacobi or Gauss-Seidel scheme, whose convergence is conditional and at most linear (Matthies and Steindorf, 2003b; Joosten et al., 2009). Cervera et al. (1996) provides an in-depth analysis of block Jacobi and Gauss-Seidel schemes applied to coupled problems, including considerations regarding efficiency, complexity, and parallelization. Matthies and Steindorf (2003a,b) suggests a block-Newton method instead, with the Jacobian of the system being approximate by a finite difference method. Under some assumptions on the subsystem solvers, this approach converges quadratically. Michler et al. (2005) propose a solution method based on the conjugation of sub-iterations via a Newton-Krylov method, which confines the GMRES acceleration to the interface degrees-of-freedom. The latter renders storage requirements for the Krylov space and computational cost of the least-squares problem low. The nesting of Newton and GMRES iterations lends itself to the reuse of Krylov vectors in subsequent linear system solutions. Küttler and Wall (2009) claims that the approach proposed by the last authors should not be regarded as a Newton-based solver but as a Krylov-based vector extrapolation scheme

One can also improve the convergence speed of the strongly coupled scheme using reduced-order models to produce a more accurate first guess and thus decrease the

number of iterations needed for the method to converge. [Vierendeels et al. \(2007\)](#) presents a technique that uses the Jacobian from reduced-order models that are built up during the coupling iterations. The reduced-order model is built for each step and approximates an arbitrary interface displacement fitting a linear regression to the previous displacement-stress points. [Degroote et al. \(2008\)](#) follows the same technique, coupling it with an Aitken delta-squared process.

[Blom \(2017\)](#) proposes a manifold mapping technique to decrease the number of sub-iterations of a high-fidelity fluid-structure interaction model. The idea is to perform many sub-iterations with a low-fidelity model instead of the high-fidelity flow and structure models.

Usage examples Regarding the use of strongly coupled schemes in the context of thermo-mechanics, there are a few contributions. [Erbts and Düster \(2012\)](#) present results concerning thermo-elasticity at finite strains, [Netz \(2013\)](#) concerning thermo-viscoelasticity, [Danowski \(2014\)](#) includes results on thermo-elasticity and thermo-elasto-plasticity. In field of fluid-structure interaction, a few examples of the use of strongly coupled schemes are [Torii et al. \(2006\)](#), [Wall et al. \(2007\)](#) and [Blom \(2017\)](#). Including more than two fields, [Erbts et al. \(2015\)](#) tackles electro-thermo-mechanical problems, as does [Wendt et al. \(2015\)](#), which also considers radiative heat transfer. In [Lenarda and Paggi \(2016\)](#), the strongly coupled scheme is used to solve coupled hygro-thermo-mechanical problems in photovoltaic laminates.

5.4 Comparison of solution techniques

According to [Felippa and Geers \(1988\)](#), the desirable properties of a time-stepping algorithm for solving coupled problems are:

- enjoys unconditional stability;
- is highly accurate;
- is easy to implement;
- is not memory intensive;
- requires low CPU time;
- satisfies software modularity constraints.

In the following, the time-stepping schemes presented above are compared with these criteria in mind. The application in view is thermomechanics.

Stability Regarding stability, the loosely coupled using an isothermal split is conditionally stable ([Armero and Simo, 1992](#)). Despite this, the limitation is not significant for metals plasticity, according to [Simo and Miehe \(1992\)](#). However, examples where the scheme diverges, can be found in [Armero and Simo \(1992\)](#). In this last contribution, the adiabatic split is introduced and shown to be unconditionally stable in the context of thermo-elasticity. [Armero and Simo \(1993\)](#) show that these properties extend to thermo-plasticity. The strongly coupled schemes are unconditionally stable because no critical time step leads to numerical instabilities in

the results. Despite this, the inner loop of the scheme may converge slowly or not at all [Matthies and Steindorf \(2003b\)](#). It depends on the spectral radius of the matrices involved ([Cervera et al., 1996](#)). There are, however, acceleration techniques that can mitigate this problem, including predictors and Aitken Δ^2 methods (see Section ??). [Danowski \(2014\)](#) presents a numerical example concerning an internal pressurized thick-walled cylinder, whose material is viscoplastic, for which the strongly coupled scheme employed diverged, despite the use of an Aitken method. On the other hand, the monolithic scheme, as long as appropriately preconditioned, is unconditionally stable ([Danowski, 2014](#)).

Accuracy Regarding accuracy, the solution found from the loosely coupled method will never exactly satisfy the fully coupled discretized equations of the problem. There will be a time lag between the thermal and the mechanical field. Loosely-coupled methods, therefore, satisfy conservation only in an asymptotic sense, i.e., for vanishing mesh width ([Michler, 2005](#)). As long as it does not diverge, the monolithic and strongly coupled satisfy the coupled discretized equations exactly.

Ease of implementation The partitioned schemes are much easier to implement as most of them can work with the field analyzers as black boxes, concerning themselves only with communication between the solvers, initial guesses, and acceleration schemes using previously computed values. The monolithic scheme requires the computation of the full stiffness matrix, including the mixed terms and appropriate preconditioning that varies widely with the specific multi-physics problem to be solved.

Memory requirements When it comes to memory requirements, the partitioned schemes often require only the diagonal blocks of the stiffness matrix found in the linearization process. Previous values also need to be saved from one iteration to the next, increasing the memory cost for some acceleration techniques. In contrast, the fully coupled monolithic scheme requires the full stiffness matrix of the coupled problem.

CPU time According to [Michler \(2005\)](#), solving a fluid-structure interaction problem with the same accuracy using a loosely and strongly coupled scheme, the latter is more efficient than the former. For the same total number of iterations, the difference in the accuracy reached ranges from one to three orders of magnitude. These results run counter to a claim in [Felippa et al. \(2001\)](#). However, this is not supported by any numerical results from the last authors. In the numerical examples presented in [Danowski \(2014\)](#), the monolithic solver is in most cases faster than a strongly coupled scheme employing an Aitken method for problems in thermomechanics. The differences range from 120% to 140% in favor of the monolithic scheme. Supporting evidence for these conclusions can also be found in [Novascone et al. \(2015\)](#). The authors report CPU time ratios between the strongly coupled and monolithic approaches, ranging from 0.635 to 3.75 on the magnitude of the coupling.

Software modularity The partitioned approaches can take full advantage of software, including closed source commercial solvers. There is little to no software reuse for the monolithic approach, save for routines that solve linear systems and the like.

Conclusions Lastly, it may be helpful to reproduce the recommendations given in [Felippa et al. \(2001\)](#) regarding the choice between partitioned and monolithic approaches. According to the authors, the circumstances that favor the partitioned approach for tackling a coupled problem are a research environment with few delivery constraints, access to existing software, localized interaction effects (e.g., surface versus volume), and widespread spatial/temporal component characteristics. The opposite circumstances: commercial environment, rigid deliverable timetable, massive software development resources, global interaction effects, and comparable length/time scales favors a monolithic approach.

Putting it all together, the most appropriate choice for the present use case is the strongly coupled schemes with appropriate acceleration techniques. They can take advantage of already existing software, provide accurate results that agree with a monolithic approach, are not memory intensive, are easy to implement, and with the use of convergence acceleration techniques, are competitive from the computational efficiency standpoint. The only drawback seems to be the possibility of divergence in the inner loop, stalling the progress of the simulation.

Table 5.1: Summary of the comparison between the FFT-Galerkin method.

	Partitioned schemes		Monolithic
	Loosely coupled	Strongly coupled	
Stability	Isothermic split: conditionally stable Adiabatic split: unconditionally stable	unconditionally stable*	unconditionally stable
Accuracy	Coupled discretized equations not satisfied exactly	Coupled discretized equations satisfied	Coupled discretized equations satisfied
Ease of implementation	Only communication between field analyzers stricly needed	Full coupling needed: • Computation of mixed terms of the Jacobian • Preconditioning needed	
Memory requirements	Only diagonal blocks of the full stiffness matrix needed	Full stiffness matrix needed	
Software modularity constraints	Full software modularity	Poor or no software modularity	

* The inner loop of the strongly coupled scheme may converge very slowly or even diverge.

Chapter 6

Strongly coupled methods for coupled fields

This chapter presents the most common strongly coupled/implicit methods employed to solve coupled field problems. This presentation seeks to provide a literature overview of the available approaches.

6.1 Equations to be solved

For the sake of clarity, the discretized equations of the thermo-mechanical problem at the next time instant, $n + 1$ are recovered here

$$\mathbf{M}\ddot{\mathbf{u}}_{n+1} + \mathbf{f}_u^{\text{int}}(\boldsymbol{\theta}_{n+1}, \mathbf{u}_{n+1}) - \mathbf{f}_{u,n+1}^{\text{ext}} = \mathbf{0}, \quad (6.1)$$

$$\mathbf{C}\dot{\boldsymbol{\theta}}_{n+1} + \mathbf{f}_\theta^{\text{int}}(\boldsymbol{\theta}_{n+1}, \mathbf{u}_{n+1}) - \mathbf{f}_{\theta,n+1}^{\text{ext}} = \mathbf{0}. \quad (6.2)$$

The complete definition of the material incremental discretized thermo-mechanical initial boundary value problem can be found in Chapter 4.

As only partitioned approaches are considered, the thermal and mechanical problems are solved separately, i.e., Equation (6.1) is solved considering a fixed temperature, and Equation (6.2) is solved assuming a fixed configuration. To ease the discussion, consider the existence of two functions \mathcal{U}_{n+1} and \mathcal{T}_{n+1} that represent these solution procedures at timestep $n + 1$. These so-called mechanical and thermal solvers satisfy

$$\mathcal{U}: \mathcal{K}_{\theta,n+1} \rightarrow \mathcal{K}_{u,n+1}, \quad \mathbf{u} = \mathcal{U}_{n+1}(\boldsymbol{\theta}), \quad (6.3)$$

$$\mathcal{T}: \mathcal{K}_{u,n+1} \rightarrow \mathcal{K}_{\theta,n+1}, \quad \boldsymbol{\theta} = \mathcal{T}_{n+1}(\mathbf{u}). \quad (6.4)$$

See Chapter 4 for detailed information on them. In the following, the subscripts on the solvers will be dropped to avoid clutter.

The goal now is to consider functions, built from \mathcal{U} and \mathcal{T} , whose roots are also the solutions to the thermo-mechanical problem (Equations (6.1) and (6.2)). Several examples can be provided. The most appropriate for the current use case are presented in what follows. They can be found in Uekermann (2016) in the context of fluid-structure interaction (FSI).

Consider the residues defined as,

$$\mathcal{R}_J: \mathcal{H}_{u,n+1} \times \mathcal{H}_{\theta,n+1} \rightarrow \mathcal{K}_{u,n+1} \times \mathcal{K}_{\theta,n+1}, \quad \mathcal{R}_J(\mathbf{u}, \boldsymbol{\theta}) = \begin{Bmatrix} \mathbf{u} - \mathcal{U}(\boldsymbol{\theta}) \\ \boldsymbol{\theta} - \mathcal{T}(\mathbf{u}) \end{Bmatrix}, \quad (6.5)$$

and

$$\mathcal{R}_{GS}: \mathcal{H}_{\theta,n+1} \rightarrow \mathcal{H}_{\theta,n+1}, \quad \mathcal{R}_{GS}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \mathcal{T} \circ \mathcal{U}(\boldsymbol{\theta}), \quad (6.6)$$

or

$$\mathcal{R}_{GS}^*: \mathcal{H}_{u,n+1} \rightarrow \mathcal{H}_{u,n+1}, \quad \mathcal{R}_{GS}^*(\mathbf{u}) = \mathbf{u} - \mathcal{U} \circ \mathcal{T}(\mathbf{u}), \quad (6.7)$$

where the subscript "J" stands for Jacobi and the subscript "GS" for Gauss-Seidel. The reason for this choice of subscripts is made clear in Section 6.3.1.

Since the methods described below for the solution of nonlinear systems of equations apply to both functions \mathcal{R}_J and \mathcal{R}_{GS} , a general function denoted as \mathcal{R} , whose variable is \mathbf{x} , is considered instead. As already stated, the solution for the thermo-mechanical problem (Equations (6.1) and (6.2)) can be abstracted as the solution of

$$\mathcal{R}(\mathbf{x}) = 0. \quad (6.8)$$

To obtain simpler expressions in what follows, consider also the function

$$\mathcal{S}(\mathbf{x}) = \mathbf{x} - \mathcal{R}(\mathbf{x}), \quad (6.9)$$

whose fixed point is the solution to the nonlinear system of equation in Equation (6.8).

6.2 A classification scheme for iterative methods

Most methods available for the solution of systems of nonlinear equations, such as the one in Equation (6.8), are iterative methods. They can be more precisely defined letting $\mathbf{x}^k, \mathbf{x}^{k-1}, \dots$, whose superscripts correspond to the loop of the iteration method, be approximants to \mathbf{x}_{n+1} , whose subscript concerns the timestep

To better understand the landscape of available methods to solve nonlinear systems of equations, the iteration functions are classified according to the information they require following the classification scheme by Traub (1982). Let \mathbf{x}^{k+1} be determined uniquely by information obtained at $\mathbf{x}^k, \mathbf{x}^{k-1}, \dots$, including the derivatives of any order of \mathcal{R} . Let the function that maps $\mathbf{x}^k, \mathbf{x}^{k-1}, \dots$ into \mathbf{x}^{k+1} be called ϕ . Thus

$$\mathbf{x}^{k+1} = \phi(\mathbf{x}^k, \mathcal{R}(\mathbf{x}^k), J_{\mathcal{R}}(\mathbf{x}^k), \dots), \quad (6.10)$$

where ϕ is called an iteration function, and $J_{\mathcal{R}}$ is the Jacobian of \mathcal{R} . To prevent cluttering \mathbf{x}^k will stand for its value as well as for the values of $\mathcal{R}(\mathbf{x}^k)$, $J_{\mathcal{R}}(\mathbf{x}^k)$ and further derivatives of higher order. Then ϕ is called a *one-point iteration function*. Most iteration functions that have been used for root-finding are one-point iteration functions. The most commonly known examples are the fixed point schemes and Newton's iteration method.

Next, let \mathbf{x}^{k+1} be determined by new information at \mathbf{x}^k and reused information at \mathbf{x}^{k-1}, \dots . Thus

$$\mathbf{x}^{k+1} = \phi(\mathbf{x}^k; \mathbf{x}^{k-1}, \dots). \quad (6.11)$$

Then ϕ is called a *one-point iteration function with memory*. The semicolon in Equation (6.11) separates the point at which new data are used from the points at which old data are reused. The secant iteration function is the best-known example of a one-point iteration function with memory.

Let \mathbf{x}^{k+1} be determined by new information at $\mathbf{x}^k, \omega_1(\mathbf{x}^k), \dots, \omega_i(\mathbf{x}^k)$, $i \geq 1$, where ω_i denote operations on \mathbf{x}^k . No old information is reused. Thus

$$\mathbf{x}^{k+1} = \phi\left[\mathbf{x}^k, \omega_1(\mathbf{x}^k), \dots, \omega_i(\mathbf{x}^k)\right]. \quad (6.12)$$

Then ϕ is called a *multipoint iterative function*. Such methods include the Aitken-Steffson method.

Finally, let \mathbf{z}_j represent the quantities $\mathbf{x}^j, \omega_1(\mathbf{x}^j), \dots, \omega_i(\mathbf{x}^j)$, $i \geq 1$. Let

$$\mathbf{x}^{k+1} = \phi(\mathbf{z}^k; \mathbf{z}^{k-1}, \dots). \quad (6.13)$$

Then ϕ is called a *multipoint iterative function with memory*. The semicolon in Equation (6.13) separates the points at which new data are used from the points at which old data are reused.

In the present work, the criteria used for the choice of the iterative method used fit roughly into the ones provided by Fang and Saad (2009) for problems in the context the electronic structure problems. They are

1. The dimensionality of the problem is large.
2. $\mathcal{R}(\mathbf{x})$ is continuously differentiable, but the analytic form of its derivative is not readily available, or it is costly to compute.
3. The cost of evaluating $\mathcal{R}(\mathbf{x})$ is computationally high.
4. The problem is noisy. In other words, the evaluated function values of $\mathcal{R}(\mathbf{x})$ usually contain errors.

Thus, the methods chosen must minimize the number of calls to \mathcal{R} , as it is expensive to compute. The amount of information saved from previous iterations must also be judiciously chosen as the problem's dimensionality is large, leading to memory limitations. Finally, the analytical form of the derivative \mathcal{R} is also not available. Thus methods that use it must be discarded.

6.2.1 Predictor

Iterative procedures are considered to solve the thermo-mechanical problem at a given timestep $n + 1$. As the first value approximating \mathbf{x}_{n+1} , one can employ the converged value of the previous timestep, \mathbf{x}_n . However, a very efficient way to increase the chances of stability and reduce computation time is to predict the optimal initial values at the beginning of every time step (Erbs and Düster, 2012; Erbs et al., 2015; Wendt et al.,

2015). The prediction of the new solution by polynomial extrapolation is based on the converged solution of the last two or three timesteps. This method is based on polynomial vector extrapolation, which is relatively easy to implement, and the extra computational input is negligible.

The maximum polynomial under consideration is of the order two, i.e., the new solution is extrapolated from the results from the last three time steps. The predictors \mathbf{x}^* for the order $p = 1$ and $p = 2$ polynomials read:

$$p = 1: \quad \mathbf{x}_{n+1}^* = 2\mathbf{x}_n - \mathbf{x}_{n-1}, \quad (6.14)$$

$$p = 2: \quad \mathbf{x}_{n+1}^* = 3\mathbf{x}_n - 3\mathbf{x}_{n-1} + \mathbf{x}_{n-2}. \quad (6.15)$$

6.2.2 Global Approaches

Following [Dennis and Schnabel \(1996\)](#), the terms "global," as in "global method" or "globally convergent algorithm," are here used to denote a method that is designed to converge to a local minimizer of a nonlinear functional or some solution of a system of nonlinear equations, from almost any starting point. The methods presented in this chapter do not qualify as global methods since if the initial trial is not close enough to the solution, they will not converge. There are, however, approaches to mitigate this problem. The ideas presented below apply with particular relevance to the Newton method (see Section 6.3.2) and related procedures. Their exposition follows [Dennis and Schnabel \(1996\)](#) where more details can be found.

Consider that the iterative solution method determines $\Delta\mathbf{x}^k$ in

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k. \quad (6.16)$$

The two global approaches here considered both come into action after $\Delta\mathbf{x}^k$ has been computed by some appropriate method (see from Section 6.3 on). At this point, one decides whether to accept the step $\Delta\mathbf{x}^k$ or to choose \mathbf{x}^{k+1} by a global strategy.

A solution to the system of equations (6.8) clearly also satisfies

$$r(\mathbf{x}) = 0, \quad \text{where } r \equiv \frac{1}{2} \|\mathcal{R}\|_2^2: \mathbb{R}^n \rightarrow \mathbb{R}, \quad (6.17)$$

so the problem can be regarded as an unconstrained minimization problem, with caveat that local minimizers of r may not be the solution to the system of equations (6.8).

The basic idea of a global method for unconstrained minimization is geometrically obvious: take steps that lead "downhill" for the function r . More precisely, one chooses a direction \mathbf{p} from the current point \mathbf{x}^k in which r decreases initially, and a new point \mathbf{x}^{k+1} in this direction from \mathbf{x}^k such that $r(\mathbf{x}^{k+1}) < r(\mathbf{x}^k)$. Such a direction is called a descent direction.

An important question to ask is, "What is a descent direction for problem (6.17)?" It is any direction \mathbf{p} for which $\nabla r(\mathbf{x}^k)^T \mathbf{p} < 0$, where

$$\nabla r(\mathbf{x}^k) = J_{\mathcal{R}}(\mathbf{x}^k)^T \mathcal{R}(\mathbf{x}^k), \quad (6.18)$$

where $J_{\mathcal{R}}(\mathbf{x}^k)$ is the Jacobian matrix of \mathcal{R} at \mathbf{x}^k . Therefore, the steepest-descent direction for (6.17) is along $-J_{\mathcal{R}}(\mathbf{x}^k)^T \mathcal{R}(\mathbf{x}^k)$.

The Newton step for the update equation (6.16) is (see Section 6.3.2)

$$\Delta \mathbf{x}_N^k = -J_{\mathcal{R}}(\mathbf{x}^k)^{-1} \mathcal{R}(\mathbf{x}^k), \quad (6.19)$$

and it is a descent direction, since

$$\nabla r(\mathbf{x}^k)^T \Delta \mathbf{x}_N^k = -\mathcal{R}(\mathbf{x}^k)^T J_{\mathcal{R}}(\mathbf{x}^k) J_{\mathcal{R}}(\mathbf{x}^k)^{-1} \mathcal{R}(\mathbf{x}^k) = -\mathcal{R}(\mathbf{x}^k)^T \mathcal{R}(\mathbf{x}^k) < 0 \quad (6.20)$$

as long as $\mathcal{R}(\mathbf{x}^k) \neq \mathbf{0}$. Hence, the appropriateness of these methods to the Newton method and related methods.

Since the Newton step yields a root of

$$M^k(\mathbf{x}^k + \Delta \mathbf{x}^k) = \mathcal{R}(\mathbf{x}^k) + J_{\mathcal{R}}(\mathbf{x}^k) \Delta \mathbf{x}^k, \quad (6.21)$$

it also goes to a minimum of the quadratic function

$$\begin{aligned} \hat{m}^k(\mathbf{x}^k + \Delta \mathbf{x}^k) &\equiv \frac{1}{2} M^k(\mathbf{x}^k + \Delta \mathbf{x}^k)^T M^k(\mathbf{x}^k + \Delta \mathbf{x}^k) \\ &= \frac{1}{2} \mathcal{R}(\mathbf{x}^k)^T \mathcal{R}(\mathbf{x}^k) + \left(J_{\mathcal{R}}(\mathbf{x}^k)^T \mathcal{R}(\mathbf{x}^k) \right)^T \Delta \mathbf{x}^k \\ &\quad + \frac{1}{2} \Delta \mathbf{x}^{kT} \left(J_{\mathcal{R}}(\mathbf{x}^k)^T J_{\mathcal{R}}(\mathbf{x}^k) \right) \Delta \mathbf{x}^k, \end{aligned} \quad (6.22)$$

because $\hat{m}^k(\mathbf{x}^k + \Delta \mathbf{x}^k) \geq 0$ for all $\Delta \mathbf{x}^k$ and $\hat{m}^k(\mathbf{x}^k + \Delta \mathbf{x}_N^k) = 0$. Therefore, $\Delta \mathbf{x}_N^k$ is a descent direction for \hat{m}^k , and since the gradients at \mathbf{x}^k of \hat{m}^k and r are the same, it is also a descent direction for r .

The above development motivates how the global methods to be described are applied, i.e., they are applied to the quadratic model $\hat{m}^k(\mathbf{x})$. Since $\nabla^2 \hat{m}^k(\mathbf{x}^k) = J_{\mathcal{R}}(\mathbf{x}^k)^T J_{\mathcal{R}}(\mathbf{x}^k)$, this model is positive definite as long as $J_{\mathcal{R}}(\mathbf{x}^k)$ is nonsingular, which is consistent with the fact that $\mathbf{x}^k + \Delta \mathbf{x}_N^k$ is the unique root of $M^k(\mathbf{x})$ and thus the unique minimizer of $\hat{m}^k(\mathbf{x})$ in this case. Thus, the model $\hat{m}^k(\mathbf{x})$ has the attractive properties that its minimizer is the Newton point for the original problem, and that all its descent directions are descent directions for $r(\mathbf{x})$ because $\nabla \hat{m}^k(\mathbf{x}^k) = \nabla r(\mathbf{x}^k)$. Therefore methods based on this model, by going downhill and trying to minimize $\hat{m}^k(\mathbf{x})$, will combine Newton's method for nonlinear equations with global methods for an associated minimization problem.

If the Jacobian of \mathcal{R} is not available and its estimate is of poor quality, the global procedure may be compromised (Kelley, 2003). However, these procedures may be unnecessary in the present use case since the initial trial is probably close enough to the solution, even without accounting for the improvements coming from more carefully chosen initial shots through predictors (see Section 6.2.1). Fang and Saad (2009) also employ a simple restarting procedure instead of a global convergence strategy. If in two consecutive values of \mathcal{R} , \mathcal{R}_{old} and \mathcal{R}_{new} , $\|\mathcal{R}_{\text{new}}\|$ is much larger than $\|\mathcal{R}_{\text{old}}\|$, the solution procedure is restarted, with the new initial trial values corresponding to \mathcal{R}_{old} . They suggest r between 0.1 and 0.3, with $\|\mathcal{R}_{\text{old}}\| < r \|\mathcal{R}_{\text{new}}\|$ leading to a restart. In their opinion, global approaches such as those suggested below are too expensive when the evaluation of \mathcal{R} is also costly to compute.

Line search The line search approach is based on the traditional idea of backtracking along the Newton direction if a complete Newton step is unsatisfactory. More precisely given a descent direction \mathbf{p}^k , a step in that direction is taken as

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \lambda^k \mathbf{p}^k, \quad (6.23)$$

for some $\lambda^k > 0$ that makes \mathbf{x}^{k+1} an acceptable iterate. The common procedure is to first try $\lambda_k = 1$ and only if this fails backtrack in a systematic way along the direction defined by that step. See [Dennis and Schnabel \(1996\)](#) for a full discussion on the choice of λ_k .

Trust region algorithms The trust region algorithm is based on estimating the region in which the local model, underlying Newton's method, can be trusted to represent the function adequately and taking a step to approximately minimize the model in this region. It drops the assumption that the step must be in the Newton direction. $\hat{m}^k(\mathbf{x}^k + \Delta\mathbf{x}^k)$ is approximately minimized subject to $\|\mathbf{x}^k\|_2 \leq \delta^k$. If $\delta^k \geq \|J_{\mathcal{R}}(\mathbf{x}^k)^{-1} \mathcal{R}(\mathbf{x}^k)\|_2$, then the step attempted is the Newton step. Otherwise, for the locally constrained optimal step, it is

$$\Delta\mathbf{x}^k = -\left(J_{\mathcal{R}}(\mathbf{x}^k)^T J_{\mathcal{R}}(\mathbf{x}^k) + \mu^k \mathbf{I}\right)^{-1} J_{\mathcal{R}}(\mathbf{x}^k)^T \mathcal{R}(\mathbf{x}^k), \quad (6.24)$$

for μ^k such that $\|\Delta\mathbf{x}^k\|_2 \cong \delta^k$. For the details on the choice of δ^k see [Dennis and Schnabel \(1996\)](#).

6.2.3 Convergence criteria

For an iterative method to be useful, there must be reasonable criteria to determine its convergence. The iteration residual is defined as

$$\mathbf{r}^k = \mathcal{R}(\mathbf{x}^k), \quad (6.25)$$

and if it is equal to zero then \mathbf{x} is the solution to the system of nonlinear equations, i.e.,

$$\mathbf{r} = \mathcal{R}(\mathbf{x}) = \mathbf{0}, \quad (6.26)$$

and hence, a reasonable convergence measure for the iteration procedure.

The discrete l^2 -norm can be used to obtain a scalar representative of the vectorial residual $\mathbf{r}^k = (r^{k,1}, \dots, r^{k,n_{\text{unknown}}})^T$ as

$$\|\mathbf{r}^k\|_{l^2} = \sqrt{\sum_i (r^{k,i})^2}. \quad (6.27)$$

Directly using (6.27) yields an absolute convergence criterion

$$\|\mathbf{r}^k\|_{l^2} < \epsilon_{\text{abs}}. \quad (6.28)$$

with $\epsilon_{\text{abs}} > 0$ as an absolute convergence tolerance, with convergence being achieved when the above condition is satisfied. However, since the absolute value of the $r^{k,i}$'s can change by orders of magnitude during one simulation, an absolute measure is

not appropriate in all situations. A relative measure solves this problem by setting the residual in relation with the current coupling iterate values as

$$\frac{\|\mathbf{r}^k\|_{l^2}}{\|\mathbf{x}^k\|_{l^2}} < \epsilon_{\text{rel}}. \quad (6.29)$$

A relative convergence measure can fail to work correctly when the coupling iterate values are close to zero, and rounding errors occur. Thus, a combination of absolute and relative measures, where the absolute measure takes care of close to zero cases, and the relative handles all other cases, is often a good choice.

6.3 One-point iteration function

6.3.1 Fixed-point approaches

The application of the fixed-point method to obtain the roots of \mathcal{R} yields

$$\mathbf{x}^{k+1} = \mathcal{S}(\mathbf{x}^k) = \mathbf{x}^k - \mathcal{R}(\mathbf{x}^k). \quad (6.30)$$

See Figure 6.1 for its geometric interpretation in one dimension.

If the particular functions defined on Equations (6.5) and (6.6) are used, one finds the two basic Schwarz procedures commonly employed in strongly coupled solution procedures. They are the additive or block Jacobi and the parallel Scharwz or Gauss-Seidel procedures. The names originate from domain decomposition, and justify the subscripts employed in Equations (6.5) and (6.6).

6.3.1.1 Block Jacobi or Schwarz additive

Applying the fixed-point approach to \mathcal{R}_J (Equation (6.5)), yields

$$\{\mathbf{u}^{k+1}, \boldsymbol{\theta}^{k+1}\}^T = \mathcal{S}_J(\mathbf{u}^k, \boldsymbol{\theta}^k) \quad (6.31)$$

$$= \{\mathbf{u}^k, \boldsymbol{\theta}^k\}^T - \mathcal{R}_J(\mathbf{u}^k, \boldsymbol{\theta}^k), \quad (6.32)$$

It is the same as solving both the mechanical (Equation (6.1)) and the thermal problem (Equation (6.2)) in parallel. Such a procedure is said to be Schwarz additive or block Jacobi, referring to the similarities with the procedure for the solution of linear systems of equations with the same name i.e.,

$$\mathbf{u}^{k+1} = \mathcal{U}(\boldsymbol{\theta}^k), \quad (6.33)$$

$$\boldsymbol{\theta}^{k+1} = \mathcal{T}(\mathbf{u}^k). \quad (6.34)$$

Box 6.1 shows the pseudo-code for the block Jacobi approach.

6.3.1.2 Block Gauss-Seidel or Schwarz multiplicative

Applying the fixed-point approach to \mathcal{R}_{GS} (Equation (6.5)), yields

$$\boldsymbol{\theta}^{k+1} = \mathcal{S}_{\text{GS}}(\boldsymbol{\theta}^k) = \boldsymbol{\theta}^k - \mathcal{R}_{\text{GS}}(\boldsymbol{\theta}^k). \quad (6.35)$$

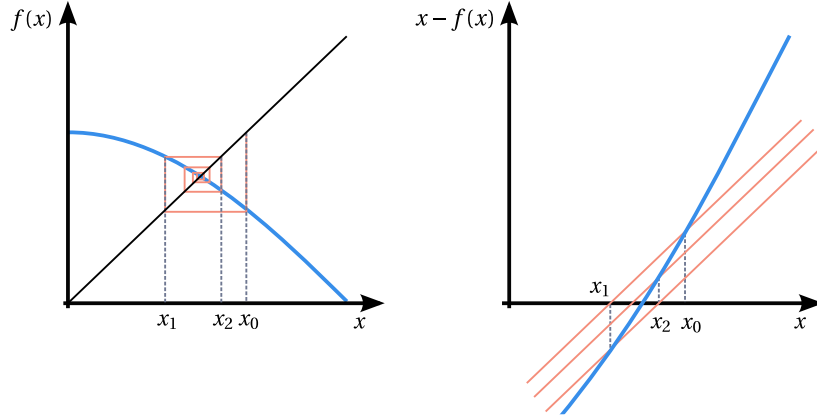


Figure 6.1: Geometric interpretation of the fixed-point iteration method in one dimension. The fixed-point of f is sought, which is equivalent to the root of $x - f(x)$.

Thus, the fields are solved sequentially, where the output of the first solver is the input of the second solver. This the solution procedure is said to be Scharwz multiplicative or block Gauss-Seidel.

$$\mathbf{u}^{k+1} = \mathcal{U}(\boldsymbol{\theta}^k), \quad (6.36)$$

$$\boldsymbol{\theta}^{k+1} = \mathcal{T}(\mathbf{u}^{k+1}). \quad (6.37)$$

One of the fields must be chosen as the first, and this may be crucial for the stability and convergence rate of the approach (Joosten et al., 2009). Here, the focus is on the sequence coinciding with the isothermic split, i.e., first, the mechanical problem is solved at a fixed temperature. Then the thermal problem is solved at a fixed configuration.

Box 6.2 shows the pseudo-code for the block Gauss-Seidel approach.

6.3.2 Newton's method

The Newton-Raphson or Newton scheme is a very popular iterative solution procedure for nonlinear systems of equations, which under appropriate conditions converges quadratically (Dennis and Schnabel, 1996; Kelley, 2003). It can be applied to Equation (6.8) yielding

$$J_{\mathcal{R}}(\mathbf{x}^k) \Delta \mathbf{x}^k = -\mathcal{R}(\mathbf{x}^k), \quad (6.38)$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k. \quad (6.39)$$

See Figure 6.2 for its geometric interpretation in one dimension.

In particular, using \mathcal{R}_J , a few simplifications can be obtained. To ease the explanation, consider, a thermal residual operator $\mathcal{R}_u(\mathbf{u}, \boldsymbol{\theta})$ and a mechanical residual operator $\mathcal{R}_\theta(\mathbf{u}, \boldsymbol{\theta})$ defined to be the first and second components in the definition of

Box 6.1: Additive Schwarz procedure, also called block Jacobi, for one timestep.

- (i) $\mathbf{u}^0 = \mathbf{u}_n$
- (ii) $\theta^0 = \theta_n$
- (iii) Set fixed-point counter to zero: $k = 0$
- (iv) Enter the fixed-point loop
 - (1) Solve the mechanical problem at fixed temperature θ^k : $\mathbf{u}^{k+1} = \mathcal{U}(\theta^k)$
 - (2) Solve the thermal problem at a fixed configuration \mathbf{u}^k : $\theta^{k+1} = \mathcal{T}(\mathbf{u}^k)$
 - (3) If the desired accuracy has not been reached, update $k = k + 1$ and go to step (1).

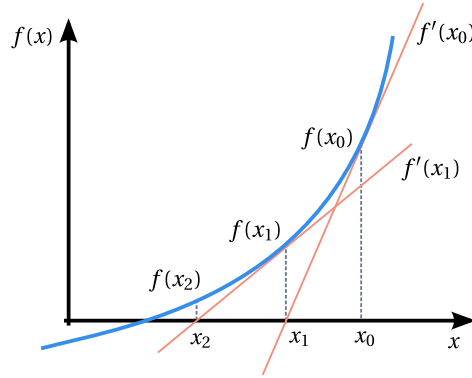


Figure 6.2: Geometric interpretation of the Newton method in one dimension for an example function f , whose derivative is denoted by f' .

\mathcal{R}_J (Equation (6.6)). Written in full

$$\mathcal{R}_u(\mathbf{u}, \theta) = \mathbf{u} - \mathcal{U}(\theta) = 0, \quad (6.40)$$

$$\mathcal{R}_\theta(\mathbf{u}, \theta) = \theta - \mathcal{T}(\mathbf{u}) = 0, \quad (6.41)$$

From this, a block Newton iteration can be written as

$$\begin{bmatrix} J_{\mathcal{R}_u}(\mathbf{u}^k, \theta^k) \\ J_{\mathcal{R}_\theta}(\mathbf{u}^k, \theta^k) \end{bmatrix} \begin{Bmatrix} \Delta \mathbf{u}^k \\ \Delta \theta^k \end{Bmatrix} = - \begin{Bmatrix} \mathcal{R}_u(\mathbf{u}^k, \theta^k) \\ \mathcal{R}_\theta(\mathbf{u}^k, \theta^k) \end{Bmatrix}, \quad (6.42)$$

and the update of the iteration variables reads

Box 6.2: Multiplicative Schwarz procedure, also called block Gauss-Seidel, for one timestep.

- (i) $\mathbf{u}^0 = \mathbf{u}_n$
- (ii) $\theta^0 = \theta_n$
- (iii) Set fixed-point counter to zero: $k = 0$
- (iv) Enter the fixed-point loop
 - (1) Solve the mechanical problem at fixed temperature θ^k : $\mathbf{u}^{k+1} = \mathcal{U}(\theta^k)$
 - (2) Solve the thermal problem at a fixed configuration \mathbf{u}^{k+1} : $\theta^{k+1} = \mathcal{T}(\mathbf{u}^{k+1})$
 - (3) If the desired accuracy has not been reached, update $k = k + 1$ and go to step (1).

$$\begin{Bmatrix} \mathbf{u}^{k+1} \\ \theta^{k+1} \end{Bmatrix} = \begin{Bmatrix} \mathbf{u}^k \\ \theta^k \end{Bmatrix} + \begin{Bmatrix} \Delta \mathbf{u}^k \\ \Delta \theta^k \end{Bmatrix}. \quad (6.43)$$

The system of equations in Equation (6.42) can be further simplified following Degroote (2010) considering the definitions of the mechanical and thermal residuals and taking their derivatives. It yields

$$\begin{bmatrix} \mathbf{I} & -J_{\mathcal{U}}(\theta^k) \\ -J_{\mathcal{T}}(\mathbf{u}^k) & \mathbf{I} \end{bmatrix} \begin{Bmatrix} \Delta \mathbf{u}^k \\ \Delta \theta^k \end{Bmatrix} = - \begin{Bmatrix} \mathcal{R}_u(\mathbf{u}^k, \theta^k) \\ \mathcal{R}_\theta(\mathbf{u}^k, \theta^k) \end{Bmatrix}, \quad (6.44)$$

Solving for $\Delta \mathbf{u}^k$ and $\Delta \theta^k$, one finds

$$(\mathbf{I} + J_{\mathcal{U}}(\theta^k) J_{\mathcal{T}}(\mathbf{u}^k)) \Delta \mathbf{u}^k = -\mathcal{R}_u(\mathbf{u}^k, \theta^k) + J_{\mathcal{U}}(\theta^k) \mathcal{R}_\theta(\mathbf{u}^k, \theta^k), \quad (6.45)$$

$$(\mathbf{I} + J_{\mathcal{T}}(\mathbf{u}^k) J_{\mathcal{U}}(\theta^k)) \Delta \theta^k = -\mathcal{R}_\theta(\mathbf{u}^k, \theta^k) + J_{\mathcal{T}}(\mathbf{u}^k) \mathcal{R}_u(\mathbf{u}^k, \theta^k). \quad (6.46)$$

Thus, the Jacobians now needed are $J_{\mathcal{U}}$ and $J_{\mathcal{T}}$. See Section 6.4.1 for the practical application of this.

Every iteration of the Newton scheme involves at least one invocation of the thermal and mechanical solvers when computing $\mathcal{R}(\mathbf{u}^k)$ or both $\mathcal{R}_u(\mathbf{u}^k, \theta^k)$ and $\mathcal{R}_\theta(\mathbf{u}^k, \theta^k)$. The critical point for black-box equation coupling is how to obtain the derivative information in the Jacobi matrices. In different ways, some of the methods presented next find approximations for the required Jacobian times vector products.

6.3.3 Constant Underrelaxation

One of the most straightforward ways to stabilize an iterative method is to use constant underrelaxation (Gatzhammer, 2014). The relaxation is performed as follows

$$\mathbf{x}^{k+1} = (1 - \omega)\mathbf{x}^k + \omega(\mathbf{x}^k - \mathcal{R}(\mathbf{x}^k)) = \mathbf{x}^k - \omega\mathcal{R}(\mathbf{x}^k), \quad (6.47)$$

where ω is the relaxation factor chosen in the range $0 < \omega < 1$, which corresponds to an underrelaxation, to achieve a stabilizing effect.

Applying to Equation (6.5)

$$\begin{Bmatrix} \mathbf{u}^{k+1} \\ \boldsymbol{\theta}^{k+1} \end{Bmatrix} = (1 - \omega) \begin{Bmatrix} \mathbf{u}^k \\ \boldsymbol{\theta}^k \end{Bmatrix} + \omega \begin{Bmatrix} \mathcal{U}(\boldsymbol{\theta}^k) \\ \mathcal{T}(\mathbf{u}^k) \end{Bmatrix} \quad (6.48)$$

Applying to Equation (6.6)

$$\boldsymbol{\theta}^{k+1} = (1 - \omega)\boldsymbol{\theta}^k + \omega\mathcal{T} \circ \mathcal{U}(\boldsymbol{\theta}^k). \quad (6.49)$$

Constant underrelaxation works well if ω is close to 1 but leads to a slow convergence if ω has to be chosen close to 0. Thus, the constant underrelaxation method creates unmanageable computational costs for severe instabilities. The optimal ω is not necessarily the largest stable one (Gatzhammer, 2014) and has to be set empirically. In what follows, alternative methods are discussed to decrease the number of iterations necessary while maintaining stability.

Box 6.3: Constant underrelaxation applied to the block Gauss-Seidel scheme.

- (i) $\boldsymbol{\theta}^0 = \boldsymbol{\theta}_{n+1}^p$
- (ii) Set fixed-point counter to zero: $k = 0$
- (iii) Enter the fixed-point loop
 - (1) Solve the mechanical problem at fixed temperature $\boldsymbol{\theta}^k$: $\mathbf{u}^{k+1} = \mathcal{U}(\boldsymbol{\theta}^k)$
 - (2) Solve the thermal problem at a fixed configuration \mathbf{u}^{k+1} : $\boldsymbol{\theta}^{k+1} = \mathcal{T}(\mathbf{u}^{k+1})$
 - (3) Compute $\boldsymbol{\theta}^{k+1}$ using constant relaxation (Equation (6.47))
 - (4) If the desired accuracy has not been reached, update $k = k + 1$ and go to step (1).

6.4 One-point iteration function with memory

6.4.0.1 Aitken relaxation

The so-called Aitken Δ^2 relaxation method was introduced by Irons and Tuck (1969) as a modified Aitken Δ^2 that does not require the computation of the function twice

per iteration as in the original method. It has been widely used in the context of FSI (Irons and Tuck, 1969; Küttler and Wall, 2008; Joosten et al., 2009; Küttler and Wall, 2009; Erbs et al., 2015; Wendt et al., 2015). It has also been used in the context of thermo-mechanics by Danowski et al. (2013).

In the one-dimensional case, this method resembles the secant method applied to the fixed point problem, which can be used to solve nonlinear equations without differentiation. Calling it an Aitken method is perhaps a misnomer since, in the Aitken-Steffensen method, the function values are computed twice per iteration (see Section 6.5.3). It is more closely related to secant methods, reusing values from previous iterations. This version of Aitken's Δ^2 method provides a dynamic under relaxation, which can be used to improve the convergence/stability properties of the coupling algorithm.

Assume that f is the function whose fixed point is sought. The linear interpolation between two points already known of the function, $(a, f(a))$ and $(b, f(b))$ is

$$y = \frac{f(b) - f(a)}{b - a}(x - a) + f(a). \quad (6.50)$$

The fixed point of this approximation is

$$c = \frac{f(b) - f(a)}{b - a}(c - a) + f(a). \quad (6.51)$$

Thus, after rearranging,

$$c = \frac{af(b) - bf(a)}{(a - f(a)) - (b - f(b))} \quad (6.52)$$

This can be rewritten as

$$c = (1 - \omega_b)b + \omega_b f(b) \quad \text{with } \omega_b = \frac{a - b}{(a - f(a)) - (b - f(b))} \quad (6.53)$$

Anticipating the next iteration step,

$$d = (1 - \omega_c)c + \omega_c f(c) \quad \text{with } \omega_c = \frac{c - b}{(b - f(b)) - (c - f(c))} \quad (6.54)$$

a convenient expression for updating the relaxation factor may be found, i.e.

$$\omega_c = -\omega_b \frac{f(b) - b}{(c - f(c)) - (f(b) - b)}. \quad (6.55)$$

See Figure 6.3 for its geometric interpretation in one dimension.

Now, for the vector case, the next step is to work out the solution to the current iteration from the outcome of the previous iteration \mathbf{x}^k plus a new increment $\Delta \mathbf{x}^k$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k. \quad (6.56)$$

The increment reads

$$\Delta \mathbf{x}^k = \omega^k \left(\mathcal{S}(\mathbf{x}^{(k)}) - \mathbf{x}^{(k)} \right) = -\omega^k \mathcal{R}(\mathbf{x}^k). \quad (6.57)$$

with ω^k being the relaxation coefficient. This coefficient is updated in every iteration cycle as a function of two previous residuals

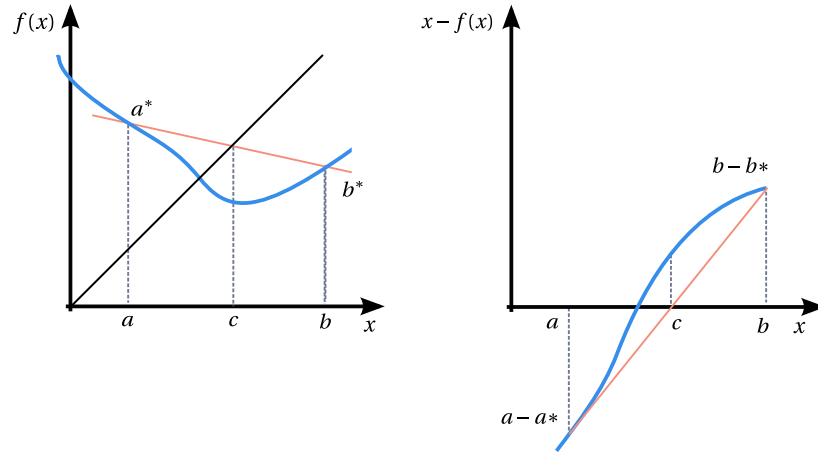


Figure 6.3: Geometric interpretation of the Aitken relaxation in one dimension for an example function f and corresponding interpretation as the secant method.

$$\omega^k = -\omega^{k-1} \frac{\left(\mathbf{r}^{(k)} - \mathbf{r}^{(k-1)}\right)^T \mathbf{r}^{(k-1)}}{\left(\mathbf{r}^{(k)} - \mathbf{r}^{(k-1)}\right)^2}. \quad (6.58)$$

Check signs!!!

Comparing with Equations (6.38) and (6.39), ω^k can be, in a sense, regarded as an approximation to the inverse of the Jacobian. Dynamic relaxation is also easy to implement, and the additional computational input is acceptable since only inner vector products must be performed. See Box 6.4 for the pseudocode.

6.4.1 Multi-secant methods

The following exposition follows closely Fang and Saad (2009). In quasi-Newton methods the Jacobian is updated in each iteration using a rank-one update. Standard quasi-Newton methods require the updated J_{k+1} to satisfy the following secant condition

$$J_{\mathcal{R}}^{k+1} \Delta \mathbf{x}^k = \Delta \mathcal{R}^k, \quad (6.59)$$

where $\Delta \mathcal{R}^k \equiv \mathcal{R}(\mathbf{x}^{k+1}) - \mathcal{R}(\mathbf{x}^k)$. Furthermore, another common requirement is the following so-called no-change condition

$$J_{\mathcal{R}}^{k+1} \mathbf{q} = J_{\mathcal{R}}^k \mathbf{q} \quad \forall \mathbf{q} \text{ such that } \mathbf{q}^T \Delta \mathbf{x}^k = 0, \quad (6.60)$$

which stipulates that there be no new information from $J_{\mathcal{R}}^k$ to $J_{\mathcal{R}}^{k+1}$ along any direction \mathbf{q} orthogonal to $\Delta \mathbf{x}^k$.

Broyden (1965) developed a method satisfying both secant condition (Equation (6.59)) and the no-change condition (Equation (6.60)). By simply imposing these conditions he arrived at the update formula

Box 6.4: Aitken relaxation for one timestep.

- (i) Set nonlinear counter to zero: $k = 0$
- (ii) $\mathbf{x}^k = \mathbf{x}_{n+1}^p$
- (iii) Enter the nonlinear loop
 - (1) Compute $\mathcal{R}(\mathbf{x}^k)$, which implies the solution of the mechanical and the thermal problems, \mathcal{U} and \mathcal{T} , respectively.
 - (2) if $k = 0$:
 - Compute \mathbf{x}^{k+1} using constant relaxation (Equation (6.47))
 - (3) else:
 - Compute \mathbf{x}^{k+1} using Aitken relaxation (Equations (6.57) and (6.58))
 - Save the current residue $\mathbf{r}^k = \mathcal{R}^k$.
 - (4) If the desired accuracy has not been reached, update $k = k + 1$ and go to step (1).

$$J_{\mathcal{R}}^{k+1} = J_{\mathcal{R}}^k + \left(\Delta \mathcal{R}^k - J_{\mathcal{R}}^k \Delta \mathbf{x}^k \right) \frac{\Delta \mathbf{x}^{kT}}{\Delta \mathbf{x}^{kT} \Delta \mathbf{x}^k}. \quad (6.61)$$

Matrix $J_{\mathcal{R}}^{k+1}$ in Equation (6.61) is the unique matrix satisfying both conditions (6.59) and (6.60). The Broyden update can also be obtained by minimizing $E(J_{\mathcal{R}}^{k+1}) = \|J_{\mathcal{R}}^{k+1} - J_{\mathcal{R}}^k\|_F^2$ with respect to terms of $J_{\mathcal{R}}^{k+1}$, subject to the secant condition (6.59).

It may seem at first that Broyden's first method can be expensive since computing the quasi-Newton step $\Delta \mathbf{x}^k$ requires solving a linear system at each iteration. However, note that, typically, the approximate Jacobian is a small rank modification of a diagonal matrix (or a matrix that is easy to invert); hence, the cost to obtain this solution is not too high as long as the number of steps is not too large.

An alternative is Broyden's second method that approximates the inverse Jacobian instead of the Jacobian itself. $G_{\mathcal{R}}^k$ is used to denote the estimated inverse Jacobian at the k th iteration. The secant condition (Equation (6.59)) now reads

$$G_{\mathcal{R}}^{k+1} \Delta \mathcal{R}^k = \Delta \mathbf{x}^k \quad (6.62)$$

By minimizing $E(G_{\mathcal{R}}^{k+1}) = \|G_{\mathcal{R}}^{k+1} - G_{\mathcal{R}}^k\|_F^2$ with respect to $G_{\mathcal{R}}^{k+1}$ subject to Equation (6.62), the following update formula is found for the inverse Jacobian

$$G_{\mathcal{R}}^{k+1} = G_{\mathcal{R}}^k + \left(\Delta \mathbf{x}_k - G_{\mathcal{R}}^k \Delta \mathcal{R}^k \right) \frac{\Delta \mathcal{R}^{kT}}{\Delta \mathcal{R}^{kT} \Delta \mathcal{R}^k} \quad (6.63)$$

which is also the only update satisfying both the secant condition (Equation (6.62)) and

the no-change condition for the inverse Jacobian

$$G_{\mathcal{R}}^k \mathbf{q} = G_{\mathcal{R}}^{k+1} \mathbf{q} \quad \forall \mathbf{q} \text{ such that } \mathbf{q}^T \Delta \mathcal{R}^k = 0. \quad (6.64)$$

The update formula in Equation (6.61) can also be obtained in terms of $G_{\mathcal{R}}^k \equiv J_{\mathcal{R}}^{k-1}$ by applying the Sherman-Morrison formula

$$G_{\mathcal{R}}^{k+1} = G_{\mathcal{R}}^k + \left(\Delta \mathbf{x}^k - G_{\mathcal{R}}^k \Delta \mathcal{R}^k \right) \frac{\Delta \mathbf{x}^{kT} G_{\mathcal{R}}^k}{\Delta \mathbf{x}^{kT} G_{\mathcal{R}}^k \Delta \mathcal{R}^k} \quad (6.65)$$

This shows, as was explained earlier, that to solve the Jacobian system associated with Broyden's first approach can be reduced to a set of update operations that are not more costly than those required by the second update. Note, however, that the above formula requires the inverse of the initial Jacobian.

From Equation (6.61) and Equation (6.63) it is possible to define Broyden's family of updates, in which an update formula takes the general form

$$G_{\mathcal{R}}^{k+1} = G_{\mathcal{R}}^k + \left(\Delta \mathbf{x}^k - G_{\mathcal{R}}^k \Delta \mathcal{R}^k \right) \mathbf{v}_k^T \quad (6.66)$$

where $\mathbf{v}_k^T \Delta \mathcal{R}^k = 1$ so that the secant condition (6.59) holds. Note that the secant condition (6.62) is equivalent to condition (6.59). Some authors called Broyden's first method Broyden's good update and Broyden's second method as Broyden's bad update. These are two particular members of Broyden's family.

6.4.1.1 Generalized Broyden

The multi-secant methods provide an approximation to the Jacobian in Equation (6.38) or Equation (6.42) using information from previous iterations. A generalized Broyden's method with a flexible rank update on the inverse Jacobian, satisfying a set of m secant equations

$$G_{\mathcal{R}}^k \Delta \mathcal{R}^i = \Delta \mathbf{x}^i \quad \text{for } i = k-m, \dots, k-1 \quad (6.67)$$

where it is assumed $\Delta \mathcal{R}^{k-m}, \dots, \Delta \mathcal{R}^{k-1}$ are linearly independent and $m \leq n$ can also be described. Aggregating Equations (6.67) in matrix form, they can be rewritten as

$$G_{\mathcal{R}}^k \mathcal{R}^k = \mathcal{X}^k. \quad (6.68)$$

where

$$\mathcal{R}^k = \left[\Delta \mathcal{R}^{k-m} \dots \Delta \mathcal{R}^{k-1} \right], \quad \mathcal{X}^k = \left[\Delta \mathbf{x}^{k-m} \dots \Delta \mathbf{x}^{k-1} \right] \in \mathbb{R}^{n \times m} \quad (6.69)$$

The no-change condition corresponding to (6.60) is

$$\left(G_{\mathcal{R}}^k - G_{\mathcal{R}}^{k-m} \right) \mathbf{q} = 0 \quad (6.70)$$

for all \mathbf{q} orthogonal to the subspace spanned by $\Delta \mathcal{R}^{k-m}, \dots, \Delta \mathcal{R}^{k-1}$, the columns of \mathcal{R}^k . In the end, this yields

$$G_{\mathcal{R}}^k = G_{\mathcal{R}}^{k-m} + \left(\mathcal{X}^k - G_{\mathcal{R}}^{k-m} \mathcal{R}^k \right) \left(\mathcal{R}^{kT} \mathcal{R}^k \right)^{-1} \mathcal{R}^{kT} \quad (6.71)$$

a rank- m update formula. Note that $\text{rank}(\mathcal{R}^k) = m$. The update formula for \mathbf{x}^{k+1} is

$$\mathbf{x}^{k+1} = \mathbf{x}^k - G_{\mathcal{R}}^k \mathcal{R}^k \quad (6.72)$$

$$= \mathbf{x}^k - G_{\mathcal{R}}^{k-m} \mathcal{R}^k - \left(\mathcal{X}^k - G_{\mathcal{R}}^{k-m} \mathcal{R}^k \right) \left(\mathcal{R}^{k\top} \mathcal{R}^k \right)^{-1} \mathcal{R}^{k\top} \mathcal{R}^k \quad (6.73)$$

$$= \mathbf{x}^k - G_{\mathcal{R}}^{k-m} \mathcal{R}^k - \left(\mathcal{X}^k - G_{\mathcal{R}}^{k-m} \mathcal{R}^k \right) \gamma_k \quad (6.74)$$

where the column vector γ_k is obtained by solving the normal equations $\left(\mathcal{R}^{k\top} \mathcal{R}^k \right) \gamma_k = \mathcal{R}^{k\top} \mathcal{R}^k$, which is equivalent to solving the least squares problem

$$\min_{\gamma} \left\| \mathcal{R}^k \gamma - \mathcal{R}^k \right\|_2. \quad (6.75)$$

Note that in Equation (6.74), if \mathcal{R}^k is square and of full rank, then for any $G_{\mathcal{R}}^{k-m}$,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \mathcal{X}^k \mathcal{R}^{k-1} \mathcal{R}^k, \quad (6.76)$$

the same form as that in the standard secant method.

6.4.1.2 Anderson mixing

The Anderson mixing scheme [5] takes the latest m steps into account to obtain a better approximation to \mathbf{x}_{n+1} without evaluating \mathcal{R} again. Consider

$$\bar{\mathbf{x}}^k = \mathbf{x}^k - \sum_{i=k-m}^{k-1} \gamma_i^k \Delta \mathbf{x}^i = \mathbf{x}^k - \mathcal{X}^k \gamma^k, \quad (6.77)$$

$$\bar{\mathcal{R}}^k = \mathcal{R}^k - \sum_{i=k-m}^{k-1} \gamma_i^k \Delta \mathcal{R}^i = \mathcal{R}^k - \mathcal{R}^k \gamma^k, \quad (6.78)$$

where $\Delta \mathbf{x}^i = \mathbf{x}^{i+1} - \mathbf{x}^i$ and $\Delta \mathcal{R}^i = \mathcal{R}^{i+1} - \mathcal{R}^i$, $\mathcal{X}^k = \left[\Delta \mathbf{x}^{k-m} \dots \Delta \mathbf{x}^{k-1} \right]$, $\mathcal{R}^k = \left[\Delta \mathcal{R}^{k-m} \dots \Delta \mathcal{R}^{k-1} \right]$, and $\gamma^k = \left[\gamma_{k-m}^k \dots \gamma_{k-1}^k \right]^\top$. Expressing the equations in the form $\bar{\mathbf{x}}^k = \sum_{j=k-m}^k w_j \mathbf{x}^j$ and $\bar{\mathcal{R}}^k = \sum_{j=k-m}^k w_j \mathcal{R}^j$, it is found that $\sum_{j=k-m}^k w_j = 1$. In other words, $\bar{\mathbf{x}}_k$ and $\bar{\mathcal{R}}_k$ are weighted averages of $\mathbf{x}_{k-m}, \dots, \mathbf{x}_k$ and $\mathcal{R}^{k-m}, \dots, \mathcal{R}^k$, respectively. The arguments $\gamma^k = \left[\gamma_{k-m}^k \dots \gamma_{k-1}^k \right]^\top$ are determined by minimizing

$$E(\gamma^k) = \langle \bar{\mathcal{R}}^k, \bar{\mathcal{R}}^k \rangle = \left\| \mathcal{R}^k - \mathcal{R}^k \gamma^k \right\|_2^2 \quad (6.79)$$

whose solution can, but should not in practice, be obtained by solving the normal equations

$$\left(\mathcal{R}^{k\top} \mathcal{R}^k \right) \gamma^k = \mathcal{R}^{k\top} \mathcal{R}^k. \quad (6.80)$$

Combining Equations (6.77), (6.78), and (6.80), one obtains

$$\mathbf{x}^{k+1} = \bar{\mathbf{x}}^k + \beta \bar{\mathcal{R}}^k \quad (6.81)$$

$$= \mathbf{x}^k + \beta \mathcal{R}^k - \left(\mathcal{X}^k + \beta \mathcal{R}^k \right) \gamma^k \quad (6.82)$$

$$= \mathbf{x}^k + \beta \mathcal{R}^k - \left(\mathcal{X}^k + \beta \mathcal{R}^k \right) \left(\mathcal{R}^{k\top} \mathcal{R}^k \right)^{-1} \mathcal{R}^{k\top} \mathcal{R}^k \quad (6.83)$$

where β is the preset mixing parameter and $\mathcal{R}^{k^T} \mathcal{R}^k$ is assumed to be nonsingular. In particular, if no previous iterate is taken into account (i.e. $m = 0$), then Equation (6.83) reads

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \beta \mathcal{R}^k \quad (6.84)$$

This scheme is referred to as simple mixing and underrelaxation if $0 < \beta < 1$ (see Section 6.3.3). The update formula (6.83) is the same as (6.74) by setting $G_{\mathcal{R}}^{k-m} = -\beta \mathbf{I}$. In this respect Anderson mixing implicitly forms an approximate inverse Jacobian $G_{\mathcal{R}}^k$ that minimizes $\|G_{\mathcal{R}}^k + \beta \mathbf{I}\|_F$ subject to (6.68). In the context of mixing, generalized Broyden's second method is equivalent to Anderson mixing. Note that if \mathcal{R}^k is square and nonsingular, then Equation (6.83) matches the formula of the standard secant method.

6.4.1.3 Generalized Broyden's family

Now we can write down the generalized Broyden family, in which an update algorithm is in the form

$$G_{\mathcal{R}}^k = G_{\mathcal{R}}^{k-m} + \left(\mathcal{X}^k - G_{\mathcal{R}}^{k-m} \mathcal{R}^k \right) V^{k^T} \quad (6.85)$$

where $V^{k^T} \mathcal{R}^k = \mathbf{I}$ so that the secant condition $G_{\mathcal{R}}^k \mathcal{R}^k = \mathcal{X}^k$ holds. The two optimal choices of $V^{k^T} = M^{k-1} N^{k^T}$ are

$$M^k = \mathcal{R}^{k^T} \mathcal{R}^k, \quad N^{k^T} = \mathcal{R}^{k^T}, \quad (6.86)$$

minimizing $\|G_{\mathcal{R}}^k - G_{\mathcal{R}}^{k-m}\|_F$ and

$$M^k = \mathcal{X}^{k^T} G_{\mathcal{R}}^k \mathcal{R}^k, \quad N^{k^T} = \mathcal{X}^{k^T} G_{\mathcal{R}}^k, \quad (6.87)$$

minimizing $\|J_{\mathcal{R}}^k - J_{\mathcal{R}}^{k-m}\|_F$. This last choice yields as the approximation for the Jacobian

$$J_{\mathcal{R}}^k = J_{\mathcal{R}}^{k-m} + \left(\mathcal{R}^k - J_{\mathcal{R}}^{k-m} \mathcal{X}^k \right) \left(\mathcal{X}^{k^T} \mathcal{X}^k \right)^{-1} \mathcal{X}^{k^T}, \quad (6.88)$$

after applying the Woodbury formula. The first choice is said to correspond to a Type-II update and the second to a Type-I update (Fang and Saad, 2009).

6.4.1.4 Anderson's family

The update formula for Anderson's family can be found from Equation (6.85) using as the approximation to the previous Jacobian the identity matrix multiplied by a constant β , i.e.,

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \beta \mathcal{R}^k - (\mathcal{X}^k + \beta \mathcal{R}^k) \mathbf{V}^{k^T} \mathcal{R}^k. \quad (6.89)$$

The two choices for \mathbf{V}^k remain the same, replacing $G_{\mathcal{R}}^{k-m}$ by $-\beta \mathbf{I}$. They now minimize $\|G_{\mathcal{R}}^k + \beta \mathbf{I}\|$ and $\|J_{\mathcal{R}}^k + (1/\beta) \mathbf{I}\|$.

6.4.1.5 The Broyden-like class

Both the generalized Broyden's family and Anderson's family can be understood as methods in the Broyden-like class as described in [Fang and Saad \(2009\)](#). Suppose the latest m iterates are available, which are denoted by $\mathbf{x}^1, \dots, \mathbf{x}^m$. Let $\Delta \mathbf{x}^i = \mathbf{x}^{i+1} - \mathbf{x}^i$ for $i = 1, \dots, m-1$. Partition $\Delta \mathbf{x}^1, \dots, \Delta \mathbf{x}^{m-1}$ into p groups,

$$\mathcal{X}^1 = [\Delta \mathbf{x}^1, \dots, \Delta \mathbf{x}^{z_1}], \quad (6.90)$$

$$\mathcal{X}^2 = [\Delta \mathbf{x}^{z_1+1}, \dots, \Delta \mathbf{x}^{z_2}], \quad (6.91)$$

$$\vdots \quad (6.92)$$

$$\mathcal{X}^p = [\Delta \mathbf{x}^{z_{p-1}+1}, \dots, \Delta \mathbf{x}^{z_p}], \quad (6.93)$$

where z_i is the index of the last entry in the i th group for $i = 1, \dots, p$; $z_0 = 0$ and $z_p = m-1$. Also partition $\Delta \mathcal{R}^1, \dots, \Delta \mathcal{R}^{m-1}$ into $\mathcal{R}^1, \dots, \mathcal{R}^p$ accordingly, where $\Delta \mathcal{R}^i = \mathcal{R}^{i+1} - \mathcal{R}^i$ with $\mathcal{R}^i = \mathcal{R}(\mathbf{x}^i)$. The sizes of the groups for $i = 1, \dots, k$ are denote by $s_i := z_i - z_{i-1}$. Note that the indexing here is different from the previous sections. The inverse of the Jacobian is iteratively approximated at the $(z_i + 1)$ st iterate for $i = 1, \dots, p$ as

$$G_{\mathcal{R}}^{i+1} = G_{\mathcal{R}}^i + (\mathcal{X}^i - G_{\mathcal{R}}^i \mathcal{R}^i) \mathbf{V}^{iT}, \quad (6.94)$$

where $\mathbf{V}^{iT} \mathcal{R}^i = \mathbf{I}$ for the secant condition. The update follows the formula of the generalized Broyden family. In the context of mixing, the base case is

$$G_{\mathcal{R}}^1 = -\beta \mathbf{I}, \quad (6.95)$$

where β is the mixing parameter. The next iterate is set as

$$\mathbf{x}^{m+1} = \mathbf{x}^m - G_{\mathcal{R}}^{k+1} \mathcal{R}^m. \quad (6.96)$$

The choice of V_i satisfying $V_i^T \mathcal{F}_i = \mathbf{I}$ is performed as described in [Section 6.4.1.3](#).

6.4.2 Practical considerations

The application of Broyden's method as described so far is unfeasible for the problems considered in this document, i.e., thermo-mechanical coupled problems with many unknowns. So far, the descriptions considered of the Broyden-like class methods require one to keep the large $G_{\mathcal{R}}^i$ matrices of size $n_{\text{unknowns}} \times n_{\text{unknowns}}$ in memory, in addition to the previous iterates. This is a significant drawback. However, [Fang and Saad \(2009\)](#) present a more memory efficient way of implementing these methods. Let

$$E^i = \mathcal{X}^i - G_{\mathcal{R}}^i \mathcal{R}^i. \quad (6.97)$$

Substituting Equation (6.97) into Equation (6.94) one obtains

$$G_{\mathcal{R}}^i = G_{\mathcal{R}}^{i-1} + E^{i-1} V^{i-1T}, \quad (6.98)$$

$$= G_{\mathcal{R}}^{i-2} + E^{i-2} V^{i-2T} + E^{i-1} V^{i-1T}, \quad (6.99)$$

$$\vdots \quad (6.100)$$

$$= G_{\mathcal{R}}^1 + \sum_{j=1}^{i-1} E^j V^{jT}, \quad (6.101)$$

for $i = 2, \dots, p+1$. Matrices $G_{\mathcal{R}}^i$ need not be explicitly stored. $G_{\mathcal{R}}^i$ is needed only to compute $G_{\mathcal{R}}^i \mathcal{R}^i$ in Equation (6.97) and $G_{\mathcal{R}}^{p+1} \mathcal{R}^m$ in Equation (6.96), and also for V^i if it depends on $G_{\mathcal{R}}^i$. Substituting Equation (6.101) into Equation (6.97) obtains

$$E^i = \mathcal{X}^i - G_{\mathcal{R}}^1 \mathcal{R}^i - \sum_{j=1}^{i-1} E^j \left(V^j{}^T \mathcal{R}^i \right), \quad (6.102)$$

for $i = 1, \dots, p$. The computation is economic for large-scale problems with $n \gg m$. The next iterate \mathbf{x}^{m+1} in (32) can also be computed in a similar manner

$$\mathbf{x}^{m+1} = \mathbf{x}^m - G_{\mathcal{R}}^{p+1} \mathcal{R}^m = \mathbf{x}^m - G_{\mathcal{R}}^1 \mathcal{R}^m - \sum_{j=1}^k E^j \left(V^j{}^T \mathcal{R}^m \right). \quad (6.103)$$

Using Type-II update, the computation of V^i is straightforward from \mathcal{R}^i . On the other hand, Type-I update involves $G_{\mathcal{R}}^i$ to compute V^i . Thus

$$N^i{}^T = \mathcal{X}^i{}^T G_{\mathcal{R}}^i = G_{\mathcal{R}}^1 \mathcal{X}^i{}^T + \sum_{j=1}^{i-1} \left(\mathcal{X}^i{}^T E^j \right) V^j{}^T. \quad (6.104)$$

After obtaining N^i , we compute $M^i = N^i{}^T \mathcal{R}^i$ and then $V^i{}^T = M^{i-1} N^i{}^T$ for $i = 1, \dots, p$.

Looking at Equations (6.102), (6.103) and (6.104), one still needs the initial approximation to the inverse of the Jacobian, $G_{\mathcal{R}}^1$, whose size is $n_{\text{unknown}} \times n_{\text{unknown}}$. In Fang and Saad (2009) the approach adopted was to follow the idea of Anderson's mixing and set

$$G_{\mathcal{R}}^1 = -\beta \mathbf{I}, \quad (6.105)$$

drastically improving the memory requirements, as only one scalar parameter, β , needs to be saved. Also, Kelley (2003), assumes in his implementation of Broyden's method, an initial approximation to $G_{\mathcal{R}}^1$ equal to the identity matrix. Information about $G_{\mathcal{R}}^1$ is applied in the preconditioning of the system instead.

To compute V^i , Fang and Saad (2009) suggest a QR decomposition with pivoting. Be it for a Type-II update, where one needs to solve a least-squares problem, or for a Type-I update, one needs to invert a generally non-symmetric matrix. This approach leads to better numerical stability when the matrices to be inverted are singular or ill-conditioned, compared with solving the normal equations. The QR decomposition has a computational cost of $\mathcal{O}(n_{\text{unknown}}^3)$ algebraic operations (Dennis and Schnabel, 1996).

If the size of the groups s_1, s_2, \dots are fixed from one Newton iteration to the next, so the E^i and V^i matrices remain the same from one iteration to the next, the computation effort to compute them is saved from one iteration to the next. Here only constant $s = s_1 = s_2 = \dots$ is considered, where $s = \infty$ corresponds to Anderson's mixing, where all previous iterates available are considered.

One question remains. How to proceed when the available memory runs out?. According to Kelley (2003), as is often the case with GMRES, the iteration can be restarted if there is no more room to store the vectors. A different approach, called limited memory in the optimization literature, is to replace the oldest of the stored steps with the most recent one.

Thus, applying the method as suggest by Fang and Saad (2009), leads to a storage need of

1. Two column vectors of size n_{unknown} for \mathbf{x}^m and \mathcal{R}^m .
2. An $n_{\text{unknown}} \times (m-1)$ matrix for $\mathcal{X}^1, \dots, \mathcal{X}^k$ (shared with E^1, \dots, E^k).
3. An $n_{\text{unknown}} \times (m-1)$ matrix for $\mathcal{R}^1, \dots, \mathcal{R}^k$ (shared with V^1, \dots, V^k).
4. For Type-I update we also store the last group N^k , since its computation involves $G_{\mathcal{R}}^k$.

and for each nonlinear iteration the computational cost is $\mathcal{O}(n_{\text{unknown}}^3)$. For $m=1$, V^i can be directly computed without needing to invert matrices, and the cost comes down to $\mathcal{O}(n_{\text{unknown}})$ per nonlinear iteration. Also, when k is different from any z_i a group is being complete, one can either use a shortened group or reuse the approximation to the Jacobian without using the new information. This save computational effort and those iteration cost only $\mathcal{O}(n_{\text{unknown}})$. See in Box 6.5 the pseudocode for the Broyden-like family with restart.

Kelley (2003) presents for the Broyden method an implementation that halves the memory requirement relative to the one present in Fang and Saad (2009). It is based on the Type-I update and to deduce it consider Equation (6.61) and Sherman-Morrison formula

$$(J_{\mathcal{R}} + \mathbf{u}\mathbf{v}^T)^{-1} = \left(\mathbf{I} - \frac{(G_{\mathcal{R}}\mathbf{u})\mathbf{v}^T}{1 + \mathbf{v}^T G_{\mathcal{R}}\mathbf{u}} \right) G_{\mathcal{R}}, \quad (6.106)$$

where as before $G_{\mathcal{R}} \equiv J_{\mathcal{R}}^{-1}$. One can rewrite Equation (6.61) as

$$J_{\mathcal{R}}^{k+1} = J_{\mathcal{R}}^k + \mathbf{u}^k \mathbf{v}^{kT}, \quad (6.107)$$

where

$$\mathbf{u}^k = (\Delta \mathcal{R}^k - J_{\mathcal{R}}^k \Delta \mathbf{x}^k) / \|\Delta \mathbf{x}^k\| \text{ and } \mathbf{v}^k = \Delta \mathbf{x}^k / \|\Delta \mathbf{x}^k\|. \quad (6.108)$$

Then, keeping in mind that $J_{\mathcal{R}}^1 = \mathbf{I}$,

$$G_{\mathcal{R}}^{k+1} = \left(\mathbf{I} - \mathbf{w}^k \mathbf{v}^{kT} \right) \left(\mathbf{I} - \mathbf{w}^{k-1} \mathbf{v}^{k-1T} \right) \dots \left(\mathbf{I} - \mathbf{w}^1 \mathbf{v}^{1T} \right) G_{\mathcal{R}}^1, \quad (6.109)$$

$$= \prod_{j=0}^k \left(\mathbf{I} - \mathbf{w}^j \mathbf{v}^{jT} \right), \quad (6.110)$$

where, for $k \geq 0$,

$$\mathbf{w}^k = \frac{G_{\mathcal{R}}^k \mathbf{u}^k}{1 + \mathbf{v}^{kT} G_{\mathcal{R}}^k \mathbf{u}^k}. \quad (6.111)$$

So, to apply $G_{\mathcal{R}}^{k+1}$ to a vector \mathbf{p} , the is cost of $\mathcal{O}(n_{\text{unknown}} k)$ floating point operations and storage of the $2k$ vectors $\{\mathbf{w}^j\}_{j=1}^k$ and $\{\Delta \mathbf{x}^j\}_{j=1}^k$. The storage can be halved with a trick (see Kelley (2003) for details)

$$\Delta \mathbf{x}^k = -G_{\mathcal{R}}^{k+1} \mathcal{R}^k, \quad (6.112)$$

$$= - \left(\mathbf{I} - \frac{\mathbf{w}^k \Delta \mathbf{x}^{kT}}{\|\Delta \mathbf{x}^k\|} \right) G_{\mathcal{R}}^k \mathcal{R}^k, \quad (6.113)$$

$$= - \frac{G_{\mathcal{R}}^k \mathcal{R}^k}{1 + \Delta \mathbf{x}^{kT} G_{\mathcal{R}}^k \mathcal{R}^k / \|\Delta \mathbf{x}^k\|^2}. \quad (6.114)$$

According to Kelley (2003), the Sherman-Morrison approach is more efficient, in terms of both time and storage, than dense matrix approaches proposed elsewhere. For example, the approach presented in Dennis and Schnabel (1996) has a $\mathcal{O}(n_{\text{unknown}})$ cost per nonlinear iteration and requires one to keep in memory the QR decomposition of the previous approximation to the Jacobian. However, the dense matrix approach can detect ill-conditioning in the approximate Jacobians. Bounded deterioration implies that the Broyden matrices will be well-conditioned if the data is sufficiently good, and superlinear convergence suggests that only a few iterates will be needed.

In the context of FSI The multi-secant quasi-Newton methods have been used in the context of FSI, although not always presented as such (Haelterman et al., 2009; Gatzhammer, 2014; Uekermann, 2016; Scheufele, 2018). Vierendeels et al. (2007) and Degroote et al. (2008) consider the system of equations (6.45) and (6.46), where recall that an estimate for the Jacobians $J_{\mathcal{U}}$ and $J_{\mathcal{T}}$ are needed. The authors achieve this by using linear reduced-order models for the fluid solver and the structure solver. These are set up from solver input and output deltas or sensitivities during the coupling iterations. The resulting method for two black-box solvers is called interface block quasi-Newton method with least-squares approximation (IBQN-LS) in Degroote (2010).

This approach can be understood in the framework of the multi-secant quasi-Newton methods presented above and originating in Fang and Saad (2009) as follows. If one looks at $\beta \mathbf{I} - (\mathcal{X}^k + \beta \mathcal{R}^k) \mathbf{V}^{kT}$ in Equation (6.89) as, in a sense, an approximation to the inverse of the Jacobian (compare with Equation (6.85)). The corresponding Jacobian is given by

$$J_{\mathcal{R}}^k = \alpha \mathbf{I} + (\mathcal{R}^k - \alpha \mathbf{I} \mathcal{X}^k) (\mathcal{X}^{kT} \mathcal{X}^k)^{-1} \mathcal{X}^{kT}, \quad (6.115)$$

where $\alpha = 1/\beta$. If one sets $\alpha = 0$, the approximation to the Jacobian obtained is

$$J_{\mathcal{R}}^k = \mathcal{R}^k (\mathcal{X}^{kT} \mathcal{X}^k)^{-1} \mathcal{X}^{kT}. \quad (6.116)$$

This corresponds to the linear reduced order models in Vierendeels et al. (2007), where \mathcal{R} is replaced by the functions corresponding to the fluid and structure solvers. If the functions considered are instead the mechanical and thermal solvers, this method can easily be applied to the thermomechanical problem. The block $(\mathcal{X}^{kT} \mathcal{X}^k)^{-1} \mathcal{X}^{kT}$ can be understood as being a part of a least-squares solution, the so-called normal equations, i.e., the equations whose solution also solve the minimization problem

$$\arg \min_{\tilde{\gamma}} \|\Delta \mathbf{x} - \mathcal{X}^k \tilde{\gamma}\|_2. \quad (6.117)$$

As such one can avoid the use of the normal equations and employ more numerically stable and efficient methods such economy size QR -decomposition. In addition, Vierendeels et al. (2007) solve the system of equation (6.45) and (6.46) in a Gauss-Seidel manner, using always the most recent values available to estimate the Jacobians.

An interface quasi-Newton method based on Equation (6.38) and Equation (6.39) is presented in Degroote (2010) for FSI. The method is called interface quasi-Newton with an approximation of the inverse of the interface Jacobian matrix by least squares (QIN-ILS). Its origin is the IBQN-LS method presented in Vierendeels et al. (2007), and it

employs only one reduced-order model for the inverse of the overall interface Jacobian matrix of the Newton system (Equation (6.38)) applied to the right-hand side vector.

If in Equation (6.83), corresponding to Anderson's mixing, one sets $\beta = -1$, the update formula comes out to be

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \mathcal{R}^k - (\mathcal{X}^k - \mathcal{R}^k) \left(\mathcal{R}^{k\top} \mathcal{R}^k \right)^{-1} \mathcal{R}^{k\top} \mathcal{R}^k. \quad (6.118)$$

Using the definition for the fixed-point function \mathcal{S} , one finds

$$\mathbf{x}^{k+1} = \mathcal{S}^k - \mathcal{S}^k \left(\mathcal{R}^{k\top} \mathcal{R}^k \right)^{-1} \mathcal{R}^{k\top} \mathcal{R}^k, \quad (6.119)$$

or

$$\Delta \mathbf{x}^k = -\mathcal{R}^k - \mathcal{S}^k \left(\mathcal{R}^{k\top} \mathcal{R}^k \right)^{-1} \mathcal{R}^{k\top} \mathcal{R}^k, \quad (6.120)$$

When no delta columns are available yet, constant relaxation is used once to ensure stability.

Box 6.5: Broyden-like class methods for one timestep with restart.

- (i) Set $\mathbf{x}^1 = \mathbf{x}_{n+1}^p$.
- (ii) Evaluate $\mathcal{R}^1 = \mathcal{R}(\mathbf{x}^1)$, which implies the solution of the mechanical and the thermal problems, \mathcal{U} and \mathcal{T} , respectively.
- (iii) Compute $\mathbf{x}^2 = \mathbf{x}^1 + \beta \mathcal{R}^1$.
- (iv) Evaluate $\mathcal{R}^2 = \mathcal{R}(\mathbf{x}^2)$, which implies the solution of the mechanical and the thermal problems, \mathcal{U} and \mathcal{T} , respectively.
- (v) Initialize counters: $k = 2$ and $i = 1$
- (vi) Enter the nonlinear loop
 - (1) If the maximum number of previous iteration m has been reached restart all \mathcal{X}^i and \mathcal{R}^i and set $i = 1$.
 - (2) Compute $\Delta \mathcal{R}^{k-1} = \mathcal{R}^k - \mathcal{R}^{k-1}$ and add it to \mathcal{R}^i .
 - (3) Compute $\Delta \mathbf{x}^{k-1} = \mathbf{x}^k - \mathbf{x}^{k-1}$ and add it to \mathcal{X}^i .
 - (4) If $s = \infty$ (Anderson mixing):
 - Compute V^{iT} according to the type of update chosen (Equation (6.87) and (6.86)).
 - Update according to $\mathbf{x}^{k+1} = \mathbf{x}^k + \beta \mathcal{R}^k - (\mathcal{X}^i + \beta \mathcal{R}^i) V^{iT} \mathcal{R}^k$.
 - (5) else:
 - If $k = z_j + 1$ for any $j \geq 1$, compute E^i (Equation (6.102)) and V^{iT} according to the type of update chosen (Equation (6.87) and (6.86)). Save them on \mathcal{X}^i and \mathcal{R}^i , respectively. Update $i = i + 1$.
 - Update according to $\mathbf{x}^{k+1} = \mathbf{x}^k + \beta \mathcal{R}^k - \sum_{j=1}^{i-1} E^j (V^j{}^T \mathcal{R}^k)$.
 - (6) Evaluate $\mathcal{R}^{k+1} = \mathcal{R}(\mathbf{x}^{k+1})$, which implies the solution of the mechanical and the thermal problems, \mathcal{U} and \mathcal{T} , respectively.
 - (7) If the desired accuracy has not been reached, update $k = k + 1$ and go to step (1).

6.5 Multipoint iteration functions

6.5.1 Finite-Difference Newton Method

This approach follows precisely the one described in Section 6.3.2 for the "standard" Newton method. The difference lies in the computation of the Jacobian. Here the Jacobian $J_{\mathcal{R}}(\mathbf{x})$ is approximated from a forward finite-difference, $J_{\mathcal{R}}^h(\mathbf{x})$, by columns. Following Kelley (2003), the j th column is

$$\left[J_{\mathcal{R}}^h(\mathbf{x}) \right]_j = \begin{cases} \frac{\mathcal{R}(\mathbf{x} + h\sigma_j \mathbf{e}_j) - \mathcal{R}(\mathbf{x})}{\sigma_j h}, & x_j \neq 0 \\ \frac{\mathcal{R}(h\mathbf{e}_j) - \mathcal{R}(\mathbf{x})}{h}, & x_j = 0 \end{cases}, \quad (6.121)$$

where \mathbf{e}_j is the unit vector in the j th coordinate direction. The difference increment h should be no smaller than the square root of the inaccuracy in \mathcal{R} (Kelley, 2003). It should, however, be scaled. Rather than simply perturbing \mathbf{x} by a difference increment h , roughly the square root of the error in \mathcal{R} , in each coordinate direction, the perturbation is multiplied to compute the j th column by

$$\sigma_j = \max(|(x)_j|, 1) \text{sign}((x)_j), \quad (6.122)$$

with a view toward varying the correct fraction of the low-order bits in $(x)_j$. The sign function is defined as

$$\text{sgn}(z) = \begin{cases} z/|z| & \text{if } z \neq 0 \\ 1 & \text{if } z = 0 \end{cases}. \quad (6.123)$$

While this scaling usually makes little difference, it can be crucial if $|(x)_j|$ is very large. Note that there is no adjustment if $|(x)_j|$ is very small because the error determines the lower limit on the size of the difference increment in \mathcal{R} . For example, if evaluations of \mathcal{R} are accurate to 16 decimal digits, the difference increment should change roughly the last eight digits of x . (Kelley, 2003)

Each column of $J_{\mathcal{R}}^h(\mathbf{x})$ requires one new function evaluation and, therefore, a finite difference Jacobian costs n_{unknown} function evaluations. If the perturbation is appropriately chosen, the method converges quadratically when the function satisfies certain conditions, and the initial attempt is close enough to the solution (Dennis and Schnabel, 1996).

The Chord and Shamanskii Methods If the computational cost of a forward difference Jacobian is high, i.e., \mathcal{R} is expensive and/or n_{unknown} is significant. If an analytic Jacobian is not available, it is wise to amortize this cost over several nonlinear iterations. The chord method does precisely that. It differs from Newton's method in that the evaluation and factorization of the Jacobian are done only once for $J_{\mathcal{R}}(\mathbf{x}^0)$. The advantages of the chord method increase as n increases, since both the n function evaluations and the $O(n^3_{\text{unknown}})$ work (in the dense matrix case) in the matrix factorization are done only once. So, while the convergence is q -linear and more nonlinear iterations will be needed than for Newton's method, the overall cost of the solution will usually be much less. A middle ground is the Shamanskii method. Here the Jacobian factorization and matrix function evaluation is done after every m computations of the step (Kelley, 2003).

Since the present use-case, the number of unknowns n is very large, and the evaluation of the function \mathcal{R} is also costly, approximating the Jacobian using a finite difference is not suitable, even utilizing the chord or Shamanskii methods.

6.5.2 Newton-Krylov methods

In the Newton-Krylov methods, the solution of the Newton system of equations in Equation (6.38) is achieved using Krylov methods, such as GMRES or BiCGSTAB. The Krylov iterative methods approximate the solution of a linear system $\mathbf{Ax} = \mathbf{b}$ using the Krylov subspace

$$\mathcal{K}_m = \text{span}\{\mathbf{r}_0, \mathbf{Ar}_0, \mathbf{A}^2\mathbf{r}_0, \dots, \mathbf{A}^{m-1}\mathbf{r}_0\}, \quad (6.124)$$

such that the m th iterate, $\mathbf{x}_m \in \mathcal{K}_m$, with $\mathbf{r}_0 = \mathbf{b} - \mathbf{Ax}_0$. The precise way the \mathbf{x}_m is built is what distinguishes the different methods.

To produce the appropriate Krylov subspace, one needs the product $J_{\mathcal{R}}(\mathbf{x}^k)\mathbf{y}$ in Equation (6.38), for some vector \mathbf{y} . It is assumed that the Jacobian is not available, so it must be approximated. Also, it would be beneficial if the full Jacobian is neither computed in its entirety nor wholly stored in memory, i.e., a matrix-free method is desirable. As in Section 6.5.1, the Jacobian-vector product is easy to approximate with a forward difference directional derivative (Kelley, 2003). The forward difference directional derivative at \mathbf{x}^k in the direction \mathbf{q} is

$$J_{\mathcal{R}}^h(\mathbf{x}^k)\mathbf{q} = \begin{cases} \mathbf{0}, & \mathbf{q} = \mathbf{0} \\ \|\mathbf{q}\| \frac{\mathcal{R}(\mathbf{x}^k + \sigma(\mathbf{x}^k, \mathbf{q})h\mathbf{q}/\|\mathbf{q}\|) - \mathcal{R}(\mathbf{x}^k)}{\sigma(\mathbf{x}^k, \mathbf{q})h}, & \mathbf{q} \neq \mathbf{0}. \end{cases} \quad (6.125)$$

The scaling is important. \mathbf{q} is scaled to be a unit vector and take a numerical directional derivative in the direction $\mathbf{q}/\|\mathbf{q}\|$. If h is roughly the square root of the error in \mathcal{R} , a difference increment in the forward difference is used to make sure that the appropriate low-order bits of \mathbf{x}^k is perturbed. So h is multiplied by

$$\sigma(\mathbf{x}^k, \mathbf{q}) = \max(|\mathbf{x}^k{}^T \mathbf{q}|, \|\mathbf{q}\|) \text{sign}(\mathbf{x}^k{}^T \mathbf{q}) / \|\mathbf{q}\|. \quad (6.126)$$

Sidi (2017) describes two different Newton-Krylov methods, the Newton-Arnoldi and the Newton-GMRES. The goal is to solve Equation (6.38) and find $\Delta\mathbf{x}^k$. The Krylov iterates are denoted by $\Delta\mathbf{x}_m^*$.

The first phase of both algorithms is identical. They both use the Arnoldi-Gram-Schmidt process to produce an orthogonal basis for the Krylov subspace, $\{\mathbf{q}_1, \dots, \mathbf{q}_m\}$ for some integer k . Consider the unitary matrix \mathbf{Q}_m

$$\mathbf{Q}_m = [\mathbf{q}_1 \cdots \mathbf{q}_m], \quad (6.127)$$

and the upper Hessenber matrix \mathbf{H}_m and related matrix $\tilde{\mathbf{H}}_m$

$$\mathbf{H}_m = \begin{bmatrix} h_{11} & h_{12} & \cdots & \cdots & h_{1m} \\ h_{21} & h_{22} & \cdots & \cdots & h_{2m} \\ & \ddots & \ddots & & \vdots \\ & & \ddots & \ddots & \vdots \\ & & & h_{m,m-1} & h_{mm} \end{bmatrix}, \quad \tilde{\mathbf{H}}_m = \begin{bmatrix} h_{11} & h_{12} & \cdots & \cdots & h_{1m} \\ h_{21} & h_{22} & \cdots & \cdots & h_{2m} \\ & \ddots & \ddots & & \vdots \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & h_{mm} \\ & & & & b_{m+1,m} \end{bmatrix}, \quad (6.128)$$

where the quantities h_{ij} are also obtained from the algorithm. \mathbf{H}_m can be interpreted as the projection of \mathbf{A} onto the Krylov subspace since $\mathbf{H}_m = \mathbf{Q}_m^T J_{\mathcal{R}}(\mathbf{x}^k) \mathbf{Q}_m$. See Box 6.6 for the pseudo-code of the Arnoldi-Gram-Schmidt process.

Box 6.6: Arnoldi process to orthonormalize the Krylov subspace

- (i) Compute $\mathbf{r}_0 = -\mathcal{R}(\mathbf{x}^k)$ and set $\beta = \|\mathbf{r}_0\|$ and $\mathbf{q}_1 = \mathbf{r}_0 / \beta$.
- (ii) $j = 1$
- (iii) $\mathbf{a}_{j+1}^{(1)} = J_{\mathcal{R}}(\mathbf{x}^k) \mathbf{q}_j$.
- (iv) Compute $h_{ij} = (\mathbf{q}_i, \mathbf{a}_{j+1}^{(i)})$ and compute $\mathbf{a}_{j+1}^{(i+1)} = \mathbf{a}_{j+1}^{(i)} - h_{ij} \mathbf{q}_i$ for $i = 1, \dots, j$.
- (v) Compute $h_{j+1,j} = \|\mathbf{a}_{j+1}^{(j+1)}\|$ and set $\mathbf{q}_{j+1} = \mathbf{a}_{j+1}^{(j+1)} / h_{j+1,j}$.
- (vi) $j = j + 1$
- (vii) If $j < m - 1$ go to Step (iii)

After obtaining an orthogonal basis for the Krylov subspace \mathcal{K}_m , the Newton-Arnoldi projects the m th residual onto the Krylov subspace and sets it to zero, similar to a weighed residual method, i.e.,

$$\mathbf{Q}_m \left(\mathcal{R}(\mathbf{x}^k) - J_{\mathcal{R}}(\mathbf{x}^k) \Delta \mathbf{x}_m^* \right) = \mathbf{0}. \quad (6.129)$$

In practice, one solves the equivalent linear system $\mathbf{H}_m \boldsymbol{\eta} = \beta \mathbf{e}_1$ for $\boldsymbol{\eta}$ and set $\Delta \mathbf{x}^k = \Delta \mathbf{x}_0^* + \mathbf{Q}_m \boldsymbol{\eta}$.

The Newton-GMRES attempts to minimize the norm of the m th residual, i.e.,

$$\Delta \mathbf{x}_m^* = \arg \min_{\mathbf{y} \in \mathcal{K}_m} \left\| \mathcal{R}(\mathbf{x}^k) - J_{\mathcal{R}}(\mathbf{x}^k) \mathbf{y} \right\|_2. \quad (6.130)$$

This is executed in practice solving the linear least-squares problem $\|\beta \mathbf{e}_1 - \tilde{\mathbf{H}}_m \boldsymbol{\eta}\|$ for $\boldsymbol{\eta}$ and set $\Delta \mathbf{x}^k = \Delta \mathbf{x}_0^* + \mathbf{Q}_m \boldsymbol{\eta}$.

Since the present use case includes many unknowns, it leads to memory concerns if the Krylov subspace is allowed to grow indefinitely. A restarted version where the maximum size of the Krylov space is restricted to m elements is preferred. Once this number is reached, the procedure is restarted. However, if m is small, the convergence can be poor.

In each nonlinear iteration of the Newton-Krylov, the number of iterations can be large, and each iteration requires an evaluation of the function. This can be a significant drawback when the intended use assumes that evaluating the function \mathcal{R} is expensive. This problem is however mitigated by the fact the Newton system is only solved until it satisfies

$$\|J_{\mathcal{R}}(\mathbf{x}^k) \Delta \mathbf{x}_m^* + \mathcal{R}(\mathbf{x}^k)\| \leq \eta \|\mathcal{R}(\mathbf{x}^k)\|, \quad (6.131)$$

where η is called the forcing term and it is chosen to avoid oversolving the Newton system (Equation (6.38)). As a simple approach, Kelley (2003) suggests $\eta = 0.1$. However,

he describes more sophisticated ways to choose this parameter. The smaller the forcing term η , the closer one gets to the "standard" Newton method. However, especially in the first nonlinear iterations, choosing a η that is too small leads to unnecessarily long computational times. The linear system is being solved with too much precision. See Box 6.7 for the pseudocode of the Newton-Arnoldi and Newton-GMRES.

Box 6.7: Timestep n of the Newton-Krylov methods with restart, Newton-Arnoldi and Newton-GMRES.

- (i) $k = 0$
- (ii) Enter the Newton loop
 - (1) Compute $\mathbf{r}_0 = -\mathcal{R}(\mathbf{x}^k)$
 - (2) Enter the Krylov loop
 - (a) Set $\beta = \|\mathbf{r}_0\|$ and $\mathbf{q}_1 = \mathbf{r}_0/\beta$.
 - (b) $j = 1$
 - (c) $\mathbf{a}_{j+1}^{(1)} = J\mathcal{R}(\mathbf{x}^k)\mathbf{q}_j$.
 - (d) Compute $h_{ij} = (\mathbf{q}_i, \mathbf{a}_{j+1}^{(i)})$ and compute $\mathbf{a}_{j+1}^{(i+1)} = \mathbf{a}_{j+1}^{(i)} - h_{ij}\mathbf{q}_i$ for $i = 1, \dots, j$.
 - (e) Compute $h_{j+1,j} = \|\mathbf{a}_{j+1}^{(j+1)}\|$ and set $\mathbf{q}_{j+1} = \mathbf{a}_{j+1}^{(j+1)}/h_{j+1,j}$.
 - (f) If using the Arnoldi method:
 - Solve the linear system $\mathbf{H}_j\boldsymbol{\eta} = \beta\mathbf{e}_1$ for $\boldsymbol{\eta}$
 - set $\Delta\mathbf{x}_j^* = \Delta\mathbf{x}_0^* + \mathbf{Q}_j\boldsymbol{\eta}$.
 - (g) Else if using the GMRES method:
 - Solve the linear least-squares problem $\|\beta\mathbf{e}_1 - \tilde{\mathbf{H}}_j\boldsymbol{\eta}\|$ for $\boldsymbol{\eta}$.
 - Set $\Delta\mathbf{x}_j^* = \Delta\mathbf{x}_0^* + \mathbf{Q}_j\boldsymbol{\eta}$.
 - (h) If $\|J\mathcal{R}(\mathbf{x}^k)\Delta\mathbf{x}_j^* + \mathcal{R}(\mathbf{x}^k)\| \leq \eta\|\mathcal{R}(\mathbf{x}^k)\|$ is not satisfied set $j = j + 1$.
 - (i) if $j > m$, restart the method going to Step (a). Else, go to Step (c).
 - (3) $\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}^k$.
 - (4) $k = k + 1$
 - (5) $\mathbf{r}_0 = -\mathcal{R}(\mathbf{x}^k)$
 - (6) If convergence has not been reached, $\|\mathbf{r}_0\| > \epsilon_1$, go to Step (2)

6.5.3 Extrapolation techniques with cycling

There is a vast literature on sequence acceleration/extrapolation methods (see [Brezinski and Zaglia \(2013\)](#) and [Sidi \(2017\)](#) for textbook treatments of this topic). One often deals with sequences and series in numerical analysis, applied mathematics, and engineering. They are produced by iterative methods, perturbation techniques, and approximation procedures depending on a parameter. Those sequences or series often converge so slowly that it is a severe drawback to their practical use. Convergence

acceleration methods present a solution and have been studied for many years and applied to various situations. They are based on the very natural idea of extrapolation. In many cases, they lead to the solution of unsolvable problems otherwise. Sequences of vectors can also be considered, with their dimension being substantial. Such sequences arise, for example, in the solution by fixed-point iterative methods of systems of linear or nonlinear algebraic equations.

An example of a scalar acceleration method is first presented to fix ideas. Let (S_n) be a sequence of numbers that converges to S . This sequence can be transformed into another, denoted (T_n) . For example, consider

$$T_n = \frac{S_n S_{n+2} - S_{n+1}^2}{S_{n+2} - 2S_{n+1} + S_n}, \quad n = 0, 1, \dots, \quad (6.132)$$

which corresponds to the Aitken Δ^2 process.

This expression can be obtained considering a transformation that would yield the limit of a geometric sequence from only three iterates, i.e., if one fits an exponential function

$$S + a\lambda^n, \quad (6.133)$$

the sequence transformation takes the horizontal asymptote of the exponential, S . For the geometrical interpretation of Aitken's Δ^2 method, see Figure 6.4.

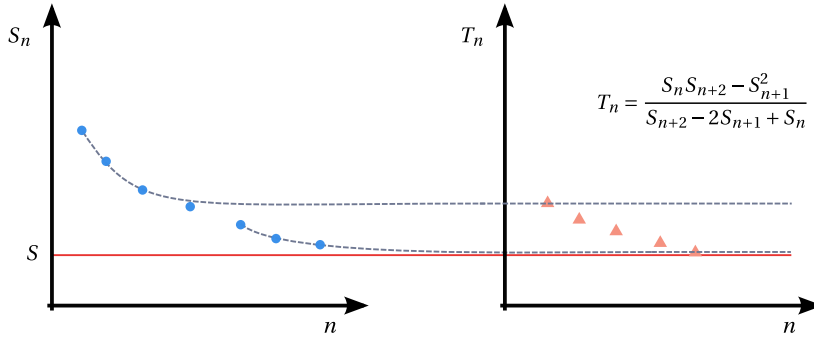


Figure 6.4: Geometrical interpretation of Aitken's Δ^2 method.

One can also show that if (S_n) goes to its limit S at a rate strictly greater than 1^1 , (T_n) does not have a better rate of convergence.

In practice, the sequence produced by Aitken's Δ^2 method tends to converge faster to the limit than (S_n) does. Very often, it is much cheaper to calculate (T_n) , which involves only the calculation of differences, one multiplication, and one division, than to calculate many more terms of the sequence (S_n) . Care must be taken, however, to avoid introducing errors due to insufficient precision when calculating the differences in the numerator and denominator of the expression.

There is, however, no universal sequence accelerator capable of accelerating all sequences. It is also the case that nonlinear transformations can even fail to converge or converge to a value other than the limit of the original sequence (Brezinski and Zaglia, 2013).

¹ (S_n) , $n \in \mathbb{N}$ converges linearly to S if there exists a number $\mu \in (0, 1)$ such that $\lim_{n \rightarrow \infty} \frac{|S_{n+1} - S|}{|S_n - S|} = \mu$.

According to [Brezinski and Zaglia \(2013\)](#), there is a very strong connection between sequence transformations and fixed point methods for solving $x = g(x)$, $g: \mathbb{R} \rightarrow \mathbb{R}$. The most well-known example of this connection is that between Aitken's Δ^2 process and Steffensen's method.

$$T_n = S_n - \frac{(S_{n+1} - S_n)^2}{S_{n+2} - 2S_{n+1} + S_n}, \quad n = 0, 1, \dots \quad \text{for Aitken's process} \quad (6.134)$$

and

$$x_{n+1} = x_n - \frac{(g(x_n) - x_n)^2}{g(g(x_n)) - 2g(x_n) + x_n}, \quad n = 0, 1, \dots \quad \text{for Steffensen's method.} \quad (6.135)$$

Turning to vector sequences and systems of nonlinear equations, let $F: (\mathbf{w}^k) \rightarrow (\mathbf{y}^k)$ be a vector extrapolation method defined by

$$\mathbf{y}^k = F(\mathbf{w}^k, \dots, \mathbf{w}^{k+m}), \quad n = 0, 1, \dots \quad (6.136)$$

For solving the fixed point problem $\mathbf{x} = \mathcal{S}(\mathbf{x})$ one can associate to it the iterative method

$$\mathbf{x}^{k+1} = F(\mathbf{x}^k, \mathcal{S}(\mathbf{x}^k), \dots, \mathcal{S}^m(\mathbf{x}^k)), \quad n = 0, 1, \dots \quad (6.137)$$

where $\mathcal{S}^{m+1}(\mathbf{x}) = \mathcal{S} \circ \mathcal{S}^m(\mathbf{x})$ and $\mathcal{S}^0(\mathbf{x}) = \mathbf{x}$. This approach is called full cycling or simply cycling. Conversely to any fixed point iteration of this form, one can associate a sequence transformation of the previous form. See [Box 6.8](#) for the general algorithm, excluding the extrapolation method.

Box 6.8: Timestep n of vector extrapolation with cycling.

- (i) Choose integers $n \geq 0$ and $k \leq 1$ and an initial vector $\mathbf{x}^k = \mathbf{x}_0^*$.
- (ii) Compute $\mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_{n+k+1}^*$ via $\mathbf{x}_{m+1}^* = \mathcal{S}(\mathbf{x}_m^*)$.
- (iii) Apply any of the four extrapolation methods, namely, MPE, RRE, MMPE, and SVD-MPE, to $\mathbf{x}_n^*, \mathbf{x}_{n+1}^*, \dots, \mathbf{x}_{n+k+1}^*$, obtaining $\mathbf{s}_{k,n}$.
- (iv) If $\mathbf{s}_{n,k}$ satisfies the accuracy test, stop.
Otherwise, set $\mathbf{x}_0^* = \mathbf{s}_{n,k}$ and go to Step (ii).

There is a variety of vector extrapolation methods, where the major two categories are polynomial methods and methods based on the ϵ -algorithm ([Brezinski and Zaglia, 2013](#); [Sidi, 2017](#)). In this presentation, only the first category is considered since the second requires a relatively large number of function evaluations per iteration, making it unsuitable for the present use-case ([Sidi, 2017](#)).

[Sidi \(2017\)](#) presents four different polynomial extrapolation methods. They all attempt to express the limit of the vector sequence as a linear combination of the p previous iterates, as follows

$$\mathbf{s} \approx \mathbf{s}_{k,m} = \sum_{j=0}^m \gamma_j \mathbf{w}^{k+j}, \quad (6.138)$$

where \mathbf{s} is the limit of the vector sequence. The methods to be presented next appear naturally when considering the vector sequence generated by

$$\mathbf{w}^{k+1} = \mathbf{T}\mathbf{w}^k + \mathbf{d}, \quad (6.139)$$

where $\mathbf{I} - \mathbf{T}$ is non-singular. It is tightly connected to the solution of linear systems of equations. Considering the minimal polynomial of \mathbf{T} with respect to $\Delta\mathbf{w}^k = \mathbf{w}^{k+1} - \mathbf{w}^k$ and $\boldsymbol{\epsilon}^k = \mathbf{w}^k - \mathbf{s}$ ², $P(\lambda)$,

$$P(\lambda) = \sum_{j=0}^i c_j \lambda^j, \quad c_i = 1, \quad (6.140)$$

where i is the degree of the polynomial, the limit of the sequence can be found exactly as

$$\mathbf{s} = \frac{\sum_{j=0}^i c_j \mathbf{w}^{k+j}}{\sum_{j=0}^i c_j}. \quad (6.141)$$

This can be derived considering the definition of $P(\lambda)$, $P(\mathbf{T})\boldsymbol{\epsilon}^k = \mathbf{0}$. Therefore,

$$\mathbf{0} = P(\mathbf{T})\boldsymbol{\epsilon}^k = \sum_{j=0}^i c_j \mathbf{T}^j \boldsymbol{\epsilon}^k = \sum_{j=0}^i c_j \boldsymbol{\epsilon}^{k+j}. \quad (6.142)$$

and so

$$\mathbf{0} = \sum_{i=0}^k c_i \boldsymbol{\epsilon}_{n+i} = \sum_{i=0}^k c_i \mathbf{x}_{n+i} - \left(\sum_{i=0}^k c_i \right) \mathbf{s}. \quad (6.143)$$

Solving this for \mathbf{s} , one obtains the desired result, provided $\sum_{j=0}^i c_j \neq 0$.

The coefficients of $P(\lambda)$ can be computed considering

$$\mathcal{W}^{i-1} \mathbf{c}' = -\Delta\mathbf{w}_{k+i}, \quad \mathbf{c}' = [c_0, c_1, \dots, c_{i-1}]^T, \quad (6.144)$$

where $\mathcal{W}^i = [\Delta\mathbf{w}^k, \dots, \Delta\mathbf{w}^{k+i}]$, since from the definition of $P(\lambda)$, one has

$$\mathbf{0} = P(\mathbf{T})\Delta\mathbf{w}^k = \sum_{j=0}^i c_j \mathbf{T}^j \Delta\mathbf{w}^k = \sum_{j=0}^i c_j \Delta\mathbf{w}^{k+j}. \quad (6.145)$$

The degree of $P(\lambda)$ can be as large as the dimension of \mathbf{w} . Hence, to be practical, the minimal polynomial extrapolation (MPE), the reduced rank extrapolation (RRE), the modified minimal extrapolation (MMPE), and the single-value decomposition, minimal polynomial extrapolation (SVD-MPE) all choose a polynomial of a lesser degree. The approximations corresponding to each extrapolation method are presented in what follows.

MPE Solve the overdetermined linear system $\mathcal{W}^{m-1} \mathbf{c}' = -\Delta\mathbf{w}_{k+m}$ in the least-squares sense for $\mathbf{c}' = [c_0, c_1, \dots, c_{m-1}]^T$. This amounts to solving the optimization problem

$$\min_{c_0, c_1, \dots, c_{m-1}} \left\| \sum_{j=0}^{m-1} c_j \Delta\mathbf{w}^{k+j} + \Delta\mathbf{w}^{k+m} \right\|_2 \quad (6.146)$$

which can also be expressed as

$$\min_{\mathbf{c}'} \left\| \mathcal{W}^{m-1} \mathbf{c}' + \mathbf{w}^{k+m} \right\|_2, \quad \mathbf{c}' = [c_0, c_1, \dots, c_{m-1}]^T. \quad (6.147)$$

With c_0, c_1, \dots, c_{k-1} available, set $c_m = 1$ and compute $\gamma_q = c_q / \sum_{j=0}^m c_j$, $q = 0, 1, \dots, m$, provided $\sum_{j=0}^m c_j \neq 0$.

²A polynomial $P(\lambda)$ is said to be minimal with respect to a vector \mathbf{a} , if $P(\mathbf{T})\mathbf{a} = \mathbf{0}$ and it is of least degree.

RRE Solve the overdetermined linear system $\mathcal{W}^m \boldsymbol{\gamma} = 0$ in the least-squares sense, subject to the constraint $\sum_{j=0}^m \gamma_j = 1$. This amounts to solving the optimization problem

$$\min_{\gamma_0, \gamma_1, \dots, \gamma_m} \left\| \sum_{j=0}^m \gamma_j \Delta \mathbf{w}^{k+j} \right\| \quad \text{subject to} \quad \sum_{j=0}^m \gamma_j = 1 \quad (6.148)$$

which can also be expressed as

$$\min_{\boldsymbol{\gamma}} \|\mathcal{W}^m \boldsymbol{\gamma}\|_2 \quad \text{subject to} \quad \sum_{j=0}^m \gamma_j = 1; \quad \boldsymbol{\gamma} = [\gamma_0, \gamma_1, \dots, \gamma_m]^T. \quad (6.149)$$

MMPE Consider a set of m linearly independent vectors \mathbf{q}_j , $j = 1, \dots, m$. Solve the linear system

$$\left(\mathbf{q}^j, \mathcal{W}^{m-1} \mathbf{c}^j \right) = - \left(\mathbf{q}^j, \Delta \mathbf{w}^{k+m} \right), \quad j = 1, \dots, m, \quad (6.150)$$

which can also be expressed as

$$\sum_{j=0}^{m-1} \left(\mathbf{q}_j, \Delta \mathbf{w}^{k+j} \right) \mathbf{c}_j = - \left(\mathbf{q}^j, \Delta \mathbf{w}^{k+p} \right), \quad j = 1, \dots, m. \quad (6.151)$$

This is, in fact, a system of m linear equations for the m unknowns c_0, c_1, \dots, c_{m-1} . With c_0, c_1, \dots, c_{m-1} available, set $c_m = 1$ and compute $\gamma_q = c_q / \sum_{j=0}^m c_j$, $i = 0, 1, \dots, m$, provided $\sum_{j=0}^m c_j \neq 0$.

SVD-MPE Solve the standard l_2 constrained minimization problem

$$\min_{\mathbf{c}} \|\mathcal{W}^m \mathbf{c}\|_2 \quad \text{subject to} \quad \|\mathbf{c}\|_2 = 1, \quad \mathbf{c} = [c_0, c_1, \dots, c_m]^T. \quad (6.152)$$

The solution \mathbf{c} is the right singular vector corresponding to the smallest singular value σ_{\min} of \mathcal{W}^m , i.e., $\mathcal{W}^{m*} \mathcal{W}^m \mathbf{c} = \sigma_{\min}^2 \mathbf{c}$, $\|\mathbf{c}\|_2 = 1$. It is assumed that σ_{\min} is simple so that \mathbf{c} is unique up to a multiplicative constant ϕ , $|\phi| = 1$.

With c_0, c_1, \dots, c_m available, compute $\gamma_q = c_q / \sum_{j=0}^m c_j$, $q = 0, 1, \dots, m$, provided $\sum_{j=0}^m c_j \neq 0$. The assumption that σ_{\min} is simple guarantees the uniqueness of the γ_i .

When $m = 1$, MPE, RRE, MMPE, and SVD-MPE can be regarded as generalizations of the Aitken Δ^2 -process to the vector case. Thus, when applied to the solution of a system of nonlinear equations using cycling

$$\mathbf{s}_{k,1} = \begin{cases} \mathbf{x}^k - \frac{(\Delta \mathbf{x}^k, \Delta \mathbf{x}^k)}{(\Delta \mathbf{x}^k, \Delta^2 \mathbf{x}^k)} \Delta \mathbf{x}^k & \text{for MPE,} \\ \mathbf{x}^k - \frac{(\Delta^2 \mathbf{x}^k, \Delta \mathbf{x}^k)}{(\Delta^2 \mathbf{x}^k, \Delta^2 \mathbf{x}^k)} \Delta \mathbf{x}^k & \text{for RRE,} \\ \mathbf{x}^k - \frac{(\mathbf{q}_1, \Delta \mathbf{x}^k)}{(\mathbf{q}_1, \Delta^2 \mathbf{x}^k)} \Delta \mathbf{x}^k & \text{for MMPE,} \\ \mathbf{x}^k - \frac{(\mathbf{g}_0, \Delta \mathbf{x}^k)}{(\mathbf{g}_0, \Delta^2 \mathbf{x}^k)} \Delta \mathbf{x}^k & \text{for SVD-MPE.} \end{cases} \quad (6.153)$$

Sidi (2017) also suggest cycling with frozen γ_i , where after some iterations the γ_i are frozen and reused henceforth. A parallel version of the full cycling procedure is also described.

Connection to Krylov subspace methods According to Sidi (2017), the so-called Krylov subspace methods are closely related to the vector extrapolation methods presented above. When the latter is applied to vector sequences obtained using fixed-point iterative methods to nonsingular linear systems of equations, they are mathematically equivalent. More precisely, the MPE and the RRE methods are mathematically equivalent to the methods of Arnoldi and generalized minimal residual (GMR).

However, Krylov subspace methods and extrapolation methods differ in their algorithmic aspects entirely: The only input of the former is a procedure that performs the matrix-vector multiplication without explicitly knowing the matrix coefficient matrix. The latter takes as their only input a vector sequence that results from a fixed-point iterative scheme without knowing the matrix coefficient matrix to know what the scheme is.

In Michler et al. (2005), a Krylov-subspace method is proposed in the context of FSI. However, as pointed out by Küttler and Wall (2009), the correct term for this approach should be instead a "Krylov-based vector extrapolation" method. The method proposed can be obtained by applying the RRE to the sequence of residuals computed as $\Delta \mathbf{r}_i^* = \mathbf{x}_i^* - \mathbf{x}^k$, where the subscript i concerns the internal loop of the vector extrapolation method, and whose limit is $\mathbf{0}$. Küttler and Wall (2009) argues that these residual differences have unfavorable numerical properties and should be avoided.

6.6 Multipoint iteration functions with memory

Multipoint iteration functions are rarer and are not thoroughly investigated in this exposition. One can mention the Airola-Nevanlinna family of methods (Fang and Saad, 2009) and Netwon-Krylov method that reuses the Krylov subspace from previous iterations (Sidi, 2017).

6.7 Conclusions

Table 6.1: Summary of the comparison between method for the solution methods of non-linear systems of equations. n here denotes the number of unknowns and m denotes depending on the context the number of previous iterates considered, the number of fixed point evaluations or the size of the Krylov subspace.

Method	Memory requirements	Nr function evaluations per iteration	Observations
Fixed-point iteration	$2 (n \times 1)$ vectors	1	<ul style="list-style-type: none"> •Often diverges. •Simplest method. •Memory efficient.
Underrelaxation	$2 (n \times 1)$ vectors	1	<ul style="list-style-type: none"> •Simple. •Improved stability over fixed-point. •Need to manually choose a relaxation parameter.
Aitken relaxation	$3 (n \times 1)$ vectors	1	<ul style="list-style-type: none"> •Very popular in FSI. •Dynamic relaxation. •Improved stability over fixed-point.
Broyden-like family (Fang and Saad, 2009)	$2 (n \times 1)$ vectors $2 (n \times (m - 1))$ matrices	1	<ul style="list-style-type: none"> •$\mathcal{O}(n^3)$ computation complexity for $m > 1$ (QR decomposition). •Low number of function evaluations •Superlinear convergence when $m = 1$.
Broyden's method (Kelley, 2003)	Up to $(m + 2)(n \times 1)$ matrices	1	<ul style="list-style-type: none"> •$\mathcal{O}(n)$ computation complexity. •Low storage. •Superlinear convergence.
Newton-Krylov	Up to $(m + 1) (n \times 1)$ vectors	$m + 1^*$	<ul style="list-style-type: none"> •Large number of iterations possible. •Popular for the solution of systems of nonlinear equations. •Quadratic convergence under appropriate conditions.
Cycling with vector extrapolation	$(m + 2) (n \times 1)$ vectors	$m + 1$	<ul style="list-style-type: none"> •Large number of function evaluations. •$\mathcal{O}(n^3)$ computational complexity (QR decomposition).

* The number of function evaluations in the Newton-Krylov methods will depend on how many iterations it will take for the inner loop to converge. There is a function evaluation per iteration of the inner loop.

Table 6.2: Summary of the update formulas for the solution methods of non-linear systems of equations.

Method	Update formula
Fixed-point	$\mathbf{x}^{k+1} = \mathbf{x}^k - \mathcal{R}^k$
Underrelaxation	$\mathbf{x}^{k+1} = \mathbf{x}^k - \omega \mathcal{R}^k,$ $0 < \omega < 1.$
Aitken relaxation	$\mathbf{x}^{k+1} = \mathbf{x}^k - \omega^{(k)} \mathcal{R}^k,$ $\omega^{(k)} = -\omega^{(k-1)} \frac{(\mathbf{r}^{(k)} - \mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}{(\mathbf{r}^{(k)} - \mathbf{r}^{(k-1)})^2}.$
Broyden's family	$\mathbf{x}^{k+1} = \mathbf{x}^k - \left(G_{\mathcal{R}}^{k-m} + (\mathcal{X}^k - G_{\mathcal{R}}^{k-m} \mathcal{Z}^k) \mathbf{V}^{kT} \right) \mathcal{R}^k,$ $\mathbf{V}^{kT} = \mathbf{M}^{k-1} \mathbf{N}^{kT},$ Type I: $M^k = \mathcal{X}^{kT} G_{\mathcal{R}}^k \mathcal{Z}^k, \quad N^{kT} = \mathcal{X}^{kT} G_{\mathcal{R}}^k,$ Type II: $M^k = \mathcal{Z}^{kT} \mathcal{Z}^k, \quad N^{kT} = \mathcal{Z}^{kT}$ \mathbf{V}^k is usually computed using <i>QR</i> -decomposition.
Anderson's family	$\mathbf{x}^{k+1} = \mathbf{x}^k - \left(-\beta \mathbf{I} + (\mathcal{X}^k + \beta \mathcal{Z}^k) \mathbf{V}^{kT} \right) \mathcal{R}^k,$ $\mathbf{V}^{kT} = \mathbf{M}^{k-1} \mathbf{N}^{kT},$ Type I: $M^k = \mathcal{X}^{kT} \mathcal{Z}^k, \quad N^{kT} = \mathcal{X}^{kT},$ Type II: $M^k = \mathcal{Z}^{kT} \mathcal{Z}^k, \quad N^{kT} = \mathcal{Z}^{kT}$ \mathbf{V}^k is usually computed using <i>QR</i> -decomposition.
Newton-Krylov	$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k$ $J_{\mathcal{R}}(\mathbf{x}^k) \Delta \mathbf{x}^k = -\mathcal{R}^k$ solved using a Krylov method to accuracy $\ J_{\mathcal{R}}(\mathbf{x}^k) \Delta \mathbf{x}_m^* + \mathcal{R}(\mathbf{x}^k)\ \leq \eta \ \mathcal{R}(\mathbf{x}^k)\ $
Cycling with vector extrapolation	$\mathbf{x}^{k+1} = \sum_{j=0}^m \gamma_j \mathcal{S}^{i+j}(\mathbf{x}^k),$ $\gamma_q = c_q / \sum_{j=0}^m c_j, \quad q = 0, 1, \dots, m,$ MPE: $\mathbf{c}' = [c_0, c_1, \dots, c_{m-1}]^T = \arg \min_{\mathbf{c}} \left\ \mathcal{S}^{m-1} \mathbf{c}' + \mathcal{S}^{i+m}(\mathbf{x}^k) \right\ _2,$ $c_m = 1,$ RRE: $\boldsymbol{\gamma} = [\gamma_0, \gamma_1, \dots, \gamma_m]^T = \arg \min_{\boldsymbol{\gamma}} \left\ \mathcal{S}^m \hat{\boldsymbol{\gamma}} \right\ _2,$ Subject to $\sum_{j=0}^m \gamma_j = 1.$ MMPE: $(\mathbf{q}^j, \mathcal{S}^{m-1} \mathbf{c}') = -(\mathbf{q}^j, \mathcal{S}^{i+m}(\mathbf{x}^k)), \quad j = 1, \dots, m,$ for a set of m linearly independent vectors $\mathbf{q}_j, j = 1, \dots, m.$ $\mathbf{c}' = [c_0, c_1, \dots, c_{m-1}]^T.$ SVD-MPE: $\mathbf{c} = \arg \min_{\mathbf{c}} \left\ \mathcal{S}^m \mathbf{c} \right\ _2 \quad \text{subject to } \ \mathbf{c}\ _2 = 1,$ $\mathbf{c} = [c_0, c_1, \dots, c_m]^T.$

Bibliography

Adam, L. and J.-P. Ponthot

2002a. Numerical simulation of viscoplastic and frictional heating during finite deformation of metal. Part I: Theory. *Journal of engineering mechanics*, 128(11):1215–1221. Publisher: American Society of Civil Engineers.

Adam, L. and J.-P. Ponthot

2002b. Numerical simulation of viscoplastic and frictional heating during finite deformation of metal. Part II: Applications. *Journal of engineering mechanics*, 128(11):1222–1232. Publisher: American Society of Civil Engineers.

Agelet de Saracibar, C.

1998. Numerical analysis of coupled thermomechanical frictional contact problems. Computational model and applications. *Archives of Computational Methods in Engineering*, 5(3):243–301.

Agelet de Saracibar, C., M. Cervera, and M. Chiumenti

1999. On the formulation of coupled thermoplastic problems with phase-change. *International Journal of Plasticity*, 15(1):1–34.

Argyris, J. H. and J. S. Doltsinis

1981. On the natural formulation and analysis of large deformation coupled thermomechanical problems. *Computer Methods in Applied Mechanics and Engineering*, 25(2):195–253.

Armero, F.

1999. Formulation and finite element implementation of a multiplicative model of coupled poro-plasticity at finite strains under fully saturated conditions. *Computer Methods in Applied Mechanics and Engineering*, 171:205–241.

Armero, F. and J. C. Simo

1992. A new unconditionally stable fractional step method for non-linear coupled thermomechanical problems. *International Journal for Numerical Methods in Engineering*, 35(4):737–766. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.1620350408>.

Armero, F. and J. C. Simo

1993. A priori stability estimates and unconditionally stable product formula algorithms for nonlinear coupled thermoplasticity. *International Journal of Plasticity*, 9(6):749–782.

- Badia, S., A. F. Martín, and R. Planas
2014. Block recursive LU preconditioners for the thermally coupled incompressible inductionless MHD problem. *Journal of Computational Physics*, 274:562–591.
- Belytschko, T. and R. Mullen
1976. Mesh partitions of explicit-implicit time integration. *Formulations and computational algorithms in finite element analysis*, Pp. 673–690. Publisher: MIT Press: New York.
- Belytschko, T. and R. Mullen
1978. Stability of explicit-implicit mesh partitions in time integration. *International Journal for Numerical Methods in Engineering*, 12(10):1575–1586. Publisher: Wiley Online Library.
- Belytschko, T., H.-J. Yen, and R. Mullen
1979. Mixed methods for time integration. *Computer Methods in Applied Mechanics and Engineering*, 17:259–275. Publisher: Elsevier.
- Blom, D.
2017. *Efficient numerical methods for partitioned fluid-structure interaction simulations*. Ph.D., Delft University of Technology, Netherlands.
- Blom, F. J.
1998. A monolithical fluid-structure interaction algorithm applied to the piston problem. *Computer methods in applied mechanics and engineering*, 167(3-4):369–391. Publisher: Elsevier.
- Borja, R., C. Tamagnini, and E. Alarcón
1998. Elastoplastic consolidation at finite strain part 2: finite element implementation and numerical examples. *Computer Methods in Applied Mechanics and Engineering*, 159:103–122.
- Brezinski, C. and M. R. Zaglia
2013. *Extrapolation methods: theory and practice*. Elsevier.
- Broyden, C. G.
1965. A class of methods for solving nonlinear simultaneous equations. *Mathematics of computation*, 19(92):577–593. Publisher: JSTOR.
- Carter, J. P. and J. R. Booker
1989. Finite element analysis of coupled thermoelasticity. *Computers & Structures*, 31(1):73–80.
- Causin, P., J.-F. Gerbeau, and F. Nobile
2005. Added-mass effect in the design of partitioned algorithms for fluid–structure problems. *Computer methods in applied mechanics and engineering*, 194(42-44):4506–4527. Publisher: Elsevier.
- Cervera, M., R. Codina, and M. Galindo
1996. On the computational efficiency and implementation of block-iterative algorithms for nonlinear coupled problems. *Engineering Computations*, 13(6):4–30. Publisher: MCB UP Ltd.

- Chen, K.
2005. *Matrix Preconditioning Techniques and Applications*, Cambridge Monographs on Applie. Cambridge University Press.
- Combescure, A. and A. Gravouil
2002. A numerical scheme to couple subdomains with different time-steps for predominantly linear transient analysis. *Computer Methods in Applied Mechanics and Engineering*, 191(11):1129–1157.
- Danowski, C.
2014. *Computational Modelling of Thermo-Structure Interaction with Application to Rocket Nozzles*. Ph.D., Technische Universität München, Germany.
- Danowski, C., V. Gravemeier, L. Yoshihara, and W. A. Wall
2013. A monolithic computational approach to thermo-structure interaction. *International Journal for Numerical Methods in Engineering*, 95(13):1053–1078. [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.4530](https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.4530).
- Degroote, J.
2010. *Development of algorithms for the partitioned simulation of strongly coupled fluid-structure interaction problems*. PhD Thesis, Ghent University. ISBN: 9789085783442.
- Degroote, J., P. Bruggeman, R. Haelterman, and J. Vierendeels
2008. Stability of a coupling technique for partitioned solvers in FSI applications. *Computers & Structures*, 86(23):2224–2234.
- Dennis, J. and R. Schnabel
1996. *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Classics in Applied Mathematics. Society for Industrial and Applied Mathematics.
- Dettmer, W. and D. Perić
2006. A computational framework for fluid–structure interaction: Finite element formulation and applications. *Computer Methods in Applied Mechanics and Engineering*, 195(41):5754–5779.
- Dittmann, M.
2017. *Isogeometric analysis and hierarchical refinement for multi-field contact problems*. Ph.D.
- Dittmann, M., M. Franke, İ. Temizer, and C. Hesch
2014. Isogeometric Analysis and thermomechanical Mortar contact problems. *Computer Methods in Applied Mechanics and Engineering*, 274:192–212.
- Erbts, P. and A. Düster
2012. Accelerated staggered coupling schemes for problems of thermoelasticity at finite strains. *Computers & Mathematics with Applications*, 64(8):2408–2430.
- Erbts, P., S. Hartmann, and A. Düster
2015. A partitioned solution approach for electro-thermo-mechanical problems. *Archive of Applied Mechanics*, 85(8):1075–1101.
- Fang, H.-r. and Y. Saad
2009. Two classes of multisecant methods for nonlinear acceleration. *Numerical linear algebra with applications*, 16(3):197–221. Publisher: Wiley Online Library.

- Farhat, C., P. Geuzaine, and G. Brown
2003. Application of a three-field nonlinear fluid–structure formulation to the prediction of the aeroelastic parameters of an F-16 fighter. *Computers & Fluids*, 32:3–29.
- Farhat, C. and M. Lesoinne
2000. Two efficient staggered algorithms for the serial and parallel solution of three-dimensional nonlinear transient aeroelastic problems. *Computer methods in applied mechanics and engineering*, 182(3-4):499–515. Publisher: Elsevier.
- Farhat, C., M. Lesoinne, P. Stern, and S. Lanteri
1997. High performance solution of three-dimensional nonlinear aeroelastic problems via parallel partitioned algorithms: methodology and preliminary results. *Advances in Engineering Software*, 28:43–61.
- Farhat, C., K. Park, and Y. Dubois-Pelerin
1991. An unconditionally stable staggered algorithm for transient finite element analysis of coupled thermoelastic problems. *Applied Mechanics and Engineering*, 85:349–365.
- Farhat, C., A. Rallu, K. G. Wang, and T. Belytschko
2010. Robust and provably second-order explicit-explicit and implicit-explicit staggered time-integrators for highly non-linear compressible fluid-structure interaction problems. *International Journal for Numerical Methods in Engineering*, 84:73–107.
- Farhat, C., K. Zee, and P. Geuzaine
2006. Provably second-order time-accurate loosely-coupled solution algorithms for transient nonlinear computational aeroelasticity. *Computer Methods in Applied Mechanics and Engineering*, 195:1973–2001.
- Felder, S., N. Kopic-Osmanovic, H. Holthusen, T. Brepols, and S. Reese
2021. Thermo-mechanically coupled gradient-extended damage-plasticity modeling of metallic materials at finite strains. *International Journal of Plasticity*, P. 103142.
- Felippa, C. and T. L. Geers
1988. Partitioned analysis for coupled mechanical systems. *Engineering Computations*, 5:123–133.
- Felippa, C. A. and K. C. Park
1980. Staggered transient analysis procedures for coupled mechanical systems: Formulation. *Computer Methods in Applied Mechanics and Engineering*, 24(1):61–111.
- Felippa, C. A., K. C. Park, and C. Farhat
2001. Partitioned analysis of coupled mechanical systems. *Computer Methods in Applied Mechanics and Engineering*, 190(24):3247–3270.
- Fernández, M. A., J.-F. Gerbeau, and C. Grandmont
2006. A projection semi-implicit scheme for the coupling of an elasticstructure with an incompressible fluid. *International Journal for Numerical Methods in Engineering*.
- Förster, C.
2007. Robust methods for fluid-structure interaction with stabilised finite elements.

- Förster, C., W. A. Wall, and E. Ramm
2007. Artificial added mass instabilities in sequential staggered coupling of nonlinear structures and incompressible viscous flows. *Computer methods in applied mechanics and engineering*, 196(7):1278–1293. Publisher: Elsevier.
- Gatzhammer, B.
2014. *Efficient and Flexible Partitioned Simulation of Fluid-Structure Interactions*. Dissertation, Technische Universität München, München.
- Gee, M. W., U. Küttler, and W. A. Wall
2011. Truly monolithic algebraic multigrid for fluid–structure interaction. *International Journal for Numerical Methods in Engineering*, 85(8):987–1016. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.3001>.
- Gillard, J.
2019. *An Efficient Partitioned Coupling Scheme for Tire Hydroplaning Analysis*. Ph.D., Technische Universität München, München.
- Gitterle, M.
2012. *A dual mortar formulation for finite deformation frictional contact problems including wear and thermal coupling*. Ph.D., Technische Universität München.
- Glaser, S.
1992. *Gekoppelte thermomechanische Berechnung duennwandiger Strukturen mit der Methode der Finiten Elemente*. Ph.D., Institute fuer Statik und Dynamik der Luft- und Raumfahrtkonstruktionen, University of Stuttgart.
- Haelterman, R., J. Degroote, D. Van Heule, and J. Vierendeels
2009. The Quasi-Newton Least Squares Method: A New and Fast Secant Method Analyzed for Linear Systems. *SIAM Journal on Numerical Analysis*, 47(3):2347–2368. _eprint: <https://doi.org/10.1137/070710469>.
- Hansen, G.
2011. A Jacobian-free Newton Krylov method for mortar-discretized thermomechanical contact problems. *Journal of Computational Physics*, 230(17):6546–6562.
- Heil, M.
2004. An efficient solver for the fully coupled solution of large-displacement fluid–structure interaction problems. *Computer Methods in Applied Mechanics and Engineering*, 193(1):1–23.
- Hesch, C. and P. Betsch
2011. Energy-momentum consistent algorithms for dynamic thermomechanical problems—Application to mortar domain decomposition problems. *International Journal for Numerical Methods in Engineering*, 86(11):1277–1302. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.3095>.
- Holt, M. and N. Yanenko
2012. *The Method of Fractional Steps: The Solution of Problems of Mathematical Physics in Several Variables*. Springer Berlin Heidelberg.

- Holzapfel, G. A. and J. C. Simo
1996. Entropy elasticity of isotropic rubber-like solids at finite strains. *Computer Methods in Applied Mechanics and Engineering*, 132(1):17–44.
- Hron, J. and S. Turek
2006. A Monolithic FEM/Multigrid Solver for an ALE Formulation of Fluid-Structure Interaction with Applications in Biomechanics. In *Fluid-Structure Interaction*., Lecture Notes in Computational Science and Engineering, vol 53.
- Hübner, B., E. Walhorn, and D. Dinkler
2004. A monolithic approach to fluid–structure interaction using space–time finite elements. *Computer Methods in Applied Mechanics and Engineering*, 193(23):2087–2104.
- Hüeber, S. and B. I. Wohlmuth
2009. Thermo-mechanical contact problems on non-matching meshes. *Computer Methods in Applied Mechanics and Engineering*, 198(15):1338–1350.
- Hughes, T. J. and W. Liu
1978. Implicit-explicit finite elements in transient analysis: stability theory.
- Ibrahimbegovic, A. and L. Chorfi
2002. Covariant principal axis formulation of associated coupled thermoplasticity at finite strains and its numerical implementation. *International Journal of Solids and Structures*, 39(2):499–528.
- Irons, B. M. and R. C. Tuck
1969. A version of the Aitken accelerator for computer iteration. *International Journal for Numerical Methods in Engineering*, 1(3):275–277. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.1620010306>.
- Jha, B. and R. Juanes
2007. A locally conservative finite element framework for the simulation of coupled flow and reservoir geomechanics. *Acta Geotechnica*, 2:139–153.
- Johansson, L. and A. Klarbring
1993. Thermoelastic frictional contact problems: Modelling, finite element approximation and numerical realization. *Computer Methods in Applied Mechanics and Engineering*, 105(2):181–210.
- Joosten, M. M., W. G. Dettmer, and D. Perić
2009. Analysis of the block Gauss–Seidel solution procedure for a strongly coupled model problem with reference to fluid–structure interaction. *International Journal for Numerical Methods in Engineering*, 78(7):757–778. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.2503>.
- Kelley, C.
2003. *Solving Nonlinear Equations with Newton's Method*, Fundamentals of Algorithms. Society for Industrial and Applied Mathematics.
- Kim, J., H. Tchelepi, and R. Juanes
2011a. Stability and convergence of sequential methods for coupled flow and geomechanics: Drained and undrained splits. *Computer Methods in Applied Mechanics and Engineering*, 200:2094–2116.

- Kim, J., H. Tchelepi, and R. Juanes
2011b. Stability and convergence of sequential methods for coupled flow and geomechanics: Fixed-stress and fixed-strain splits. *Computer Methods in Applied Mechanics and Engineering*, 200:1591–1606.
- Klöppel, T., A. Popp, U. Küttler, and W. A. Wall
2011. Fluid–structure interaction for non-conforming interfaces based on a dual mortar formulation. *Computer Methods in Applied Mechanics and Engineering*, 200(45):3111–3126.
- Küttler, U., M. Gee, C. Förster, A. Comerford, and W. A. Wall
2010. Coupling strategies for biomedical fluid–structure interaction problems. *International Journal for Numerical Methods in Biomedical Engineering*, 26(3-4):305–321. [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/cnm.1281](https://onlinelibrary.wiley.com/doi/pdf/10.1002/cnm.1281).
- Küttler, U. and W. A. Wall
2008. Fixed-point fluid–structure interaction solvers with dynamic relaxation. *Computational Mechanics*, 43(1):61–72.
- Küttler, U. and W. A. Wall
2009. Vector Extrapolation for Strong Coupling Fluid-Structure Interaction Solvers. *Journal of Applied Mechanics*, 76(2).
- Lenarda, P. and M. Paggi
2016. A geometrical multi-scale numerical method for coupled hygro-thermo-mechanical problems in photovoltaic laminates. *Computational Mechanics*, 57(6):947–963.
- Lewis, R. and Y. Sukirman
1993. Finite element modelling of three-phase flow in deforming saturated oil reservoirs. *International Journal for Numerical and Analytical Methods in Geomechanics*, 17(8):577–598. Publisher: Wiley Online Library.
- Lin, P. T., J. N. Shadid, R. S. Tuminaro, M. Sala, G. L. Hennigan, and R. P. Pawlowski
2010. A parallel fully coupled algebraic multilevel preconditioner applied to multiphysics PDE applications: Drift-diffusion, flow/transport/reaction, resistive MHD. *International Journal for Numerical Methods in Fluids*, 64(10-12):1148–1179. [tex.eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/fld.2402](https://onlinelibrary.wiley.com/doi/pdf/10.1002/fld.2402).
- Matthies, H. and J. Steindorf
2003a. Partitioned Strong Coupling Algorithms for Fluid-Structure-Interaction. *Computers & Structures*, 81:805–812.
- Matthies, H. and J. Steindorf
2003b. Strong Coupling Methods.
- Mayr, M., T. Klöppel, W. A. Wall, and M. W. Gee
2015. A Temporal Consistent Monolithic Approach to Fluid-Structure Interaction Enabling Single Field Predictors. *SIAM Journal on Scientific Computing*, 37(1):B30–B59. Publisher: Society for Industrial and Applied Mathematics.
- Mayr, M., M. H. Noll, and M. W. Gee
2020. A hybrid interface preconditioner for monolithic fluid–structure interaction solvers. *Advanced Modeling and Simulation in Engineering Sciences*, 7(1):15.

- Michler, C.
2005. *Efficient numerical methods for fluid-structure interaction*. Ph.D., Delft University of Technology, Netherlands. ISBN: 90-9019533-5.
- Michler, C., E. H. v. Brummelen, S. J. Hulshoff, and R. d. Borst
2003. The relevance of conservation for stability and accuracy of numerical methods for fluid–structure interaction. *Computer Methods in Applied Mechanics and Engineering*, 192(37):4195–4215.
- Michler, C., S. Hulshoff, E. Van Brummelen, and R. De Borst
2004. A monolithic approach to fluid–structure interaction. *Computers & fluids*, 33(5-6):839–848. Publisher: Elsevier.
- Michler, C., E. H. van Brummelen, and R. de Borst
2005. An interface Newton–Krylov solver for fluid–structure interaction. *International Journal for Numerical Methods in Fluids*, 47(10-11):1189–1195. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/flf.850>.
- Miehe, C.
1995a. Entropic thermoelasticity at finite strains. Aspects of the formulation and numerical implementation. *Computer Methods in Applied Mechanics and Engineering*, 120(3):243–269.
- Miehe, C.
1995b. A theory of large-strain isotropic thermoplasticity based on metric transformation tensors. *Archive of Applied Mechanics*, 66(1):45–64.
- Mikelić, A. and M. F. Wheeler
2013. Convergence of iterative coupling for coupled flow and geomechanics. *Computational Geosciences*, 17(3):455–461. Publisher: Springer.
- Miller, B. A.
2015. *Loosely Coupled Time Integration of Fluid- Thermal-Structural Interactions in Hypersonic Flows*. Ph.D., Ohio State University.
- Neishlos, H.
1983. Finite-Element Mesh Partitioning for Time Integration of Transient Problems. In *Numerical Solution of Partial Differential Equations: Theory, Tools and Case Studies: Summer Seminar Series Held at CSIR, Pretoria, February 8–10, 1982*, D. P. Laurie, ed., Pp. 225–245. Basel: Birkhäuser Basel.
- Netz, T.
2013. *High-order space and time discretization scheme applied to problems of finite thermo-viscoelasticity*. Ph.D., Institute of Applied Mechanics, Clausthal University of Technology.
- Novascone, S. R., B. W. Spencer, J. D. Hales, and R. L. Williamson
2015. Evaluation of coupling approaches for thermomechanical simulations. *Nuclear Engineering and Design*, 295:910–921.
- Oancea, V. G. and T. A. Laursen
1997. A finite element formulation of thermomechanical rate-dependent frictional sliding. *International Journal for Numerical Methods in Engineering*, 40(23):4275–4311. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/%28SICI%291097-0207%2819971215%2940%3A23%3C4275%3A%3AAID-NME257%3E3.0.CO%3B2-K>.

- Pantuso, D., K.-J. Bathe, and P. A. Bouzinov
2000. A finite element procedure for the analysis of thermo-mechanical solids in contact. *Computers & Structures*, 75(6):551–573.
- Park, K.
1983. Stabilization of partitioned solution procedure for pore fluid-soil interaction analysis. *International Journal for Numerical Methods in Engineering*, 19:1669–1673.
- Park, K., C. Felippa, and J. DeRuntz
1977. Stabilization of staggered solution procedures for fluid-structure interaction analysis. *Computational methods for fluid-structure interaction problems*, 26(94-124):51. Publisher: ASME New York.
- Piperno, S.
1997. Explicit/implicit fluid/structure staggered procedures with a structural predictor and fluid subcycling for 2D inviscid aeroelastic simulations. *International Journal for Numerical Methods in Fluids*, 25:1207–1226.
- Piperno, S. and C. Farhat
2001. Partitioned procedures for the transient solution of coupled aeroelastic problems—Part II: energy transfer analysis and three-dimensional applications. *Computer methods in applied mechanics and engineering*, 190(24-25):3147–3170. Publisher: Elsevier.
- Piperno, S., C. Farhat, and B. Larrouturou
1995. Partitioned procedures for the transient solution of coupled aroelastic problems Part I: Model problem, theory and two-dimensional application. *Computer Methods in Applied Mechanics and Engineering*, 124:79–112.
- Rothe, S., P. Erbs, A. Düster, and S. Hartmann
2015. Monolithic and partitioned coupling schemes for thermo-viscoplasticity. *Computer Methods in Applied Mechanics and Engineering*, 293:375–410.
- Saetta, A. and R. Vitaliani
1992. Unconditionally convergent partitioned solution procedure for dynamic coupled mechanical systems. *International Journal for Numerical Methods in Engineering*, 33:1975–1996.
- Scheufele, K.
2018. *Coupling schemes and inexact Newton for multi-physics and coupled optimization problems*. PhD Thesis, Universität Stuttgart, Stuttgart.
- Seitz, A.
2019. *Computational Methods for Thermo-Elasto-Plastic Contact*. Ph.D., Technische Universität München, Germany.
- Seitz, A., W. A. Wall, and A. Popp
2018. A computational approach for thermo-elasto-plastic frictional contact based on a monolithic formulation using non-smooth nonlinear complementarity functions. *Advanced Modeling and Simulation in Engineering Sciences*, 5(1):5.
- Shadid, J., R. Pawlowski, J. Banks, L. Chacón, P. Lin, and R. Tuminaro
2010. Towards a scalable fully-implicit fully-coupled resistive MHD formulation with stabilized FE methods. *Journal of Computational Physics*, 229(20):7649–7671.

- Sidi, A.
2017. *Vector extrapolation methods with applications*. SIAM.
- Simo, J. C. and F. Armero
1992. Recent Advances in the Numerical Analysis and Simulation of Thermoplasticity at Finite Strains.
- Simo, J. C. and C. Miehe
1992. Associative coupled thermoplasticity at finite strains: Formulation, numerical analysis and implementation. *Computer Methods in Applied Mechanics and Engineering*, 98(1):41–104.
- Smith, B., P. Bjorstad, and W. Gropp
2004. *Domain Decomposition: Parallel Multilevel Methods for Elliptic Partial Differential Equations*. Cambridge University Press.
- Tezduyar, T. E., S. Sathe, R. Keedy, and K. Stein
2006. Space-time finite element techniques for computation of fluid–structure interactions. *Computer methods in applied mechanics and engineering*, 195(17-18):2002–2027. Publisher: Elsevier.
- Torii, R., M. Oshima, T. Kobayashi, K. Takagi, and T. E. Tezduyar
2006. Computer modeling of cardiovascular fluid–structure interactions with the deforming-spatial-domain/stabilized space-time formulation. *Computer Methods in Applied Mechanics and Engineering*, 195(13-16):1885–1895. Publisher: Elsevier.
- Traub, J. F.
1982. *Iterative methods for the solution of equations*, volume 312. American Mathematical Soc.
- Uekermann, B. W.
2016. *Partitioned fluid-structure interaction on massively parallel systems*. PhD Thesis, Technische Universität München.
- Verdugo, F. and W. A. Wall
2016. Unified computational framework for the efficient solution of n-field coupled problems with monolithic schemes. *Computer Methods in Applied Mechanics and Engineering*, 310:335–366.
- Vierendeels, J., L. Lanoye, J. Degroote, and P. Verdonck
2007. Implicit coupling of partitioned fluid–structure interaction problems with reduced order models. *Computers & structures*, 85(11-14):970–976. Publisher: Elsevier.
- Wall, W. A., S. Genkinger, and E. Ramm
2007. A strong coupling partitioned approach for fluid–structure interaction with free surfaces. *Computers & Fluids*, 36(1):169–183.
- Wendt, G., P. Erbs, and A. Düster
2015. Partitioned coupling strategies for multi-physically coupled radiative heat transfer problems. *Journal of Computational Physics*, 300:327–351.

White, J. A. and R. Borja

2008. Stabilized low-order finite elements for coupled solid-deformation/fluid-diffusion and their application to fault zone transients. *Computer Methods in Applied Mechanics and Engineering*, 197:4353–4366.

Wriggers, P., C. Miehe, M. Kleiber, and J. C. Simo

1992. On the coupled thermomechanical treatment of necking problems via finite element methods. *International Journal for Numerical Methods in Engineering*, 33(4):869–883. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.1620330413>.

Zavarise, G., P. Wriggers, E. Stein, and B. A. Schrefler

1992. Real contact mechanisms and finite element formulation—a coupled thermomechanical approach. *International Journal for Numerical Methods in Engineering*, 35(4):767–785. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.1620350409>.

Zhang, Q. and T. Hisada

2004. Studies of the strong coupling and weak coupling methods in FSI analysis. *International Journal for Numerical Methods in Engineering*, 60(12):2013–2029. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.1034>.

Zienkiewicz, O., D. K. Paul, and A. Chan

1988. Unconditionally stable staggered solution procedure for soil-pore fluid interaction problems. *International Journal for Numerical Methods in Engineering*, 26:1039–1055.