



Thesis submitted to obtain the title of
Doctor of Philosophy

Doctoral School of Engineering Science
Field: Computer Science

Real-time Soft Tissue Modelling on GPU for Medical Simulation

Prepared by Olivier COMAS at
INRIA Lille, SHAMAN Team and CSIRO ICT Brisbane, AEHRC

DRAFT - Mon 30th of August 2010 at 21:53

Jury:

<i>Reviewers:</i>	Thomas Sangild SØRENSEN	- University of Aarhus, Denmark
	Bernard	- INRIA (Shaman)
<i>Advisor:</i>	Stéphane COTIN	- INRIA (Shaman), France
<i>President:</i>	Bernard	- INRIA (Shaman)
<i>Examinators:</i>	Bernard	- INRIA (Shaman)
	Bernard	- INRIA (Shaman)

Contents

Contents	i
I Introduction	1
1 Medical simulation	3
1.1 General context	4
1.1.1 Simulators for everything: we started by simulating what is simple	4
1.1.2 As complexity increases, applications to medical training, patient-specific planning and per-operative guidance	4
1.2 Challenges	4
1.2.1 Photorealistic rendering	4
1.2.2 Accurate haptic feedback	4
1.2.3 Physically realist organ modelling in real-time	4
2 Background in continuum mechanics for soft-tissue modelling	5
2.1 Introduction	6
2.2 Description of motion	6
2.2.1 Lagrangian description	7
2.2.2 Eulerian description	7
2.2.3 Displacement field	8
2.3 Analysis of deformation	8
2.3.1 Deformation gradient tensor	8
2.3.2 Change of volume	9
2.3.3 Change of surface	10
2.3.4 Volumetric and isochoric components	10
2.4 Strain measures	11
2.4.1 Cauchy-Green deformation tensors	12
2.4.2 Green-Lagrange strain tensor	12
2.4.3 Cauchy and Euler tensor	13
2.4.4 Principal strains and invariants	14
2.4.5 Infinitesimal strain tensor	14
2.5 Stress	15
2.5.1 Cauchy stress	15
2.5.2 First Piola-Kirchhoff stress tensor	17
2.5.3 Second Piola-Kirchhoff stress tensor	18
2.5.4 Principal stresses and invariants	18
2.6 Constitutive equations	18
2.6.1 Elasticity	19

2.6.2	Linear materials: generalised Hooke's law	20
2.6.3	Orthotropic materials	20
2.6.4	Isotropic materials	21
2.6.5	Non-linear materials	22
2.6.6	Viscoelastic materials	23
2.7	Tissue characterisation	24
3	A practical approach of the finite element method	25
3.1	Introduction	26
3.1.1	A numerical method	26
3.1.2	The basic ideas of FEM	26
3.2	Discretisation	27
3.2.1	Meshing process	27
3.2.2	Solution interpolation	28
3.2.3	Natural coordinates	30
3.2.4	Geometry interpolation	31
3.2.5	A particular case: isoparametric elements	32
3.3	Derivation of element equations	33
3.3.1	Strong and weak forms	33
3.3.2	Time dependence	34
3.3.3	Dynamic system of equations	34
3.3.4	Static system of equations	35
3.3.5	A few words on the matrices involved	35
3.4	Assembly of element equations	37
3.5	Solution of global problem	38
3.5.1	Explicit time integration	38
3.5.2	Implicit time integration	40
3.5.3	Static solutions	41
3.5.4	Solvers	41
II	Solid organs modelling	47
4	Modelling the deformation of solid objects in real-time	49
4.1	Introduction: the problematic	50
4.2	Techniques based on geometry	51
4.2.1	Free-form deformation	51
4.2.2	Shape matching	51
4.3	Techniques relying on physics	53
4.3.1	3D Chainmail algorithm	53
4.3.2	Modal analysis	54
4.3.3	Mass-spring model	57
4.4	Techniques based on continuum mechanics	59
4.4.1	The finite element method	59

4.4.2	Meshless methods	67
5	The total Lagrangian explicit dynamics (TLED) algorithm	71
5.1	Description of the TLED algorithm	72
5.1.1	Key ideas	72
5.1.2	Computation of element nodal forces	73
5.1.3	Computation of node displacements	76
5.2	Anisotropic and viscoelastic constitutive equations	78
5.2.1	An extension to the TLED algorithm	78
5.2.2	Visco-hyperelasticity	79
5.2.3	Hyperelastic response	80
5.2.4	Recapitulation	82
5.3	Constitutive update procedure for explicit analyses	82
5.3.1	Stress update equations	82
5.3.2	State variable update equations	83
5.3.3	Summary	84
6	GPU implementation of TLED	85
6.1	Summary of the TLED formulation	86
6.2	General-purpose computation on graphics processing units	87
6.2.1	Goal and motivation	87
6.2.2	Programming languages for GPUs	88
6.3	Implementation into SOFA	91
6.3.1	SOFA, an open source simulation framework	91
6.3.2	CUDA description	92
6.3.3	First implementation: scatter as a gather	94
6.3.4	Second implementation: a better use of shared memory	97
6.3.5	Third implementation: atomic writes	98
6.4	Results	98
6.4.1	Pure shear of a cube	99
6.4.2	Compression of a cube	100
6.4.3	GPU performance	101
6.4.4	Simulation of liver deformation	103
6.5	Discussion	103
6.5.1	Critical time step	103
6.5.2	Handling contacts	104
III	Hollow organs modelling	111
7	State of art: hollow structures	113
7.1	Expose the problem: models presented for solid organs are insufficient. How to model hollow organs?	114
7.1.1	Colonoscopy simulator project: needs for colon	114

7.1.2	Cataract surgery, stenosis: other needs for implants/blood vessels	114
7.2	Non-physic approaches (computer graphics stuff)	114
7.3	Physically accurate approches (plates/shells)	114
8	Mass-spring model for colon implemented on GPU	115
8.1	Colon is big, we need a fast model	116
8.2	Limitations (difficult to link it to actual physical properties)	116
9	A co-rotational triangular shell model (ISBMS 2010)	117
9.1	Model description	118
9.2	Validation	118
9.3	Application to implant deployment simulation in cataract surgery . .	118
10	'Shell meshing' technique (MICCAI 2010)	119
10.1	State of art: reconstruction/simplification	120
10.2	Our method	120
IV	Conclusion	121
11	Summary	123
12	Discussion and perspectives (Interaction solid/hollow organs)	125
A	Tensors	127
B	The weighted residual method	129
	References	131

Part I

Introduction

Draft Version

CHAPTER 1

MEDICAL SIMULATION

A short abstract for the upcoming chapter

Draft Version

1.1 General context

- 1.1.1 Simulators for everything: we started by simulating what is simple
- 1.1.2 As complexity increases, applications to medical training, patient-specific planning and per-operative guidance

1.2 Challenges

- 1.2.1 Photorealistic rendering
- 1.2.2 Accurate haptic feedback
- 1.2.3 Physically realist organ modelling in real-time

CHAPTER 2

BACKGROUND IN CONTINUUM MECHANICS FOR SOFT-TISSUE MODELLING

As seen in the previous chapter, realistic modelling of organ deformation is a challenging research field that opens the door to new clinical applications including: medical training and rehearsal systems, patient-specific planning of surgical procedures and per-operative guidance based on simulation. In all these cases the clinician needs fast updates of the deformation model to obtain a real-time display of the computed deformations. If for medical training devices the haptic feedback from touching organs merely needs to feel real, the accuracy of the information provided to the clinician in the cases of planning or per-operative guidance is crucial. Therefore, a substantial comprehension of the mechanics involved and a knowledge of the physical properties of anatomical structures are both mandatory in our quest to realistically model the deformation of organs. This chapter will start by introducing the main concepts of continuum mechanics that are fundamental to study the mechanical response of organs. It will then present mathematical models able to describe the different mechanical aspects of materials and briefly discuss the mechanical characterisation of tissues.

2.1 Introduction

In our everyday life, matter appears smooth and continuous: from the wood used to build your desk to the water you drink. But this is just illusion. The concept that matter is composed of discrete units has been around for millennia. In fact, we know with certainty that our world is composed of microscopic atoms and molecules separated by empty space since the beginning of the twentieth century (Lautrup, 2005). However, certain physical phenomena can be predicted with theories that pay no attention to the molecular structure of materials. Consider for instance the deformation of the horizontal board of a bookshelf under the weight of books. The bending of the shelf can be modelled without considering its molecular composition. The branch of physics in which materials are treated as continuous is known as *continuum mechanics*. Continuum mechanics studies the response of materials to different loading conditions. In this theory, matter is assumed to exist as a continuum, meaning that the matter in the body is continuously distributed and fills the entire region of space it occupies (Lai et al., 1996). This assumption is generally valid if the length scales of interest are large compared with the length scales of discrete molecular structure, but whether the approximation of continuum mechanics is actually justified in a given situation is merely a matter of experimental test.

Modelling anatomical structures requires an understanding of the deformation and stresses caused by the different interactions that occur during medical procedures. A sufficient knowledge of continuum mechanics is therefore essential to follow the rest of this manuscript. Continuum mechanics can be divided into two main parts: general principles common to all media (analysis of deformation, strain and stress concepts) and constitutive equations defining idealised materials. This chapter will not only deal with those two aspects but will also discuss the different aspects of the mechanical behaviour that we find in anatomical structures. Note that this chapter will mostly follow the notation used by Bonnet and Wood (1997) and Reddy (2007). The interested reader may refer to these books for more details.

2.2 Description of motion

Let us consider a body B of known geometry in a three-dimensional Euclidian space \mathbb{R}^3 . For a given geometry and loading, B will undergo a set of macroscopic changes which is called *deformation*. The region of space occupied by the body at a given time t is termed a *configuration*. A change in the configuration of a continuum body results in a displacement. The displacement of a body has two components: a rigid-body displacement and a deformation. A rigid-body displacement consists of a simultaneous translation and rotation of the body without changing its shape or size. Deformation implies the change in shape and/or size of the body from an initial configuration noted κ_0 to a new configuration κ called the *current* or *deformed configuration*.

Let us now consider a given particle of the body that we call X . What we will call particle in the following is in fact an infinitesimal volume of material. We denote

the position it occupies in the initial configuration \mathbf{X} and note its position in the deformed configuration \mathbf{x} , both expressed in the chosen frame of reference. The mapping χ defined as the following:

$$\chi : \mathcal{B}_{\kappa_0} \rightarrow \mathcal{B}_\kappa \quad (2.1)$$

is called the *deforming mapping* of the body \mathcal{B} from κ_0 to κ . When analysing the deformation of a continuous body, it is necessary to describe the evolution of configurations through time. Its mathematical description follows one of the two approaches: the material description or the spatial description. The material description is known as *Lagrangian description* whereas the spatial description is called *Eulerian description*. These two approaches are detailed next.

2.2.1 Lagrangian description

In the Lagrangian description, the position and physical properties of the particles are referred to a reference configuration κ_R , often chosen to be the undeformed configuration κ_0 . Thus, the current coordinates ($\mathbf{x} \in \kappa$) are expressed in terms of the reference coordinates ($\mathbf{X} \in \kappa_0$):

$$\mathbf{x} = \chi(\mathbf{X}, t), \quad \chi(\mathbf{X}, 0) = \mathbf{X}, \quad (2.2)$$

and the variation of a typical variable ϕ over the body is described with respect to the coordinates \mathbf{X} and time t :

$$\phi = \phi(\mathbf{X}, t). \quad (2.3)$$

For a fixed value of $\mathbf{X} \in \kappa_0$, $\phi(\mathbf{X}, t)$ gives the value of ϕ at time t associated with the fixed particle X whose position in the reference configuration is \mathbf{X} . We note that a change in time t implies that the same particle X has a different value ϕ . Thus, the Lagrangian description focuses its attention on the particles of the continuous body and it is usually used in solid mechanics.

2.2.2 Eulerian description

Rather than following the particles, the Eulerian description observes the changes at fixed locations. In the Eulerian description, the motion is this time referred to the current configuration κ and ϕ is described with respect to the current position ($\mathbf{x} \in \kappa$):

$$\phi = \phi(\mathbf{x}, t), \quad \mathbf{X} = \mathbf{X}(\mathbf{x}, t). \quad (2.4)$$

For a fixed value of $\mathbf{x} \in \kappa$, $\phi(\mathbf{x}, t)$ gives the value of ϕ associated with the particle occupying the position $\mathbf{x} \in \kappa$, which may very well be a different particle for each new time t . Because a change in time t implies that a different value ϕ is observed at the same spatial location $\mathbf{x} \in \kappa$, now probably occupied by a different particle, the Eulerian description is focused on a spatial position. This approach is convenient for the study of fluid flow where the kinematic property of greatest interest is the

rate at which change is taking place rather than the shape of the body of fluid at a reference time (Spencer, 1980). Because we are only interested in the study of solid bodies, the Lagrangian description will be used in the rest of the text.

2.2.3 Displacement field

The displacement \mathbf{u} of a particle X is called *displacement vector* and its expression in the Lagrangian description is given by the following:

$$\mathbf{u}(\mathbf{X}, t) = \mathbf{x}(\mathbf{X}, t) - \mathbf{X}. \quad (2.5)$$

A *displacement field* is a vector field of all displacement vectors for all particles in the body, which relates the deformed configuration κ with the undeformed configuration κ_0 . Indeed if the displacement field is known, we can construct the current configuration κ from the undeformed configuration κ_0 : $\chi(\mathbf{X}, t) = \mathbf{X} + \mathbf{u}(\mathbf{X})$.

2.3 Analysis of deformation

2.3.1 Deformation gradient tensor

The displacement field tells us how a point displaces from the reference to the deformed configuration. However, we would also like to know how a piece of material is stretched and rotated as the body moves from the reference to the deformed configuration. Since the length of a material line $d\mathbf{X}$ can change when going to the deformed configuration as well as its orientation, we can say that $d\mathbf{X}$ deforms into $d\mathbf{x}$. The question then becomes how to relate $d\mathbf{x}$ in the deformed configuration with $d\mathbf{X}$ of the reference configuration.

Consider two particles P_1 and P_2 in a continuous body separated by an infinitesimal distance $d\mathbf{X}$:

$$d\mathbf{X} = \mathbf{X}_{P_2} - \mathbf{X}_{P_1}. \quad (2.6)$$

After deformation, the two particles have deformed to their current positions given by the mapping χ (2.1) as:

$$\mathbf{x}_{P_1} = \chi(\mathbf{X}_{P_1}, t) \quad \text{and} \quad \mathbf{x}_{P_2} = \chi(\mathbf{X}_{P_2}, t). \quad (2.7)$$

Using (2.6) the distance $d\mathbf{x}$ between P_1 and P_2 can then be expressed as:

$$d\mathbf{x} = \mathbf{x}_{P_2} - \mathbf{x}_{P_1} = \chi(\mathbf{X}_{P_1} + d\mathbf{X}, t) - \chi(\mathbf{X}_{P_1}, t). \quad (2.8)$$

The *deformation gradient* \mathbf{F} can be defined as:

$$\mathbf{F} = \frac{\partial \chi}{\partial \mathbf{X}} \quad (2.9)$$

and the vector $d\mathbf{x}$ can then be obtained in terms of $d\mathbf{X}$ as:

$$d\mathbf{x} = \mathbf{F} d\mathbf{X}. \quad (2.10)$$

Note that \mathbf{F} transforms vectors from the reference configuration into vectors in the current configuration and is therefore a second-order tensor.

Knowing that $\chi(\mathbf{X}, t)$ is of course \mathbf{x} , the deformation gradient may also be written as:

$$\mathbf{F} = \frac{\partial \mathbf{x}}{\partial \mathbf{X}} = \nabla_0 \mathbf{x} = \nabla_0 \mathbf{u} + \mathbf{I}, \quad (2.11)$$

where ∇_0 is the gradient operator with respect to \mathbf{X} and \mathbf{u} the displacement vector. In indicial notation in a Cartesian coordinate system, (2.11) can be explicated as:

$$[\mathbf{F}] = \begin{bmatrix} \frac{\partial \mathbf{x}_1}{\partial \mathbf{X}_1} & \frac{\partial \mathbf{x}_1}{\partial \mathbf{X}_2} & \frac{\partial \mathbf{x}_1}{\partial \mathbf{X}_3} \\ \frac{\partial \mathbf{x}_2}{\partial \mathbf{X}_1} & \frac{\partial \mathbf{x}_2}{\partial \mathbf{X}_2} & \frac{\partial \mathbf{x}_2}{\partial \mathbf{X}_3} \\ \frac{\partial \mathbf{x}_3}{\partial \mathbf{X}_1} & \frac{\partial \mathbf{x}_3}{\partial \mathbf{X}_2} & \frac{\partial \mathbf{x}_3}{\partial \mathbf{X}_3} \end{bmatrix}. \quad (2.12)$$

2.3.2 Change of volume

At this point, we have seen how a deformation can affect a vector. We will now look into its effect on volumes. Our motivation comes from the need to write global equilibrium statements that involve integrals over volumes. We can define volume elements in the reference and deformed configurations. Consider three non-coplanar line elements $d\mathbf{X}^{(1)}$, $d\mathbf{X}^{(2)}$ and $d\mathbf{X}^{(3)}$ forming a parallelepiped in the reference configuration as can be seen on Fig. 2.1.

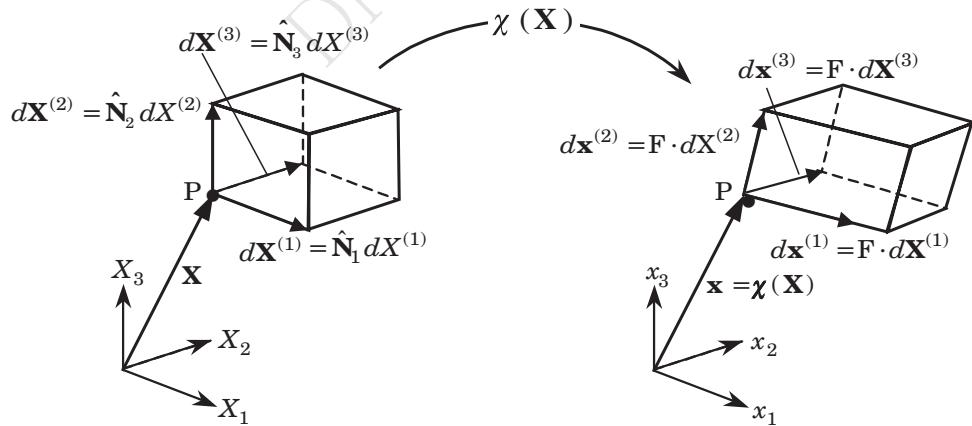


Figure 2.1: Transformation of a volume element under a deformation mapping.
TODO: CHANGE TO A BETTER FIGURE.

The three vectors after deformation $d\mathbf{x}^{(1)}$, $d\mathbf{x}^{(2)}$ and $d\mathbf{x}^{(3)}$ can be obtained with:

$$d\mathbf{x}^{(i)} = \mathbf{F} d\mathbf{X}^{(i)}, \quad i = 1, 2, 3. \quad (2.13)$$

The volume of the parallelepiped that we will note dV can be calculated using the triple product between the three vectors:

$$\begin{aligned} dV &= d\mathbf{X}^{(1)} \cdot d\mathbf{X}^{(2)} \times d\mathbf{X}^{(3)} = (\hat{\mathbf{N}}_1 \cdot \hat{\mathbf{N}}_2 \times \hat{\mathbf{N}}_3) dX^{(1)} dX^{(2)} dX^{(3)} \\ &= dX^{(1)} dX^{(2)} dX^{(3)}, \end{aligned} \quad (2.14)$$

where $\hat{\mathbf{N}}_i$ denote the unit vector along $d\mathbf{X}^i$. The corresponding volume in the deformed configuration is given by:

$$\begin{aligned} dv &= d\mathbf{x}^{(1)} \cdot d\mathbf{x}^{(2)} \times d\mathbf{x}^{(3)} \\ &= (\mathbf{F} \cdot \hat{\mathbf{N}}_1) \cdot (\mathbf{F} \cdot \hat{\mathbf{N}}_2) \times (\mathbf{F} \cdot \hat{\mathbf{N}}_3) dX^{(1)} dX^{(2)} dX^{(3)} \quad \text{by (2.13)} \\ &= \det \mathbf{F} dX^{(1)} dX^{(2)} dX^{(3)}. \end{aligned} \quad (2.15)$$

The determinant of \mathbf{F} is called the *Jacobian* and it is denoted by $J = \det \mathbf{F}$. And we have:

$$dv = J dV. \quad (2.16)$$

Thus, J has the physical meaning of being the local ratio of current to reference volume of a material volume element.

2.3.3 Change of surface

For similar reasons, let us find the relation between an element of area in the reference configuration $d\mathbf{A}$ which becomes $d\mathbf{a}$ in the deformed configuration. Considering an arbitrary material line $d\mathbf{L}$ in the reference configuration and noting $d\mathbf{l}$ the same line after deformation (see Fig. 2.2, the corresponding reference and current volumes are given respectively by:

$$d\mathbf{V} = d\mathbf{L} \cdot d\mathbf{A} \quad (2.17)$$

$$\text{and } d\mathbf{v} = d\mathbf{l} \cdot d\mathbf{a}. \quad (2.18)$$

Using (2.16) which relates the reference and deformed volumes and recalling that $d\mathbf{l} = \mathbf{F} d\mathbf{L}$, we have:

$$J d\mathbf{L} \cdot d\mathbf{A} = (\mathbf{F} d\mathbf{L}) \cdot d\mathbf{a}. \quad (2.19)$$

The fact that this expression is valid for any vector $d\mathbf{L}$ and using the property $\mathbf{a} \cdot \mathbf{b} = \mathbf{a}^T \mathbf{b}$ yields:

$$d\mathbf{a} = J \mathbf{F}^{-T} d\mathbf{A}. \quad (2.20)$$

2.3.4 Volumetric and isochoric components

We recall that the Jacobian J has the physical meaning of being the local ratio of current to reference volume of a material volume element. Therefore, if $J = 1$ the volume does not change locally during the deformation and the latter is qualified as

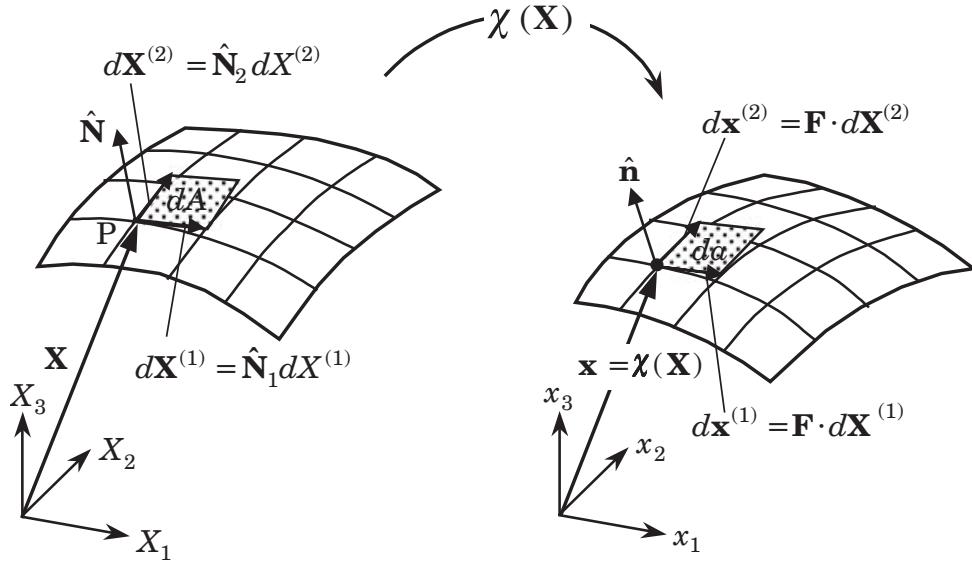


Figure 2.2: Transformation of a surface element under a deformation mapping.
TODO: CHANGE TO A BETTER FIGURE.

isochoric at this given particle P . If $J = 1$ everywhere in the body, the deformation of the body is isochoric.

When dealing with incompressible and nearly incompressible materials it is necessary to separate the volumetric from the isochoric components of the deformation. Such a separation must ensure that the isochoric component $\hat{\mathbf{F}}$ does not imply any change in volume. This condition can be achieved by choosing $\hat{\mathbf{F}}$ as:

$$\hat{\mathbf{F}} = J^{-1/3} \mathbf{F}. \quad (2.21)$$

Indeed,

$$\begin{aligned} \det \hat{\mathbf{F}} &= \det(J^{-1/3} \mathbf{F}) \\ &= (J^{-1/3})^3 \det \mathbf{F} \\ &= J^{-1} J \\ &= 1. \end{aligned} \quad (2.22)$$

The deformation gradient \mathbf{F} can now be expressed in terms of the volumetric and isochoric component J and $\hat{\mathbf{F}}$ as:

$$\mathbf{F} = J^{1/3} \hat{\mathbf{F}}. \quad (2.23)$$

2.4 Strain measures

The length of a material curve from the reference configuration can change when displaced to a curve in the deformed configuration. If all the curves do not change

length, it is said that a rigid body displacement occurred. The concept of strain is used to evaluate how much a given displacement differs locally from a rigid body displacement. Therefore, although we know how to transform vectors from the reference configuration into vectors in the current configuration using the deformation gradient, it is more useful to find a measure of the change in length of $d\mathbf{X}$. Many measures of strains can be defined and the most common ones will now be introduced.

2.4.1 Cauchy-Green deformation tensors

Let us consider two particles P_1 and P_2 in the neighbourhood of each other, separated by $d\mathbf{X}$ in the reference configuration. In the deformed configuration P_1 and P_2 occupy the positions \tilde{P}_1 and \tilde{P}_2 and they are separated by $d\mathbf{x}$. We are interested in the change of distance between the two points P_1 and P_2 as the body deforms.

The squared distances between P_1 and P_2 and \tilde{P}_1 and \tilde{P}_2 are respectively given by:

$$(dS)^2 = d\mathbf{X} \cdot d\mathbf{X} \quad (2.24)$$

$$\begin{aligned} (ds)^2 &= d\mathbf{x} \cdot d\mathbf{x} = (\mathbf{F} d\mathbf{X}) \cdot (\mathbf{F} d\mathbf{X}) \\ &= (\mathbf{F} d\mathbf{X})^T (\mathbf{F} d\mathbf{X}) = (d\mathbf{X}^T \mathbf{F}^T) (\mathbf{F} d\mathbf{X}) = d\mathbf{X}^T (\mathbf{F}^T \mathbf{F} d\mathbf{X}) \\ &= d\mathbf{X} \cdot (\mathbf{F}^T \mathbf{F} d\mathbf{X}). \end{aligned} \quad (2.25)$$

using the property that the dot product between two vectors a and b can also be expressed as the simple product between the transpose of a and b ($a \cdot b = a^T b$). We define the *right Cauchy-Green deformation tensor* \mathbf{C} as:

$$\mathbf{C} = \mathbf{F}^T \mathbf{F}. \quad (2.26)$$

Thus, the change of distance between the two points P_1 and P_2 after deformation of the continuous body may be written as:

$$(ds)^2 = d\mathbf{X} \cdot (\mathbf{C} d\mathbf{X}). \quad (2.27)$$

By definition, \mathbf{C} is a symmetric second-order tensor. The transpose of \mathbf{C} is denoted \mathbf{B} and is called the *left Cauchy-Green deformation tensor*:

$$\mathbf{B} = \mathbf{F} \mathbf{F}^T. \quad (2.28)$$

2.4.2 Green-Lagrange strain tensor

Using (2.24) and (2.27), the change in the squared lengths due to the body deformation between the reference and the current configuration can be expressed as:

$$\begin{aligned} (ds)^2 - (dS)^2 &= d\mathbf{X} \cdot (\mathbf{C} d\mathbf{X}) - d\mathbf{X} \cdot d\mathbf{X} \\ &= d\mathbf{X} \cdot (\mathbf{C} - \mathbf{I}) d\mathbf{X}. \end{aligned} \quad (2.29)$$

Let us define the *Green-Lagrange strain tensor* \mathbf{E} as:

$$\mathbf{E} = \frac{1}{2}(\mathbf{C} - \mathbf{I}) \quad (2.30)$$

so we can write:

$$(ds)^2 - (dS)^2 = 2 d\mathbf{X} \cdot \mathbf{E} d\mathbf{X}. \quad (2.31)$$

By definition, the Green-Lagrange strain tensor is a symmetric second-order tensor. Also, the change in squared lengths is zero if and only if $\mathbf{E} = \mathbf{0}$. Using (2.11), the Green strain tensor may be expanded as the following:

$$\begin{aligned} \mathbf{E} &= \frac{1}{2}(\mathbf{F}^T \mathbf{F} - \mathbf{I}) \\ &= \frac{1}{2} [(\nabla_0 \mathbf{u} + \mathbf{I})^T (\nabla_0 \mathbf{u} + \mathbf{I}) - \mathbf{I}] \\ &= \frac{1}{2} [(\nabla_0 \mathbf{u})^T + \nabla_0 \mathbf{u} + (\nabla_0 \mathbf{u})^T (\nabla_0 \mathbf{u})]. \end{aligned} \quad (2.32)$$

The Green strain tensor can be expressed in terms of its components in any coordinate system. In particular, in the Cartesian coordinate system (X_1, X_2, X_3) the components E_{ij} of \mathbf{E} are the following:

$$E_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial X_j} + \frac{\partial u_j}{\partial X_i} + \frac{\partial u_k}{\partial X_i} \frac{\partial u_k}{\partial X_j} \right), \quad i = 1, 2, 3. \quad (2.33)$$

The components E_{11} , E_{22} and E_{33} are called *normal strains* and it can be shown that they are in fact the ratio of the change in length to the original length along each of the three unit vectors. The components E_{12} , E_{23} and E_{13} are called *shear strains* and they can be interpreted as a measure of the change in angle between line elements that were perpendicular to each other in the undeformed configuration.

2.4.3 Cauchy and Euler tensor

The change in the squared lengths during the body deformation can also be expressed relative to the current length. The length dS can be written in terms of $d\mathbf{x}$ as:

$$(dS)^2 = d\mathbf{X} \cdot d\mathbf{X} = d\mathbf{x} \cdot (\mathbf{F}^{-T} \mathbf{F}^{-1}) d\mathbf{x} = d\mathbf{x} \cdot \tilde{\mathbf{B}} d\mathbf{x} \quad (2.34)$$

where $\tilde{\mathbf{B}} = \mathbf{F}^{-T} \mathbf{F}^{-1}$ is called the *Cauchy strain tensor*. $\tilde{\mathbf{B}}$ is in fact the inverse of the left Cauchy-Green tensor \mathbf{B} introduced previously.

In a similar way we defined the Green strain tensor, we can write the change in the squared lengths but relative to the current length:

$$(ds)^2 - (dS)^2 = 2 d\mathbf{x} \cdot \mathbf{e} d\mathbf{x}. \quad (2.35)$$

where \mathbf{e} is called *Almansi-Hamel strain tensor* or simply *Euler strain tensor*.

2.4.4 Principal strains and invariants

The components ε_{ij} of the strain tensor depend on the coordinate system at the point under consideration. However, the strain tensor itself is a physical quantity and as such, it is independent of the coordinate system chosen to represent it. Certain operations on strain tensors give the same result independently of the coordinate system chosen to represent the components of strain. As an example, a vector is a simple tensor of rank one. In three dimensions, it has three components. The value of these components will depend on the coordinate system chosen to represent the vector, but the length of the vector is a physical quantity (a scalar) and is independent of the coordinate system chosen to represent the vector. Similarly, every second rank tensor (such as the strain tensors) has three independent invariant quantities associated with it. One set of such *invariants* are the principal strains of the strain tensor, which are just the eigenvalues of the strain tensor. Because they are independent of any coordinate system, they are very convenient to define strain energy density functions (see section 2.6.5). The most commonly used invariants are:

$$\begin{aligned} I_1 &= \text{tr}(\varepsilon) \\ I_2 &= \frac{1}{2} \left\{ \text{tr}(\varepsilon^2) - [\text{tr}(\varepsilon)]^2 \right\} \\ I_3 &= \det(\varepsilon). \end{aligned} \quad (2.36)$$

It can be shown that in the coordinate system $(\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3)$ formed by the three eigenvectors, the expression of the strain tensor is:

$$\boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 & 0 & 0 \\ 0 & \varepsilon_2 & 0 \\ 0 & 0 & \varepsilon_3 \end{bmatrix}. \quad (2.37)$$

Since there are no shear strain components in this particular coordinate system, the principal strains represent the maximum and minimum stretches of an elemental volume. Because of the obvious simplicity of the strain tensor's expression, this coordinate system is often used in mechanics to derive and solve equations.

2.4.5 Infinitesimal strain tensor

In some cases, it is possible to simplify the expression of the Green strain tensor defined in (2.32). Indeed, when the displacement gradients are small (that is, $|\nabla \mathbf{u}| \ll 1$) we can neglect the non-linear terms in the definition. In the case of infinitesimal strains, the Green-Lagrange strain tensor and the Eulerian strain tensor are approximately the same and can be approximated by the infinitesimal strain tensor denoted $\boldsymbol{\varepsilon}$ and is given by:

$$\boldsymbol{\varepsilon} = \frac{1}{2} [\nabla \mathbf{u} + (\nabla \mathbf{u})^T]. \quad (2.38)$$

Its Cartesian components ε_{ij} are the following:

$$\varepsilon_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial X_j} + \frac{\partial u_j}{\partial X_i} \right). \quad (2.39)$$

The strain quantities defined in the previous sections are non-linear expressions in terms of the mapping χ and will lead to non-linear governing equations. In solid mechanics, whenever the hypothesis is acceptable, it is common practice to assume that the displacements are small and the infinitesimal strain tensor is then used as a measure of the deformation.

2.5 Stress

Stress is a measure of the intensity of the internal forces acting between particles of a deformable body across imaginary internal surfaces. These internal forces are produced between the particles in the body as a reaction to external forces applied on the body. The *SI* unit for stress is pascal (symbol *Pa*), which is equivalent to one newton (force) per square metre (unit area). As with strain, different measures of stress can be defined. For instance, stress can be measured per unit of deformed area (section 2.5.1) or undeformed area (sections 2.5.2 and 2.5.3). In general, the stress is not uniformly distributed over the cross section of a material body, and consequently the stress at a point on a given area is different from the average stress over this entire area. Therefore, it is necessary to define the stress not over a given area but at a specific point in the body. Stress at a point in a three-dimensional continuum can be measured in terms of nine quantities, three per plane, on three mutually perpendicular planes at the point. These nine quantities may be viewed as the components of a second-order tensor that we call *stress tensor*.

2.5.1 Cauchy stress

The Euler-Cauchy's stress principle states that on any surface (real or imaginary) that divides the body, the action of one part of the body on the other is equivalent to the system of distributed forces and couples on the surface dividing the body. With that consideration in mind, let us consider a plane S that passes through an arbitrary internal point P which has a unit normal vector \mathbf{n} . This plane separates the body into two regions and we are interested in the forces applied by one region onto the other. First we will introduce the true stress, defined as the stress in the deformed configuration χ measured per unit area of the deformed configuration. The force $\Delta\mathbf{f}(\mathbf{n})$ acting on a small element of area Δa in a continuous medium depends not only on the size of the area but also depends on the orientation of this area \mathbf{n} . The *Cauchy stress vector* at P on S can be defined as:

$$\mathbf{t}(\mathbf{n}) = \lim_{\Delta a \rightarrow 0} \frac{\Delta\mathbf{f}(\mathbf{n})}{\Delta a}. \quad (2.40)$$

This equation means that the stress vector depends on its location in the body and the orientation of the plane on which it is acting. In general, the stress vector is

not perpendicular to that plane and it can be separated into two components: one component normal to the plane called *normal stress* and one component tangent to the plane named *shear stress*. Since \mathbf{t} depends on \mathbf{n} but is not in the direction of \mathbf{n} , it may be interesting to look into the relationship between \mathbf{t} and \mathbf{n} .

To establish this relationship, let us consider an infinitesimal tetrahedron in a Cartesian coordinate system as represented Fig. 2.3. We note $-\mathbf{t}_1$, $-\mathbf{t}_2$, $-\mathbf{t}_3$

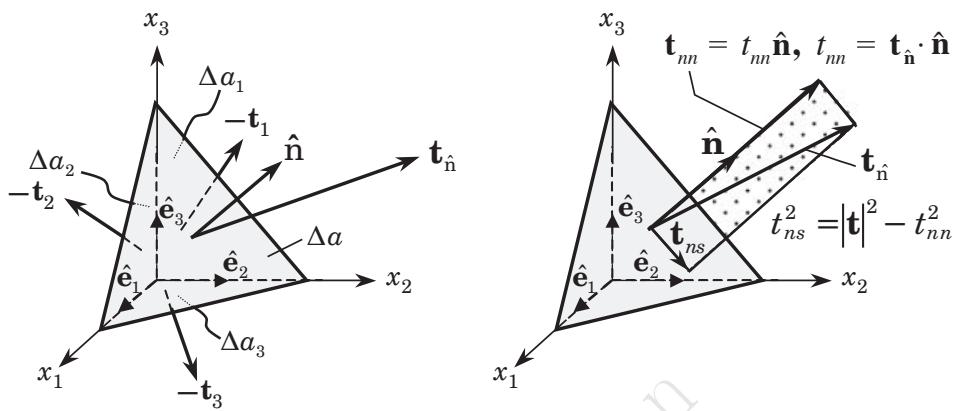


Figure 2.3: Tetrahedral element in Cartesian coordinates. **TODO: CHANGE TO A BETTER FIGURE.**

and \mathbf{t} the stress vectors in the outwards directions on the faces of the tetrahedron whose areas are respectively Δa_1 , Δa_2 , Δa_3 and Δa . For the mass m inside the tetrahedron, we have by Newton's second law:

$$\sum \mathbf{F} = \mathbf{t}\Delta a - \mathbf{t}_1\Delta a_1 - \mathbf{t}_2\Delta a_2 - \mathbf{t}_3\Delta a_3 + \Delta v\mathbf{f} = m\mathbf{a}, \quad (2.41)$$

where \mathbf{f} is the force per unit volume acting on the body, Δv the volume of the tetrahedron and \mathbf{a} the acceleration. The areas Δa_1 , Δa_2 and Δa_3 being the projections of Δa , they are related to Δa in the Cartesian coordinate system ($\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$) by:

$$\Delta a_i = (\mathbf{n} \cdot \mathbf{e}_i)\Delta a, \quad i = 1, 2, 3. \quad (2.42)$$

Dividing (2.41) by Δa , using (2.42) and noting that $\Delta v/\Delta a \rightarrow 0$ when the tetrahedron shrinks to a point yields:

$$\begin{aligned} \mathbf{t} &= (\mathbf{n} \cdot \mathbf{e}_1)\mathbf{t}_1 + (\mathbf{n} \cdot \mathbf{e}_2)\mathbf{t}_2 + (\mathbf{n} \cdot \mathbf{e}_3)\mathbf{t}_3 \\ &= \mathbf{n} \cdot (\mathbf{e}_1\mathbf{t}_1 + \mathbf{e}_2\mathbf{t}_2 + \mathbf{e}_3\mathbf{t}_3). \end{aligned} \quad (2.43)$$

We define the *Cauchy stress tensor* $\sigma = \mathbf{e}_1\mathbf{t}_1 + \mathbf{e}_2\mathbf{t}_2 + \mathbf{e}_3\mathbf{t}_3$ and we have:

$$\mathbf{t}(\mathbf{n}) = \mathbf{n} \cdot \sigma. \quad (2.44)$$

The stress vector \mathbf{t} represents the vectorial stress on a plane whose normal is \mathbf{n} . As we demonstrated, the stress tensor is a property of the medium that is independent of \mathbf{n} . It represents the current force per unit of deformed area: $d\mathbf{f} = \mathbf{t}da = \sigma \cdot da$.

The matrix form of (2.44) in a Cartesian coordinate system is given by:

$$\begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} = \begin{bmatrix} \sigma_{11} & \sigma_{21} & \sigma_{31} \\ \sigma_{12} & \sigma_{22} & \sigma_{32} \\ \sigma_{13} & \sigma_{23} & \sigma_{33} \end{bmatrix} \begin{bmatrix} n_1 \\ n_2 \\ n_3 \end{bmatrix}. \quad (2.45)$$

The Cauchy stress tensor is used for stress analysis of material bodies experiencing small deformations. This stress is basically defined as force/(unit deformed area). The strain measure that is appropriate to use with the Cauchy stress tensor is the infinitesimal strain tensor (see section 2.4.5). The problem with using the Cauchy stress tensor for analysing materials undergoing large deformation is that we generally do not know the area in the deformed configuration. Therefore, we need to develop alternative stress tensors that we can use in the reference configuration. We will introduce two new stress measures: the first Piola-Kirchhoff stress tensor and the second Piola-Kirchhoff stress tensor.

2.5.2 First Piola-Kirchhoff stress tensor

We know that the Cauchy stress tensor represents the current force per unit of deformed area: $d\mathbf{f} = \mathbf{t}(\mathbf{n})da = \boldsymbol{\sigma} \cdot d\mathbf{a}$. We define a stress vector \mathbf{T} over the area element dA with normal \mathbf{N} in the underdeformed configuration such that it results in the same total force

$$d\mathbf{f} = \mathbf{t}(\mathbf{n})da = \mathbf{T}(\mathbf{N})dA. \quad (2.46)$$

The two stress vectors have the obviously the same direction but different magnitudes since they are not applied to the same area.

In a similar way we defined the Cauchy stress $\boldsymbol{\sigma}$ with $\mathbf{t}(\mathbf{n}) = \boldsymbol{\sigma} \cdot \mathbf{n}$, we introduce a stress tensor \mathbf{P} called the *first Piola-Kirchhoff stress tensor* such that $\mathbf{T}(\mathbf{N}) = \mathbf{P} \cdot \mathbf{N}$. The first Piola-Kirchhoff stress tensor gives the current force per unit undeformed area.

The relation between the first Piola-Kirchhoff stress tensor and the Cauchy stress tensor can be obtained as follows. From (2.46) we have:

$$\mathbf{T} = \frac{da}{dA} \mathbf{t}. \quad (2.47)$$

Using $d\mathbf{f} = \mathbf{t}da$ and the relation (2.20) between $d\mathbf{A}$ and $d\mathbf{a}$, we have:

$$\mathbf{T} = J \mathbf{t} \mathbf{F}^{-T}. \quad (2.48)$$

Finally we can express the first Piola-Kirchhoff stress tensor \mathbf{P} as

$$\mathbf{P} = J \boldsymbol{\sigma} \mathbf{F}^{-T}. \quad (2.49)$$

2.5.3 Second Piola-Kirchhoff stress tensor

The *second Piola-Kirchhoff stress tensor* \mathbf{S} is introduced as the stress tensor associated with the force $d\mathbf{F}$ in the undeformed area $d\mathbf{A}$ that corresponds to the force $d\mathbf{f}$ in the deformed area $d\mathbf{a}$:

$$d\mathbf{F} = \mathbf{S} d\mathbf{A}. \quad (2.50)$$

Thus, the second Piola-Kirchhoff stress tensor gives the transformed current force per unit undeformed area. Similar to the relationship (2.13) between $d\mathbf{x}$ and $d\mathbf{X}$, the force $d\mathbf{f}$ is related to the force $d\mathbf{F}$ by:

$$\begin{aligned} d\mathbf{F} &= \mathbf{F}^{-1} d\mathbf{f} \\ &= \mathbf{F}^{-1} (\mathbf{P} d\mathbf{A}) \\ &= \mathbf{S} d\mathbf{A}. \end{aligned} \quad (2.51)$$

Therefore, the second Piola-Kirchhoff stress tensor is related to the first Piola-Kirchhoff stress tensor as follows:

$$\mathbf{S} = \mathbf{F}^{-1} \mathbf{P} = J \mathbf{F}^{-1} \sigma \mathbf{F}^{-T}. \quad (2.52)$$

2.5.4 Principal stresses and invariants

In the same way that we defined principal strains in section 2.4.4 and like any second-rank tensor, the stress tensor has at least three invariants and we can find a coordinate system $(\mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3)$ in which the expression of the stress tensor is:

$$\sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix}. \quad (2.53)$$

This result means that at every point in a stressed body there are at least three planes called *principal planes*, with normal vectors \mathbf{n}_1 , \mathbf{n}_2 and \mathbf{n}_3 called *principal directions* (in fact, the eigenvectors), where the corresponding stress vector is perpendicular to the plane and where there are no normal shear stresses. The three stresses normal σ_1 , σ_2 and σ_3 to these principal planes are called *principal stresses*.

2.6 Constitutive equations

So far, all the relations that we have established are valid for every continuum. Indeed, no mention was made of any material in the derivations. Therefore, these relations are not sufficient to describe the response of a material to a given loading. We know that the stresses result from the deformation of the material and it is now necessary to express them in terms of some measure of this deformation such as, for instance, the strain. These relationships, known as *constitutive equations*, obviously depend on the type of material under consideration and may be dependent on time or not. Constitutive equations are not derived from any physical principles, although

they are subject to obeying certain rules and the entropy inequality. In essence, constitutive equations are mathematical models of the behaviour of materials that are validated against experimental results.

The relationship between stress and strain depends on the nature, or constitution of the particular material that we are studying. Of course, any material mechanics, and especially tissue mechanics, can be quite complex, and many assumptions must be made when deriving a constitutive equation. It is basically impossible to derive a particular constitutive equation that would accurately model all aspects of tissue behaviour under any type of loading. This is because tissues may behave very differently under small loads as compared to large loads. For instance, we may expect damage under large loads. Also, under different loading rates different tissues may behave differently. Therefore, whenever we develop a constitutive equation to model a tissue, we need to balance the need to accurately model tissue behaviour under the range of loadings with the need to have a constitutive equation that is simple enough to:

1. be computed in a numerical model
2. experimentally measure all the material constants used in the constitutive equation

This section will now introduce the most common constitutive models relevant to soft-tissue modelling.

2.6.1 Elasticity

If under applied loads a material stores but does not dissipate energy, and it returns to its original shape when the loads are removed, we call that material *elastic*. The work done by the stress is generally dependent on the deformation path. However, for elastic materials the stress does not depend on the path of deformation and therefore the state of stress in the deformed configuration is determined only by the state of deformation. It means that the stress does not depend on the time taken to achieve the deformation or the rate at which the state of deformation is reached.

One special case of elastic materials are *hyperelastic* materials. A material is said to be hyperelastic if there exists a *strain energy density* function $U_0(\varepsilon)$ such that

$$\boldsymbol{\sigma} = \frac{\partial U_0}{\partial \varepsilon}. \quad (2.54)$$

However, for *incompressible* materials the stress tensor is not completely determined by the deformation. The hydrostatic pressure¹ affects the stress and the above equation is then written as:

$$\boldsymbol{\sigma} = -p\mathbf{I} + \frac{\partial U_0}{\partial \varepsilon}, \quad (2.55)$$

¹pressure exerted by a fluid at equilibrium due to the force of gravity

where p is the hydrostatic pressure. In order to develop a mathematical model of an hyperelastic material, U_0 is expanded in Taylor's series about $\varepsilon = 0$:

$$U_0 = C_0 + C_{ij}\varepsilon_{ij} + \frac{1}{2!}\hat{C}_{ijkl}\varepsilon_{ij}\varepsilon_{kl} + \frac{1}{3!}\hat{C}_{ijklmn}\varepsilon_{ij}\varepsilon_{kl}\varepsilon_{mn} + \dots, \quad (2.56)$$

where C_0 , C_{ij} and \hat{C} are material stiffnesses. The form of U_0 may varied for different materials. A material is said to be *linear* if the relationship between strain and stress is also linear. Based on this definition, U_0 is a cubic or higher-order function of the strains for non-linear materials whereas U_0 is a quadratic function of strain for linear materials. We will start by discussing linear materials for the case of infinitesimal deformations (sections 2.6.2 to 2.6.4). Consequently, we will not make any distinction between the various measures of stress and strain and use σ and ε for stress and strain tensors respectively.

2.6.2 Linear materials: generalised Hooke's law

The linear constitutive model for infinitesimal deformations is called the generalised Hooke's law. To derive the relation between stress and strain for a linear material, we assume the quadratic form of U_0 :

$$U_0 = C_0 + C_{ij}\varepsilon_{ij} + \frac{1}{2!}\hat{C}_{ijkl}\varepsilon_{ij}\varepsilon_{kl}. \quad (2.57)$$

From there and using (2.54) it can be shown that there exists some coefficients C_{ijkl} such as:

$$\sigma_{ij} = C_{ijkl}\varepsilon_{kl}. \quad (2.58)$$

The C_{ijkl} are the coefficients of the fourth-order tensor \mathbf{C} that relates the second-order tensors of strain and stress. The 81 scalar components of this fourth-order tensor ($= 3^4$) can be reduced to 21 independent coefficients using existing symmetries. Recalling that both strain and stress tensors are symmetric yields:

$$\begin{bmatrix} \sigma_{11} \\ \sigma_{22} \\ \sigma_{33} \\ \sigma_{23} \\ \sigma_{13} \\ \sigma_{12} \end{bmatrix} = \begin{bmatrix} C_{1111} & C_{1122} & C_{1133} & C_{1123} & C_{1113} & C_{1112} \\ C_{1122} & C_{2222} & C_{2233} & C_{2223} & C_{2213} & C_{2212} \\ C_{1133} & C_{2233} & C_{3333} & C_{3323} & C_{3313} & C_{3312} \\ C_{1123} & C_{2223} & C_{3323} & C_{2323} & C_{2313} & C_{2312} \\ C_{1113} & C_{2213} & C_{3313} & C_{2313} & C_{1313} & C_{1312} \\ C_{1112} & C_{2212} & C_{3312} & C_{2312} & C_{1312} & C_{1212} \end{bmatrix} \begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{22} \\ \varepsilon_{33} \\ \varepsilon_{23} \\ \varepsilon_{13} \\ \varepsilon_{12} \end{bmatrix}. \quad (2.59)$$

2.6.3 Orthotropic materials

Further reduction in the number of independent stiffness parameters can be achieved using the so-called material symmetry. This symmetry is due to the internal structure of the material, that is the arrangement of molecules (crystallographic form). We note that the symmetry under discussion is a directional property and not a positional property. Thus, a material may have certain elastic symmetry at every

point of a material body and the properties may vary from point to point. Positional dependence of material properties is what we called the inhomogeneity of the material.

Orthotropic materials are characterised by three mutually orthogonal planes of material symmetry, that is, three mutually orthogonal preferred directions. In this case, the number of coefficients is reduced to 9 and the stress-strain relation for an orthotropic material is the following:

$$\begin{bmatrix} \sigma_{11} \\ \sigma_{22} \\ \sigma_{33} \\ \sigma_{23} \\ \sigma_{13} \\ \sigma_{12} \end{bmatrix} = \begin{bmatrix} C_{1111} & C_{1122} & C_{1133} & 0 & 0 & 0 \\ C_{1122} & C_{2222} & C_{2233} & 0 & 0 & 0 \\ C_{1133} & C_{2233} & C_{3333} & 0 & 0 & 0 \\ 0 & 0 & 0 & C_{2323} & 0 & 0 \\ 0 & 0 & 0 & 0 & C_{1313} & 0 \\ 0 & 0 & 0 & 0 & 0 & C_{1212} \end{bmatrix} \begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{22} \\ \varepsilon_{33} \\ \varepsilon_{23} \\ \varepsilon_{13} \\ \varepsilon_{12} \end{bmatrix}. \quad (2.60)$$

Most often, the material properties are determined in a laboratory in terms of the engineering constants such as Young's modulus, shear modulus, and so on. These constants are measured using simple tests like uniaxial tension test or pure shear test. Because of their direct and obvious physical meaning, engineering constants are used in place of the more abstract stiffness coefficients C_{ijkl} . Therefore, we will now reformulate (2.60) with respect to those engineering constants. Noting E_i the Young's modulus in the direction x_i and G_{ij} the shear moduli in the $x_i - x_j$ plane, it can be shown that:

$$\begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{22} \\ \varepsilon_{33} \\ \varepsilon_{23} \\ \varepsilon_{13} \\ \varepsilon_{12} \end{bmatrix} = \begin{bmatrix} \frac{1}{E_1} & -\frac{\nu_{21}}{E_2} & -\frac{\nu_{31}}{E_3} & 0 & 0 & 0 \\ -\frac{\nu_{12}}{E_1} & \frac{1}{E_2} & -\frac{\nu_{32}}{E_3} & 0 & 0 & 0 \\ -\frac{\nu_{13}}{E_1} & -\frac{\nu_{23}}{E_2} & \frac{1}{E_3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{G_{23}} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{G_{13}} & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{G_{12}} \end{bmatrix} \begin{bmatrix} \sigma_{11} \\ \sigma_{22} \\ \sigma_{33} \\ \sigma_{23} \\ \sigma_{13} \\ \sigma_{12} \end{bmatrix}. \quad (2.61)$$

2.6.4 Isotropic materials

A material is said *isotropic* if its properties are independent of directions. Therefore, we have the following relations:

$$E_1 = E_2 = E_3 = E, \quad G_{12} = G_{13} = G_{23} = G, \quad \nu_{12} = \nu_{23} = \nu_{13} = \nu. \quad (2.62)$$

The stress-strain relations may be written:

$$\sigma_{ij} = \frac{E}{1 + \nu} \varepsilon_{ij} + \frac{\nu E}{(1 + \nu)(1 - 2\nu)} \varepsilon_{kk} \delta_{ij}. \quad (2.63)$$

where the summation of repeated indices is implied (Einstein's notation) and δ_{ij} is the Kronecker delta². The inverse relations are:

$$\varepsilon_{ij} = \frac{1 + \nu}{E} \sigma_{ij} - \frac{\nu}{E} \sigma_{kk} \delta_{ij}, \quad (2.64)$$

recalling that $G = E/[2(1 + \nu)]$.

2.6.5 Non-linear materials

For most materials, the linear relation between stress and strain is only valid for small deformation. Beyond a certain threshold, Hooke's law is no longer valid. Note that it does mean that the material is not elastic, it may well recover all its deformation after removal of the loads. However, materials often exhibit an elastic threshold in the non-linear range after which permanent deformation occurs. This range of deformation is qualified of *plastic*.

The strain energy density function takes different forms for different materials. It is often expressed as a linear combination of unknown coefficients (determined experimentally) and principal invariants of Green strain tensor \mathbf{E} , left Cauchy-Green strain tensor \mathbf{B} and deformation gradient tensor \mathbf{F} (see section 2.4.4 for more information on invariants). The constitutive equations of two well-known non-linear elastic materials (Mooney-Rivlin and neo-Hookean) will now be introduced.

Mooney-Rivlin

The Mooney-Rivlin model was proposed by Melvin Mooney and Ronald Rivlin in two independent papers in 1952. For an incompressible Mooney-Rivlin material, the strain energy density function U_0 is taken as a linear function of the principal invariants of the left Cauchy-Green strain tensor \mathbf{B} and is written as:

$$U_0 = C_1(I_B - 3) + C_2(II_B - 3), \quad (2.65)$$

where C_1 and C_2 are two constants and I_B and II_B are the two principal invariants of \mathbf{B} . For those materials, it can be shown that the stress tensor has the form

$$\sigma = -p\mathbf{I} + \alpha\mathbf{B} + \beta\mathbf{B}^{-1}, \quad (2.66)$$

where ρ is the material density and α and β are given by

$$\alpha = 2\rho \frac{\partial U_0}{\partial I_B} = 2\rho C_1, \quad \beta = -2\rho \frac{\partial U_0}{\partial II_B} = -2\rho C_2. \quad (2.67)$$

The Mooney-Rivlin incompressible model is often used to represent the behaviour of rubber-like materials.

² δ_{ij} is equal to 1 if $i = j$ and to 0 if $i \neq j$

Neo-Hookean

The neo-Hookean model was proposed by Ronald Rivlin in 1948. It can be derived from the Mooney-Rivlin model in the special case where $C_2 = 0$ in the strain energy density function. Therefore, we have:

$$U_0 = C_1(I_B - 3). \quad (2.68)$$

and the constitutive equation for the stress may be written as:

$$\sigma = -p\mathbf{I} + 2\rho C_1 \mathbf{B}. \quad (2.69)$$

The neo-Hookean model provides a reasonable approximation of rubber-like material for moderate strains. **TODO: Do I keep this range of validity, which seems somewhat limited...** It is typically accurate for strains less than 20 % (Gent, 2001).

2.6.6 Viscoelastic materials

Viscoelasticity is the property of materials that exhibit both viscous and elastic characteristics when undergoing deformation. Viscous materials, like honey, resist shear flow and strain linearly with time when a stress is applied. Elastic materials strain instantaneously when stretched and just as quickly return to their original state once the stress is removed. Viscoelastic materials have elements of both of these properties and, as such, exhibit time dependent strain. Depending on the change of strain rate versus stress inside a material the viscosity can be categorised as having a linear or non-linear response. When a material exhibits a linear response it is categorised as a Newtonian material. In this case the stress is linearly proportional to the strain rate. If the material exhibits a non-linear response to the strain rate, it is categorised as a Non-Newtonian material.

Viscoelastic materials present a few characteristics:

1. if the stress is held constant, the strain increases with time (creep)
2. if the strain is held constant, the stress decreases with time (relaxation)
3. the effective stiffness depends on the rate of application of the load
4. if cyclic loading is applied, hysteresis occurs, leading to a dissipation of mechanical energy

In other words, the resistance of a viscoelastic material to a given force depends on the velocity to which this force is applied.

TODO: Add details on constitutive equations? Level of details to be determined later (depending on TLED needs)

2.7 Tissue characterisation

One of the main obstacles in developing realistic organ models is the lack of data on the material properties of live organ tissues. Measuring and characterising *in vivo* organ properties is a highly challenging task, but is a requirement for realistic organ modelling. Indeed, as we have seen in sections 2.6.0 and 2.6.5, constitutive equations require the experimental determination of material constants. Organ models with incorrect material properties will affect training in surgical simulator systems (Sedef et al., 2006).

The constitutive equations are used to define idealised materials, which represent different aspects of the mechanical behaviour of natural materials. However, no real material is known to behave exactly as one of these mathematical models. In particular, few biological tissues obey Hooke's law. If because of their high stiffness (and therefore small deformations) cortical bones are accurately described by a linear model, Hooke's law becomes insufficient to describe the mechanical behaviour of the liver with precision for instance. In fact, *in vivo* experiments conducted by Melvin et al. (1973) suggest that mechanical behaviour of liver is non-linear and viscoelastic. The brain also appears to be non-linear and viscoelastic (Miller and Chinzei, 1997). The liver has the additional property of anisotropy, more precisely, it is transversely isotropic according to Chui et al. (2007).

These various mechanical aspects of materials make realistic modelling of anatomical structures fairly challenging. In medical simulation, where modelling must occur in real-time or close to real-time, very often simplifications have to be made. As we will see in chapter 4 and 7, numerous models are available in the literature to describe physics of objects, from fairly simple and naive approaches to more complex and thorough representations. However, the most common technique to solve equations of continuum mechanics is the Finite Element Method. Moreover, since the two main contributions of this thesis deal with finite element methods formulations, the general principles of this method will be detailed in the next section.

CHAPTER 3

A PRACTICAL APPROACH OF THE FINITE ELEMENT METHOD

The previous chapter explained in details how physics of anatomical structures could be described by continuum mechanics. When an organ is submitted to a force, the first effect of this application is a change in the organ's shape. This deformation can be measured using different strain tensors. By modelling the mechanical aspects of the material with an appropriate constitutive law, we can then deduce the stress tensor (or internal forces) at each point within the organ. The next step is to apply Newton's second law (law of conservation of momentum), which yields a partial differential equation at every point of the volume. The finite element method is a numerical procedure that allows to turn an infinite number of partial differential equations into a finite number of algebraic equations. The method follows different steps. We will first begin by describing the subdivision of the domain into smaller sub-domains into which local simpler equations are derived. We will then discuss how all the local equations are summed over the whole domain which eventually leads to the global solution.

3.1 Introduction

3.1.1 A numerical method

One of the most important things engineers and scientists do is to model physical phenomena. Using assumptions concerning how the phenomena works and using the appropriate laws of physics governing the process, they can derive a mathematical model, often characterised by complex differential or integral equations relating various quantities of interest. However, because of the complexities in the geometry and complex boundary conditions found in real life problems, they cannot analytically solve these equations. A few decades ago, the only possible approach was to drastically simplify them, which was not always sufficient to find an approximate solution. Nowadays, in practice, most of the problems are solved using numerical methods. Indeed, with suitable mathematical models and numerical methods, computers can help solving many practical problems of engineering. Numerical methods typically transform differential equations governing a continuum to a set of algebraic equations of a discrete model of the continuum that are to be solved using computers (Reddy, 1993). The *finite element method* (FEM) is the most popular numerical procedure that is used to approximately solve differential equations, especially in continuum mechanics.

3.1.2 The basic ideas of FEM

The finite element method begins by dividing the structure into small pieces, manageable regions, called *elements*. The collection of all these elements makes up a *mesh* which approximates the problem geometry. Why is this a good idea? If we can expect the solution for an engineering problem to be very complex, it is mostly due to the complex geometry, on which a global solution is difficult to figure out. If the problem domain can be divided (*meshed*) into small elements, the solution within an element is easier to approximate (MacDonald, 2007). Over each finite element, the unknown variables are approximated using known functions. These functions can be linear or higher-order polynomial expressions that depend on the geometrical locations (*nodes*) used to define the finite element shape. The governing equations are integrated over each finite element using linear algebra techniques and the solution over the entire domain problem is obtained by summing (*assembling*) the solution of each element. Thus, the finite element method transforms an infinite number of differential equations (one can be defined at any point of the continuum) into a finite number of algebraic equations (depending on the chosen number of elements).

This approach may be compared to trying to find the area under a curve. We know that we can find the exact solution for the area under the curve by integration. However, sometimes the function describing the curve is not known, or is difficult to integrate. One method to obtain an approximate solution is to break up the area into a series of rectangles and add the areas of all rectangles (see Fig. 3.1). It is worth noting that the solution accuracy can be increased by reducing the width of

the rectangles to better follow the curve.

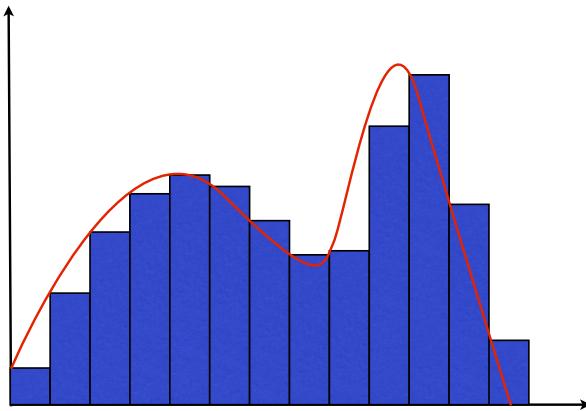


Figure 3.1: Approximation of the area under a curve by adding the areas of all rectangles (left Riemann sum for this example).

It is crucial to keep in mind that approximations occur at different stages during finite element analysis. The division of the whole domain into finite elements may not be exact (see Fig. 3.2), introducing error in the domain being modelled. The second stage is when element equations are derived. As mentioned earlier, the unknowns of the problem are approximated using the idea that any continuous function can be represented by a linear combination of known functions and undetermined coefficients. Algebraic relations between the undetermined coefficients are then obtained by satisfying the governing equations over each element. There are a few types of approach for establishing these equations but, without going into details, the mathematical foundations of all these approaches are energy principles or the *weighted residual methods* (see Appendix B), which both lead to integrals during the process. Therefore, the second stage creates two sources of error: the representation of the solution by a linear combination of functions and the evaluation of the integrals. Finally, errors are introduced when solving the assembled system of algebraic equations.

3.2 Discretisation

3.2.1 Meshing process

The first step in the finite element method is to create a mesh of the domain to study. Mesh generation is a very important task and can be very time consuming. The domain has to be meshed properly into elements of specific shapes. All the elements together form the entire domain of the problem without any gap or overlapping. For example, triangles or quadrilaterals can be used in two dimensions, and tetrahedra and hexahedra in three dimensions. Information, such as *element connectivity*, must be created during the meshing process for later use in the formulation of the FEM

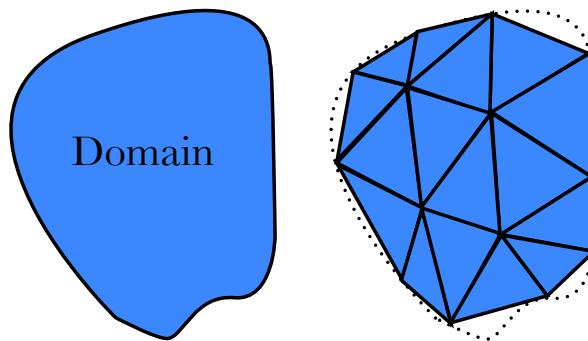


Figure 3.2: The division of the domain into elements is not always perfect.

equations. The number of elements into which the domain is divided in a problem depends mainly on the geometry of the domain and on the desired accuracy of the solution. Usually, the number of elements increases with the complexity of the geometry. For instance, if one part of the domain is thinner (see Fig. 3.3), the size of the elements must be reduced in order to tile this particular part, which increases the total number of elements, and hence the number of algebraic equations to solve. However, adding elements is sometimes desirable. Indeed, increasing the number of elements tends to get an approximate solution closer to the exact analytical solution (as reducing the width of the rectangles allows for a better approximation of the area under the curve). Elements are classically added in the particular regions of interest for a given problem, if this information is known *a priori*. A trade-off between accuracy and computational time must be found. This is why the meshing process on the problem domain must be carried out carefully. One needs to create a mesh which gives an accurate enough solution for the desired application while restraining the computational time according to the time constraints of the problem. As an example, finding the maximum load that can sustain a bridge in structural mechanics demands a high degree of accuracy, no matter how much time it is needed to compute the solution. Conversely, organ deformation in medical training simulators must be computed at an interactive rate so that no apparent delay can be observed between the manipulation of a given organ and its deformation. It does not mean that precision is not required, simply that the time constraints will necessarily limit the accuracy of the solution.

3.2.2 Solution interpolation

As we have seen, the finite element method is based on finding an approximate solution over each simple element rather than the whole domain. Any continuous function f may be approximated by a linear combination of known functions ϕ_i and



Figure 3.3: By adding details on the tree's geometry (**right**), the surface delimited by the tree becomes more difficult to mesh and would require smaller elements to tile all additional small branches.

undetermined coefficients c_i :

$$f \approx \tilde{f} = \sum_{i=1}^n c_i \phi_i \quad (3.1)$$

Moreover, the approximation solution u^e within the element Ω_e must fulfill certain conditions of continuity in order to be convergent to the actual solution u as the number of elements is increased. The finite element method approximates the solution by the following polynomial expression:

$$u^e(x) = \sum_{i=1}^n u_i^e \psi_i^e(x) \quad (3.2)$$

where Ω_e is a one-dimensional element¹, x a position within this element, u_i^e are the values of the solution u at the nodes and $\psi_i^e(x)$ the approximation functions over the element. This particular form will be assumed for brevity but the interested reader may refer to Reddy (1993) for demonstration. Note that u_i^e plays the role of the undetermined coefficients and ψ_i^e the role of approximation functions. Writing the approximation in terms of the nodal values of the solution is necessitated by the fact that we require the solution to be the same at points common to the elements in order to connect the approximate solution from each element and form a continuous solution over the whole domain. The approximation ψ_i^e can be linear or higher-order polynomial expressions and are called interpolation functions. Depending on the degree of polynomial approximation used to represent the solution, additional nodes may be identified inside the element. It is worth noting that the type of interpolation directly affects the accuracy of the solution. In other words, finding

¹For the sake of simplicity, all derivations are carried out in 1D but remain valid for each component in higher dimensions.

the solution to the problem consists of merely figuring out the values of the sought variable at every node of the mesh. Its value at any other point of the domain may then be deducted by interpolation within elements.

The next logical question is how to derive the interpolation functions ψ_i^e for a given element. Their derivation depends only on the geometry of the element and the number and location of the nodes. The number of nodes must be equal to the number of terms in the polynomial expansion. Therefore each element contains a single interpolation function for each of its nodes. As stated by interpolation theory, each interpolation function is required to be equal to 1 at its corresponding node, and 0 at all other nodes:

$$\psi_i^e(x_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (3.3)$$

where x_j is the position of node j . From there, we can state that the interpolation function at node i may be written as:

$$\psi_i^e(x) = C_i \prod_{j \neq i}^{n-1} (x - x_j) \quad (3.4)$$

where j are the indices of the other nodes and C_i is a constant to be determined such as:

$$\psi_i^e(x_i) = 1. \quad (3.5)$$

This function is zero at all nodes except the i^{th} node. It is worth noting that such a definition for the interpolation functions ψ_i^e used in (3.2) yields:

$$u^e(x_j) = \sum_{i=1}^n u_i^e \psi_i^e(x_j) = u_j^e \quad (3.6)$$

as expected.

Since elements over the whole domain are generally not identical (different shapes), this process of determining all interpolation functions can become fastidious. We will now introduce a general method for defining interpolation functions which allows for arbitrary element shapes (subject to a given topology). To this end, we discuss the concept of element *natural coordinates*.

3.2.3 Natural coordinates

The natural coordinate system allows us to map every element into a typical and simpler element. As an example, an 8-node hexahedron element is shown in Fig. 3.4 as it appears in the global (x_1, x_2, x_3) and natural (ξ_1, ξ_2, ξ_3) coordinate systems. While in its natural coordinate system the element is a regular aligned cube, in the global coordinate system the element may assume any admissible arbitrary form. Essentially, admissible means that the element must not be too distorted, and must

certainly not fold back on itself. Another example is given Fig. 3.5 with a tetrahedral element. The components of each node has a simple expression in natural coordinates, usually 0, 1 or -1 and the centre of the coordinate system is taken at one of the nodes or at the centre of the element.

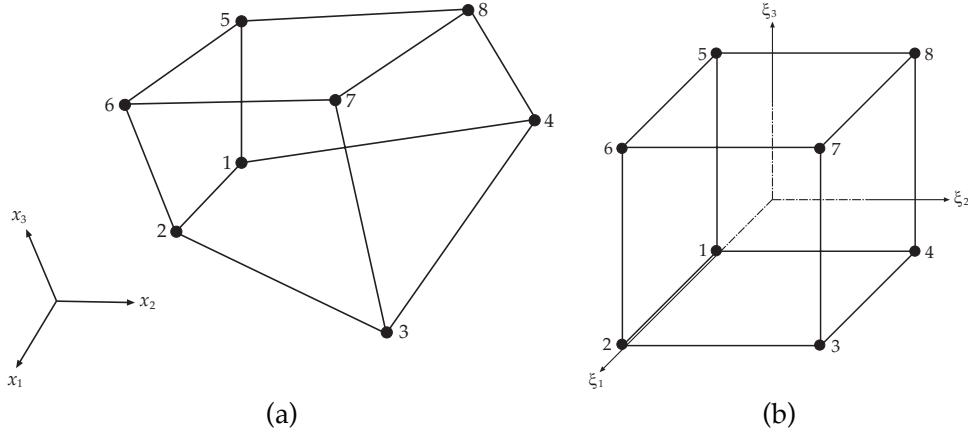


Figure 3.4: Hexahedral element as it appears in the global coordinate system (a) and its natural coordinate system (b).

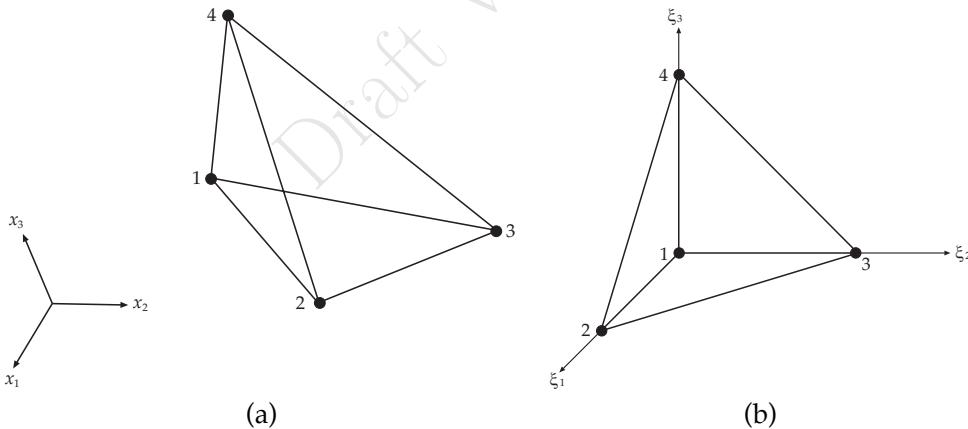


Figure 3.5: Tetrahedral element as it appears in the global coordinate system (a) and its natural coordinate system (b).

The use of the natural coordinate system has several advantages (Biswas et al., 1976). Not only the interpolation functions can be derived only once per type of element in the mesh, but their expression is also much simpler, regardless of the actual element shape. Consequently, the element equations and the derived element matrices get to be simplified too. However, we now need to switch between the natural and global coordinate systems before solving the equations on the whole

domain. Indeed, it is necessary to account for the diversity of element shapes within the mesh. Solving the element equations expressed into the natural coordinate system would not yield the expected result for the whole domain.

3.2.4 Geometry interpolation

Let us assume a relation between the global coordinate x and the natural coordinate ξ in the following form:

$$x = f(\xi) \quad (3.7)$$

where f is assumed to be a one-to-one correspondence (that is, a bijective function). This function may be seen as a transformation between the natural shape of the element and its arbitrary shape within the actual mesh, it describes a change in geometry. It is natural to think of approximating the geometry in a similar way that we approximated the solution in section 3.2.2. Hence, the transformation defined by (3.7) may also be written as

$$x = \sum_{i=1}^m x_i^e \hat{\psi}_i^e(\xi) \quad (3.8)$$

where x_i^e is the global coordinate of the i^{th} node of the element Ω_e , m the number of nodes for the element and $\hat{\psi}_i^e(\xi)$ are the interpolation functions of degree $m - 1$. Thus, we have a linear transformation when $m = 2$ and the relation between x and ξ is quadratic when $m = 3$. The interpolation functions $\hat{\psi}_i^e(\xi)$ are called *shape functions* because they are used to express the geometry or shape of the element.

Remark. It is worth noting that (3.8) can be easily extended to three-dimensional problems. Let us consider a Cartesian coordinate system and we denote the position of node i $\mathbf{x}_i = [x_{i1}, x_{i2}, x_{i3}]^T$. The position of a point $\mathbf{x} = [x_1, x_2, x_3]^T$ within the element may then be interpolated by applying (3.8) to each component of the position:

$$\begin{aligned} x_1 &= \sum_{i=1}^m x_{i1} \hat{\psi}_i^e(\xi) \\ x_2 &= \sum_{i=1}^m x_{i2} \hat{\psi}_i^e(\xi) \\ x_3 &= \sum_{i=1}^m x_{i3} \hat{\psi}_i^e(\xi) \end{aligned} \quad (3.9)$$

If we construct the vector $\hat{\mathbf{x}}^e$ of all nodal positions \mathbf{x}_i of the element, we can define a shape function matrix \mathbf{H} such as:

$$\mathbf{x} = \mathbf{H} \hat{\mathbf{x}}^e \quad (3.10)$$

where \mathbf{H} is a $3 \times m$ matrix built from submatrices \mathbf{H}_i

$$\mathbf{H}_i = \begin{bmatrix} \hat{\psi}_i^e & 0 & 0 \\ 0 & \hat{\psi}_i^e & 0 \\ 0 & 0 & \hat{\psi}_i^e \end{bmatrix} \quad i = 1, 2, \dots, m. \quad (3.11)$$

3.2.5 A particular case: isoparametric elements

Most of the time, we choose to interpolate the solution and the element geometry with the same interpolation functions. In this case, the elements are said to be *isoparametric*. Under this configuration, natural coordinates appear as parameters that define the shape functions. Note that the shape functions describe the geometry and the solution with the same degree of interpolation. Because isoparametric elements are very common in practice, the use of the phrase *shape functions* is often extended to denote the interpolation functions employed to approximate the solution.

3.3 Derivation of element equations

Let us sum up where we are into the finite element analysis process so far. The whole domain has been divided into sub-domains, smaller, that we call elements, for which the solution to be sought will be simpler to find. We know how to approximate the solution within each of these elements using a linear combination of shape functions. And we even have a method to simplify the expression of these shape functions using natural coordinates.

3.3.1 Strong and weak forms

The next step is therefore to derive the element equations themselves. Obviously these equations depend on the type of problem that we seek to solve. From now on, we will only consider the context of continuum mechanics since our overall goal is the modelling of organ deformation. However, a similar approach can be used in the other fields of physics. In continuum mechanics, applying Newton's second law yields a partial differential system of equations. Such equations are called *strong forms*. The strong form, in contrast to a *weak form*, requires strong continuity on the dependent field variables, such as displacements (Liu and Quek, 2003). The functions defining these field variables have to be differentiable up to the order of the partial differential equations that exist in the strong form. Obtaining the exact solution for a strong form is usually very difficult for practical engineering problems. We know that the finite element method can be used to find an approximated solution but the method usually works well only with regular geometry and boundary conditions.

A weak form is often created using energy principles or weighted residual methods. The weak form is often an integral form and requires a weaker continuity on the field variables. Due to the weaker requirement on the field variables, and the

integral operation, a formulation based on a weak form usually produces a set of discretised system equations that give much more accurate results, especially for problems of complex geometry. Hence, the weak form is usually preferred for obtaining an approximate solution. Using the weak form usually leads to a set of well-behaved algebraic system equations, if the problem domain is discretised properly into elements.

3.3.2 Time dependence

In all the equations of the finite element method that we derived so far, we have not taken time into account. For some problem, this is not an issue, the motion is sufficiently slow to assume that the dynamic effects (such as damping effect) are neglectable. However, this assumption is not always valid, particulary in our case where we want to model organ deformation. The deformation obviously depends on time. Finite element models of time-dependent problems can be developed in two alternatives ways (Reddy, 1993):

- (a) coupled formulation in which the time t is treated as an additional coordinate along with the spatial coordinate x
- (b) decoupled formulation where time and spatial variations are assumed to be separable.

Thus, the approximation of the solution u takes one of these two forms:

$$u(x, t) \approx u^e(x, t) = \sum_{i=1}^n \hat{u}_i^e \hat{\psi}_i^e(x, t) \quad (\text{coupled formulation}) \quad (3.12)$$

$$u(x, t) \approx u^e(x, t) = \sum_{i=1}^n u_i^e(t) \psi_i^e(x) \quad (\text{decoupled formulation}) \quad (3.13)$$

where $\hat{\psi}_i^e(x, t)$ are time-space interpolation functions, \hat{u}_i^e are the nodal values that are independent of x and t , $\psi_i^e(x)$ are the usual interpolation functions in spatial coordinate x only and the nodal values $u_i^e(t)$ are functions of time t only. Space-time coupled finite element formulations are not common and they have not been adequately studied. Hence, we consider the decoupled formulation only. Of course, the assumption that the time and spatial variations are separable is not valid in general. However, with sufficiently small time steps, it is possible to obtain accurate solutions nevertheless.

3.3.3 Dynamic system of equations

The space-time decoupled finite element formulation of time-dependent problems involves two steps:

1. **Spatial approximation.** The solution u is first approximated by expressions of the form (3.13) while keeping the time-dependent term during the finite

element model derivation. This step leads to a set of ordinary differential equations in time.

- 2. Temporal approximation.** The system of ordinary differential equations is then approximated in time where different schemes may be used to assess the time derivatives. This step yields a set of algebraic equations for u_i^e at discrete time $t_{k+1} = (k + 1)\Delta t$ where Δt is the time increment and $k > 0$ an integer.

At the end of this two-step approximation, we have a continuous spatial solution at discrete intervals in time:

$$u(x, t_k) \approx u^e(x, t_k) = \sum_{i=1}^n u_i^e(t_k) \psi_i^e(x) \quad \text{with } k = 0, 1, \dots \quad (3.14)$$

The construction of the weak form using either energy principles or weighted residual methods will not be detailed here. The reader may easily find it in dozens of textbooks. By substituting (3.14) in the weak form, we obtain the finite element equations of equilibrium in matrix form:

$$\mathbf{M}^e \ddot{\mathbf{U}}^e + \mathbf{D}^e \dot{\mathbf{U}}^e + \mathbf{K}^e(\mathbf{U}^e) \cdot \mathbf{U}^e = \mathbf{R}^e, \quad (3.15)$$

where \mathbf{M}^e is a constant mass matrix, \mathbf{D}^e is a constant damping matrix, $\mathbf{K}^e(\mathbf{U}^e)$ is the stiffness matrix, which is a function of nodal displacements \mathbf{U}^e , and \mathbf{R}^e are externally applied loads. All matrices are expressed at the level of individual elements (which is the reason for the superscript e).

3.3.4 Static system of equations

Under some assumption, (3.15) may be simplified. For instance, let us consider the problem of finding the maximum load that a desk can sustain. It is fair to say that strain and stress within the structure does not fluctuate over time: the system is in an equilibrium state. At all times, the sum of all forces applied on the system is equal to zero. Therefore, we can assume that the equations describing the structure do not depend on time: the problem is said to be *static*. Static analyses assume that loading of a body occurs slowly enough that inertial and damping terms may be neglected. In this case, the static system of equations can be easily obtained by merely dropping out the dynamics terms in (3.15). Indeed, the static case may just be seen as a special case of the dynamic equations and the static equations are:

$$\mathbf{K}^e(\mathbf{U}^e) \cdot \mathbf{U}^e = \mathbf{R}^e. \quad (3.16)$$

3.3.5 A few words on the matrices involved

Mass matrix

The mass matrix \mathbf{M} for the whole domain may be computed by summing the mass contributions from every elements:

$$\mathbf{M} = \sum_e \mathbf{M}^e, \quad (3.17)$$

where \mathbf{M}^e is the mass of element Ω_e given by:

$$\mathbf{M}^e = \int_{V^e} \rho^e \mathbf{H}^T \mathbf{H} dV \quad (3.18)$$

where ρ the mass density, V^e the volume of the element and \mathbf{H} is the shape function matrix defined by (3.10). This integral may be evaluated using a procedure of numerical integration and generally results in a band mass matrix (non-zero entries are confined to a diagonal band, comprising the main diagonal and more diagonals on either side). This method for computing the mass matrix is called *variational mass lumping*.

However, it is very common in finite element analysis to use a diagonal mass matrix. The most common way of achieving this is to define the diagonal entries with the sum of the corresponding rows and then set all other values to 0. The physical interpretation of this is lumping all of the mass in the volume surrounding a node at the node itself. A diagonal mass matrix may offer significant computational and storage advantages. First, it is easily inverted, since the inverse of a diagonal matrix is also diagonal. And a diagonal mass matrix can also be stored simply as a vector. Therefore, such a matrix may be built by directly computing the mass of each element and assigning an equal proportion of this to each of the element's nodes. This method is named *direct mass lumping*. The mass m^e of element Ω_e is given by:

$$m^e = \rho^e \int_{V^e} dV, \quad (3.19)$$

which is integrated over the element using a numerical method. Each diagonal component of the mass matrix then receives a contribution of m^e/m .

Damping matrix

In a similar way to the mass matrix, a damping matrix may be defined by summing the damping contributions from every elements:

$$\mathbf{D} = \sum_e \mathbf{D}^e, \quad (3.20)$$

where \mathbf{D}^e is the damping of element Ω_e given by:

$$\mathbf{D}^e = \int_{V^e} \kappa^e \mathbf{H}^T \mathbf{H} dV \quad (3.21)$$

where κ^e is a damping parameter. In practice, the components of damping matrices are difficult to determine experimentally. Therefore, it is common to employ the *Rayleigh damping*. The Rayleigh damping assumes that the damping may be written in the form of:

$$\mathbf{D}^e = \alpha \mathbf{M}^e + \beta \mathbf{K}^e, \quad (3.22)$$

where α and β are coefficients defining mass and stiffness proportional damping components. However, for systems with large degrees of freedom, it is difficult to

Elements	1	2
1	1	2
2	2	3
3	3	4

Table 3.1: Element connectivity information

guess meaningful values of α and β . Chowdhury and Dasgupta (2003) have proposed a method to calculate these coefficients. In practice, further assumptions may only lead to include the mass proportional component ($\beta = 0$), which yields a diagonal damping matrix if a diagonal mass matrix is used.

Stiffness matrix

The expression of the stiffness matrix depends on the chosen strain tensor as well as the appropriate constitutive model for the material. Its general form cannot be explicated. However, it can be demonstrated that \mathbf{K}^e may be written as

$$\mathbf{K}^e = \int_{V^e} \mathbf{B}^T \mathbf{E} \mathbf{B} dV, \quad (3.23)$$

where \mathbf{E} is a material matrix, which only depends on the constitutive model chosen, and \mathbf{B} is the *strain-displacement matrix* defined as the derivatives of the shape functions (see Liu and Quek (2003); Belytschko et al. (2000) for details).

3.4 Assembly of element equations

We have now derived the equations for each element. However, since the element is physically connected to its neighbours, the resulting algebraic equations will contain more unknowns than the number of algebraic equations. In other words, elements share vertices. The deformation of an element induces the deformation of its neighbours. Thus, it becomes necessary to put the elements together to eliminate the extra unknowns, to model this interaction between elements. The assembly of the global stiffness matrix should be carried out as soon as the element matrices are computed, rather than waiting until the element coefficients of all elements are computed. This would require storage of the coefficients for each element. Instead, we can perform the assembly while we calculate the element matrices.

The assembly operation is actually quite simple. Let us consider an example with stiffness matrices for an one-dimensional problem. We assume that we have three two-node elements. And we suppose the element connectivity information given by Table 3.1, which provides the global node number for each local node number within an element. We also assume the following stiffness matrices for the three elements:

$$\mathbf{K}^1 = \begin{bmatrix} -1 & 2 \\ 3 & 1 \end{bmatrix}, \quad \mathbf{K}^2 = \begin{bmatrix} 5 & 2 \\ -2 & 1 \end{bmatrix} \text{ and } \mathbf{K}^3 = \begin{bmatrix} 3 & 1 \\ 2 & -1 \end{bmatrix} \quad (3.24)$$

The assembled matrix is a square matrix of dimension number of nodes \times number of degrees of freedom. In our case, the global matrix is therefore a 4×4 and is assembled in the following manner:

$$\mathbf{K} = \begin{bmatrix} -1 & 2 & 0 & 0 \\ 3 & 1+5 & 2 & 0 \\ 0 & -2 & 1+3 & 1 \\ 0 & 0 & 2 & -1 \end{bmatrix} \quad (3.25)$$

In other words, the global matrix is merely a rewriting on the whole domain of the element stiffness matrices, which are expressed with respect to local node numbers, in terms of global node numbers. In addition to constructing the global mass and damping matrices (as shown in section 3.3.5), this process leads to the following global system of equations:

$$\mathbf{M}\ddot{\mathbf{U}} + \mathbf{D}\dot{\mathbf{U}} + \mathbf{K}(\mathbf{U}) \cdot \mathbf{U} = \mathbf{R}. \quad (3.26)$$

The same set of equations was established at the element level in section 3.3.3 but we may note that the same system is also valid for the whole domain.

3.5 Solution of global problem

The complete and global system of equations has now been established. In this section, we will examine the techniques commonly used for solving these equations. In the general case of a dynamic system of equations, the equations are not yet algebraic and still depends on derivatives of the unknown displacements with respect to time. Therefore, an additional step needs to be performed before using a solver by integrating the equations over time. As we will see shortly, the different methods of time integration are classified into two main schemes: *explicit* or *implicit* integration. We begin by introducing the simplest method, the Euler method to explain the concept of an explicit method. Then the central difference method will be detailed. We then describe a few implicit methods, allowing us to highlight the relative advantages of the two types of methods. Finally, the static case will be solved by emphasizing the similarities with implicit systems through the critical step of linearisation of the governing equations (Belytschko et al., 2000).

3.5.1 Explicit time integration

The Euler method

We want to approximate the solution of the following differential equation:

$$\dot{\mathbf{U}}^t = f(t, \mathbf{U}^t). \quad (3.27)$$

Note that we restrict ourselves to first-order differential equations (meaning that only the first derivative of \mathbf{U}^t appears in the equation, and no higher derivatives).

However, a higher-order equation can easily be converted to a system of first-order equations by introducing extra variables. For example, the second-order equation $\ddot{\mathbf{U}}^t = -\mathbf{U}^t$ can be rewritten as two first-order equations: $\dot{\mathbf{U}}^t = \mathbf{V}^t$ and $\dot{\mathbf{V}}^t = -\mathbf{U}^t$. By solving equations like (3.27), we will therefore be able to integrate the global system of equations over time.

We start by replacing the derivative $\dot{\mathbf{U}}^t$ by the finite difference approximation:

$$\dot{\mathbf{U}}^t \approx \frac{\mathbf{U}^{t+\Delta t} - \mathbf{U}^t}{\Delta t}, \quad (3.28)$$

which gives:

$$\mathbf{U}^{t+\Delta t} \approx \mathbf{U}^t + \Delta t \dot{\mathbf{U}}^t \quad (3.29)$$

And substituting the last relation in (3.27) yields:

$$\mathbf{U}^{t+\Delta t} \approx \mathbf{U}^t + \Delta t f(t, \mathbf{U}^t) \quad (3.30)$$

This formula is usually applied in a recursive scheme. We choose a time step Δt and we construct the sequence of time $t^n = n\Delta t$. Using (3.30), we can approximate the exact solution \mathbf{U}^{t^n} at each time step by \mathbf{U}^n :

$$\mathbf{U}^{n+1} = \mathbf{U}^n + \Delta t f(t^n, \mathbf{U}^n). \quad (3.31)$$

This method is called *Euler method* and is said to be explicit because the new value \mathbf{U}^{n+1} is defined in terms of variables from the previous time step \mathbf{U}^n .

The central difference method

The central difference method is one of the most popular explicit methods in computational mechanics. This method is based on central difference formulas for the velocity and acceleration. Let us the time t of the simulation be subdivided into time steps Δt . Let us denote by \mathbf{U}^t , $\dot{\mathbf{U}}^t$ and $\ddot{\mathbf{U}}^t$ the matrices of nodal displacements, velocities and accelerations, respectively, at time t . The central difference method makes use of the following finite differences approximations for acceleration and velocity:

$$\ddot{\mathbf{U}}^t = \frac{1}{(\Delta t)^2} (\mathbf{U}^{t-\Delta t} - 2\mathbf{U}^t + \mathbf{U}^{t+\Delta t}) \quad (3.32)$$

$$\dot{\mathbf{U}}^t = \frac{1}{2\Delta t} (\mathbf{U}^{t+\Delta t} - \mathbf{U}^{t-\Delta t}). \quad (3.33)$$

One of the key advantage of the method is that the stiffness term in (3.26) may be computed from

$$\mathbf{K}(\mathbf{U}) \cdot \mathbf{U} = \mathbf{F}(\mathbf{U}) = \sum_e \tilde{\mathbf{F}}^e, \quad (3.34)$$

where $\tilde{\mathbf{F}}^e$ are the global nodal force contributions due to stresses in element Ω_e . $\tilde{\mathbf{F}}$ may be computed using strain-displacement equations, the constitutive equation

and the element connectivity. Therefore, let us assume that all nodal forces \mathbf{F}^t at time t are known.

By substituting (3.32) and (3.33) into (3.26) we obtain:

$$\left(\frac{\mathbf{M}}{(\Delta t)^2} + \frac{\mathbf{D}}{2\Delta t} \right) \mathbf{U}^{t+\Delta t} = \mathbf{R}^t - \mathbf{F}^t + \frac{2\mathbf{M}}{(\Delta t)^2} \mathbf{U}^t + \left(\frac{\mathbf{D}}{2\Delta t} - \frac{\mathbf{M}}{(\Delta t)^2} \right) \mathbf{U}^{t-\Delta t}. \quad (3.35)$$

Both the right-hand side of the equation and the factor in front $\mathbf{U}^{t+\Delta t}$ may be evaluated, which gives the nodal displacements for the next time step. The entire update can be accomplished without solving any system of equations provided that the mass matrix \mathbf{M} and the damping matrix \mathbf{D} are diagonal. As we have seen in section 3.3.5, these approximations are often used in practice. And this is the characteristic of an explicit method: the time integration of the general system of equations for finite element model does not require the solution of any equations. This is a key point because it gives explicit methods a tremendous computational advantage. However, explicit methods are only conditionally stable. If the time step exceeds a critical value, the solution may be erroneous. This critical time step depends on the mesh and the material properties. It will decrease with mesh refinement and increasing stiffness of the material.

3.5.2 Implicit time integration

The backward Euler method

The problem that we wish to solve is the same than the one we introduced with the Euler method, which we reproduce here for convenience:

$$\dot{\mathbf{U}}^t = f(t, \mathbf{U}^t). \quad (3.36)$$

This time, instead of using (3.28), we approximate the derivative $\dot{\mathbf{U}}^t$ by:

$$\dot{\mathbf{U}}^t \approx \frac{\mathbf{U}^t - \mathbf{U}^{t-\Delta t}}{\Delta t}. \quad (3.37)$$

The update expression for y^n then becomes:

$$\mathbf{U}^{n+1} = \mathbf{U}^n + \Delta t f(t^{n+1}, \mathbf{U}^{n+1}). \quad (3.38)$$

The *backward Euler method* is an implicit method, meaning that we have to solve an equation to find \mathbf{U}^{n+1} . Of course, it costs time to solve this equation and this cost must be taken into consideration when we select the time integration method to use. The advantage of implicit methods is that they are usually more stable for solving a stiff equation, meaning that larger steps can be used.

The Runge-Kutta method

The accuracy of backward Euler method can be improved by making the solution depends on more values. The 4th order Runge-Kutta method gives the following

update equation:

$$\mathbf{U}^{n+1} = \mathbf{U}^n + \frac{1}{6} \Delta t (k_1 + 2k_2 + 2k_3 + k_4), \quad (3.39)$$

where :

$$\begin{aligned} k_1 &= f(t^n, \mathbf{U}^n) \\ k_2 &= f\left(t^n + \frac{1}{2}\Delta t, \mathbf{U}^n + \frac{1}{2}\Delta t k_1\right) \\ k_3 &= f\left(t^n + \frac{1}{2}\Delta t, \mathbf{U}^n + \frac{1}{2}\Delta t k_2\right) \\ k_4 &= f(t^n + \Delta t, \mathbf{U}^n + \Delta t k_3) \end{aligned}$$

This method is a fourth-order method because the total accumulated error has order Δt^4 . It may be seen as a generalisation of the backward Euler method (Press et al., 1992). Again, we have to solve an equation to find \mathbf{U}^{n+1} .

3.5.3 Static solutions

First, we recall the general system of equations in static introduced in section 3.3.4, which is also valid for the whole system:

$$\mathbf{K}(\mathbf{U}) \cdot \mathbf{U} = \mathbf{R}. \quad (3.40)$$

The situation is similar to implicit time integration methods for dynamic systems by the fact that we also have to solve a set of algebraic equations to find the nodal displacements. In the general case, the equations are non-linear and therefore quite difficult to solve. A suitable solver for this kind of system is, for instance, the *Newton-Raphson* method. In essence, this method consists of obtaining linear equations instead. This process is called *linearisation*. One may refer to (Belytschko et al., 2000) for more details.

However, if the problem is linear, the stiffness matrix does not depend on the displacement and the resolution is facilitated. The unknown displacements may be computed by:

$$\mathbf{U} = \mathbf{K}^{-1} \mathbf{R}. \quad (3.41)$$

Many techniques are available in the literature to solve this type of system, and a few of them are introduced in the next section.

3.5.4 Solvers

We will now examine the techniques commonly used for solving linear systems of algebraic equations. Note that this is only the most classical approaches and the list is obviously not exhaustive. Two categories of techniques are in opposition: *direct* or *iterative* solvers (Press et al., 2002). In any case, we seek the solution for the unknown solution \mathbf{x} of the following system of linear equations:

$$\mathbf{A}\mathbf{x} = \mathbf{b}. \quad (3.42)$$

The system may already be linear by nature, or been linearised using techniques such as Newton-Raphson methods. At this point in the discussion, we are therefore only interested in solving a linear system of equations. If the matrix \mathbf{A} is squared, which is the case in finite element methods, the problem may be written as follows:

$$\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}. \quad (3.43)$$

However, only solvers relevant to the context of finite element methods will be discussed. Indeed the matrix \mathbf{A} takes a specific form in problems solved by finite element methods. The matrix \mathbf{A} is *sparse*. A sparse matrix is merely a matrix populated mainly with zeros. An example is given by Fig. 3.6 for a finite element problem in two dimensions. Because of the peculiar structure of the matrix involved in finite element computations, methods have been specifically designed to take advantage of this structure. Indeed, it is wasteful to use general methods of linear algebra on such problems, because most of the arithmetic operations devoted to solving the set of equations or inverting the matrix involve zero operands. Furthermore, it is wasteful to reserve storage for zero elements.

Direct solvers

As their name is suggesting, direct solvers make an attempt into finding a solution directly by inverting the matrix \mathbf{A} . They execute in a predictable number of operations. We will now introduce a few direct methods in the context of finite element methods.

LU decomposition. The LU decomposition is a matrix decomposition which writes a matrix as the product of a lower triangular matrix and an upper triangular matrix. In other words, \mathbf{A} may be written as:

$$\mathbf{A} = \mathbf{L}\mathbf{U}, \quad (3.44)$$

where \mathbf{L} is lower triangular (has elements only on the diagonal and below) and \mathbf{U} is upper triangular (has elements only on the diagonal and above). As an example, for the case of a 4×4 matrix \mathbf{A} , we have:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} = \begin{bmatrix} \alpha_{11} & 0 & 0 & 0 \\ \alpha_{21} & \alpha_{22} & 0 & 0 \\ \alpha_{31} & \alpha_{32} & \alpha_{33} & 0 \\ \alpha_{41} & \alpha_{42} & \alpha_{43} & \alpha_{44} \end{bmatrix} \begin{bmatrix} \beta_{11} & \beta_{12} & \beta_{13} & \beta_{14} \\ 0 & \beta_{22} & \beta_{23} & \beta_{24} \\ 0 & 0 & \beta_{33} & \beta_{34} \\ 0 & 0 & 0 & \beta_{44} \end{bmatrix} \quad (3.45)$$

By substituting (3.44) into (3.42) we obtain:

$$\mathbf{Ax} = (\mathbf{Lu})\mathbf{x} = \mathbf{L}(\mathbf{Ux}) = \mathbf{b}. \quad (3.46)$$

Therefore, we can solve (3.42) by first solving the vector \mathbf{y} such that:

$$\mathbf{Ly} = \mathbf{b} \quad (3.47)$$

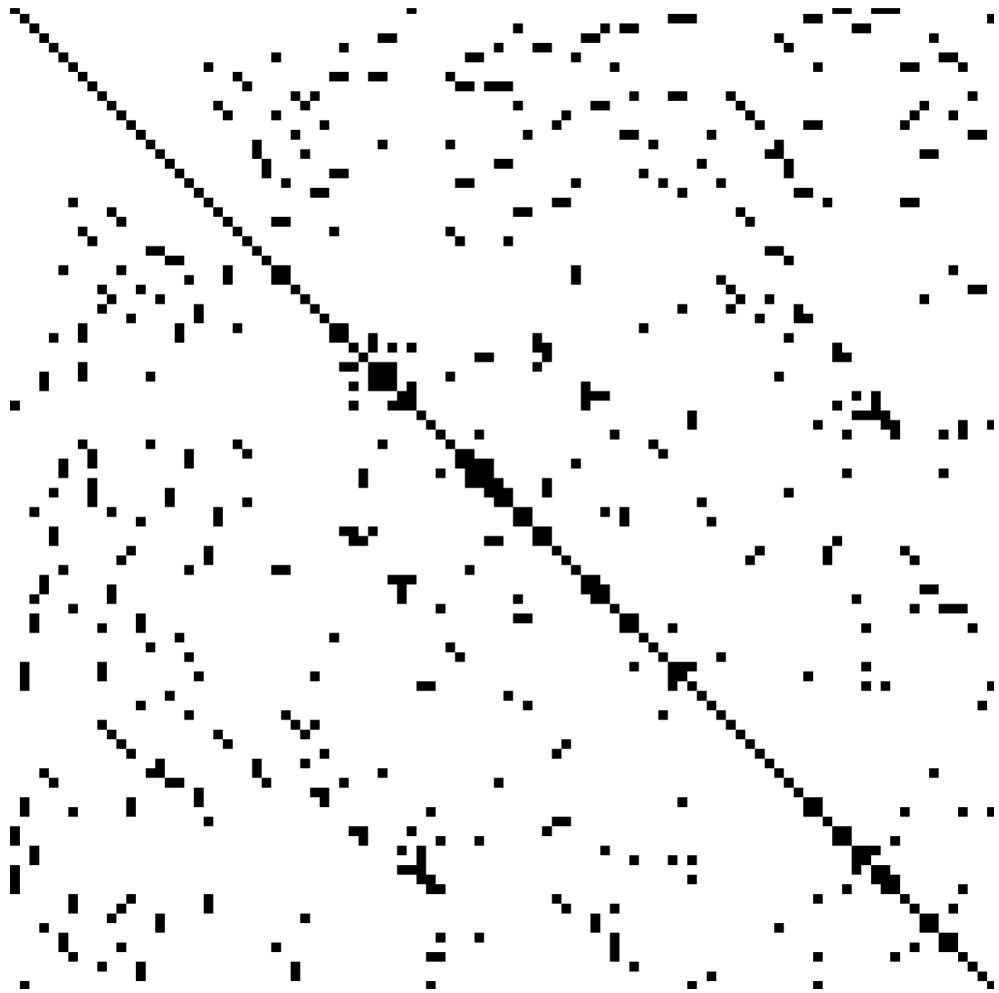


Figure 3.6: Illustration of a sparse matrix for a two-dimensional finite element problem. Each non-zero elements of the matrix is represented in black. Image courtesy of Oleg Alexandrov.

and then solving

$$\mathbf{U}\mathbf{x} = \mathbf{y}. \quad (3.48)$$

The advantage of breaking up one linear set of equations into two successive ones is that the solution of a triangular set of equations is quite trivial. Indeed, (3.47) can be solved by:

$$\begin{cases} y_1 = \frac{b_1}{\alpha_{11}} \\ y_i = \frac{1}{\alpha_{ii}} \left[b_i - \sum_{j=1}^{i-1} \alpha_{ij} y_j \right] \quad i = 2, 3, \dots, N. \end{cases} \quad (3.49)$$

for a $N \times N$ matrix \mathbf{A} . We then solve (3.48) in a similar way:

$$\begin{cases} x_N = \frac{y_N}{\beta_{NN}} \\ x_i = \frac{1}{\beta_{ii}} \left[y_i - \sum_{j=i+1}^N \beta_{ij} x_j \right] \quad i = N-1, N-2, \dots, 1. \end{cases} \quad (3.50)$$

TODO: Cholesky decomposition?

Iterative solvers

An iterative method attempts to solve the system of equations by finding successive approximations to the solution starting from an initial guess. In the case of a system of linear equations, the two main classes of iterative methods are the stationary iterative methods, and the more general Krylov subspace methods. Stationary iterative methods solve a linear system with an operator approximating the original one and based on a measurement of the error in the result (called the residual), form an equation of correction for which this process is repeated. While these methods are simple to derive, convergence is only guaranteed for a limited class of matrices. Examples of stationary iterative methods are the Jacobi method and Gauss-Seidel method. Krylov subspace methods work by forming an orthogonal basis of the sequence of successive matrix powers times the initial residual (the Krylov sequence). The approximations to the solution are then formed by minimising the residual over the subspace formed. The typical method in this class is the conjugate gradient method (CG). In practice, for large system of equations (as often encountered in finite element analysis), direct methods would take too much time and an iterative solvers is commonly used.

The Gauss-Seidel method. For this method, convergence is only guaranteed if the matrix is either diagonally dominant, or symmetric and positive-definite. In a similar way than the LU decomposition, we start by decomposing the matrix \mathbf{A} into a lower triangular component \mathbf{L}_* and a strictly upper triangular component \mathbf{U} :

$$\mathbf{A} = \mathbf{L}_* + \mathbf{U}. \quad (3.51)$$

The system of linear equations (3.42) may be rewritten as:

$$\mathbf{L}_* \mathbf{x} = \mathbf{b} - \mathbf{U} \mathbf{x} \quad (3.52)$$

The Gauss-Seidel method is an iterative technique that solves the left-hand side of the expression for \mathbf{x} using previous values for \mathbf{x} on the right-hand side. That is:

$$\mathbf{x}^{n+1} = \mathbf{L}_*^{-1}(\mathbf{b} - \mathbf{U} \mathbf{x}^n). \quad (3.53)$$

Similarly to the LU decomposition method, we take advantage of the triangular form of \mathbf{L}_* . The components of \mathbf{x}^{n+1} can be computed sequentially in the following manner:

$$x_i^{n+1} = \frac{1}{a_{ii}} \left[b_i - \sum_{j>i} a_{ij} x_j^n - \sum_{j<i} a_{ij} x_j^{n+1} \right] \quad i = 1, 2, \dots, N. \quad (3.54)$$

The procedure is usually continued until the changes made by an iteration are below a given threshold.

The conjugate gradient method. The simplest class of conjugate gradient method is based on the idea of minimising the quadratic function:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} + c. \quad (3.55)$$

Indeed, if the matrix \mathbf{A} is symmetric and positive-definite, $f(\mathbf{x})$ is minimised by the solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$ (Shewchuk, 1994). The minimisation is carried out by generating a succession of search directions \mathbf{p}_k and improved minimisers \mathbf{x}_k . At each stage a quantity α_k is found that minimises $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)$, and \mathbf{x}_{k+1} is set equal to the new point $\mathbf{x}_k + \alpha_k \mathbf{p}_k$. The \mathbf{p}_k and \mathbf{x}_k are built up in such a way that \mathbf{x}_{k+1} is also the minimiser of f over the whole vector space of directions already taken, that is the space formed by the basis of vector $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k$. After N iterations, you arrive at the minimiser over the entire vector space, that is, the solution to our linear system of equations.

Later, the *biconjugate gradient method* was introduced as a generalisation of the ordinary conjugate gradient to solve systems where the matrix \mathbf{A} was not necessarily symmetric and positive-definite (see Belytschko et al. (2000) for details).

Part II

Solid organs modelling

Draft Version

CHAPTER 4

MODELLING THE DEFORMATION OF SOLID OBJECTS IN REAL-TIME

Modelling the deformation of solid organs is a crucial problem in medical simulation. If the finite element method is classically used in mechanical engineering for its ability to solve the equations of continuum mechanics fairly accurately, its development in medical simulation was more progressive. The main reason for this is that the finite element method is computationally demanding and not easily compatible with time constraints required in interactive simulation. Therefore, other methods were developed to allow faster computations, at the cost of larger approximations of course. Over time, along with the increase in computational power, methods became more elaborate and less approximative. These different approaches may be grouped into three main categories: (1) geometrically based techniques, (2) approaches physically motivated and (3) methods actually based on the equations of continuum mechanics. In this chapter, we will review the most classic methods applied to modelling of solid objects with a particular emphasis on medical simulation.

4.1 Introduction: the problematic

At the beginning of computer simulation, users were immersed into environments that were merely a decor. If it could virtually take the user somewhere else, a static scene quickly showed its limits. Indeed, in a real-life environment many objects are in motion, interacting with each others. Therefore, in order to increase the fidelity of simulations, elements of physics were progressively added. It could be as simple as adding motion to clouds in a flight simulator for instance. But by improving the dynamism of the scene, the simulation appears more natural to users. Over time, with the increase of computers' computational power, objects stopped being all rigid and simulators started taking laws of physics into account. In other words, objects became deformable. If this took the realism a step further, the demands in computational power were multiplied. Simulators eventually faced a challenge: they had to weight the desired degree of realism against the computational power at their disposal.

In the field of medical simulation, the correctness of deformation is often crucial. This is particularly true in the cases of per-operative guidance. As an example, removing a brain tumor relies heavily on knowing the spatial relation between the patient's brain and the images acquired prior to surgery. Unfortunately, drilling a hole into the patient's skull releases pressure and the brain deforms. This phenomenon is called brain shift. The overall deformation is quite complex and the displacements depend on local elasticity of the brain, the size and the location of the tumor. A large brain shift, if not corrected, will result in inaccuracies in the surgical procedure and has the potential to cause damage to normal tissue. A solution is to accurately model the mechanical response of the patient's brain to predict the actual deformation and hence keep the knowledge of the tumor's location.

Conversely, when medical simulation is applied to training, such a precision in terms of distance is not required. The physician would not even visually notice reasonable errors as soon as it looks realistic enough. However, in some medical procedures, the sense of touch and the forces that the physician can feel are a key factor for the success of the operation. For instance, during a colonoscopy procedure, force feedback gives precious information about what type of loop is forming and the physician can then react accordingly to prevent loop formation. Consequently, a good colonoscopy simulator is required to provide accurate force feedbacks. If the contact forces were not reproduced to the operator, he would never be able to learn how to detect loop forming. Obviously, the computation of a realistic force to be returned demands an appropriate mechanical modelling of all structures.

If the reasons differ, the need for modelling the mechanical response of organs with precision remains. However, strong time constraints often limit the complexity of the modelling and affects the overall accuracy. Because of this trade-off, various approaches were proposed over the years to fit into real-time constraints and may be grouped into three main categories: (1) geometrically based techniques, (2) approaches physically motivated (usually relying on Newton's second law) and (3) methods actually based on the equations of continuum mechanics.

4.2 Techniques based on geometry

4.2.1 Free-form deformation

The first needs to deforming solid objects came from for the field of solid modelling. The goal was to allow the designer to shape an object into the form he wishes in a similar way that clay is manipulated by a sculptor's hands. Sederberg and Parry (1986) introduced a technique called *free-form deformation* based on the paper of Barr (1984). It involves a mapping from \mathbb{R}^3 to \mathbb{R}^3 through a tensor product trivariate Bernstein polynomial. The deformed shape of the object is interpolated using control points as in the case of Bézier curves and surface patches (see Fig. 4.1). The

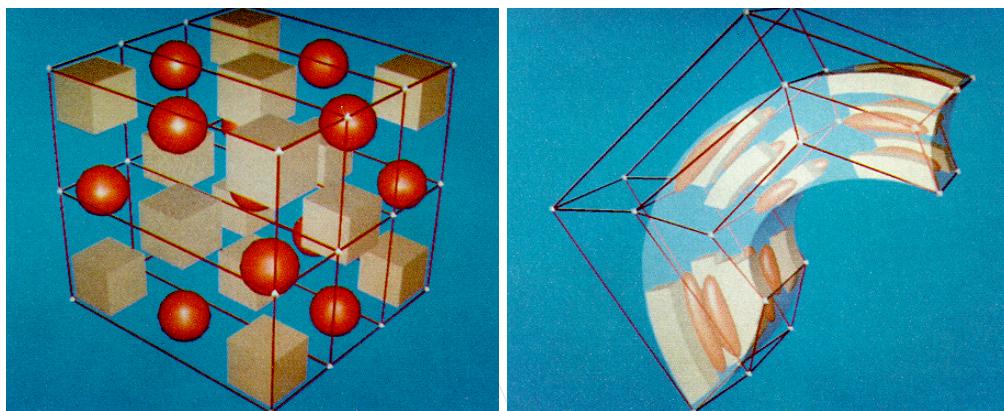


Figure 4.1: Let us imagine that the objects that we want to deform (spheres and cubes) are placed into a transparent parallelepiped (**left**). Control points are indicated by small white diamonds. The free-form deformation technique offers to compute the deformation of the objects from the positions of the control points (**right**).

capability of locally applying deformations makes the technique strongly analogous to sculpting with clay, which was the sought target. Although the authors do not give details on computational efficiency, their work constitutes the beginnings of deformable solid objects.

Coquillard later developed the idea of Extended free-form deformation (1990) and Animated free-form technique (1991). The idea of free-form deformation with lattices of arbitrary topology was introduced by MacCracken and Joy (1996). Finally, Schein and Elber (2004) presented a technique to facilitate the incorporation of discontinuities in a model while deforming it.

4.2.2 Shape matching

One field particularly interested in modelling deformable objects in a very efficient way is the game industry. Of course, for games we are more attracted by computational efficiency and extreme stability features than accuracy to physics laws.

Most of the time, we can tolerate a degradation of realism as long as the result looks realistic. Müller et al. (2005) developed a technique with this idea in mind called *shape matching*. We start with a set of particules with masses m_i in an initial configuration where we denote by \mathbf{x}_i^0 the initial positions of the particules. No connectivity information is required. The particles are simulated as a simple particle system without particle-particle interactions, but including response to collisions with the environment and including external forces such as gravity. The positions in the deformed configuration are noted \mathbf{x}_i . The method consists of finding a set of points \mathbf{g}_i which minimises the difference between the two sets \mathbf{x}_i^0 and \mathbf{x}_i . The first step is to find the rotation matrix \mathbf{R} and the translation vectors \mathbf{t} and \mathbf{t}_0 such that

$$\sum_i m_i (\mathbf{R}(\mathbf{x}_i^0 - \mathbf{t}_0) + \mathbf{t} - \mathbf{x}_i)^2 \quad (4.1)$$

is minimal. The optimal translation vectors \mathbf{t} and \mathbf{t}_0 turn out to be the centre of mass of the initial shape and actual shape, that we will note \mathbf{x}_c^0 and \mathbf{x}_c , respectively. The optimal rotation matrix is found by first finding the optimal linear transformation \mathbf{A} . The optimal rotation is eventually obtained through the rotational part of \mathbf{A} after a polar decomposition. Finally, the goal positions \mathbf{g}_i can be computed as follows:

$$\mathbf{g}_i = \mathbf{R}(\mathbf{x}_i^0 - \mathbf{x}_c^0) + \mathbf{x}_c. \quad (4.2)$$

Once the goal positions are known, they are used to integrate the positions of the particules:

$$\begin{cases} \mathbf{v}_i(t+h) = \mathbf{v}_i(t) + \alpha \frac{\mathbf{g}_i(t) - \mathbf{x}_i(t)}{h} + h \frac{f_{\text{ext}}(t)}{m_i} \\ \mathbf{x}_i(t+h) = \mathbf{x}_i(t) + h \mathbf{v}_i(t+h) \end{cases} \quad (4.3)$$

where α is a parameter between 0 and 1 which simulates stiffness. Because the method in this basic form is not suitable for objects undergoing large deformations, it is then extended to linear and quadratic deformations by replacing \mathbf{R} in (4.2) with combinations of the form $\beta\mathbf{A} + (1-\beta)\mathbf{R}$ where β is an additional parameter. Linear transformation can represent shear and stretch while quadratic transformation add twist and bending. Of course, this technique is not physically realistic. But this is beyond the point. This method allows the authors to simulate hundreds of deformable objects at an interactive rate with an unconditional stability (see Fig. 4.2). Another illustration of this stability is shown Fig. 4.3.

Later, Rivers and James (2007) extented the technique to regular lattices with embedded geometry. From a surface mesh to be deformed, they start by voxelising the model to construct a lattice of cubic cells containing the mesh. The embedded mesh is then deformed using trilinear interpolation of lattice vertex positions. A particule is placed at each vertex of the lattice and associated with a shape matching region. Their technique of *fast lattice shape matching* allows them to compute per-region transformations in an efficient manner (see Fig. 4.4).

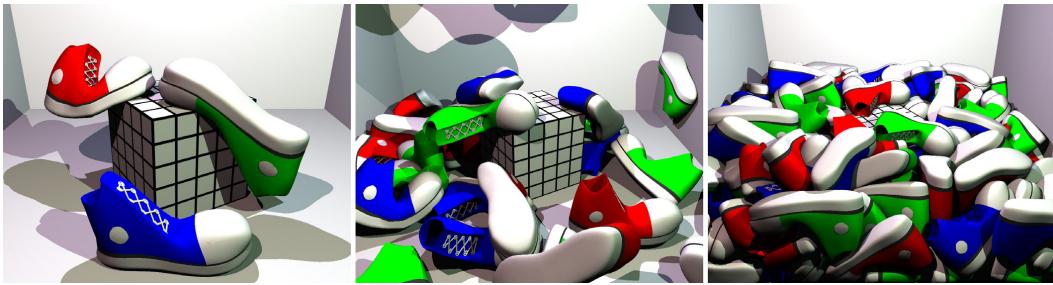


Figure 4.2: Shape matching allows to simulate hundreds of deformable objects in real-time.

One limitation of fast lattice shape matching is that small features of a geometry yield an explosion of the runtime cost. Indeed, a small surface feature may require fine sampling in order to be deformed independently from non-adjacent material, but this fine sampling must be applied to the whole object. Therefore, Steinemann et al. (2008) further improved the technique by allowing dynamic adaptive selection of levels of detail through the use of an octree-based hierarchical sampling. This adaptative sampling is illustrate with Fig. 4.5. Their technique also handles efficient topological changes.

Although realism is not a crucial feature for game-like environments, other models based on laws of physics were created for other applications.

4.3 Techniques relying on physics

4.3.1 3D Chainmail algorithm

The *3D Chainmail* was introduced by Gibson (1997) as a fast algorithm for deforming solid objects. Rather than carrying out complex computations on a low resolution geometry (as opposed to FEM techniques), the main idea is to take advantage of the original data resolution (of a CT scan for instance) but perform simple deformation calculation for each element. In this algorithm, each volume element is linked to its six nearest neighbours. If a link between two nodes is stretched to its limit, displacements are transferred to its neighbouring links. It works in a very similar way to a chain and is illustrated by Fig. 4.6. The chain mail algorithm is then followed by an elastic relaxation, which relaxes the shape of the deformed object approximated by the chain mail algorithm. This step consists of minimising an elastic energy measure within the object to insure that the calculated deformation is a minimum energy configuration.

This method was extented by Schill et al. (1998) to integrate inhomogeneous and anisotropic behaviour. This extension of the chain mail algorithm allows the authors to simulate the vitreous humor in the eye during a vitrectomy procedure. However, the main drawback of this algorithm is that it does not model volume preservation, an important characteristic of many tissues in the human body.

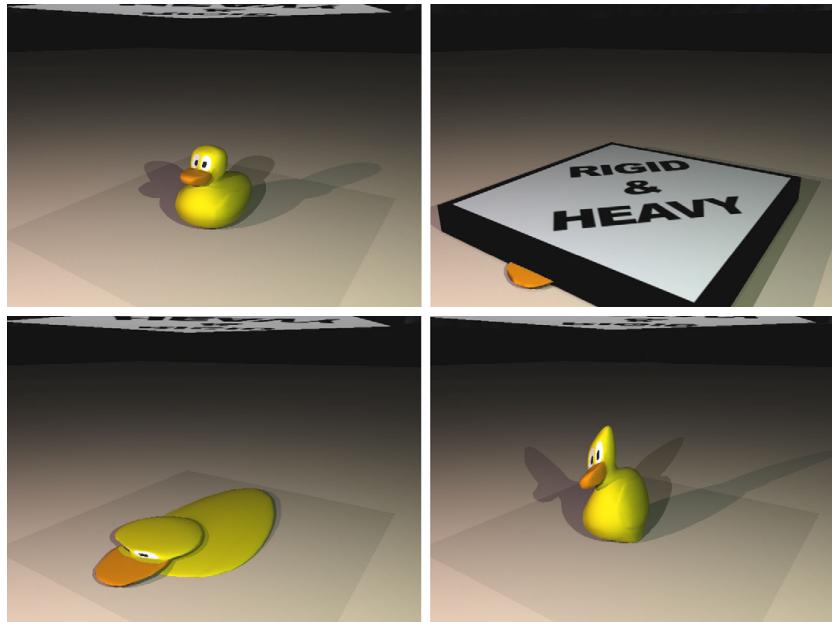


Figure 4.3: A duck model is squeezed to demonstrate the stability of the shape matching technique. Although the flat shape of the duck is not physically plausible, this example exposes the impressive ability of the approach to recover from highly deformed or inverted configurations.

4.3.2 Modal analysis

Modal analysis was first introduced to the graphics community by Pentland and Williams (1989) as a fast method for approximating deformation. Modal analysis is the process of taking the non-linear description of a system, finding a good linear approximation, and then finding a coordinate system that diagonalises the linear approximation. This process transforms a complicated system of non-linear equations into a simple set of decoupled linear equations that may be individually solved analytically (Hauser et al., 2003). Moreover, because each decoupled equation may be solved analytically, stability concerns due to the choice of a time integration procedure are eliminated. The linearised system of equations may be written as:

$$\mathbf{M}\ddot{\mathbf{U}} + \mathbf{D}\dot{\mathbf{U}} + \mathbf{K}\mathbf{U} = \mathbf{R}. \quad (4.4)$$

where \mathbf{M} is the mass matrix, \mathbf{D} is the damping matrix, \mathbf{K} is the stiffness matrix and \mathbf{R} are the externally applied loads. Modal decomposition refers to the process of diagonalising equation (4.4). By using Rayleigh damping (see discussion on damping in section 3.3.5), (4.4) may be rewritten as:

$$\mathbf{K}(\mathbf{U} + \alpha_1\dot{\mathbf{U}}) + \mathbf{M}(\alpha_2\dot{\mathbf{U}} + \ddot{\mathbf{U}}) = \mathbf{R}, \quad (4.5)$$

where α_1 and α_2 are the Rayleigh coefficients. We denote \mathbf{W} the solution to the generalised eigenproblem $\mathbf{Kx} + \lambda\mathbf{Mx} = 0$ and $\mathbf{\Lambda}$ the diagonal matrix of eigenvalues,

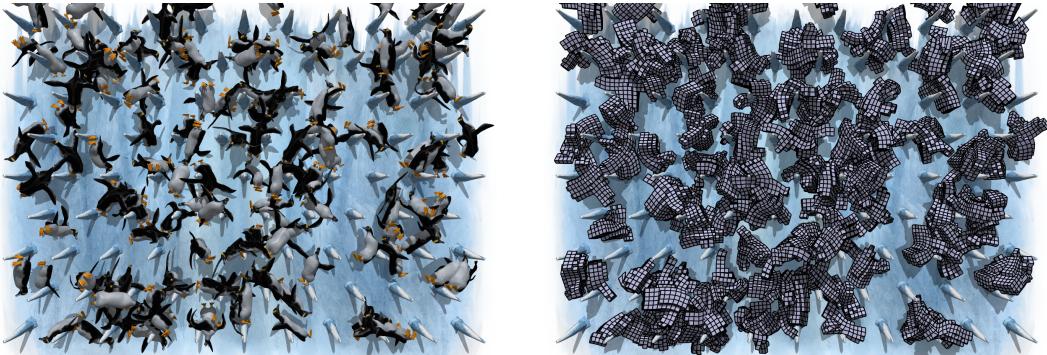


Figure 4.4: (**Left**) 150 cartoon penguins deforming dynamically using fast lattice shape matching(**Right**). Deformed lattices consisting of 150 particles per penguin (22 500 particles). The penguins can be deformed at 25 frames per second on a Pentium4 3.4 GHz.



Figure 4.5: Complex scene with 40 deforming flowers. A coarse sampling is employed when the flowers are moved by the wind, for a total of 5 680 nodes in the scene. However, the models are dynamically refined when touched by the user, as shown in the right image.

then (4.5) may be expressed as:

$$\boldsymbol{\Lambda}(\mathbf{Z} + \alpha_1 \dot{\mathbf{Z}}) + (\alpha_2 \dot{\mathbf{Z}} + \ddot{\mathbf{Z}}) = \mathbf{G}, \quad (4.6)$$

where $\mathbf{Z} = \mathbf{W}^{-1}\mathbf{U}$ is the vector of modal coordinates and $\mathbf{G} = \mathbf{W}^T\mathbf{U}$ is the external force vector in the modal coordinate system. Each row of this equation now corresponds to a single scalar differential equation:

$$\ddot{z}_i + (\alpha_1 \lambda_i + \alpha_2) \dot{z}_i + \lambda_i z_i = g_i. \quad (4.7)$$

The analytical solutions to each equations are well known and are the following:

$$z_i = c_1 e^{t\omega_i^+} + c_2 e^{t\omega_i^-} \quad (4.8)$$

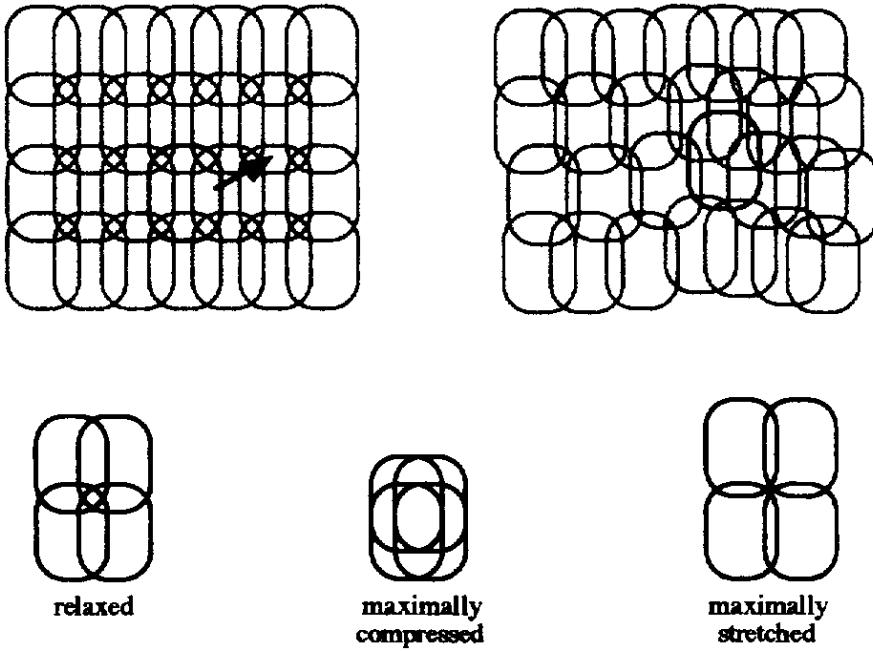


Figure 4.6: Deformation of a 2D chain mail when the link indicated by an arrow is moved.

where c_1 and c_2 are complex constants and ω_i is the complex frequency given by:

$$\omega_i^{\pm} = \frac{-(\alpha_1 \lambda_i + \alpha_2) \pm \sqrt{(\alpha_1 \lambda_i + \alpha_2)^2 - 4\lambda_i}}{2}. \quad (4.9)$$

The main advantage of modal analysis is the possibility to model only the necessary modes. Indeed if the eigenvalue λ_i associated with a particular mode is large, then the force required to cause a discernible displacement of that mode will also be large. Similarly, some displacements may be too small to be detected. Removing those modes from the computation will not change the appearance of the resulting simulation while substantially reducing the computational cost.

James and Pai (2002) went even further by using graphics hardware to reduce CPU costs to a minimum. In particular, they use modal analysis to model human skin and secondary tissues in a laparoscopic surgical simulation to allow the CPU to focus on more complex tissue models and user contact interactions (see Fig. 4.7). The technique is fairly efficient as it allows precomputations of the vibration modes.

Although modal analysis significantly accelerates the simulation, it generates noticeable artifacts when applied to large deformations due to the step of linearisation. It can produce unnatural results with deformations of large amplitude. Choi and Ko (2005) proposed to overcome this limitation by taking the rotational component of the deformation into account at each node. Barbić and James (2005) make use of



Figure 4.7: (**Left**) A jumping motion that leads to significant thigh and belly vibrations of human skin when rendered with modal analysis. (**Right**) The same hardware accelerated technique is applied to secondary tissues in a laparoscopic simulation to ease the CPU of computations.

modal analysis to build a deformation basis from modal derivatives. It allows them to simulate large deformations of solid objects with non-linear material.

4.3.3 Mass-spring model

A classic technique for modelling a deformable object is the mass-spring system. In this method, the geometry is described by a network of masses connected together by springs. The force F exerted onto each mass by the spring may be computed with:

$$F_{\text{spring}} = -kx, \quad (4.10)$$

where k is the spring stiffness constant, x represents the stretch of the spring from its rest position. From our experience, we know that a mass attached to a single spring has a tendency to oscillate. If no appropriate measure is taken, a network of springs will also oscillate as well. To prevent the system from oscillating, a damper is also added between each couple of masses. The force created by the damper acts to reduce the velocity and may be expressed as:

$$F_{\text{damper}} = -dv \quad (4.11)$$

where d is the damping constant and v the derivative of x with respect to time. The damper added between each couple of masses allows to model friction and reduce the amplitude of the oscillations. Applying Newton's second law $\sum \mathbf{F} = m\mathbf{a}$ yields the standard differential equation for a mass-spring system:

$$m\ddot{x} = -d\dot{x} - kx. \quad (4.12)$$

For complex systems with more than one spring attached to each mass, the right-hand side of the equation (4.12) becomes more complex as it needs to be repeated and adjusted for each spring. The easiest way to solve a complex mass-spring system iteratively, is to first calculate the total force currently acting on each mass by summing all forces currently exerted by all springs attached to this mass (using (4.10) and (4.11)). Dividing the total force by the mass gives the current acceleration of each mass (by (4.12)). This equation must be numerically solved for each mass, either by explicit or implicit methods. A discussion of integration methods for mass-spring models can be found in Shinya (2005).

Mass-spring systems have been extensively used and are still very popular. And this is not surprising because they are easy to implement and computationally efficient. Of course, they have drawbacks. But as we will see shortly, various recent works attempt to circumvent them. One limitation often mentioned is the difficulty to derive spring stiffnesses from elastic properties (Young's modulus and Poisson's ratio). To overcome this deficiency, Lloyd et al. (2007) mention two ways of obtaining the parameters of a mass-spring model. The first one consists in varying the parameters until the behaviour of the system is similar to the one obtained through the experiments or the one obtained with the Finite Element Method (Bianchi et al., 2004). The second way consists in establishing an analytical reasoning for calculating the constants of the mass-spring model. For instance, starting from the definition of Young's modulus, Poisson's ration, shear and bulk modulus, Baudet et al. (2007) apply simple tests to their mass-spring model to find the most appropriate stiffnesses. Lloyd et al. (2007) proposes a linearisation of the mass-spring equations with the aim of equating the linearised stiffness matrix to the stiffness matrix of a linear finite element method. By following the same method, San Vicente et al. (2009) derived a mass-spring model equivalent to a linear finite element model for maxillofacial surgery simulation.

Because they do not derive from the equations of continuum mechanics, mass-spring models have a limited capacity to model the various aspects of a material like anisotropy, viscoelasticity etc. Even worse, according to Bourguignon and Cani (2000), if all springs are set to the same stiffness, the mesh geometry may generate undesired anisotropy as shown in Fig. 4.8. However, some works tried to take advantage of this feature by giving directions of interest to their model by specifically designing the mesh to align springs on specific directions. For instance, it was used by Miller (1988) to animate the motion of snakes and worms. It was also used by Ng-Thow-Hing and Fiume (1997) in their muscle model where some of the springs were aligned with the muscle fibers and more recently by Magjarevic et al. (2009) for modelling the mechanical behavior of the myocardial tissue. More details on controlling anisotropy in mass-spring systems may be found in Bourguignon and Cani (2000).

A mass-spring model may also be enhanced to feature viscoelasticity. Indeed, Stiles and Alexander (1972) proposed a viscoelastic mass-spring system to test hypotheses on the mechanical response of muscles. Much more recently, Tamura et al.

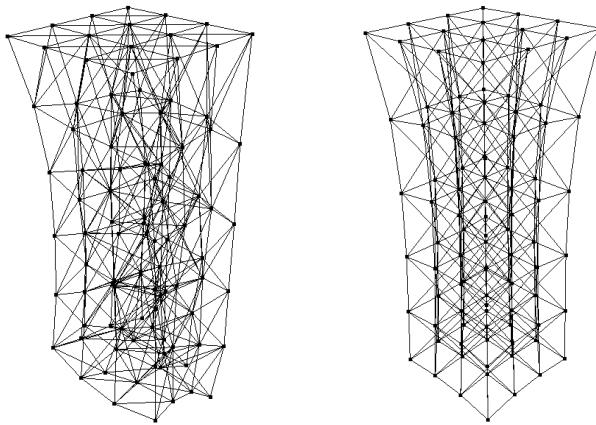


Figure 4.8: The two meshes are undergoing a downward pull at their bottom. While the tetrahedral mass-spring system (**left**) shows signs of anisotropy, the hexahedral mesh (**right**) with springs aligned along the gravity does not.

(2005) introduced a method to simulate viscoelastic material using a mass-spring system. They use a large number of particles to create a randomly connected mesh to mimic the structure of polymeric materials and hence their viscoelasticity characteristics. Basafa and Farahmand (2010) went even further and extended a mass-spring system to simulate non-linear viscoelastic deformations of soft tissue for laparoscopic surgery. They tune their parameters by a simple optimisation procedure to fit the mechanical response obtained on a set of experimental data.

Terzopoulos et al. (1991) even designed a mass-spring system that experience transitions from solid to liquid. Nodes are connected by springs in a hexahedral lattice and each of them has an associated temperature. Spring stiffnesses decreased with the increase in temperature. The diffusion of heat through the material is computed by using a discretised form of the heat equation. When the melting point is reached, the stiffness is set to zero and the node is detached from this particular spring. Once a node is freed from all his neighbouring springs, it becomes an independent glop of fluid and its behaviour is then modelled using a discrete fluid model.

If mass-spring systems are quite versatile and easy to implement, the most complete mathematical formalism available to describe the mechanical behaviour of a solid is provided by continuum mechanics. Consequently, it feels natural to derive computational models from the equations of continuum mechanics and a few techniques based on these equations will now be introduced.

4.4 Techniques based on continuum mechanics

4.4.1 The finite element method

The finite element method brings accuracy a step further. The method will not be explained here since it was already discussed with great details in chapter 3. This is the method of first choice for modelling deformable objects with precision as it is directly derived from the most complete theory available for describing the mechanical response of a continuum. Different kind of analysis may be carried out with the finite element method. As seen in section 3.3, an analysis may be either static (that is, the transient response is ignored) or dynamic (more general). In all analyses, the deformation of the body at a point is described by some measure of strain. If this measure is chosen to be linear (see section 2.4.5 for details), the analysis is said to be *geometrically* linear. In general, we consider this assumption to be valid for strain less than 10 %. In contrast, geometrically non-linear analyses make use of a non-linear strain tensor to measure the deformation and thus may handle large deformations. Another characteristic of an analysis is whether the relationship between strain and stress (the constitutive model) is assumed to be linear or non-linear. The analysis is qualified as *materially* linear or non-linear, respectively.

Because of its complexity, the finite element method is substantially more computationally expensive than the ones introduced so far. And this has been a problem in computer graphics and interactive simulation for many years. So researchers came up with simplifying assumptions and ideas to make the computations more efficient. Over time, with the increase in computational power, the simplifications were progressively relaxed and the observed tendency for real-time simulation is now towards more and more complete finite element methods. In this section, we will describe this evolution in the use of finite element methods for interactive simulation.

The use of the finite element method in computer graphics was pioneered by Terzopoulos et al. (1987) for simulations of elastic deformations. They derived a geometrically non-linear strain tensor from differential geometry theory. Obviously, the computations were not carried out in real-time but this technique allows Terzopoulos et al. to calculate the interaction between solid and deformable models in a physically realistic manner (see Fig. 4.9). In order to speed up the computations, linear strain measures started to be used (see section 2.4.5). For instance, Gourret et al. (1989) make use of a linear finite element method to model the deformation of grasped objects. The equations are however solved in static for simplicity. An approach to solve the system of equations in dynamic in a reasonable amount of time is to decrease the number of elements. This reduction must be compensated by the use of high-order elements, that is the degree of the shape functions used in the interpolation of the solution is higher to maintain the overall accuracy. This approach was applied to an elaborate model of a muscle by Chen and Zeltzer (1992) and to modelling skin in a simulator of plastic surgery by Pieper et al. (1992) where they both use 20-node brick elements. Yet, finite element modelling remained very

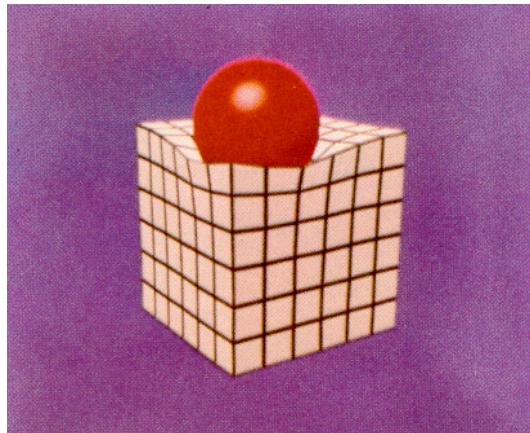


Figure 4.9: A ball resting on a supporting elastic solid as computed by Terzopoulos et al..

computationally expensive and the interactivity required for medical simulators was still to be reached.

In 1996, Bro-Nielsen and Cotin developed a real-time finite element method formulation by introducing three improvements. The first one is to compress the stiffness matrix of the volumetric system into a matrix with the same complexity as a surface model of the same object. This technique is called *condensation*. The deformation is solved only at surface points while still taking the volumetric nature of the object into account. The second improvement concerns the fact that they explicitly precompute the inverse of the stiffness matrix and use matrix vector multiplication with this matrix to achieve a low calculation time. Finally, they exploit the sparse structure of the force vector. They showed the efficiency of their method by modelling a leg with 700 nodes at 20 frames per second.

Later, James and Pai (1999) had another idea to improve the speed of the computation. They based their model on boundary integrals and the *boundary element method*, used for the first time in computer graphics. The method consists of calculating only boundary displacements and forces. Unfortunately, the authors do not provide the performance of their framework.

Cotin et al. (1999) took advantage of the superposition principle which applies to linear problems. Their idea was to approximate the deformation of a solid object with a linear combination of pre-computed states. In addition, they enhanced the formulation by adding a non-linear radial component. This method was applied to compute the deformation of liver in a hepatic surgery simulator. Their approach substantially improved the computational time for a finite element model as it allowed them to compute the positions and the forces of a 1 500 node liver at 300 Hz. An important drawback emphasised by the authors themselves is the impossibility to simulate tissue cutting. Indeed, the large precomputations associated with this method are dependent of the geometry of the mesh and therefore cannot be updated

in real-time.

Besides real-time computation of soft tissue, another major concern for medical simulators is the ability to cutting through soft tissue models. Mor and Kanade (2000) proposed a technique to generate a minimal set of new elements to replace intersected tetrahedra along the path formed by the cut. The method do not wait for the cut to be completed for splitting an element, instead, a minimal subdivision of a partially cut tetrahedron is generated. The authors used a very small time step to insure the stability of the soft tissue model and their simulations were not real-time. In addition, their method can create very small elements and the model can become unstable. This occurs because the stiffness matrices for the small elements are quite large, and even small displacements can generate very large forces.

Cotin et al. (2000) proposed a hybrid model allowing real-time cutting. They combined a quasi-static precomputed model, which is very efficient but does not allow topology changes, with a tensor-mass model which requires more computation but authorises cutting. The key idea is to separate a same structure into two parts: one susceptible to be cut during the surgery procedure and one that will not be submitted to any topological change. While the former will be simulated with the tensor-mass model, the deformation of the latter will be computed using the precomputed model. Since both linearly elastic models follow the same physical law, the combination of these two models should behave exactly as a global, linearly elastic model. To achieve this goal, Cotin et al. describe how to impose the additional boundary conditions at the connection nodes for the global model to be consistent in terms of forces and displacements. The efficiency of this hybrid model was demonstrated by simulating an hepatectomy (removing of one of the eight anatomical segments of the liver). They reported the simulation (allowing deformation and cutting) of a mesh with more than 8 000 tetrahedra in real-time (see Fig. 4.10).

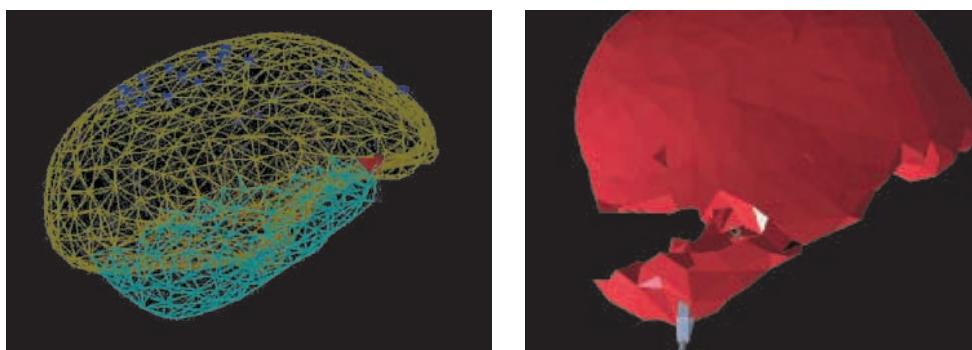


Figure 4.10: Simulation of a hepatectomy using a hybrid model. 18% of the mesh was modelled with the tensor-mass model (green) and therefore could be cut out. The rest of the mesh was simulated with the precomputed finite element model.

Serby et al. (2001) described a new approach to cutting with finite element

models. The central idea is not to introduce new nodes/elements but to displace the existing ones to account for the topological changes introduced by a cut. They also added a step of homogenisation of the mesh to avoid tiny elements. Thus, the problem of decreasing element size is minimised and consequently the stability of the solution of the equations of motion is increased.

Cutting is a problem for all techniques precomputing the inverse of the stiffness matrix. Indeed, topological changes require updates of this matrix, which is too computationally heavy to be done in real-time. Nienhuys and van der Stappen (2001) avoid the precomputations by approximating the inverse at runtime with an iterative algorithm (Conjugate Gradient). Moreover, instead of subdividing each tetrahedron to make the cut appear where the user performed it, they adapt the mesh locally so that there are always triangles on the scalpel path, and perform the cut along those triangles. If the approach is interesting, the authors reported a lag between the scalpel and the realised cut. In addition, repositioning the nodes within the mesh can easily generate degeneracies, which have to be eliminated.

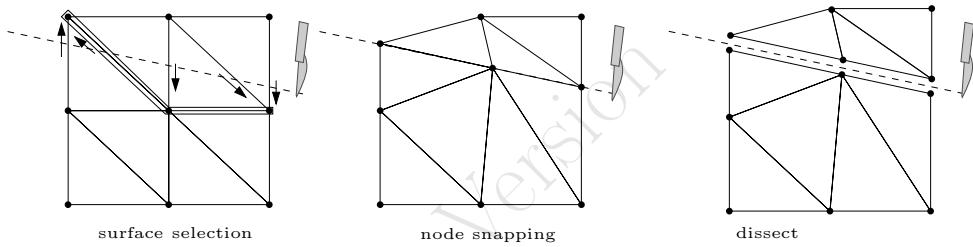


Figure 4.11: The three steps in performing a cut (represented here in 2D).

As we know, the assumption of linearity is usually only valid for very small deformations and strains. If this is often an acceptable trade-off between the visual deformation result and the speed of the system, research was carried out to handle large deformations. However, since geometrically non-linear analyses are highly computationally demanding, another approach was found which allows large deformation while still relying on a linear strain measure: the *co-rotational* methods. A very complete description of the co-rotational approach is given by Felippa (2000). As explained in section 2.4, an appropriate strain measure is a tensor which only measures the deformation applied to the object. This is the case of the Cauchy-Green or Green tensors for instance. Those strain metrics only measures deformation, they are invariant with respect to rigid-body transformations applied. However, geometrically linear analyses make use of the infinitesimal strain tensor, which is only an approximation of the Green tensor. Consequently, this linearised strain measure is not invariant with respect to rigid-body transformations. In particular, the error brought by the approximated (linear) strain measure increases with the rotation of the element. This error in the measure creates ghost forces which overdeforms the simulated object. The idea behind co-rotational methods is the decomposition

of the motion into a rigid-body and deformational components. In other words, the deformation of each element may be seen as measured in a coordinate system which follows the rigid-body motion of the element, hence cancelling the part of the error due to the rotation in the strain measure. Consequently, co-rotational formulations allow large displacements and rotations (since the rigid-body motion is not considered) as long as the strain remains small (the strain measure is still linear).

Among the firsts to apply this technique to static analyses are Argyris et al. (1964) and Wempner (1969). Although the birth of co-rotational methods is not very clear (Felippa and Haugen, 2005), Belytschko and Hsieh (1973) had a similar thinking and introduced the so-called *convected coordinates* for solving dynamic equations. However, while the spirit is the same than co-rotational methods, convected coordinates forms a curvilinear system that fits the change of metric as the body deforms, therefore the convected metric necessarily encompasses deformations. In other words, the decomposition in rigid-body and deformational components is not exact. During the last decade, the co-rotational approach has regained interest in the computer graphics community.

Hauth and Strasser (2004) made use of the polar decomposition of the deformation gradient to extract the rotation. They also applied this formulation in a hierarchical finite element setting where the coarse level is computed with a non-linear strain tensor and the finer levels are computed with a co-rotational method. In addition, they examined the error that is arising from using the co-rotational approximation instead of a non-linear strain tensor. They measured that the co-rotational approximation induces an error of about 5 – 20 % but can give a speedup of about a factor seven or more.

Pouliquen et al. (2005) addressed the simulation of human fingers with haptic feedback. During the grasping task, the deformable finger pads undergo large rotations while their deformation remains small. Using a co-rotational approach for modelling the finger pads allows them to efficiently compute the contact forces between the deformable pads and the rigid object to be grasped.

If co-rotational approaches increase accuracy in the strain measure, all the examples we have seen so far make use of a linear constitutive law. Yet, the mechanical behaviour of soft tissue is known to be non-linear, viscoelastic and often anisotropic (Fung, 1993). Moreover, Misra et al. (2007) demonstrated a noticeable difference between the force feedbacks provided by linear and non-linear tissue models. One of the first materially non-linear finite element formulation used in medical simulation was probably the work of Sagar et al. (1994). While developing a teleoperated micro-surgical robot for eye surgery, they wanted to provide a virtual environment has part of the system for training. Since the cornea is anisotropic and non-linearly elastic and undergoes large deformations, they have used large deformation elasticity theory with orthotropic and non-linear material properties in the finite element model. Coupled with a detailed visual model of the eye, the simulation produces a fair level of realism with a performance greater than 10 Hz (see Fig. 4.12).

Keeve et al. (1998) compared a fully (geometrically and materially) non-linear

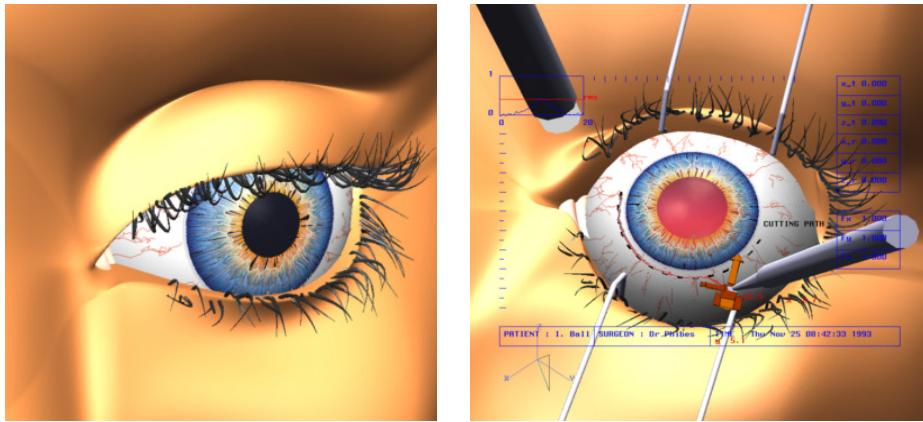


Figure 4.12: **(Left)** Exterior view of the model showing the eye and the eyelashes. **(Right)** Surgical virtual environment including micro-tools and guidance information.

finite element model with a mass-spring model to compute the deformation of skin after bone realignment in a craniofacial surgery simulator. The efficiency of the mass-spring model gives the ability to the surgeon to realise interactive simulation of the resulting tissue changes and to improve his planning process. The best precision was given by the finite element model which could be used off-line to verify the chosen surgical procedure.

The first geometrically non-linear finite element model running in real-time was reported by DeBunne et al. (2001). They provided an adaptative method for animating viscoelastic deformable objects in a guaranteed frame rate. As the object moves and deforms, the sampling is refined to concentrate the computational load into the regions that deform the most (see Fig. 4.13). They reported animating a few hundred points in real-time at a frequency higher than 300 Hz.

In the same spirit, Wu et al. (2001) applied an adaptative meshing technique to provide sufficient detail where required while minimising unnecessary computation. They included both Mooney-Rivlin and Neo-Hookean material models and used the non-linear Cauchy-Green tensor. In addition to the local mesh refinement technique, they used mass lumping and an explicit time integration to deform a liver model of 1 210 tetrahedra in real-time.

Still in 2001, Picinbono et al. proposed to use a non-linear model only where the displacements are larger than a given threshold, the remaining part of the object still uses linear elasticity.

Recently, research on non-linear finite element and more elaborate constitutive models increased significantly. For instance, Pras Pathmanathan and Rajagopal (2004) used a fully non-linear model to predict deformation of breast under compression during mammography. However, they were not interested in real-time. Zhong et al. (2005) claimed achieving real-time with their non-linear finite element method using an interpolation approach but do not give any figure. Sedef et al.

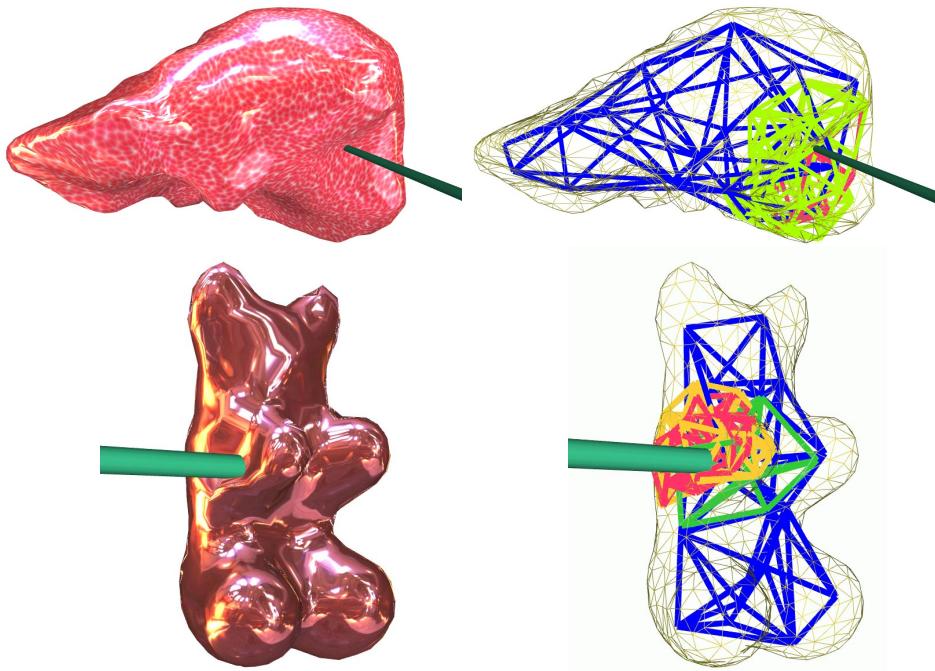


Figure 4.13: The approach of Debunne et al. (2001) makes use of a local refinement technique to ensure high physical fidelity while bounding the global computation load to guarantee real-time simulation with a geometrically non-linear and visco-elastic finite element model.

(2006) proposed a viscoelastic model to simulate the deformation of a liver. They calculate the deformation in real-time by carrying out the computations in static and using the superposition principle. However, their approach used a linear strain measure, which limits their system to small deformations. Yan et al. (2007) claimed to compute their non-linear finite element model in real-time using a technique of *graded mesh* but do not give any information on their performance.

Finally, Miller et al. (2007) proposed the total Lagrangian explicit dynamics (TLED) algorithm, a very efficient fully non-linear formulation. The algorithm is based on the finite element method using the total Lagrangian formulation, where stresses and strains are measured with respect to the original configuration. This choice allows for precomputing of most spatial derivatives before the commencement of the time-stepping procedure. The authors reported that the average number of floating-point operations per element per time step is 35 % lower than for the similar implementation of the algorithm based on updated Lagrangian formulation. They were able to compute the deformation of a hexahedral mesh of 6 000 elements in about 16 ms for a time step of 1 ms. If they could not achieve real-time computations, their work constituted a step towards the simulation of entire organs in real-time.

The finite element method is very computationally demanding. However, by

nature, very similar computations are carried out for each element. Consequently, the finite element method is a good candidate for being implemented on parallel architectures. Thus, Székely et al. (2000) implemented a fully non-linear formulation for an 8-processor machine. However, the dedicated hardware was not yet assembled and functional at the time so they did not provide any measure of performance.

Early 2000s, the rapid increase in the performance of graphics hardware, coupled with recent improvements in its programmability, have made graphics hardware a compelling platform for computationally demanding tasks in a wide variety of application domains. Researchers and developers have become interested in harnessing this power for general-purpose computing, an effort known collectively as GPGPU (for General-Purpose computing on Graphics Processing Units). A survey of GPGPU algorithms and techniques was written by Owens et al. (2007).

To the best of our knowledge, Rumpf and Strzodka (2001) were the firsts to leverage the power of GPUs to solve partial differential equations. They established a correspondence between common mathematical operations and OpenGL operations and structures. Thus, as an example, a vector is represented as a RGBA color vector in a color image on the graphics card. They decided to implement an iterative solver for a linear system of equations, which we find at the core of most finite element codes. Using this hardware accelerated solver, they solved the linear heat equations and presented an anisotropic diffusion model for image processing. They obtained promising results but noted the hardware restrictions that forced them to approximate all non-linear functions by linear ones in the implementation of the anisotropic diffusion. Although a complete hardware accelerated finite element method implementation was yet to be done, Rumpf and Strzodka definitely stepped towards the right direction.

Wu and Heng (2004) proposed the first application of GPU in medical simulation . They were able to achieve an interactive frame rate with topology changes in surgical simulation. Not only they used the condensation technique for merely calculating the surface nodes in the non-operation parts, but they implemented the conjugate gradient solver onto GPU. Indeed, entering 2003, the full IEEE floating point precision data type supported in shaders and texture stages allows the GPU to be applied into more general applications. The conjugate gradient solver includes three primary operations: multiplication of the sparse matrix and the vector, addition of two vectors and the sum-reduction operation. The core component is the multiplication of the sparse matrix and the vector. Therefore, they migrated this part of calculation from the CPU into the fragment processor of the GPU, to take advantage of the fragment processor in its efficient manipulation of the local texture memory on the mathematical calculation. The performance of their method on the GPU is up to 2.15 times faster than the CPU implementation.

In 2007, Taylor et al. introduced the first GPU implementation of a fully non-linear finite element method (more details in Taylor et al. (2008)). They reformulated the TLED algorithm presented by Miller et al. (2007) to accomodate with the limited hardware possibilities of GPUs. In particular, the main restriction when

using graphics-based GPU execution is the inability to scatter, which necessitated reformulation of force summation as a gather in their implementation. It is worth noting that the entire finite element method implementation was ported to GPU. They reported significant solution speed gains up to $16.8\times$ over the CPU implementation. Thus, they were able to compute a simple cube model with up to 16 000 tetrahedral elements in real-time.

4.4.2 Meshless methods

If the finite element method is a robust method and widely used in engineering, the method has a few shortcomings (Liu and Gu, 2005). First, the creation of a mesh is compulsory and producing a good quality mesh is both difficult and time consuming. Second, the stresses obtained in FEM are often discontinuous at the interfaces of the elements. Special techniques are required in a post-processing stage to recover accurate stresses. In addition, under large deformations, considerable loss in accuracy in FEM results can arise from the element distortions. The origin of these problems is the use of elements. Consequently, the idea of getting rid of the elements and the mesh has emerged and the concept of *meshless* methods was born. Instead, meshless methods use a set of nodes scattered within the problem domain. They do not form a mesh, meaning it does not require any a priori information on the relationship between the nodes for the interpolation of the unknown variables. In theory, this should facilitate topological changes, a feature fairly desirable in surgery simulation. It is worth noting though that meshless methods are still in their developing stage and are not as mature as finite element formulations. We will now briefly review the research carried out in soft tissue modelling for medical simulation that makes use of meshless method. A complete review and details on meshless methods are beyond the scope of this thesis.

A meshless method was introduced for the first time for modelling soft tissue in medical simulation by De et al. (2001). They use the so-called method of finite spheres (MFS). Nodal points are only sprinkled locally around the surgical tool tip and not over the whole domain. The interpolation is performed by functions that are nonzero only on spheres surrounding the nodes. A point collocation technique is used to generate the discrete equations. The authors qualitatively compared the method of finite spheres with the finite element technique from a commercial package and the displacement profile was similar in both cases. They were able to achieve a computational rate of about 100 Hz when 34 spheres where used. Real-time haptic rates of about 1 kHz was then obtained using a force extrapolation technique.

Lim and De (2004) used the same technique of finite spheres introduced by De et al. to model soft tissue but they added progressive cutting, without the generation of new primitives. As the cut progresses, the nearest vertex to the intersection point of the tool is snapped to the tool path. Then, they observed that the prescribed boundary condition changes on only a very small portion as the tool interacts with the organ. Rather than solving the entire problem over and over again, they

proposed instead to make incremental corrections, which results in an accelerated solution procedure.

Meshless methods were then introduced to computer graphics by Müller et al. (2004). In each step, they compute the spatial derivatives of the displacement field using a moving least squares (MLS) procedure. From these derivatives, they obtain strains, stresses and elastic forces at each simulated point. In addition, they proposed techniques for modelling and animating a point-sampled surface that dynamically adapts to deformations of the underlying volumetric model. They used a non-linear measure of strain with a linear constitutive model. They reported real-time simulation of models exhibiting elastic, plastic, melting and flowing effects.

Horton et al. (2006) tested meshless methods for surgical simulation. They proposed a total Lagrangian explicit dynamics algorithm using a non-linear material formulation. They applied the method on an example of brain shift and compared their results with a commercial finite element code. The shape functions are created from a cloud of unconnected nodes using the element free Galerkin method. They were interested in the displacement of the centre of mass for the brain and found the absolute displacement differences between the two algorithm are less than 0.85 mm which is the resolution of the MRIs in their study. The same team kept investigating and testing the technique. However, they used moving least squares procedure for creating the shape functions instead of the element free Galerkin method. They also realised experimental indentations (Horton et al., 2007) and provided more extensive comparisons with a finite element method (Horton et al., 2010).

Another example of application is the use of the technique called moving particule semi-implicit method (MPS) by Chhatkuli et al. (2009). They applied this method to simulate the deformation of lung and a tumor inside the left lung during inspiration. The lung tissues were considered to be homogeneous, isotropic and viscoelastic. They compared their results with the experimental CT taken at the end of inspiration and found out that the deformation predicted by numerical simulation matches reasonably well with the experimental results.

TODO: add drawbacks of meshless methods: supposedly good for handling topological changes but other issues. Couldn't really find drawbacks in meshless' literature.

CHAPTER 5

THE TOTAL LAGRANGIAN EXPLICIT DYNAMICS (TLED) ALGORITHM

*TLED proposed by Miller almost real-time GPU implementation to reach real-time
But first: TLED description
(MICCAI 2008; MedIA 2009)*

Draft Version

5.1 Description of the TLED algorithm

5.1.1 Key ideas

The total Lagrangian explicit dynamics (TLED) algorithm was first introduced to the field of medical simulation by Miller et al. (2007). More details on this finite element method formulation may also be found in Bathe (1995). It is worth noting that the finite element analysis carried out by the TLED algorithm is dynamic and fully non-linear (geometrically and materially). Let us begin by exposing the key ideas of the approach.

The deformation of a body may be described by two different kind of description: Lagrangian or Eulerian (see section 2.2). Because the Lagrangian description focuses its attention on the particles of the continuous body, it is usually used in solid mechanics. In the Lagrangian description, the position and physical properties of the particles are referred to a reference configuration. For the great majority of commercial finite element programs, the reference configuration is the previous configuration, that is the one at the end of the previous time step. In this case, where all variables are referred to the previous configuration, the formulation is called *updated Lagrangian*. The advantage of this approach is the simplicity of the incremental strain description. The disadvantage is that all derivatives with respect to spatial coordinates must be recomputed in each time step, because the reference configuration is changing. The reason for this choice is historical, at the time of solver development the memory was expensive and caused more problems than actual speed of computations (Miller et al., 2007).

The first key idea of the TLED algorithm is to refer all variables to the undeformed configuration. This type of formulation is said to be *total Lagrangian*. Because of this, the choice of the second Piola-Kirchhoff stress tensor as a measure of stress is required. The strain measure that is work-conjugate with the second Piola-Kirchhoff stress is the Green-St.Venant strain tensor. That is, the work of stress increments on the strain increments gives an accurate expression for work (Ji et al., 2010). We should note that all derivatives in the definition of the Green tensor are with respect to the initial and undeformed configuration. The decisive advantage is that all derivatives are calculated with respect to the undeformed configuration and therefore can be precomputed. This is at the cost to a more complex strain-displacement matrix due an initial displacement effect in the incremental strain (Bathe, 1995). However, the TLED algorithm performs significantly fewer mathematical operations in each time step.

The second crucial idea is to use the central difference method for an explicit time integration of dynamic equilibrium equations. The advantage is that the stiffness term $\mathbf{K}(\mathbf{U}) \cdot \mathbf{U}$ of the system of equations may be computed from:

$$\mathbf{K}(\mathbf{U}) \cdot \mathbf{U} = \mathbf{F}(\mathbf{U}) = \sum_e \tilde{\mathbf{F}}^e, \quad (5.1)$$

where $\tilde{\mathbf{F}}^e$ are the global nodal force contributions due to stresses in element Ω_e . This implies that the stiffness matrix does not need to be assembled since element

nodal force contributions can be computed at the element level instead, which is a decisive computational advantage.

We will first explain how to compute element nodal forces based on the deformation gradient and a neo-Hookean constitutive model. We will then show how to use these element nodal forces to calculate the unknown displacements for each node of the model (Taylor, 2006).

5.1.2 Computation of element nodal forces

We adopt the notation of Bathe with respect to indication of the relevant configuration of the body: a left superscript indicates the configuration in which a quantity occurs and, when applicable, a left subscript indicates the configuration with respect to which the quantity is measured. We note the coordinates of a point at time t ${}^t\mathbf{x}$.

The deformation gradient

As we know, a fundamental measure of deformation is the deformation gradient tensor ${}_0^t\mathbf{X}$ and we recall that it may be written as:

$${}_0^t\mathbf{X} = \frac{\partial {}^t\mathbf{x}}{\partial {}^0\mathbf{x}}. \quad (5.2)$$

This tensor describes the stretches and rotations that the material fibres have undergone from time 0 to time t . In order to compute the element deformation gradients, we first compute derivatives of displacements $u_{i,j}$ with respect to global coordinates. Since we use a total Lagrangian framework, the derivatives are referred to the undeformed configuration and we have:

$${}_0^t u_{i,j} = \frac{\partial {}^t u_i}{\partial {}^0 x_j} = \sum_{a=1}^N \frac{\partial h_a}{\partial {}^0 x_j} {}^t u_{ai}, \quad (5.3)$$

where h_a is the shape function associated with node a and N the number of nodes of the element. If we define a matrix ${}_0^t \partial \mathbf{u}_x$ of displacement derivatives

$${}_0^t \partial \mathbf{u}_x = \begin{bmatrix} {}_0^t u_{1,1} & {}_0^t u_{1,2} & {}_0^t u_{1,3} \\ {}_0^t u_{2,1} & {}_0^t u_{2,2} & {}_0^t u_{2,3} \\ {}_0^t u_{3,1} & {}_0^t u_{3,2} & {}_0^t u_{3,3} \end{bmatrix} \quad (5.4)$$

The deformation gradient may be obtained from the displacement derivatives with:

$${}^t\mathbf{X} = {}_0^t \partial \mathbf{u}_x + \mathbf{I}. \quad (5.5)$$

The Green-Lagrange strain tensor

From the deformation gradient we may obtain the right Cauchy-Green deformation tensor ${}^t_0\mathbf{C}$:

$${}^t_0\mathbf{C} = {}^t_0\mathbf{X}^T {}^t_0\mathbf{X}, \quad (5.6)$$

and then yields the Green-Lagrange strain tensor ${}^t_0\mathbf{E}$:

$${}^t_0\mathbf{E} = \frac{1}{2}({}^t_0\mathbf{C} - \mathbf{I}), \quad (5.7)$$

where \mathbf{I} is the rank 2 identity tensor.

Constitutive equations and evaluation of stress

We know that the stresses result from the deformation of the material and they may be expressed in terms of some measure of this deformation such as the strain. The constitutive equations, which depends on the material under consideration, relate the stresses to the strain. As we have seen in section 2.6.1, hyperelastic materials are a general class of materials in which the constitutive relationship is expressed in the form of a strain energy density function t_0W . In such materials, the stresses are obtained by differentiating this energy density function with respect to the appropriate strain measure (that is energically conjugate). In a total Lagrangian framework, the appropriate stress and strain measures are second Piola-Kirchhoff stress ${}^t_0\mathbf{S}$ and Green-Lagrange strain ${}^t_0\mathbf{E}$. Consequently, the stress may be computed from:

$${}^t_0\mathbf{S} = \frac{\partial {}^t_0\mathbf{W}}{\partial {}^t_0\mathbf{E}}. \quad (5.8)$$

The form of the strain energy function depends on the chosen material. For the TLED algorithm we choose a neo-Hookean material already described section 2.6.5. For a compressible neo-Hookean material, the strain energy density is given by:

$${}^t_0W = \frac{1}{2}\mu({}^t_0I_1 - 3 - 2\ln {}^t J) + \frac{1}{2}\lambda({}^t J - 1)^2, \quad (5.9)$$

where t_0I_1 is the first invariant of the right Cauchy-Green deformation tensor ${}^t_0\mathbf{C}$ given by ${}^t_0I_1 = \text{tr}({}^t_0\mathbf{C})$ and ${}^t J = \det({}^t_0\mathbf{X}) = \det({}^t_0\mathbf{C})^2$ is the Jacobian. Then, using the chain rule we have:

$$\frac{\partial {}^t_0W}{\partial {}^t_0\mathbf{E}} = \frac{\partial {}^t_0W}{\partial {}^t_0\mathbf{C}} \frac{\partial {}^t_0\mathbf{C}}{\partial {}^t_0\mathbf{E}}. \quad (5.10)$$

Noting that

$$\frac{\partial {}^t_0\mathbf{C}}{\partial {}^t_0\mathbf{E}} = 2 \quad \text{by (5.7)}, \quad (5.11)$$

Therefore, we may evaluate the stress from:

$${}^t_0\mathbf{S} = 2 \frac{\partial {}^t_0W}{\partial {}^t_0\mathbf{C}}, \quad (5.12)$$

which we can then express as the following:

$${}^t_0\mathbf{S} = 2 \left(\frac{\partial {}^t_0W}{\partial {}^t_0\mathbf{I}_1} \frac{\partial {}^t_0\mathbf{I}_1}{\partial {}^t_0\mathbf{E}} + \frac{\partial {}^t_0W}{\partial {}^t J} \frac{\partial {}^t J}{\partial {}^t_0\mathbf{E}} \right). \quad (5.13)$$

It may be shown that this yields:

$${}^t_0S_{ij} = \mu(\delta_{ij} - {}^t_0C_{ij}^{-1}) + \lambda {}^tJ({}^tJ - 1){}^t_0C_{ij}^{-1}. \quad (5.14)$$

Element nodal forces

For a given element at time t , the global nodal force contributions $\tilde{\mathbf{F}}$ due to stresses in the element may be computed from:

$${}^t\tilde{\mathbf{F}} = \int_{^0V} {}^t_0\mathbf{B}_L^t {}^t_0\hat{\mathbf{S}} d^0V, \quad (5.15)$$

where ${}^t_0\hat{\mathbf{S}}$ is the vector form of the second Piola-Kirchhoff stress given by:

$${}^t_0\hat{\mathbf{S}} = [{}^t_0S_{11} \quad {}^t_0S_{22} \quad {}^t_0S_{33} \quad {}^t_0S_{12} \quad {}^t_0S_{23} \quad {}^t_0S_{13}]^T, \quad (5.16)$$

${}^t_0\mathbf{B}_L^T$ is the strain-displacement matrix and d^0V is the initial (undeformed) volume of the element. The strain-displacement matrix \mathbf{B}_L relates the strains in an element to the element's nodal displacements. If an implicit analysis is performed, the strain-displacement matrix is used in the assembly of the stiffness matrix. In explicit analyses it is used to compute element nodal force contributions according to (5.15). In geometrically linear analyses, we assume that the displacements are infinitesimally small so that the geometry of an element does not change over time and the strain-displacement matrix is constant. Conversely, because geometrically non-linear analyses (such as carried out by the TLED algorithm) take the change of geometry of the elements into account, the strain-displacement matrix $\mathbf{B}_L = {}^t_0\mathbf{B}_L$ varies. However, the strain-displacement matrix at time t may be computed by transforming a linear matrix ${}^t_0\mathbf{B}_{L0}$ using the deformation gradient. We begin by defining the linear matrix as:

$${}^t_0\mathbf{B}_{L0} = \begin{bmatrix} {}^t_0\mathbf{B}_{L0}^{(1)} & {}^t_0\mathbf{B}_{L0}^{(2)} & \dots & {}^t_0\mathbf{B}_{L0}^{(N)} \end{bmatrix}, \quad (5.17)$$

where N is the number of nodes per element and the submatrix ${}^t_0\mathbf{B}_{L0}^{(a)}$ are given by:

$${}^t_0\mathbf{B}_{L0}^{(a)} = \begin{bmatrix} {}^t_0h_{a,1} & 0 & 0 \\ 0 & {}^t_0h_{a,2} & 0 \\ 0 & 0 & {}^t_0h_{a,3} \\ {}^t_0h_{a,2} & {}^t_0h_{a,1} & 0 \\ 0 & {}^t_0h_{a,3} & {}^t_0h_{a,2} \\ {}^t_0h_{a,3} & 0 & {}^t_0h_{a,1} \end{bmatrix} \quad a = 1, 2, \dots, N. \quad (5.18)$$

The subscripted comma denotes partial differentiation such as:

$${}^t_0h_{a,i} = \frac{\partial h_a}{\partial {}^0x_i}. \quad (5.19)$$

The full strain-displacement matrix (accounting for initial displacement effect) is then computed via

$${}^t_0\mathbf{B}_L^{(a)} = {}^t_0\mathbf{B}_{L0}^{(a)} {}^t_0\mathbf{X}^T. \quad (5.20)$$

We note that the linear strain-displacement matrix is composed of derivatives of shape functions with respect to the original undeformed coordinates of the body. These derivatives are therefore constant and may be precomputed, affording the total Lagrangian formulation a significant computational advantage.

5.1.3 Computation of node displacements

We recall that the dynamic system of equations describing a deformable solid is the following (see section 3.26):

$$\mathbf{M}\ddot{\mathbf{U}} + \mathbf{D}\dot{\mathbf{U}} + \mathbf{K}(\mathbf{U}) \cdot \mathbf{U} = \mathbf{R}. \quad (5.21)$$

where \mathbf{M} is a constant mass matrix, \mathbf{D} is a constant damping matrix, $\mathbf{K}^e(\mathbf{U}^e)$ is the stiffness matrix, which is a function of nodal displacements \mathbf{U} , and \mathbf{R} are externally applied loads. And we seek the unknown displacements \mathbf{U} .

Mass and damping matrices

The mass matrix for the system may be computed by summing contributions from individual elements:

$$\mathbf{M} = \sum_e \mathbf{M}^e, \quad (5.22)$$

where \mathbf{M}^e is the mass contribution from element Ω_e given by:

$$\mathbf{M}^e = \int_{^0V_e} \rho_e \mathbf{H}^T \mathbf{H} d^0V. \quad (5.23)$$

Following the total Lagrangian framework, we observe that the expression is integrated over the undeformed volume and the mass density ρ_e of the undeformed configuration is also used. We then diagonalise the mass matrix by applying the technique of mass lumping as explained in section 3.3.5. Therefore, the mass matrix may be built by directly computing the mass of each element and assigning an equal proportion of this to each of the element's nodes. In other words, if the mass of element Ω_e is m_e then for every node a which is attached to this element, the a^{th} diagonal component M_{aa} of the global mass matrix will receive a contribution of m_e/N where N is the number of nodes per element. The complete matrix is compiled by summing contributions from all elements.

We also employ a Rayleigh damping from which we only consider the mass-proportional component to obtain a diagonal damping matrix. Consequently, \mathbf{D} is computed from:

$$\mathbf{D} = \alpha \mathbf{M}. \quad (5.24)$$

Gathering of nodal forces

As stated early in this chapter, the stiffness matrix of (5.21) may be obtained from

$$\mathbf{K}(\mathbf{U}) \cdot \mathbf{U} = \mathbf{F}(\mathbf{U}) = \sum_e \tilde{\mathbf{F}}^e. \quad (5.25)$$

The previous section detailed how to calculate element nodal forces $\tilde{\mathbf{F}}^e$. For each node we can now add the element force contributions up from all elements attached to this node. Gathering those nodal force contributions allows us to compute the stiffness term of the equilibrium equations ${}^t\mathbf{F}$. Therefore, (5.21) may be re-written as the following:

$$\mathbf{M}^t \ddot{\mathbf{U}} + \mathbf{D}^t \dot{\mathbf{U}} + {}^t\mathbf{F} = {}^t\mathbf{R}. \quad (5.26)$$

Explicit time integration

The remaining step is the time integration of (5.26) to find out the expression of ${}^t\mathbf{U}$. For this, we employ the explicit central difference method described page 39. We assume that all displacements, velocities and accelerations for the current time step is known and we seek a formula for computing displacements at time $t + \Delta t$. Let us remind (3.35) for convenience:

$$\left(\frac{\mathbf{M}}{(\Delta t)^2} + \frac{\mathbf{D}}{2\Delta t} \right) \mathbf{U}^{t+\Delta t} = \mathbf{R}^t - \mathbf{F}^t + \frac{2\mathbf{M}}{(\Delta t)^2} \mathbf{U}^t + \left(\frac{\mathbf{D}}{2\Delta t} - \frac{\mathbf{M}}{(\Delta t)^2} \right) \mathbf{U}^{t-\Delta t}. \quad (5.27)$$

By using diagonalised mass and damping matrices, this expression may be rearranged to give a formula for updating displacements component-wise:

$${}^{t+\Delta t} \mathbf{U}_i = \frac{{}^t R_i - {}^t F_i + \frac{2M_{ii}}{\Delta t^2} {}^t U_i + \left(\frac{D_{ii}}{2\Delta t} - \frac{M_{ii}}{\Delta t^2} \right) {}^{t-\Delta t} U_i}{\frac{D_{ii}}{2\Delta t} + \frac{M_{ii}}{\Delta t^2}} \quad (5.28)$$

Let us define the following vectors:

$$A_i = \frac{1}{\frac{D_{ii}}{2\Delta t} + \frac{M_{ii}}{\Delta t^2}} \quad (5.29)$$

$$B_i = \frac{\frac{2M_{ii}}{\Delta t^2}}{\frac{D_{ii}}{2\Delta t} + \frac{M_{ii}}{\Delta t^2}} = \frac{2M_{ii}}{\Delta t^2} A_i \quad (5.30)$$

$$C_i = \frac{\frac{D_{ii}}{2\Delta t} - \frac{M_{ii}}{\Delta t^2}}{\frac{D_{ii}}{2\Delta t} + \frac{M_{ii}}{\Delta t^2}} = \frac{D_{ii}}{2\Delta t} A_i - \frac{B_i}{2}, \quad (5.31)$$

Equation (5.28) may be rewritten as:

$${}^{t+\Delta t} \mathbf{U}_i = A_i ({}^t R_i - {}^t F_i) + B_i {}^t U_i + C_i {}^{t-\Delta t} U_i. \quad (5.32)$$

It is worth noting that the coefficient vectors \mathbf{A} , \mathbf{B} and \mathbf{C} can be precomputed.

Of crucial importance in explicit analyses is the restriction on time step sizes imposed by stability constraints . As mentioned, the central difference method

is only stable if $\Delta t < \Delta t_{cr}$. Essentially, stability requires that the time step be not larger than the time required for a dilatational wave to propagate through the smallest element. In linear analyses this translates into the following formula for Δt_{cr} :

$$\Delta t_{cr} = \frac{L_e}{c}, \quad (5.33)$$

where L_e is the smallest characteristic element length (roughly interpreted as the smallest edge length) in the assembly, and c is the dilatational wave speed of the material, given by

$$c = \sqrt{\frac{E(1-\nu)}{\rho(1+\nu)(1-2\nu)}}, \quad (5.34)$$

where E is Young's modulus and ν is Poisson's ratio. If nonlinear analyses (either geometric or material) are conducted, a somewhat smaller time step must be employed since c will change with deformation, and in particular will generally increase. Additionally, the presence of damping necessitates a further decrease in the limit. Nonetheless, (5.34) may be used to estimate Δt_{cr} .

Some important observations arise from (5.33) and (5.34). Firstly, since the stiffness of soft tissues is very much smaller than that of common engineering materials ($E \approx 3 \times 10^3$ Pa for brain tissue versus $\approx 2 \times 10^{11}$ Pa for steel), the allowable time step is much larger for analysis of the former. This is a key reason for the expediency of explicit time integration for analysis of soft tissues. Secondly, soft tissues are generally considered to be incompressible, leading to $\nu \approx 0.49$ commonly being employed. However, if this value can be relaxed even further, significant increases in Δt_{cr} are possible. As an example, for a given L_e a material with $E = 3\,000$ Pa and $\nu = 0.45$ allows a time step of more than double that of a material with $\nu = 0.49$. Of course, lowering the Poisson ratio introduces inaccuracies, but for some applications (interactive medical simulation for instance) it may be that this is acceptable. Finally, it must be highlighted that Δt_{cr} is proportional to L_e . If the element size is decreased (by using a finer mesh, or by analysing very small objects), inexorably Δt_{cr} is also reduced.

5.2 Anisotropic and viscoelastic constitutive equations

5.2.1 An extension to the TLED algorithm

An attractive feature of explicit analyses is the relative ease with which arbitrarily complicated constitutive models may be incorporated. This arises from the fact that element stresses are computed directly from strains in the course of the procedure. Additionally, there is no requirement for computation of tangent matrices as in implicit or quasi-static procedures, since there is no involvement of Newtonian iterations. While the developments above incorporated a hyperelastic constitutive formulation (neo-Hookean), thus accommodating nonlinearity of the strain-related tissue stress response, the formulation was intentionally simple nonetheless since

the focus was on validation of the computational framework. Two further key features of the response of most biological tissues are time- (and rate-)dependence and anisotropy (Fung, 1993).

Time-dependence manifests itself in many aspects of the mechanical response. Soft tissues under constant load generally exhibit creep, while those under constant deformation exhibit stress relaxation. Additionally, most tissues appear stiffer at higher loading velocities. In particular, so-called visco-hyperelastic models based on strain energy functions with time-dependent parameters have been shown to reproduce both the time-dependent and large strain aspects of the response (Miller and Chinzei, 1997; Miller, 2000; Nava et al., 2008).

Anisotropic mechanical response may arise, for example, from the presence of a highly organised microstructure such as those of connective tissues. These are predominantly composed of collagen or elastin fibres embedded in an amorphous matrix (Fung, 1993), and may be considered as fibre reinforced composites in some cases. Alternatively, the presence of vasculature and other functional components means even non-load bearing organs may exhibit directional dependence (Picinbono et al., 2003; Prange and Margulies, 2002). Whereas isotropic constitutive models may be formulated in terms of the usual principal strain invariants (as used in (5.9) for a neo-Hookean material), directional dependence requires inclusion of so-called pseudo-invariants of strain and material direction.

In order to model such soft-tissue features with the explicit TLED formulation without significant performance penalties, an efficient constitutive update procedure involving time integration of the relevant hereditary equations is required. We address this problem by presenting a procedure similar to that developed by Poon and Ahmad (1998) for analysis of anisotropic linear viscoelastic models. This was also adapted by Taylor et al. (2007b) for solution of their fibre composite-based microstructural model. In these cases the rate independent responses were based on linear elasticity, rather than a hyperelastic formulation suitable for large deformations. In the present work we begin with a class of anisotropic visco-hyperelastic models, and develop a constitutive update procedure for explicit analyses based on these.

5.2.2 Visco-hyperelasticity

Large recoverable deformations and the time- and rate-dependence of the mechanical response of soft-tissues have led to the formulation of visco-hyperelastic constitutive models in which an underlying hyperelastic formulation is augmented by time dependent (viscoelastic) material parameters. Models of this type are well known in the continuum mechanics community, and provide a kinematically consistent basis for modelling non-linear materials at large deformations.

For such materials the constitutive response is defined in terms of a time-dependent Helmholtz free energy (strain energy) function $\hat{\Psi}$, expressed in the form

of a convolution integral:

$$\hat{\Psi}(\Psi, t) = \int_0^t \alpha(t-s) \frac{\partial \Psi}{\partial s} ds, \quad (5.35)$$

where t is time and Ψ is the underlying hyperelastic strain energy function. The relaxation functions $\alpha(t)$ commonly assume the form of a Prony series:

$$\alpha(t) = \alpha_\infty + \sum_{i=1}^N \alpha_i e^{-t/\tau_i}, \quad (5.36)$$

where α_∞ , α_i and τ_i are positive real constants. Such forms for the relaxation functions have a physical interpretation, namely that of a generalised Maxwell model (Holzapfel, 1996). If we impose the condition

$$\left(\alpha_\infty + \sum_{i=1}^N \alpha_i \right) = 1 \quad (5.37)$$

we may rewrite(5.36) as

$$\alpha(t) = 1 - \sum_{i=1}^N \alpha_i (1 - e^{-t/\tau_i}), \quad (5.38)$$

which will be of use in subsequent sections.

The required stress \mathbf{S} may be obtained via differentiation with respect to strain:

$$\mathbf{S} = 2 \frac{\partial \hat{\Psi}(\Psi, t)}{\partial \mathbf{C}} = \int_0^t \alpha(t-s) \left(2 \frac{\partial}{\partial s} \frac{\partial \Psi}{\partial \mathbf{C}} \right) ds = \int_0^t \alpha(t-s) \frac{\partial \Phi}{\partial s} ds, \quad (5.39)$$

where \mathbf{C} is the right Cauchy-Green deformation tensor, and we have introduced $\Phi \stackrel{\text{def}}{=} 2\partial\Psi/\partial\mathbf{C}$ as the instantaneous hyperelastic stress response.

Models of this form have been presented by Miller and co-workers for analysis of brain tissue (Miller and Chinzei, 1997, 2002) and of liver and kidney (Miller, 2000). They were shown to model the tissue responses to compressive (and in the case of brain, tensile) loading at strain rates varying over two orders of magnitude very well. We next consider evaluation of the hyperelastic stress response Φ .

5.2.3 Hyperelastic response

We firstly consider the case of isotropic materials, from which anisotropic formulations follow. We then consider the standard cases of transverse isotropy and orthotropy, which may be viewed as arising from, for example, the presence of mutually orthogonal reinforcing fibre phases. In considering these cases we encompass constitutive equations which have been proposed for a wide variety of biological tissues.

Isotropic materials

For isotropic materials the strain energy is a function of strain only, hence $\Psi = \Psi(\mathbf{C})$. In this work we consider strain energy functions with separated isochoric (volume preserving) and volumetric components (Holzapfel et al., 2000):

$$\Psi(\mathbf{C}) = \Psi^{iso}(\bar{\mathbf{C}}) + \Psi^{vol}(J) = \Psi^{iso}(\bar{I}_1, \bar{I}_2) + \Psi^{vol}(J), \quad (5.40)$$

where J is the Jacobian determinant, $\bar{\mathbf{C}} = J^{2/3}\mathbf{C}$ is the modified right Cauchy-Green deformation, and $\bar{I}_1 = \text{tr}\bar{\mathbf{C}}$ and $\bar{I}_2 = [(tr\bar{\mathbf{C}})^2 - tr(\bar{\mathbf{C}}^2)]/2$ are invariants of $\bar{\mathbf{C}}$. The hyperelastic stress Φ then also consists of isochoric and volumetric components:

$$\Phi = \Phi^{iso} + \Phi^{vol}, \quad (5.41)$$

where

$$\Phi^{vol} = 2 \frac{\partial \Psi^{vol}(J)}{\partial \mathbf{C}} = J \frac{d\Psi^{vol}(J)}{dJ} \mathbf{C}^{-1}, \quad (5.42)$$

and

$$\Phi^{iso} = 2 \frac{\partial \Psi^{iso}(\bar{\mathbf{C}})}{\partial \bar{\mathbf{C}}} = J^{-2/3} \text{Dev}\bar{\Phi}, \quad (5.43)$$

where $\text{Dev}(\bullet) = (\bullet) - (1/3)[(\bullet) : \mathbf{C}] \mathbf{C}^{-1}$ is the referential configuration deviatoric operator for a second order tensor (Holzapfel et al., 2000), and

$$\bar{\Phi} = 2 \frac{\partial \Psi^{iso}(\bar{\mathbf{C}})}{\partial \bar{\mathbf{C}}} = \bar{\gamma}_1 \mathbf{I} + \bar{\gamma}_2 \bar{\mathbf{C}}, \quad (5.44)$$

with

$$\bar{\gamma}_1 = 2 \left(\frac{\partial \Psi^{iso}}{\partial \bar{I}_1} + \bar{I}_1 \frac{\partial \Psi^{iso}}{\partial \bar{I}_2} \right) \quad \text{and} \quad \bar{\gamma}_2 = -2 \frac{\partial \Psi^{iso}}{\partial \bar{I}_2}. \quad (5.45)$$

Transversely isotropic materials

Transversely isotropic materials are characterised by a single preferred direction \mathbf{a}_0 in the reference configuration; the mechanical response is isotropic in the plane orthogonal to this direction. The strain energy is then a function of both \mathbf{C} and a structure tensor $\mathbf{A}_0 \stackrel{\text{def}}{=} \mathbf{a}_0 \otimes \mathbf{a}_0$, where \otimes denotes a tensor product. Analogous to the isotropic case, we consider strain energy functions of the form

$$\Psi(\mathbf{C}, \mathbf{A}_0) = \Psi^{iso}(\bar{I}_1, \bar{I}_2, \bar{I}_4, \bar{I}_5) + \Psi^{vol}(J), \quad (5.46)$$

where $\bar{I}_4 = \mathbf{a}_0 \cdot \bar{\mathbf{C}} \mathbf{a}_0$ and $\bar{I}_5 = \mathbf{a}_0 \cdot \bar{\mathbf{C}}^2 \mathbf{a}_0$ are pseudo-invariants of $\bar{\mathbf{C}}$ and \mathbf{A}_0 .

As can be seen transversely isotropic strain energy functions (5.46) differ from isotropic ones (5.40) only in the form of the isochoric term. Therefore the volumetric and isochoric stresses Φ^{vol} and Φ^{iso} remain as in (5.42) and (5.43) but with $\bar{\Phi}$ given by

$$\bar{\Phi} = \bar{\gamma}_1 \mathbf{I} + \bar{\gamma}_2 \bar{\mathbf{C}} + \bar{\gamma}_4 \mathbf{A}_0 + \bar{\gamma}_5 (\mathbf{a}_0 \otimes \bar{\mathbf{C}} \mathbf{a}_0 + \bar{\mathbf{C}} \mathbf{a}_0 \otimes \mathbf{a}_0), \quad (5.47)$$

with

$$\bar{\gamma}_a = 2 \frac{\partial \Psi^{iso}}{\partial \bar{I}_a}, \quad a = 4, 5. \quad (5.48)$$

Orthotropic materials

Orthotropic materials are characterised by three mutually orthogonal preferred directions, which we identify with unit vectors \mathbf{a}_0 and \mathbf{b}_0 (and corresponding structure tensors \mathbf{A}_0 and \mathbf{B}_0), in the reference configuration. We need specify only two vectors, since the direction orthogonal to these naturally emerges as a preferred direction also. Orthotropic strain energy functions are then of the form

$$\Psi(\mathbf{C}, \mathbf{A}_0, \mathbf{B}_0) = \Psi^{iso}(\bar{I}_1, \bar{I}_2, \bar{I}_4, \bar{I}_5, \bar{I}_6, \bar{I}_7) + \Psi^{vol}(J), \quad (5.49)$$

where $\bar{I}_6 = \mathbf{b}_0 \cdot \bar{\mathbf{C}}\mathbf{b}_0$ and $\bar{I}_7 = \mathbf{b}_0 \cdot \bar{\mathbf{C}}^2\mathbf{b}_0$ are pseudo-invariants of $\bar{\mathbf{C}}$ and \mathbf{B}_0 . In a similar manner to the transversely isotropic case we obtain Φ^{vol} and Φ^{iso} from (5.42) and (5.43), with $\bar{\Phi}$ now given by

$$\begin{aligned} \bar{\Phi} = & \bar{\gamma}_1 \mathbf{I} + \bar{\gamma}_2 \bar{\mathbf{C}} + \bar{\gamma}_4 \mathbf{A}_0 + \bar{\gamma}_5 (\mathbf{a}_0 \otimes \bar{\mathbf{C}}\mathbf{a}_0 + \bar{\mathbf{C}}\mathbf{a}_0 \otimes \mathbf{a}_0) \\ & + \bar{\gamma}_6 \mathbf{B}_0 + \bar{\gamma}_7 (\mathbf{b}_0 \otimes \bar{\mathbf{C}}\mathbf{b}_0 + \bar{\mathbf{C}}\mathbf{b}_0 \otimes \mathbf{b}_0) \end{aligned} \quad (5.50)$$

where

$$\bar{\gamma}_a = 2 \frac{\partial \Psi^{iso}}{\partial \bar{I}_a}, \quad a = 6, 7. \quad (5.51)$$

5.2.4 Recapitulation

For a visco-hyperelastic material stress may be obtained from (5.39). This form is general in the sense that any underlying hyperelastic response may be used, including the anisotropic formulations described (5.40), supplemented with (5.42), (5.43), and (5.44) is the general form of an isotropic hyperelastic stress response, defined in terms of invariants. Transversely isotropic or orthotropic models may be produced by substituting (5.44) for (5.47) or (5.50), respectively. The specification of particular forms of Ψ^{vol} and Ψ^{iso} would be motivated by the particular tissue/material under analysis, and may stem from phenomenological or microstructural considerations. Finally we note that for separated isochoric and volumetric hyperelastic functions as used here, viscoelastic terms may be applied to either or both independently.

5.3 Constitutive update procedure for explicit analyses

Use of the above visco-hyperelastic models within the TLED algorithm (or any other explicit dynamic finite element procedure) requires a constitutive update scheme involving time integration of (5.39).

5.3.1 Stress update equations

We proceed by restating (5.39) and using (5.38):

$$\mathbf{S} = \int_0^t \left[1 - \sum_{i=1}^N \alpha_i (1 - e^{(s-t)/\tau_i}) \right] \frac{\partial \Phi}{\partial s} ds. \quad (5.52)$$

This may be separated into rate-dependent and -independent terms as

$$\mathbf{S} = \boldsymbol{\Phi} - \sum_{i=1}^N \boldsymbol{\Upsilon}_i, \quad (5.53)$$

where

$$\boldsymbol{\Upsilon}_i = \int_0^t \alpha_i \left(1 - e^{(s-t)/\tau_i}\right) \frac{\partial \boldsymbol{\Phi}}{\partial s} ds, \quad i \in [1, N] \quad (5.54)$$

are rate-dependent terms associated with each term in the Prony series.

In an incremental analysis we require the stress at the current increment given the deformation state and history of the material. Adding superscripts to indicate time increments the stress may be updated using

$$\mathbf{S}^n = \boldsymbol{\Phi}^n - \sum_{i=1}^N \boldsymbol{\Upsilon}_i^n. \quad (5.55)$$

The instantaneous terms $\boldsymbol{\Phi}^n$ may be computed directly from the (known) current deformation \mathbf{C}^n . The main difficulty is then computation of the incremental rate-dependent terms $\boldsymbol{\Upsilon}_i^n$. Following Poon and Ahmad (1998), our strategy is to maintain each $\boldsymbol{\Upsilon}_i^n$ as a separate state variable to be updated at each increment also.

5.3.2 State variable update equations

Our approach is to convert the integral equation (5.54) into a rate form which may then be numerically integrated to produce an incremental update formula for $\boldsymbol{\Upsilon}_i^n$. We note that (5.54) is of the form

$$y(t) = \int_0^t f(t, s) ds. \quad (5.56)$$

Poon and Ahmad (1998) provide the following formula for differentiating such equations with respect to t :

$$\dot{y} = f(t, t) + \int_0^t \dot{f}(t, s) ds. \quad (5.57)$$

Applying this formula to (5.54) we obtain

$$\begin{aligned} f(t, t) &= \alpha_i \left(1 - e^{(t-t)/\tau_i}\right) \frac{\partial \boldsymbol{\Phi}}{\partial s} \\ &= 0 \end{aligned} \quad (5.58)$$

and

$$\begin{aligned} \dot{f}(t, s) &= \frac{d}{dt} \left[\alpha_i \left(1 - e^{(s-t)/\tau_i}\right) \frac{\partial \boldsymbol{\Phi}}{\partial s} \right] \\ &= \alpha_i e^{(s-t)/\tau_i} \frac{\partial \boldsymbol{\Phi}}{\partial s} \frac{1}{\tau_i}, \end{aligned} \quad (5.59)$$

leading to the required rate form for Υ_i :

$$\begin{aligned}\dot{\Upsilon}_i &= \frac{1}{\tau_i} \int_0^t \alpha_i e^{(s-t)/\tau_i} \frac{\partial \Phi}{\partial s} ds \\ &= \frac{1}{\tau_i} (\alpha_i \Phi - \Upsilon_i).\end{aligned}\quad (5.60)$$

Equation (5.60) may be integrated using a convenient numerical method. In particular the unconditionally stable backward Euler method suggested by Poon and Ahmad (1998) yields the following formula for Υ_i^n :

$$\begin{aligned}\Upsilon_i^n &= \left(\frac{\alpha_i \Phi^n}{\tau_i} + \frac{\Upsilon_i^{n-1}}{\Delta t} \right) / \left(\frac{1}{\Delta t} + \frac{1}{\tau_i} \right) \\ &= A \Phi^n + B \Upsilon_i^{n-1},\end{aligned}\quad (5.61)$$

where Δt is the time step size, and $A = \Delta t \alpha_i / (\Delta t + \tau_i)$ and $B = \tau_i / (\Delta t + \tau_i)$ are constant coefficients.

5.3.3 Summary

The constitutive update procedure consists of

1. Updating state variables (one for each Prony term) via (5.61).
2. Updating stresses via (5.55).

Assuming a constant time step size Δt , the coefficients A and B may be precomputed.

CHAPTER **6**

GPU IMPLEMENTATION OF TLED

(ISBMS 2008a)

Draft Version

6.1 Summary of the TLED formulation

A complete description of the TLED algorithm for soft tissue simulation was given in the previous chapter. However, let us remind the main steps of the algorithm before tackling the implementation. Briefly, the algorithm consists of a precomputation phase in which element shape function derivatives $\partial\mathbf{h}$ (and other quantities) and the system mass matrix \mathbf{M} are calculated, followed by a time-loop in which incremental solutions for the node displacements \mathbf{U} are found. During each step of the time-loop we:

1. Apply loads (displacements and/or forces) and boundary conditions to relevant nodal degrees of freedom
2. For each element compute
 - (a) deformation gradient \mathbf{X} and right Cauchy-Green deformation tensor \mathbf{C}
 - (b) linear strain-displacement matrix \mathbf{B}_L
 - (c) second Piola-Kirchhoff stress \mathbf{S}
 - (d) element nodal forces $\tilde{\mathbf{F}}$, and add these forces to the total nodal forces \mathbf{F}
3. For each node compute new displacements \mathbf{U} using the central difference method.

The nodal force contributions $\tilde{\mathbf{F}}$ from each element are obtained

$$\tilde{\mathbf{F}} = \int_{^0V} \mathbf{B}_L^t \hat{\mathbf{S}} d^0V, \quad (6.1)$$

where 0V is the initial volume of the element and $\hat{\mathbf{S}}$ is the vector form of the stress tensor \mathbf{S} . This integral is generally evaluated numerically, for example using Gaussian quadrature. For reduced integration 8-node hexahedral elements we obtain

$$\tilde{\mathbf{F}} = 8 \mathbf{B}_L^t \hat{\mathbf{S}} J \quad (6.2)$$

where J is the Jacobian determinant. For 4-node tetrahedral elements we obtain

$$\tilde{\mathbf{F}} = V \mathbf{B}_L^t \hat{\mathbf{S}}. \quad (6.3)$$

The above equations make no assumption concerning the constitutive model employed. The deformation state \mathbf{F} in each element is known, allowing stresses \mathbf{S} to be computed from any valid constitutive equation.

6.2 General-purpose computation on graphics processing units

6.2.1 Goal and motivation

Graphics processing units (GPU) functionality has, traditionally, been very limited. In fact, for many years the GPU was only used to accelerate certain parts of the graphics pipeline. A GPU is essentially a special purpose hardware designed to accelerate each stage of the geometry pipeline, the process of matching image data or a computer model to the pixels on the computer's screen. Initially, GPU could only run two kinds of program: vertex and pixel shaders. Vertex shaders are run once for each vertex given to the graphics processor. The purpose is to transform each vertex's 3D position in virtual space to the 2D coordinate at which it appears on the screen. Vertex shaders can manipulate properties such as position, color, texture coordinate and normal vector. Pixel shaders are functions that compute color and other attributes of each pixel. They range from always outputting the same color, to applying a lighting value, to doing shadows, specular highlights or translucency for instance.

Eventually, vertex and pixel shaders became programmable to enable game programmers to generate even more realistic effects. Programmable pixel shaders allow the programmer to substitute, for example, a lighting model other than those provided by default by the graphics card. Shaders have enabled graphics programmers to create lens effects, displacement mapping and depth of field. This evolution of GPU's hardware and the increasing programmable capability naturally lead to use GPUs for non-graphics applications. The term GPGPU (General-purpose computation on graphics processing units) was coined by Mark Harris in 2002 when he recognised an early trend of using GPUs for non-graphics applications. However, capabilities of GPU's were still fairly limited for non-graphics applications at this time.

Things dramatically changed in 2007 when the two types of shaders were unified. While early shader models used very different instruction sets for vertex and pixel shaders, unified shader models have almost the same capabilities. An unified shader architecture allows more flexible use of the graphics rendering hardware. The computing units of the GPU can run vertex or pixel shaders according to work loads. Along with unified shader models, a new type of shader was created: geometry shaders. They are executed after vertex shaders and can generate new graphics primitives, such as points, lines and triangles.

GPUs may be seen as high-performance many-core processors that can be used to accelerate a wide range of applications. An example (given by Sanford Russell from NVIDIA) to think about to illustrate the difference between a traditional CPU and a GPU is this: if you were looking for a word in a book, and handed the task to a CPU, it would start at page 1 and read it all the way to the end, because it's a serial processor. It would be fast, but would take time because it has to go in order.

A GPU, which is a parallel processor, would tear the book into a thousand pieces and read it all at the same time. Even if each individual word is read more slowly, the book may be read in its entirety quicker, because words are read simultaneously.

In addition to an increasing programmability, it is easy to understand why the development of GPGPU is soaring. GPU is now used in imaging, finance, signal processing, simulation, audio and video processing, astronomy, weather forecasting, molecular modelling, cryptography, quantum mechanics and many other fields. A fairly recent review of GPGPU algorithms may be found in Owens et al. (2007). This is therefore not surprising that research has been carried out to leverage the power of GPUs to solve partial differential equations of continuum mechanics via the computationally expensive finite element method. Indeed, we have already insisted on the strong time constraints demanded by the field of medical simulation and the path towards GPGPU is natural. A non-exhaustive review of algorithms implemented on GPU in the context of medical simulation was given chapter 4 (page 67).

Before introducing our own GPU implementation of the TLED formulation, we will discuss the different GPU programming languages at our disposal.

6.2.2 Programming languages for GPUs

Interface with a GPU has traditionally been via a graphic application programming interface (API) such as OpenGL or DirectX. OpenGL is an open standard API that provides a number of functions for the rendering of 2D and 3D graphics and is available on most modern operating systems including but not limited to Windows, Mac OS X and Linux. It was initially designed by Silicon Graphics Inc. and now managed by the technology consortium Khronos group. DirectX is a proprietary API that provides functions to render three dimensional graphics, and uses hardware acceleration if it is available on the graphics card. It was designed by Microsoft for use on the Windows platform.

Because OpenGL and DirectX are both low-level APIs and include many irrelevant functionality from a GPGPU perspective, higher level languages were designed to give developers more direct control of the graphics pipeline without having to use assembly language or hardware-specific languages.

Graphics API

OpenGL (version 1.5 and newer) provides shader language based on the C programming language called OpenGL shading language (GLSL). In the DirectX API (DirectX 9 and newer), shaders are programmed with the high level shader language (HLSL). Cg is another high-level shading language developed by NVIDIA in close collaboration with Microsoft and it is very similar to Microsoft's HLSL. However, the Cg compiler is not specific to a graphics API and can output both OpenGL and DirectX shader programs.

These languages all share basic data types (float, integer, bool etc.) and also feature vector and matrix data types that are based on the basic data types (like

float3 or float4x4 in Cg and vec3 or mat4 in GLSL) and standard library functions are provided (vector and matrix multiplication, cross product etc.). In addition, they support loops and branching like if/else, while, and for. However, because these languages are just an abstraction on the top of graphics API, there are significant restrictions on what they can do. Specifically, shader programs are unable to determine the destination of their outputs. They merely operate on fragment data (color, texture coordinates, normal etc.) and return the modified fragment at the same location. Each shader is also executed independently of its neighbours since no communication is allowed between processors. Note that while this was made possible on recent hardware, these graphics API do not take advantage of this feature and therefore this statement remains valid.

Memory storage is allowed through the use of textures, which essentially are structured collections of floating point values arranged as arrays. Shaders can read textures from any location and as many times as they like. As such it is a concept very similar to regular arrays. However, the individual components of textures consist of up to four floating point values in the form of a vector. This comes from the fact that graphics processing units were initially designed to manipulate color images, so each component of a texture is in fact a vector holding the red, green, blue and alpha (transparency) values for the given pixel. Textures are therefore optimised for retrieving up to four floating point values in a single read.

When it comes to writing to GPU memory though, things are a bit more limited. Indeed, in contrast to texture reads, shaders can only write once into textures (at the end of the set of instructions). Moreover, as stated earlier, shaders can only write in predefined locations corresponding to their potential screen location. Therefore, graphics API does not allow writing to random memory locations (this feature is called *scattering*). Nevertheless, shaders can have multiple render targets. The write locations within each texture are still fixed though.

Non-graphics API

Eventually, non-graphics API were released to support the development of GPGPU. The two major graphics card manufacturers first released their own API. While the objective of hiding graphics aspects is the same, the approaches of NVIDIA and ATI (now owned by AMD) are quite different.

In November 2006, NVIDIA released CUDA which consists of C language extensions. CUDA is fairly high level. Programs (called kernels) may be coded as C functions and invoked directly without going through a render pass as with graphics APIs. When called, kernels are executed N times in parallel by N threads as opposed to only once like regular C functions. Memory management functions analogous to standard C are also provided (memcpy, free etc.).

At the same time, ATI released CTM released a software development layer aimed for low-level access (assembly-style language). Memory management functions similar to graphics API are also provided. A year later, in December 2007, ATI released Stream SDK and added a new high-level language derived from C,

called ATI Brook+ which later evolved to ATI CAL (Compute Abstraction Layer). In both cases, these APIs work solely on each manufacturer's devices, CUDA with NVIDIA cards and Stream SDK on ATI cards.

OpenCL is a more general framework for writing programs that execute across heterogeneous platforms consisting of CPUs, GPUs, and other processors. It is not specific to a single type of hardware. OpenCL was initially developed by Apple and refined into an initial proposal in collaboration AMD (ATI), IBM, Intel, and NVIDIA. Apple submitted this initial proposal to the Khronos Group and the first release occurred in December 2008. AMD decided to support OpenCL instead of the now deprecated CTM in its Stream SDK.

The latest API to be released was DirectCompute from Microsoft. DirectCompute was released as part of DirectX11 in October 2009 and allows access to general-purpose computing with a fairly low level language.

The choice of API

The first GPU implementation of the TLED algorithm was carried out by Taylor et al. (2007a, 2008). They used the Cg language and had to reformulate the TLED algorithm presented to accomodate with the limited possibilities of this graphics API. Indeed, the major difference between this parallel implementation and a serial one is the writing of element nodal forces to memory at the end of a first kernel and their subsequent retrieval and summation during a second kernel. In a serial implementation nodal force contributions would most likely be added to a global sum as they are computed, rather than stored and summed in a second loop. But as stated earlier, this graphics API does not allow scattering, that is writing to random memory locations. Because of this limitation, element nodal forces cannot be directly added to the total nodal forces. Instead, the authors had to reformulate force summation as a gather operation in their implementation. They reported significant solution speed gains up to 16.8× over the CPU implementation. Thus, they were able to compute a simple cube model with up to 16 000 tetrahedral elements in real-time.

A major improvement was introduced with the release of CUDA. In theory, CUDA allows scattered writes and this new feature could dramatically change the GPU implementation of the TLED as the reformulation of the scatter as a gather would no longer be useful. A single kernel would be sufficient and an increase in performance is expected. Along with a more readable and simple code, this is the main feature that made us choose the CUDA API to reimplement the TLED algorithm.

In addition, we decided to develop our CUDA-based re-implementation of the TLED algorithm within the international and open source framework SOFA. We will show that this integration has a very limited cost in terms of performance by comparing the SOFA version with a standalone implementation. By providing an efficient and accurate non-linear FEM for soft tissue modelling to worldwide researchers, we thus hope to assist in enhancing the realism of medical simulators.

6.3 Implementation into SOFA

6.3.1 SOFA, an open source simulation framework

Objectives

The multi-disciplinary aspect of medical simulation requires the integration within a single environment of solutions in areas as diverse as visualisation, biomechanical modelling, haptic feedback and contact modelling. Although their interaction is essential to design a realistic simulator, only few teams have the sufficient resources to build such frameworks. This diversity of problems creates challenges for researchers to advance specific areas, and leads rather often to duplication of effort. The open source SOFA framework (Allard et al., 2007) was created to overcome this issue by providing researchers with an advanced software architecture that facilitates the development of new algorithms and simulators. SOFA has been mostly developed by INRIA (the French national institute for research in computer science and control) and CIMIT (Center for Integration of Medicine and Innovative Technology) and is primarily targeted at real-time simulation with an emphasis on medical simulation. SOFA is highly modular and flexible: it allows independently developed algorithms to interact together within a common simulation while minimising the development time required for integration. The overall goal is to develop a flexible framework while minimising the impact of this flexibility on the computation overhead. To achieve these objectives, SOFA proposes a new architecture that implements a series of concepts described below.

SOFA architecture

High-level modularity. The SOFA architecture relies on the innovative notion of multi-model representation where an object is explicitly decomposed into various representations: Behaviour Model, Collision Model, Collision Response Model, Visual Model and Haptic Model. Each representation can then be optimised for a particular task (biomechanics, collision detection, visualisation, haptics) while at the same time improving interoperability by creating a clear separation between the functional aspects of the simulation components. These representations are then connected together via a mechanism called *mapping*. Various mapping functions can be defined, and each mapping associates a set of primitives of a representation to a set of primitives in the other representation (Fig. 6.1). For instance, a mapping can connect degrees of freedom in a Behaviour Model to vertices in a Visual Model.

Finer level modularity. In order to easily compare algorithms within SOFA, more flexibility was added to the Behaviour Model by introducing an even finer level of granularity. A series of generic primitives common to most physics-based simulations have been defined: DoF, Mass, Force Field and Solver. The DoF component describes the degrees of freedom, and their derivatives, of the object. The Mass component represents its mass. The Force Field describes both internal and

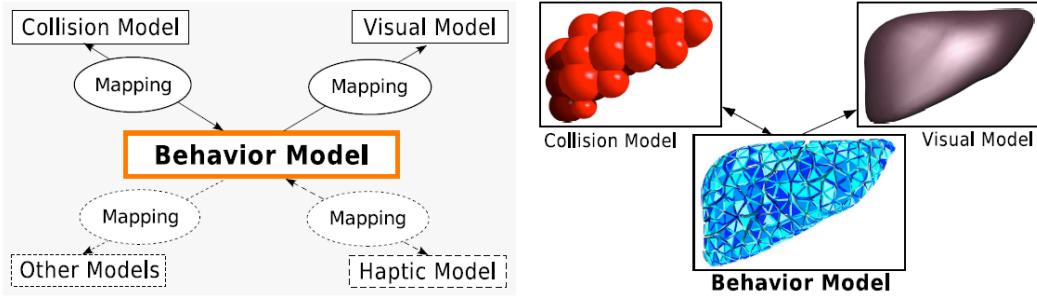


Figure 6.1: Multi-model representation in SOFA. **Left:** a Behaviour Model controls the other representations via a series of mappings. **Right:** examples of representations with a liver model.

external forces that can be applied to this object. The Solver component handles the time step integration, i.e. advancing the state of the system from time t to time $t + \Delta t$.

Scene-graph. Finally, another key aspect of SOFA is the use of a scene-graph to organise and process the elements of a simulation. Each component is attached to a node of a tree structure. This simple structure makes it easy to visit all or a subset of the components in a scene, and dependencies between components are handled by retrieving sibling components attached to the same node. During the simulation loop, most computations can be expressed as a traversal of the scene-graph. For instance, at each time step, the simulation state is updated by processing all Solver components, which will then forward requests to the appropriate components by recursively sending actions within its sub-tree.

These different functionalities and levels of abstraction allow the user to switch from one component to another by simply editing an XML file, without having to recompile. In particular this permits testing of different computational models of soft tissue deformation, and to assess the pros and cons of various algorithms within the same context.

6.3.2 CUDA description

GPUs achieve a high floating point capacity by distributing computation across a high number of parallel execution threads. They perform optimally as single instruction, multiple data devices. CUDA is a relatively new C API for compatible NVIDIA GPUs. CUDA organises threads in two hierarchical levels: blocks, which are groups of threads executed on one of the GPU's multiprocessors, and grids, which are groups of blocks launched concurrently on the device, and which all execute the same kernel. Figure 6.2 represents this thread organisation.

CUDA allows developers to specify the number of threads per block in each execution (the so-called execution configuration), effectively defining the distribution of computational load across all processors. For a given kernel the block dimensions

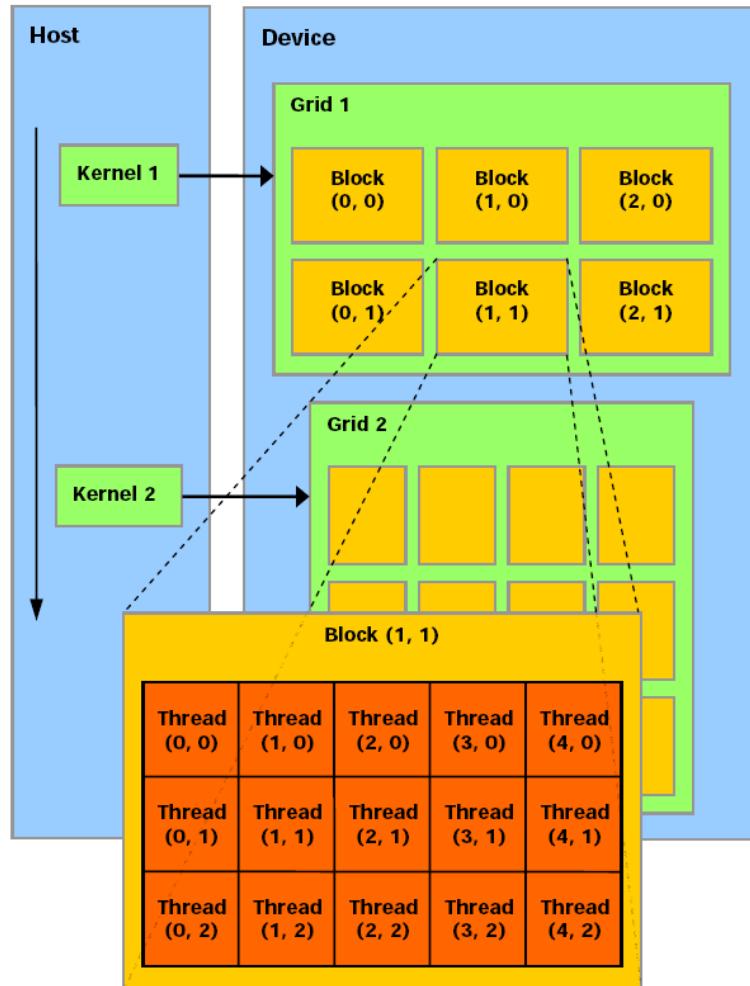


Figure 6.2: Each kernel is executed by CUDA as a group of threads within a grid.
Image courtesy of NVIDIA.

are chosen to optimise the utilisation of the available computational resources. Care should be taken at the multiprocessor level in balancing the available memory required by the kernels with the ability to hide global memory latency. Since a finite amount of memory is available on a multiprocessor, the memory requirements of a kernel will determine how many threads can run concurrently on each. Importantly, CUDA's use of time slicing allows more than one block to be executed concurrently on a single multiprocessor, which has important implications for hiding memory latency. If more than one block is executing, the multiprocessor is able to switch processing between blocks while others are stalled on memory accesses, whereas it has no option but to wait for these if only one block is executing. Therefore for memory bandwidth bound kernels it may be preferable to launch several smaller blocks on each multiprocessor rather than a single larger one if both configurations make the same use of multiprocessor memory resources. While tools are available from NVIDIA for estimating the optimal execution configuration, it has proved necessary to fine tune the configuration experimentally for each kernel.

6.3.3 First implementation: scatter as a gather

SOFA integration.

Implementing a biomechanical model in SOFA translates essentially into writing a new Force Field, i.e. describing the algorithm used to compute internal forces in the model. It merely comes down to creating a single C++ class and changing the position reads and force writes to integrate the algorithm into SOFA's design. The precomputation phase takes place in the initialisation method where relevant variables are computed and passed to an external C function that allocates memory on the GPU and binds textures to it. During the simulation loop, the Solver requests the computation of the forces by launching the appropriate kernels on the GPU.

Kernel organisation.

Although CUDA allows scattered writes in theory, it offers no write conflict management between threads. Why is this a problem? The major steps of the TLED algorithm are the following: (1) we compute the element nodal forces, that is for each node of a given element, we calculate the force contribution that this element has on this node; (2) we obtain the global force for each node by summing up the force contributions of all elements to which this node is attached; (3) we compute the displacement for each node by integration using the central difference method. In a serial implementation (CPU) we would do a loop on all elements and for each element we compute all element nodal force contributions. We would then directly add the element force contributions computed for each node into a global node array at the corresponding location for each node of the current processed element. In a parallel implementation, many elements are processed concurrently, in parallel. Inevitably, two elements sharing a node may be processed in the same time and may therefore try to both write their computed nodal force contribution into the exact

same location in GPU's memory. At the time of this first implementation (2007), the result of concurrent writes was undefined. In fact, only one write was guaranteed to happen and there was no way to know which one succeeded. Potential measures to address this have proved extremely detrimental to performance. For this reason, the kernel arrangement in our first CUDA implementation eventually ends up to be essentially the same as the one used by Taylor et al. (2007a).

Consequently, the TLED CUDA implementation also relies on 2 kernels. The first kernel operates over elements in the model and computes the element stresses based on the current model configuration. It then converts these into nodal force contributions, which are written to global memory. The second kernel operates over nodes and reads the previously calculated element force contributions and sums them for each node. The SOFA central difference solver computes the displacements from the nodal forces. Therefore, due to the impracticability of scattered writes, the sum operation is reformulated as a gather and the second kernel is needed to sum the nodal forces. The use of the developed viscoelastic constitutive update scheme necessitates storage of an additional array of state variables Υ_i .

Memory usage.

One efficient method for reading global memory data within kernels is texture fetching. Textures may be bound to either *cudaArrays* or regions of linear memory. CudaArrays have been designed to achieve optimal fetching when the access pattern has a high level of 2D locality. In the present application, the access pattern among threads is essentially random (since unstructured meshes are used) and our experiments have shown that texture fetching from linear memory is in fact fastest. Therefore all global memory precomputed variables were accessed using this method.

In SOFA, the forces are stored on the GPU in global memory. Since this memory space is not cached, it is important to follow the appropriate access pattern to obtain maximum memory bandwidth, especially given how costly accesses to device memory are. A multiprocessor takes 4 clock cycles to issue one memory instruction for a set of threads. When accessing global memory, there are, in addition, 400 to 600 clock cycles of memory latency. A suboptimal access pattern would yield incoherent writes. The memory bandwidth would then be an order of magnitude lower. In order to prevent this, a key feature of CUDA has been used: *shared memory*. This is a very fast memory shared by all the processors of a multiprocessor. Hence, the results of the second kernel are first copied to shared memory and then moved to global memory. If the copies are well organised, it is possible to re-order the access to fulfil all the memory requirements (for both shared and global) and thus reach the maximum bandwidth.

CPU-GPU interaction.

CPU-GPU interaction is generally a significant bottleneck in General Purpose GPU applications due to the relatively low interface bandwidth and it is desirable to

minimise such interaction. However, interaction cannot be entirely removed from the present implementation since, for example, the solver requires inputs in the form of loaded nodes (which may change due to the interaction with the user) and their displacements, and may need to provide outputs in the form of reaction forces for haptic feedback. CUDA alleviates the problem somewhat by allowing allocation of areas of page-locked host memory which are directly accessible by the GPU and therefore offer much higher bandwidth. In SOFA, all transfers between CPU and GPU are made via this mechanism.

Element technology.

We used both reduced integration 8-node hexahedral and 4-node (linear) tetrahedral elements. Tetrahedral meshes are easily generated and therefore widely used in simulations. However, hexahedra are preferable both in terms of solution accuracy and computational efficiency. Indeed, 4-node linear tetrahedra are known to be susceptible to volumetric locking (Hughes, 2000). Yet, a disadvantage of hexahedra is the difficulty in automatically generating hexahedral meshes over arbitrary domains. Construction of hexahedral meshes is time consuming and laborious, which fact is of even greater significance when patient-specific simulations are considered. For this reason tetrahedral meshes are widely used in simulations despite their short comings.

Another drawback of using hexahedra is the existence of so-called Hourglass modes that have to be addressed to avoid deterioration of the solution (Flanagan and Belytschko, 1981). Techniques for suppressing these modes exist, but naturally involve additional computations. However, for a given number of degrees of freedom (DOF), a hexahedral mesh can be built with far fewer elements than a tetrahedral one. Since the majority of the calculations in explicit dynamic analyses are performed per element, this results in reduced overall computation time.

From a GPU perspective, hexahedral element computations are substantially heavier and demand more memory resources than linear tetrahedral elements. Most of the matrices (such as shape function derivatives and nodal displacements) are twice as large for hexahedra, which necessitates twice as many texture fetches per element, and use of twice as many registers per thread. Similarly, twice as many nodal forces per element are written to global memory. Additional variables associated with hourglass control are required also. Therefore, on a per element basis hexahedra are significantly less efficient than tetrahedra, especially for GPU execution where memory efficiency is crucial. Thus, we observe that the occupancy (GPU percentage usage) drops from 25 % to only 8 % when using hexahedral elements. However, from the point of view of an entire model the lower number of hexahedra required for a given number of degrees of freedom still outweighs this element-wise inefficiency.

6.3.4 Second implementation: a better use of shared memory

Limitations of the previous implementation

The previous TLED algorithm has been extensively tested on regular meshes like cubes. The model was purposely simple to test our new CUDA implementation. However, the TLED algorithm was developed to compute the deformation of organs. Thus, we started to assess our TLED implementation on an irregular mesh: a liver. From a segmented liver obtained from IRCAD (Research institute against digestive system cancer, France), we meshed the surface using a marching cube technique (Lorensen and Cline, 1987). After smoothing it with Blender, a tetrahedral mesh has been generated with Tetgen, a free mesh generator. From there, we realised that the TLED was highly non efficient with such a mesh. After an extensive analysis with the profiler provided by NVIDIA, we became able to explain it. The liver mesh has a higher valency, meaning that each node has more elements attached than within a regular mesh like a cube. And on GPU, nodes are processed in parallel by group of 32 threads (called a warp). But if the instructions are different within a warp, the warp is divergent and the computations are serialised. With its higher valency, the liver mesh involves much more divergent warps when the forces are added up. To give an idea of the difference, a cube (343 nodes and 1 296 elements) is computed in 0.16 ms when a liver mesh (399 nodes and 1 250 elements) is processed in 0.27 ms. Based on the limitations discovered when the TLED has been applied to this irregular mesh, we came up with a new design, that we implemented and tested.

A single kernel approach

The main bottleneck of the previous implementation was adding up the forces on each node in a second kernel after the computation of element force contributions in a first kernel. The latter were stored in 4 different textures (8 for the hexahedral formulation) and in order to add up the forces the second kernel had to switch all the time between the correct textures within the warps, causing them to diverge and slowing down substantially the overall process. Our idea was therefore to use shared memory to avoid writing in global memory at the end of the first kernel and prevent the warps to diverge. As explained before, shared memory is a very fast memory shared by all the processors of a multiprocessor. And each block of threads is sent to a multiprocessor to be executed. Obviously the amount of shared memory is finite, which limits the validity of this new approach and the results in performance will depend on the hardware (which may have different amount of shared memory).

In the precomputation phase, nodes and elements are sorted by block to leave enough space for storage in shared memory of the element force contributions initially calculated at the end of the first kernel. We add nodes to a block along with their unique elements until the shared memory is saturated. Of course, in this approach some elements will be computed several times (because being into different blocks). However, storing the force contributions in shared memory instead of global

memory to read them later on via texture fetches widely outperforms the need of computing several elements several times.

TODO: If time, try this implementation on Geforce GTX 480 as it features 3 times the amount of shared memory.

6.3.5 Third implementation: atomic writes

In April 2010, NVIDIA released a family of graphics cards based on a new GPU architecture called Fermi. Among many new features and improvements over previous generation, these cards added the capability of so-called *atomic writes* for floats. The write operation is said to be atomic in the sense that it is guaranteed to be performed without interference from other threads. In other words, no other thread can access this address until the operation is complete. The concurrent writes are simply serialised in hardware. This feature allows us to implement the TLED algorithm as we would on a CPU, using a single kernel for all computations, and without the need to re-order the nodes and elements organisation as with the second implementation.

TODO: If time, try simple implementation using atomic writes with floats on Geforce GTX 480

6.4 Results

We present a series of examples based on pure shear and compression of cube models to demonstrate the validity of the constitutive update procedure and the performance of the GPU-based implementation. We conclude with an example of simulation of liver deformation, since as mentioned there is experimental evidence that liver exhibits an anisotropic response (Chui et al., 2007). In each case we used a transversely isotropic visco-hyperelastic model with elastic strain energy components defined by

$$\begin{aligned}\Psi^{iso} &= \Psi^I(\bar{I}_1) + \Psi^{TI}(\bar{I}_4) \\ &= \frac{\mu}{2}(\bar{I}_1 - 3) + \frac{\eta}{2}(\bar{I}_4 - 1)^2,\end{aligned}\tag{6.4}$$

$$\Psi^{vol} = \frac{\kappa}{2}(J - 1)^2,\tag{6.5}$$

where Ψ^I and Ψ^{TI} are isotropic and transversely isotropic components, respectively, μ is the small strain shear modulus, κ is the bulk modulus, and η is a material parameter with units of Pa. The isotropic component Ψ^I represents the well known neo-Hookean model, while the anisotropic component Ψ^{TI} is the same form used by Picinbono et al. (2003). Here the dependence of Ψ^{TI} on \bar{I}_5 is omitted and we include only \bar{I}_4 terms. This was done foremost for simplicity, but is also a common practical feature of anisotropic models for soft tissues. The main reason is the clear physical interpretation of \bar{I}_4 as the square of the stretch λ ¹ along the preferred

¹Strictly, $\bar{I}_4 = J^{-2/3}\lambda^2$

direction \mathbf{a}_0 Holzapfel et al. (2000), whereas \bar{I}_5 has a less clear physical meaning. \bar{I}_4 terms are therefore more amenable to experimental investigation, and of more use in modelling structural elements (e.g. fibres) responsible for the anisotropy. It must be emphasised that this model is used as an example only, and we make no claim concerning its appropriateness for any particular tissue.

The instantaneous isochoric and volumetric responses are then given by

$$\Phi^{iso} = J^{-2/3} \left[\mu \left(\mathbf{I} + \frac{I_1}{3} \mathbf{C}^{-1} \right) + \eta (\bar{I}_4 - 1) \left(\mathbf{A}_0 + \frac{I_4}{3} \mathbf{C}^{-1} \right) \right], \quad (6.6)$$

$$\Phi^{vol} = \kappa J (J - 1) \mathbf{C}^{-1}. \quad (6.7)$$

We used viscoelastic isochoric terms only, with a single Prony series term for simplicity. The complete stress response is then

$$\begin{aligned} \mathbf{S} &= \Phi^{vol} + \int_0^t \left[1 - \alpha_1 \left(1 - e^{(s-t)/\tau_1} \right) \right] \frac{\partial \Phi^{iso}}{\partial \tau} ds \\ &= \Phi^{vol} + \Phi^{iso} (1 - \alpha_1) + \alpha_1 \int_0^t e^{(s-t)/\tau_1} \frac{\partial \Phi^{iso}}{\partial \tau} ds. \end{aligned} \quad (6.8)$$

Note that in this case the developed constitutive update procedure need only to be employed for isochoric terms.

Material parameters were chosen based on recent results for the viscoelastic response of human liver *in vivo*. Nava et al. (2008) used an aspiration device to obtain an estimate of $\mu = 19\,700$ Pa for the combined liver/liver capsule. However, they noted that an earlier study suggested the presence of the liver capsule meant this value could overestimate liver parenchyma stiffness by up to $3\times$. Since it is the parenchyma we are concerned with, we employed a value of $\mu = 19\,700/3 = 6\,567$ Pa. Assuming near incompressibility we specified $\kappa = 326\,210$ Pa (corresponding to a Poisson ratio of 0.49). Nava et al. employed two Prony terms, and we use parameters from the first of these: $\alpha_1 = 0.5$, $\tau_1 = 0.58$. The short time constant τ_1 means this Prony term most affects faster loadings, which are of interest in surgical simulations. Finally, since Nava et al. assumed an isotropic response their data provide no basis for estimating η . Therefore we chose a conservative value of $\eta = 2\mu$, conservative since Picinbono et al. (2003) used a value closer to $\eta = 16\mu$.

6.4.1 Pure shear of a cube

The first example involved shearing of a unit cube model. We imposed a ramped loading such that the displacements u of the loaded nodes at time t were given by $u = rt$, where r is the loading speed, and we chose the x -direction to be the preferred material direction ($\mathbf{a}_0 = [1\ 0\ 0]$). Since the deformation is homogenous only a single element was required to reproduce the stresses exactly. Moreover an analytical solution is available with which numerical results may be compared (Taylor et al., 2009). Four independent nonzero stress terms are produced.

In Fig. 6.3 we demonstrate the effects of material anisotropy independent of viscoelastic effects by plotting the limiting instantaneous (hyperelastic) stress curves, produced by setting the relaxation parameter $\alpha_1 = 0$. Fig. 6.3(a) shows curves for the anisotropic model, while Fig. 6.3(b) shows curves for an equivalent isotropic model (i.e. with $\eta = 0$). Firstly we confirm that the FE solutions (data marks) match the analytical solutions (solid curves) very accurately. Secondly we observe that the anisotropic term significantly affects both the magnitude and shape of the stress curves. In particular the S_{11} curve which is coincident with the S_{33} curve in the isotropic model actually changes sign when anisotropy is introduced. The other curves exhibit a more pronounced stiffening for the anisotropic case.

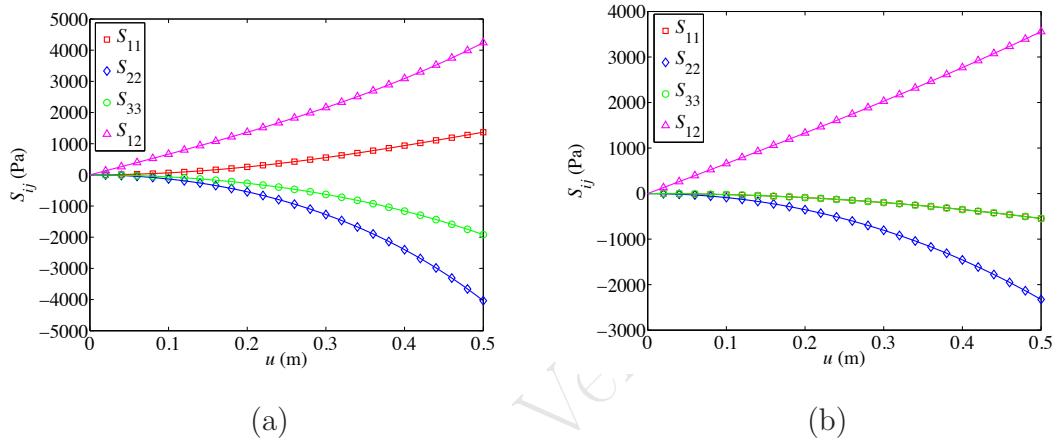


Figure 6.3: Instantaneous stress curves for pure shear deformation of (a) the anisotropic model compared with those of (b) an isotropic model. Solid lines correspond to the analytical solution, while markers indicate the FE solution.

In Fig. 6.4 we demonstrate the effect of loading rate on the complete viscoelastic model by plotting the S_{11} component for strain rates of 2.5 s^{-1} , 0.25 s^{-1} , and 0.025 s^{-1} , plus the bounding instantaneous and equilibrium responses. The strain rate-dependence introduced by the model (and commonly observed in biological tissues) is clearly shown. Additionally the close match between the analytic and numerical solutions demonstrates the validity of the developed constitutive update scheme.

6.4.2 Compression of a cube

The second example involved compression of a cube of edge length 0.1 m (comparable to some human organs). The cube was imagined to be fixed to opposing load platens and free on the remaining four faces. Compression was applied along the x -axis. A stress-relaxation type test protocol was simulated, in which the cube was compressed by 30 % in 0.5 seconds and the compression was held for a further 4.5 seconds. We demonstrate the effects of anisotropy on the deformation response by

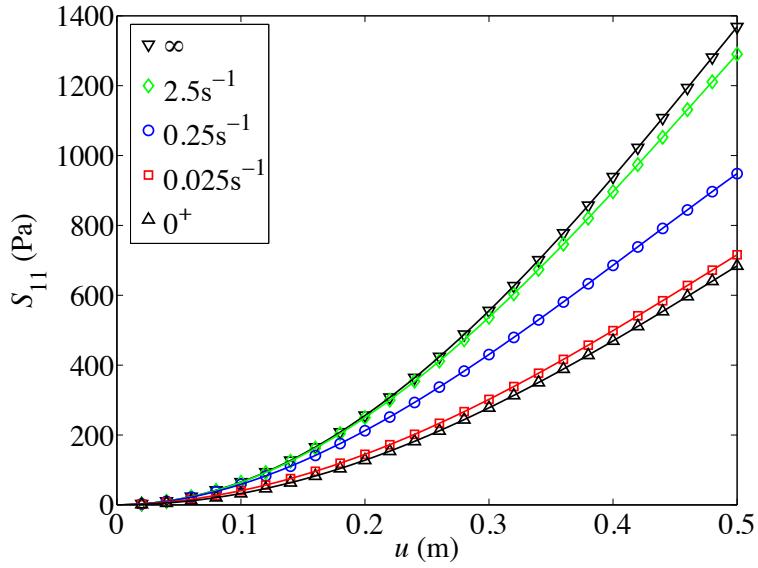


Figure 6.4: S_{11} curves for pure shear deformation of the anisotropic viscoelastic model at varying strain rates. Curves for strain rates of $2.5\text{ }s^{-1}$, $0.25\text{ }s^{-1}$, and $0.025\text{ }s^{-1}$ are given (as labelled), along with the bounding instantaneous and equilibrium responses. The latter two are labelled ∞ and 0^+ , respectively, indicating that they correspond to strain rates approaching these values. Solid lines correspond to the analytical solution, while markers indicate the FE solution.

comparing the deformed shapes of anisotropic models with $\mathbf{a}_0 = [0 \ 1 \ 0] \stackrel{\text{def}}{=} \mathbf{a}_0^y$ and $\mathbf{a}_0 = [0 \ \frac{1}{\sqrt{2}} \ \frac{1}{\sqrt{2}}] \stackrel{\text{def}}{=} \mathbf{a}_0^{yz}$ with that of an isotropic model (with $\eta = 0$).

Figure 6.5(a) shows the undeformed cube, while Fig. 6.5(b) shows the deformed isotropic model. With no preferred direction the lateral expansion was uniform. In Fig. 6.5(c) the increased y -direction stiffness of the first anisotropic model lead to a much reduced expansion in this direction (vertical in the image), and an accompanying increase in the orthogonal z -direction. For the second anisotropic model the direction defined by $y = z$ was stiffened, and Fig. 6.5(d) shows the resulting reduced expansion along this axis and the increased orthogonal expansion.

The end face reaction force history for each model is shown in Fig. 6.6. A distinct decay curve resulting from stress relaxation, and similar to that commonly noted in biological tissues (Fung, 1993), was observed. The increased stiffness afforded by the anisotropic models (albeit in an orthogonal direction to the loading) result in greater reaction forces than the isotropic model.

6.4.3 GPU performance

General GPU performance

Next, we assessed the computational performance of the GPU implementation and efficiency of the constitutive update scheme by measuring computation times for a range of mesh densities. Using the cube geometry as above we generated models with mesh sizes ranging from 3 993 DOF (1 331 nodes, 1 000 hexahedral elements) to 177 957 DOF (59 319 nodes, 54 872 hexahedral elements). We compared the GPU solution times for a single time step with those of a CPU-based implementation developed using C++. The test machine included an Intel Core2Duo 2.4GHz CPU, 2GB RAM, and an NVIDIA GeForce 8800GTX GPU. Three constitutive models were considered: the transversely isotropic viscoelastic model (TIV), TIV minus the viscoelastic terms (TIE), and TIE minus the anisotropy terms (NHE). Comparison of times for TIV and TIE indicates the computational load introduced by the developed viscoelastic constitutive update scheme. Comparison of times for TIE and NHE indicates the computational load introduced by material anisotropy. Figure 6.7 shows the results of the experiments.

In Fig. 6.7(a) and Fig. 6.7(b) it can be seen that GPU solution affords significant speed improvements over CPU solution for all constitutive models. Additionally it appears that solution times are little affected by the introduction of the more complex constitutive models, and importantly by the use of the developed constitutive update scheme. This is borne out in Fig. 6.7(d), where we observe that the maximum solution time ratio for model TIE to model NHE (anisotropy vs isotropy) was 1.013, and that of model TIV to model TIE (viscoelastic vs elastic) was 1.043. The largest total solution time increase for an anisotropic viscoelastic model compared with an isotropic hyperelastic one was 5.1% We conclude that the key features of anisotropy and viscoelasticity may be included in simulations at very little additional computational cost. Referring to Fig. 6.7(c) we observe a maximum speed improvement for GPU solution over CPU solution of 56.3×.

Fig. 6.7(b) also shows the decreasing critical time steps Δt_{cr} with increasing model sizes (which results from the decreasing element sizes), using liver-like and brain-like materials. We define liver-like materials to be those described by model TIV with material parameters quoted at the beginning of the section, whereas brain-like materials are those described by model NHE with material parameters $\mu = 1006.71$ Pa and Poisson's ratio $\nu = 0.45$. The intersections of these curves with the solution time curves provide an estimate of the real-time capacity of the implementation. This is discussed further in section 6.5.1.

TODO: Add comparison of different implementations

Performance within SOFA

The efficiency and the side-effects of porting the algorithm into the flexible SOFA framework need to be measured. Therefore the performance has been assessed by comparing the computational time of the algorithm running within and outside

SOFA. NVIDIA provides a tool to check the GPU implementation by evaluating many variables during the execution like for instance timings, counts of inconsistent reads and writes or GPU occupancy. We used this tool to carry out two measures:

1. GPU time only estimates the GPU computational time.
2. CPU time allows the evaluation of the execution time with the additional overhead due to the framework.

The tests were performed on a simple test scene featuring a cube under gravity. Hexahedral meshes with different resolutions from 1 331 to 29 791 nodes were used and the results are presented in figure 6.8. In our standalone implementation, the second kernel not only accumulates nodal forces but also adds gravity and updates positions based on the central difference integration scheme. These operations are split into separated components in SOFA, in order to introduce more flexibility (such as applying additional forces or changing the integration algorithm). While this introduces no noticeable difference on a CPU-based simulation, when using the GPU it is more costly due to overheads in the CUDA API for the additional kernel launches. Although it reduces the performance by 8.4 % for large meshes, this could be optimised away by adding a kernel specific to a given combination of components.

6.4.4 Simulation of liver deformation

Finally, we performed simulations of manipulation of a liver model in order to demonstrate the feasibility of the procedure with more realistic organ geometries. The model geometry was based on the liver mesh available in the SOFA framework, which we smoothed and remeshed to improve uniformity of element sizes and node valencies. A mesh of 748 nodes and 2 428 tetrahedral elements was solvable in real-time. For anisotropic models a preferred material direction $\mathbf{a}_0 = [0\ 1\ 0]$ was used, corresponding to the vertical direction in Fig. 6.9. We note that in a real organ the preferred direction likely varies depending on the local vasculature and orientation of lobules (Chui et al., 2007), but in the absence of experimental data for the present specimen the mentioned direction was assumed throughout. Moreover, on this point we note that different preferred directions may be specified for each element in the mesh, if such data are available.

The Von Mises equivalent stress distributions resulting from stretching the model with and without inclusion of viscoelasticity are compared in Fig. 6.9. For the viscoelastic case (Fig. 6.9, right image) the model is depicted several seconds after the completion of the stretch, so that stresses had reached an approximately steady state. The stress dissipation introduced by the viscoelastic model is clearly evident, and for example would result in significantly altered virtual tool reaction forces.

TODO: Add movies online? If possible...

6.5 Discussion

6.5.1 Critical time step

As already emphasised in section 5.1.3 page 77, the use of an explicit integration scheme imposes that the solution time remains below the critical time step. This critical time step depends on two kind of criterion: the size of the elements (due to the choice of the mesh) and the material's properties (Young's modulus and Poisson's ratio).

The first factor leads to a strong constraint on the quality of the mesh. Indeed, meshing a complex shape with good and uniform elements is not a straightforward task and rather often bad and very small elements are created. In such a case, the critical time step for the whole simulation would be small even if there is only a single small element in the mesh. As an example, the liver mesh used in the section 6.4.4 had to be smoothed and remeshed properly to allow a decent simulation.

Material's properties are also a very influent factor for the determination of a critical time step. If the critical time step may be fairly large for very soft tissue (like brain), it is already divided by 2 when modelling liver. As we have seen, the value of the Poisson's ratio may be relaxed from 0.49 to 0.45 for instance. If it allows a time step twice as large, this relaxation introduces inaccuracies. If this is certainly acceptable for medical training simulators, the demand in precision for surgical planning or per-operating assistance would probably prevent such an approximation.

In a nutshell, additional care must be taken to produce quality meshes and the range of materials that may be used while maintaining interactive computations is limited when using an explicit time integration scheme. However, it allows a huge gain in speed for very soft tissue when the compromises on material's properties are acceptable for the desired application, as with medical simulators used for training. For the interested reader, a more thorough discussion on the critical time step of explicit analyses may be found in Taylor et al. (2008).

6.5.2 Handling contacts

The present work was entirely concerned with numerical solution of constitutive models in explicit FE analyses, however an equally important issue from the point of view of interactive simulation is the modelling of interaction between organs and virtual tools, and its possible impact on Δt_{cr} and the stability of the simulation. Explicit time integration schemes are known to be more sensitive to contacts than implicit schemes. While this is beyond the scope of the present work, we note that a range of contact formulations are available in industry standard commercial explicit FE codes (e.g. LS-Dyna and ABAQUS). The methods are well tested and understood in the computational mechanics community, with two main types in common use: penalty methods and kinematic constraint methods. Penalty methods involve use of fictitious contact springs at interfaces, effectively to penalise interpenetration of opposing surfaces. If excessively stiff springs are used Δt_{cr} may indeed be reduced,

however except in extreme conditions (e.g. from explosive or high speed impacts, or imposition of very large forces on highly constrained objects) spring stiffnesses of the same order of magnitude as the material stiffness are suitable. The effect on Δt_{cr} is therefore minimal (Hallquist, 2006). Kinematic constraint methods involve (loosely) direct relocation of penetrating nodes back to the outside surface of the penetrated body, and have no effect on Δt_{cr} . These may be especially useful for modelling interaction between soft tissues and stiff (effectively rigid) surgical tools. The computational cost of either variety has not been investigated here.

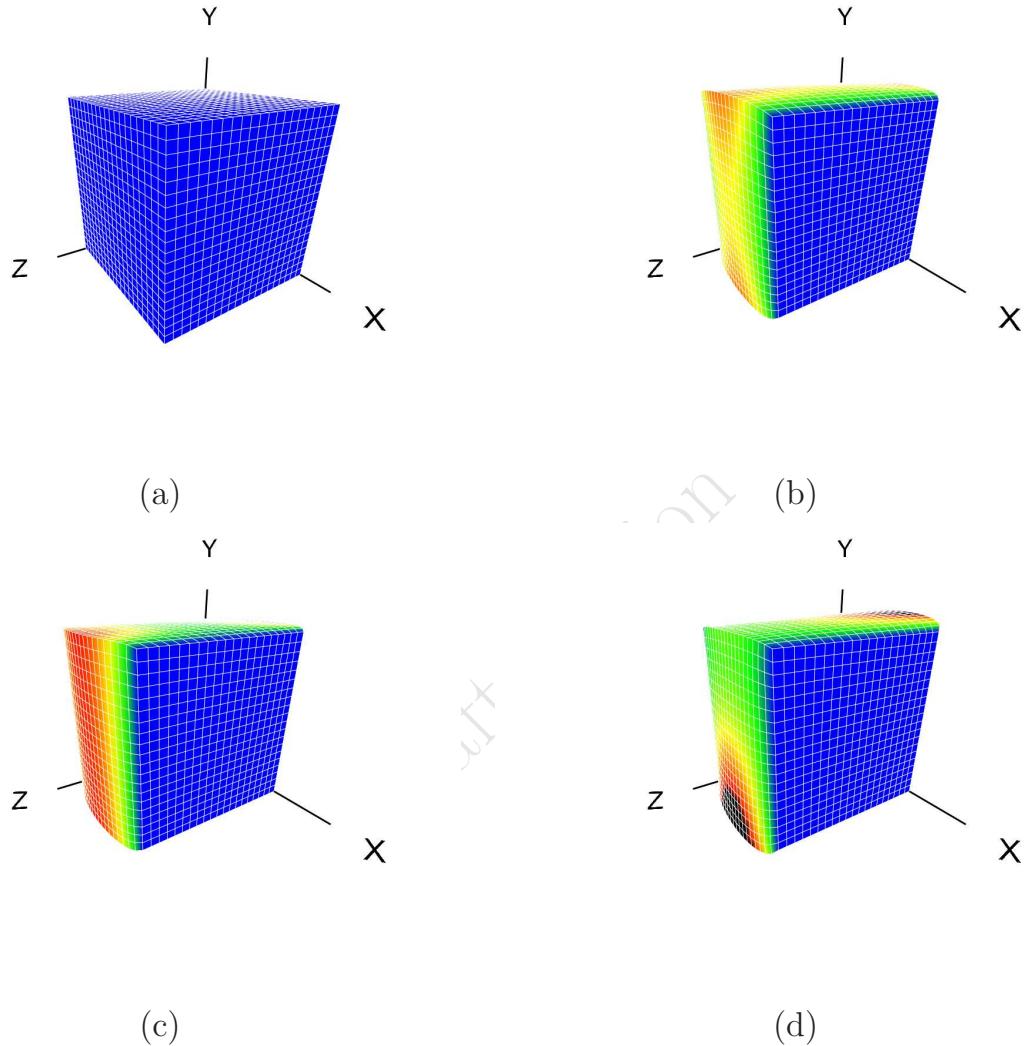


Figure 6.5: Deformation patterns of transversely isotropic models compared with that of an isotropic model: (a) the undeformed cube, (b) the deformed isotropic model, (c) the deformed anisotropic model with $\mathbf{a}_0 = \mathbf{a}_0^y$, (d) the deformed anisotropic model with \mathbf{a}_0^{yz} . Colour maps indicate relative magnitude of lateral displacement $((u_y^2 + u_z^2)^{1/2})$.

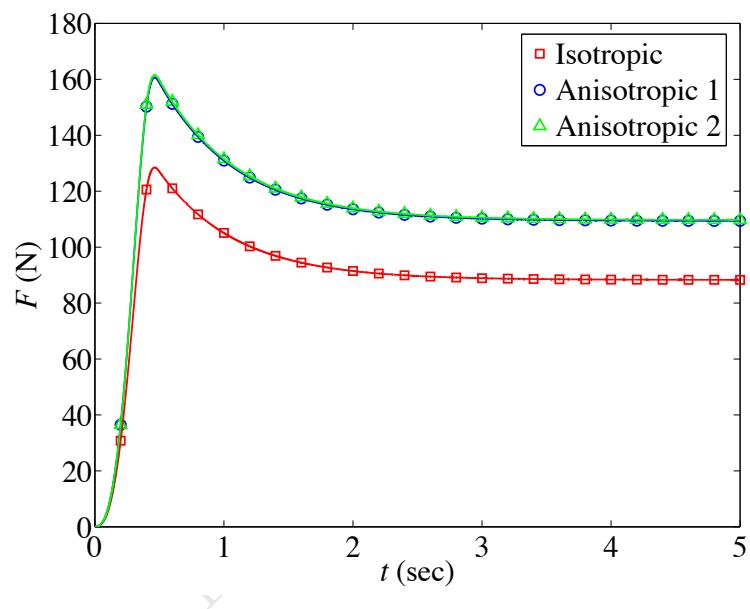


Figure 6.6: FE computed reaction forces on the end faces of the cube models over time. Curves are labelled according to the constitutive model used, where *Anisotropic 1* refers to the transversely isotropic viscoelastic model with $\mathbf{a}_0 = \mathbf{a}_0^y$ and *Anisotropic 2* refers to the case of $\mathbf{a}_0 = \mathbf{a}_0^{yz}$. Note that in this figure both markers and solid lines correspond to FE solutions.

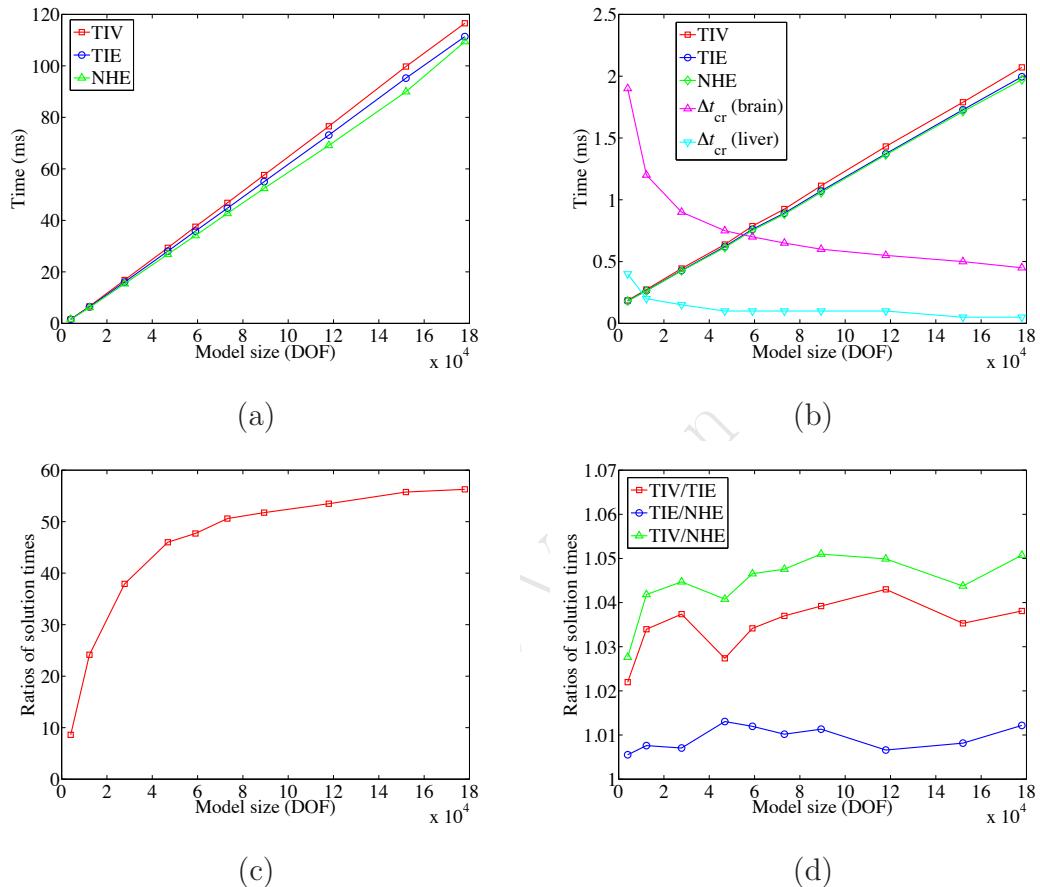


Figure 6.7: (a), (b) show CPU and GPU solution times, respectively, for a single time step for each constitutive model as a function of model size, (c) shows the ratio of CPU to GPU solution times for model TIV, and (d) shows ratios of GPU solution times for the constitutive models: TIV to TIE, TIE to NHE, and TIV to NHE. (b) also shows the critical time steps Δt_{cr} for each model size using liver-like and brain-like properties.

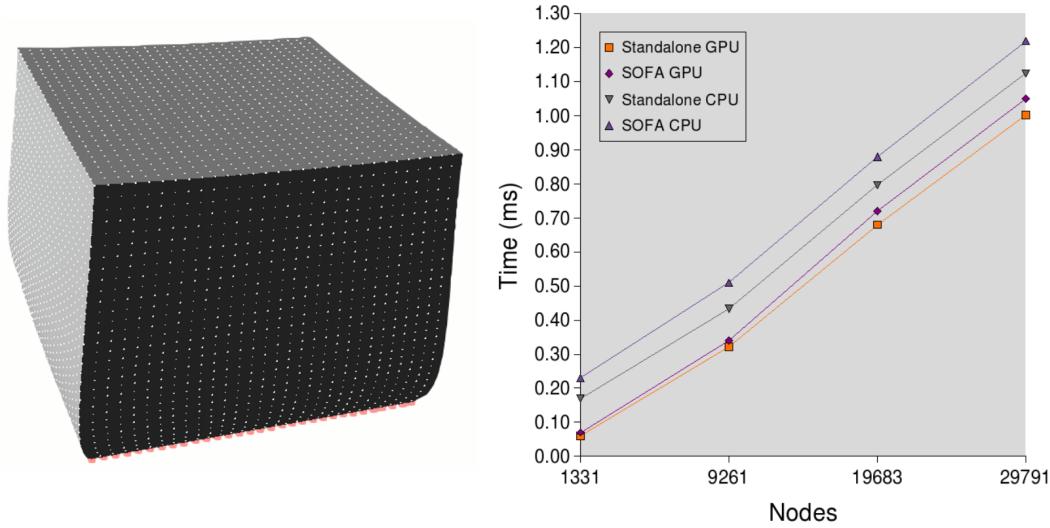


Figure 6.8: **Left:** FEM mesh of cube with 29 791 nodes deformed under a uniformed load. **Right:** Comparison of GPU computational timings and CPU overheads between the SOFA and standalone implementations for different mesh sizes.

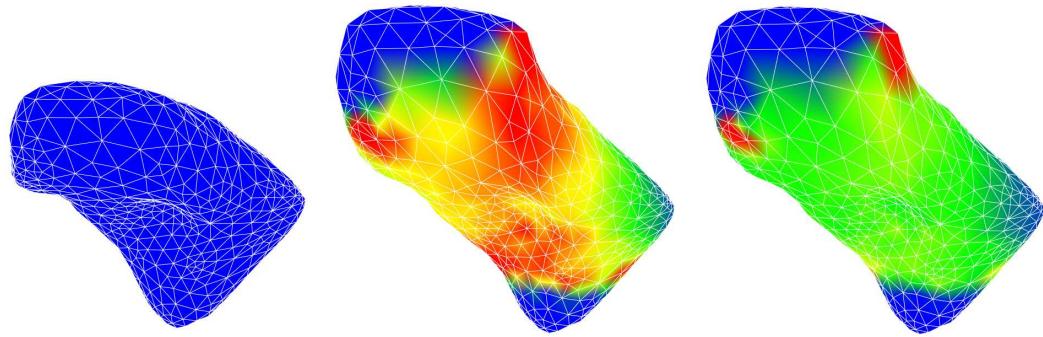


Figure 6.9: Deformation of a liver model using the TI model (middle) and the TIV model (right). The undeformed mesh is shown on the left. Colour maps indicate the relative Von Mises equivalent stress magnitude.

Part III

Hollow organs modelling

Draft Version 1.0

CHAPTER **7**

STATE OF ART: HOLLOW STRUCTURES

A short abstract for the upcoming chapter

Draft Version

7.1 Expose the problem: models presented for solid organs are insufficient. How to model hollow organs?

- 7.1.1 Colonoscopy simulator project: needs for colon**
- 7.1.2 Cataract surgery, stenosis: other needs for implants/blood vessels**
- 7.2 Non-physic approaches (computer graphics stuff)**
- 7.3 Physically accurate approaches (plates/shells)**

CHAPTER **8**

MASS-SPRING MODEL FOR COLON IMPLEMENTED ON GPU

A short abstract for the upcoming chapter

Draft Version

- 8.1 Colon is big, we need a fast model**
- 8.2 Limitations (difficult to link it to actual physical properties)**

Draft Version

CHAPTER 9

A CO-ROTATIONAL TRIANGULAR SHELL MODEL (ISBMS 2010)

A short abstract for the upcoming chapter

Draft Version

9.1 Model description

9.2 Validation

9.3 Application to implant deployment simulation in cataract surgery

Draft Version

CHAPTER 10

'SHELL MESHING' TECHNIQUE (MICCAI 2010)

A short abstract for the upcoming chapter

Draft Version

10.1 State of art: reconstruction/simplification

10.2 Our method

Draft Version

Part IV

Conclusion

CHAPTER 11

SUMMARY

A short abstract for the upcoming chapter

Draft Version

CHAPTER 12

DISCUSSION AND PERSPECTIVES (INTERACTION SOLID/HOLLOW ORGANS)

A short abstract for the upcoming chapter

Draft Version

APPENDIX A

TENSORS

Draft Version

APPENDIX B

THE WEIGHTED RESIDUAL METHOD

The *method of weighted residuals* is an approximation technique for solving differential equations. A subclass of this method, the *Galerkin method of weighted residuals*, is often used to derive the element equations for the finite element method.

Let us suppose that we have a linear differential operator D acting on a function u to produce a function p :

$$D(u(x)) = p(x). \quad (\text{B.1})$$

We wish to approximate u by a function \tilde{u} , which is a linear combination of basis functions chosen from a linearly independent set, that is:

$$u \approx \tilde{u} = \sum_{i=1}^n a_i \phi_i. \quad (\text{B.2})$$

Now, when substituted into the differential operator D , the result of the operation is not, in general, $p(x)$. Hence, an error or *residual* will exist:

$$R(x) = D(\tilde{u}(x) - p(x)) \neq 0. \quad (\text{B.3})$$

The idea behind the weighted residual method is to force the residual to zero in some average sense over the domain:

$$\int_X R(x) W_i(x) dx = 0, \quad i = 1, 2, \dots, n \quad (\text{B.4})$$

where the number of *weight functions* W_i is exactly equal to the number of unknown constants a_i in \tilde{u} . The result is a set of n algebraic equations for the unknown constants a_i . Depending on the choice for W_i , several sub-methods can be derived: collocation, sub-domain, least squares, Ritz, Galerkin etc. For instance, the weight function used by the Galerkin method TODO: need to check this: Wi are the shape functions? is:

$$W_i = \frac{\partial \tilde{u}}{\partial a_i}. \quad (\text{B.5})$$

References

- [Dir, 2010] 2010.
URL [http://msdn.microsoft.com/library/ee663275\(VS.85\).aspx](http://msdn.microsoft.com/library/ee663275(VS.85).aspx)
- [Allard et al., 2007] J. Allard, S. Cotin, F. Faure, P. j. Bensoussan, F. Poyer, C. Duriez, H. Delingette and L. Grisoni. *SOFA – an Open Source Framework for Medical Simulation*. In *In Medicine Meets Virtual Reality (MMVR 15)*, 2007.
- [Argyris et al., 1964] J. Argyris, S. Kelsey and H. Kamel. *Matrix Methods of Structural Analysis, A Precis of Recent Developments*. Pergamon Press, 1964.
- [Barbič and James, 2005] J. Barbič and D. L. James. *Real-Time Subspace Integration for St. Venant-Kirchhoff Deformable Models*. ACM Transactions on Graphics (SIGGRAPH 2005), vol. 24, no. 3, pages 982–990, 2005.
- [Barr, 1984] A. H. Barr. *Global and local deformations of solid primitives*. SIGGRAPH Comput. Graph., vol. 18, no. 3, pages 21–30, 1984.
- [Basafa and Farahmand, 2010] E. Basafa and F. Farahmand. *Real-time simulation of the nonlinear visco-elastic deformations of soft tissues*. International Journal of Computer Assisted Radiology and Surgery, pages 1–11, 2010.
- [Bathe, 1995] K.-J. Bathe. *Finite Element Procedures*. Prentice Hall, second ed., 1995.
- [Baudet et al., 2007] V. Baudet, M. Beuve, F. Jaillet, B. Shariat and F. Zara. *Integrating Tensile Parameters in 3D Mass-Spring System*. Tech. rep., LIRIS - CNRS, 2007.
- [Belytschko and Hsieh, 1973] T. Belytschko and B. J. Hsieh. *Non-linear transient finite element analysis with convected co-ordinates*. International Journal for Numerical Methods in Engineering, vol. 7, no. 3, pages 255–271, 1973.
- [Belytschko et al., 2000] T. Belytschko, W. K. Liu and B. Moran. *Nonlinear Finite Elements for Continua and Structures*. Wiley, first ed., 2000.
- [Bianchi et al., 2004] G. Bianchi, B. Solenthaler, G. Székely and M. Harders. *Simultaneous Topology and Stiffness Identification for Mass-Spring Models Based on FEM Reference Deformations*. In *MICCAI 2004*, vol. 3217, pages 293–301, 2004.
- [Biswas et al., 1976] D. Biswas, K. S. Ram and S. S. Rao. *Application of ‘natural coordinate system’ in the finite element solution of multigroup neutron diffusion equation*. Annals of Nuclear Energy, vol. 3, no. 11-12, pages 465–469, 1976.
- [Board et al., 2005] O. A. R. Board, D. Shreiner, M. Woo, J. Neider and T. Davis. *OpenGL Programming Guide: The Official Guide to Learning OpenGL*. Addison-Wesley Professional, fifth ed., 2005.

- [Bonnet and Wood, 1997] J. Bonnet and R. D. Wood. *Nonlinear Continuum Mechanics For Finite Element Analysis*. Cambridge University Press, 1997.
- [Bourguignon and Cani, 2000] D. Bourguignon and M.-P. Cani. *Controlling Anisotropy in Mass-Spring Systems*. In *Eurographics Workshop on Computer Animation and Simulation*, pages 113–123, Springer-Verlag, 2000.
- [Bro-Nielsen and Cotin, 1996] M. Bro-Nielsen and S. Cotin. *Real-time Volumetric Deformable Models for Surgery Simulation using Finite Elements and Condensation*. In *Computer Graphics Forum*, pages 57–66, 1996.
- [Chen and Zeltzer, 1992] D. T. Chen and D. Zeltzer. *Pump it up: computer animation of a biomechanically based model of muscle using the finite element method*. SIGGRAPH Comput. Graph., vol. 26, no. 2, pages 89–98, 1992.
- [Chhatkuli et al., 2009] S. Chhatkuli, H. Kojima, S. Koshizuka, and M. Uesaka. *Mesh-free Simulation of Lung Deformation during Inspiration*. In *World Congress on Engineering and Computer Science*, vol. 2, 2009.
- [Choi and Ko, 2005] M. G. Choi and H.-S. Ko. *Modal Warping: Real-Time Simulation of Large Rotational Deformation and Manipulation*. IEEE Transactions on Visualization and Computer Graphics, vol. 11, no. 1, pages 91–101, 2005.
- [Chowdhury and Dasgupta, 2003] I. Chowdhury and S. Dasgupta. *Computation of Rayleigh damping coefficients for large systems*. The Electronic Journal of Geotechnical Engineering, vol. 8, no. C, 2003.
- [Chui et al., 2007] C. Chui, E. Kobayashi, X. Chen, T. Hisada and I. Sakuma. *Transversely isotropic properties of porcine liver tissue: experiments and constitutive modelling*. Medical and Biological Engineering and Computing, vol. 45, no. 1, pages 99–106, 2007.
- [Coquillart, 1990] S. Coquillart. *Extended free-form deformation: a sculpturing tool for 3D geometric modeling*. In *SIGGRAPH '90: Proceedings of the 17th annual conference on Computer graphics and interactive techniques*, pages 187–196, ACM, New York, NY, USA, 1990.
- [Coquillart and Jancéne, 1991] S. Coquillart and P. Jancéne. *Animated free-form deformation: an interactive animation technique*. In *SIGGRAPH '91: Proceedings of the 18th annual conference on Computer graphics and interactive techniques*, pages 23–26, ACM, New York, NY, USA, 1991.
- [Cotin et al., 1999] S. Cotin, H. Delingette and N. Ayache. *Real-time elastic deformations of soft tissues for surgery simulation*. IEEE Transactions On Visualization and Computer Graphics, vol. 5, no. 1, pages 62–73, 1999.
- [Cotin et al., 2000] S. Cotin, H. Delingette and N. Ayache. *A hybrid elastic model for real-time cutting, deformations, and force feedback for surgery training and simulation*. The Visual Computer, vol. 16, no. 8, pages 437–452, 2000.

- [De et al., 2001] S. De, J. Kim and M. A. Srinivasan. *A meshless numerical technique for physically based real time medical simulations*. In *Medicine Meets Virtual Reality (MMVR)*, vol. 81, pages 113–118, 2001.
- [Debunne et al., 2001] G. Debunne, M. Desbrun, M.-P. Cani and A. H. Barr. *Dynamic real-time deformations using space & time adaptive sampling*. In *SIGGRAPH '01: Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 31–36, ACM, New York, NY, USA, 2001.
- [Felippa, 2000] C. A. Felippa. *A systematic approach to the element independent corotational dynamics of finite elements*. Tech. rep., Department of Aerospace Engineering Sciences and Center for Aerospace Structures. University of Colorado, 2000.
- [Felippa and Haugen, 2005] C. A. Felippa and B. Haugen. *A unified formulation of small-strain corotational finite elements: I. Theory*. Computer Methods in Applied Mechanics and Engineering, vol. 194, no. 21-24, pages 2285–2335, 2005.
- [Flanagan and Belytschko, 1981] D. P. Flanagan and T. Belytschko. *A uniform strain hexahedron and quadrilateral with orthogonal hourglass control*. International Journal for Numerical Methods in Engineering, vol. 17, no. 5, pages 679–706, 1981.
- [Fung, 1993] Y. C. Fung. *Biomechanics: Mechanical Properties of Living Tissues*. Springer, second ed., 1993.
- [Gent, 2001] A. N. Gent. *Engineering With Rubber: How to Design Rubber Components*. Hanser Gardner Publications, 2001.
- [Gibson, 1997] S. F. Gibson. *3D chainmail: a fast algorithm for deforming volumetric objects*. In *I3D '97: Proceedings of the 1997 symposium on Interactive 3D graphics*, pages 149–154, ACM, New York, NY, USA, 1997.
- [Gourret et al., 1989] J.-P. Gourret, N. M. Thalmann and D. Thalmann. *Simulation of object and human skin formations in a grasping task*. SIGGRAPH Comput. Graph., vol. 23, no. 3, pages 21–30, 1989.
- [Hallquist, 2006] J. O. Hallquist. *LS-DYNA Theory Manual*. Livermore Software Technology Corporation, Livermore, California, 2006.
- [Hauser et al., 2003] K. K. Hauser, C. Shen and J. F. O'Brien. *Interactive Deformation Using Modal Analysis with Constraints*. In *Graphics Interface*, pages 247–256, CIPS, Canadian Human-Computer Communication Society, 2003.
- [Hauth and Strasser, 2004] M. Hauth and W. Strasser. *Corotational Simulation of Deformable Solids*. Journal of WSCG, pages 137–145, 2004.
- [Holzapfel et al., 2000] G. Holzapfel, T. Gasser and R. Ogden. *A New Constitutive Framework for Arterial Wall Mechanics and a Comparative Study of Material Models*. Journal of Elasticity, vol. 61, no. 1, pages 1–48, 2000.

- [Holzapfel, 1996] G. A. Holzapfel. *On large strain viscoelasticity: continuum formulation and finite element applications to elastomeric structures*. International Journal for Numerical Methods in Engineering, vol. 39, no. 22, pages 3903–3926, 1996.
- [Horton et al., 2006] A. Horton, A. Wittek and K. Miller. *Towards Meshless Methods for Surgical Simulation*. 2006.
- [Horton et al., 2007] A. Horton, A. Wittek and K. Miller. *Subject-specific biomechanical simulation of brain indentation using a meshless method*. In *MICCAI'07: Proceedings of the 10th international conference on Medical image computing and computer-assisted intervention*, pages 541–548, Springer-Verlag, Berlin, Heidelberg, 2007.
- [Horton et al., 2010] A. Horton, A. Wittek, G. R. Joldes and K. Miller. *A meshless Total Lagrangian explicit dynamics algorithm for surgical simulation*. International Journal for Numerical Methods in Biomedical Engineering, vol. 26, no. 8, pages 977–998, 2010.
- [Hughes, 2000] T. J. R. Hughes. *The Finite Element Method: Linear Static and Dynamic Finite Element Analysis*. Dover Publications, 2000.
- [James and Pai, 1999] D. L. James and D. K. Pai. *ArtDefo: accurate real time deformable objects*. In *SIGGRAPH '99: Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 65–72, ACM Press/Addison-Wesley Publishing Co., New York, NY, USA, 1999.
- [James and Pai, 2002] D. L. James and D. K. Pai. *DyRT: dynamic response textures for real time deformation simulation with graphics hardware*. In *SIGGRAPH '02: Proceedings of the 29th annual conference on Computer graphics and interactive techniques*, pages 582–585, ACM, New York, NY, USA, 2002.
- [Ji et al., 2010] W. Ji, A. M. Waas and Z. P. Bazant. *Errors Caused by Non-Work-Conjugate Stress and Strain Measures and Necessary Corrections in Finite Element Programs*. Journal of Applied Mechanics, vol. 77, no. 4, pages 044,504–5, 2010.
- [Keeve et al., 1998] E. Keeve, S. Girod, R. Kikinis and B. Girod. *Deformable Modeling of Facial Tissue for Craniofacial Surgery Simulation*. Computer Aided Surgery, vol. 3, pages 3–228, 1998.
- [Lai et al., 1996] W. Lai, D. Rubin and E. Krempl. *Introduction to Continuum Mechanics*. Butterworth-Heinemann, third ed., 1996.
- [Lautrup, 2005] B. Lautrup. *Physics of continuous matter: exotic and everyday phenomena in the macroscopic world*. Institute of Physics Publishing, 2005.

- [Lim and De, 2004] Y. J. Lim and S. De. *On the use of meshfree methods and a geometry based surgical cutting algorithm in multimodal medical simulations*. Haptic Interfaces for Virtual Environment and Teleoperator Systems, 2004. HAPTICS '04. Proceedings. 12th International Symposium on, pages 295–301, 2004.
- [Liu and Gu, 2005] G. Liu and Y. Gu. *An Introduction to Meshfree Methods and Their Programming*. Springer, 2005.
- [Liu and Quek, 2003] G. Liu and S. S. Quek. *Finite Element Method: A Practical Course*. Butterworth-Heinemann, 2003.
- [Lloyd et al., 2007] B. A. Lloyd, G. Székely and M. Harders. *Identification of spring parameters for deformable object simulation*. IEEE Transactions on Vizualisation and Computer Graphics, pages 1081–1094, 2007.
- [Lorensen and Cline, 1987] W. E. Lorensen and H. E. Cline. *Marching cubes: A high resolution 3D surface construction algorithm*. SIGGRAPH Comput. Graph., vol. 21, no. 4, pages 163–169, 1987.
- [MacCracken and Joy, 1996] R. MacCracken and K. I. Joy. *Free-form deformations with lattices of arbitrary topology*. In *SIGGRAPH '96: Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 181–188, ACM, New York, NY, USA, 1996.
- [MacDonald, 2007] B. J. MacDonald. *Practical Stress Analysis with Finite Elements*. Glasnevin Publishing, first ed., 2007.
- [Magjarevic et al., 2009] R. Magjarevic, J. H. Nagel, J. Sloten, P. Verdonck, M. Nyssen, J. Haueisen, O. Jarrousse, T. Fritz and O. Dössel. *A modified Mass-Spring system for myocardial mechanics modeling*. In *4th European Conference of the International Federation for Medical and Biological Engineering*, vol. 22, pages 1943–1946, Springer Berlin Heidelberg, 2009.
- [Melvin et al., 1973] J. Melvin, R. Stalnaker and V. Roberts. *Impact injury mechanisms in abdominal organs*. SAE Transactions, , no. 730968, pages 115–126, 1973.
- [Miller, 1988] G. S. P. Miller. *The motion dynamics of snakes and worms*. In *SIGGRAPH '88: Proceedings of the 15th annual conference on Computer graphics and interactive techniques*, pages 169–173, ACM, New York, NY, USA, 1988.
- [Miller, 2000] K. Miller. *Constitutive modelling of abdominal organs*. Journal of Biomechanics, vol. 33, no. 3, pages 367 – 373, 2000, ISSN 0021-9290.
URL <http://www.sciencedirect.com/science/article/B6T82-3Y6Y6WS-D/2/d62ff415a78f2aed9c18401fa2a35cda>
- [Miller and Chinzei, 1997] K. Miller and K. Chinzei. *Constitutive modelling of brain tissue: Experiment and theory*. Journal of Biomechanics, vol. 30, no. 11-12, pages 1115 – 1121, 1997.

- [Miller and Chinzei, 2002] K. Miller and K. Chinzei. *Mechanical properties of brain tissue in tension*. Journal of Biomechanics, vol. 35, no. 4, pages 483–490, 2002.
- [Miller et al., 2007] K. Miller, G. Joldes, D. Lance and A. Wittek. *Total Lagrangian explicit dynamics finite element algorithm for computing soft tissue deformation*. Communications in Numerical Methods in Engineering, vol. 23, no. 2, pages 121–134, 2007.
- [Misra et al., 2007] S. Misra, A. M. Okamura and K. T. Ramesh. *Force Feedback is Noticeably Different for Linear versus Nonlinear Elastic Tissue Models*. In *WHC '07: Proceedings of the Second Joint EuroHaptics Conference and Symposium on Haptic Interfaces for Virtual Environment and Teleoperator Systems*, pages 519–524, IEEE Computer Society, Washington, DC, USA, 2007.
- [Mor and Kanade, 2000] A. B. Mor and T. Kanade. *Modifying Soft Tissue Models: Progressive Cutting with Minimal New Element Creation*. In *MICCAI '00: Proceedings of the Third International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 598–607, Springer-Verlag, London, UK, 2000.
- [Müller et al., 2004] M. Müller, R. Keiser, A. Nealen, M. Pauly, M. Gross and M. Alexa. *Point based animation of elastic, plastic and melting objects*. In *SCA '04: Proceedings of the 2004 ACM SIGGRAPH/Eurographics symposium on Computer animation*, pages 141–151, Eurographics Association, Aire-la-Ville, Switzerland, Switzerland, 2004.
- [Müller et al., 2005] M. Müller, B. Heidelberger, M. Teschner and M. Gross. *Meshless deformations based on shape matching*. ACM Trans. Graph., vol. 24, no. 3, pages 471–478, 2005.
- [Nava et al., 2008] A. Nava, E. Mazza, M. Furrer, P. Villiger and W. H. Reinhart. *In vivo mechanical characterization of human liver*. Medical Image Analysis, vol. 12, no. 2, pages 203–216, 2008.
- [Ng-Thow-Hing and Fiume, 1997] V. Ng-Thow-Hing and E. Fiume. *Interactive display and animation of B-spline solids as muscle shape primitives*. In *Eurographics Workshop on Animation and Simulation*, pages 81–97, 1997.
- [Nienhuys and van der Stappen, 2001] H.-W. Nienhuys and A. F. van der Stappen. *A Surgery Simulation Supporting Cuts and Finite Element Deformation*. In *Medical Image Computing and Computer-Assisted Intervention –MICCAI 2001*, vol. 2208, pages 145–152, Springer Berlin / Heidelberg, 2001.
- [Owens et al., 2007] J. D. Owens, D. Luebke, N. Govindaraju, M. Harris, J. Krüger, A. E. Lefohn and T. J. Purcell. *A Survey of General-Purpose Computation on Graphics Hardware*. Computer Graphics Forum, vol. 26, no. 1, pages 80–113, 2007.

- [Pentland and Williams, 1989] A. Pentland and J. Williams. *Good vibrations: modal dynamics for graphics and animation*. SIGGRAPH Comput. Graph., vol. 23, no. 3, pages 207–214, 1989.
- [Picinbono et al., 2003] G. Picinbono, H. Delingette and N. Ayache. *Non-linear anisotropic elasticity for real-time surgery simulation*. Graph. Models, vol. 65, no. 5, pages 305–321, 2003.
- [Pieper et al., 1992] S. Pieper, J. Rosen and D. Zeltzer. *Interactive graphics for plastic surgery: a task-level analysis and implementation*. In *I3D '92: Proceedings of the 1992 symposium on Interactive 3D graphics*, pages 127–134, ACM, New York, NY, USA, 1992.
- [Poon and Ahmad, 1998] H. Poon and M. F. Ahmad. *A material point time integration procedure for anisotropic, thermo rheologically simple, viscoelastic solids*. Computational Mechanics, vol. 21, no. 3, pages 236–242, 1998.
- [Pouliquen et al., 2005] M. Pouliquen, C. Duriez, C. Andriot, A. Bernard, L. Chodorge and F. Gosselin. *Real-Time Finite Element Finger Pinch Grasp Simulation*. In *WHC '05: Proceedings of the First Joint Eurohaptics Conference and Symposium on Haptic Interfaces for Virtual Environment and Teleoperator Systems*, pages 323–328, IEEE Computer Society, Washington, DC, USA, 2005.
- [Prange and Margulies, 2002] M. T. Prange and S. S. Margulies. *Regional, Directional, and Age-Dependent Properties of the Brain Undergoing Large Deformation*. Journal of Biomechanical Engineering, vol. 124, no. 2, pages 244–252, 2002.
- [Pras Pathmanathan and Rajagopal, 2004] J. W. S. M. B. M. N. P. N. Pras Pathmanathan, David Gavaghan and V. Rajagopal. *Predicting Tumour Location by Simulating Large Deformations of the Breast Using a 3D Finite Element Model and Nonlinear Elasticity*. pages 217–224, Springer Berlin / Heidelberg, 2004.
- [Press et al., 1992] W. H. Press, B. P. Flannery, S. A. Teukolsky and W. T. Vetterling. *Numerical Recipes in FORTRAN 77: The Art of Scientific Computing*. Cambridge University Press, second ed., 1992.
- [Press et al., 2002] W. H. Press, S. A. Teukolsky, W. T. Vetterling and B. P. Flannery. *Numerical Recipes in C++: The Art of Scientific Computing*. Cambridge University Press, second ed., 2002.
- [Reddy, 1993] J. N. Reddy. *Introduction to the Finite Element Method*. McGraw-Hill, second ed., 1993.
- [Reddy, 2007] J. N. Reddy. *An Introduction to Continuum Mechanics*. Cambridge University Press, 2007.
- [Rivers and James, 2007] A. R. Rivers and D. L. James. *FastLSM: fast lattice shape matching for robust real-time deformation*. In *SIGGRAPH '07: ACM SIGGRAPH 2007 papers*, page 82, ACM, New York, NY, USA, 2007.

- [Rumpf and Strzodka, 2001] M. Rumpf and R. Strzodka. *Using Graphics Cards for Quantized FEM Computations*. In *Proceedings of IASTED Visualization, Imaging and Image Processing Conference (VIIP'01)*, pages 193–202, 2001.
- [Sagar et al., 1994] M. A. Sagar, D. Bullivant, G. D. Mallinson and P. J. Hunter. *A virtual environment and model of the eye for surgical simulation*. In *SIGGRAPH '94: Proceedings of the 21st annual conference on Computer graphics and interactive techniques*, pages 205–212, ACM, New York, NY, USA, 1994.
- [San Vicente et al., 2009] G. San Vicente, C. Buchart, D. Borro and J. Celigüeta. *Maxillofacial surgery simulation using a mass-spring model derived from continuum and the scaled displacement method*. International journal of computer assisted radiology and surgery, vol. 4, no. 1, pages 89–98, 2009.
- [Schein and Elber, 2004] S. Schein and G. Elber. *Discontinuous Free Form Deformations*. In *PG '04: Proceedings of the Computer Graphics and Applications, 12th Pacific Conference*, pages 227–236, IEEE Computer Society, Washington, DC, USA, 2004.
- [Schill et al., 1998] M. A. Schill, S. F. F. Gibson, H.-J. Bender and R. Männer. *Biomechanical Simulation of the Vitreous Humor in the Eye Using and Enhanced Chain-Mail Algorithm*. In *MICCAI '98: Proceedings of the First International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 679–687, Springer-Verlag, London, UK, 1998.
- [Sedef et al., 2006] M. Sedef, E. Samur and C. Basdogan. *Real-Time Finite-Element Simulation of Linear Viscoelastic Tissue Behavior Based on Experimental Data*. IEEE Comput. Graph. Appl., vol. 26, no. 6, pages 58–68, 2006, ISSN 0272-1716.
- [Sederberg and Parry, 1986] T. W. Sederberg and S. R. Parry. *Free-form deformation of solid geometric models*. SIGGRAPH Comput. Graph., vol. 20, no. 4, pages 151–160, 1986.
- [Serby et al., 2001] D. Serby, M. Harders and G. Székely. *A New Approach to Cutting into Finite Element Models*. In *MICCAI '01: Proceedings of the 4th International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 425–433, Springer-Verlag, London, UK, 2001.
- [Shewchuk, 1994] J. R. Shewchuk. *An Introduction to the Conjugate Gradient Method Without the Agonizing Pain*. Tech. rep., Shool of Computer Science - Carnegie Mellon University, 1994.
- [Shinya, 2005] M. Shinya. *Theories for Mass-Spring Simulation in Computer Graphics: Stability, Costs and Improvements*. IEICE - Trans. Inf. Syst., vol. E88-D, no. 4, pages 767–774, 2005.
- [Spencer, 1980] A. Spencer. *Continuum Mechanics*. Longman Group Limited, 1980.

- [Steinemann et al., 2008] D. Steinemann, M. A. Otaduy and M. Gross. *Fast adaptive shape matching deformations*. In *SCA '08: Proceedings of the 2008 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 87–94, Eurographics Association, 2008.
- [Stiles and Alexander, 1972] R. N. Stiles and D. M. Alexander. *A viscoelastic-mass model for muscle*. Mathematical Biosciences, vol. 14, no. 3-4, pages 343–354, 1972.
- [Székely et al., 2000] G. Székely, C. Brechbühler, J. Dual, R. Enzler, J. Hug, R. Hutter, N. Ironmonger, M. Kauer, V. Meier, P. Niederer, A. Rhomberg, P. Schmid, G. Schweitzer, M. Thaler, V. Vuskovic, G. Tröster, U. Haller and M. Bajka. *Virtual Reality-Based Simulation of Endoscopic Surgery*. Presence: Teleoper. Virtual Environ., vol. 9, no. 3, pages 310–333, 2000.
- [Tamura et al., 2005] N. Tamura, N. Tsumura, T. Nakaguchi and Y. Miyake. *Spring-bead animation of viscoelastic materials*. In *SIGGRAPH '05: ACM SIGGRAPH 2005 Sketches*, page 64, ACM, New York, NY, USA, 2005.
- [Taylor et al., 2007a] Z. Taylor, M. Cheng and S. Ourselin. *Real-Time Nonlinear Finite Element Analysis for Surgical Simulation Using Graphics Processing Units*. Medical Image Computing and Computer-Assisted Intervention –MICCAI 2007, pages 701–708, 2007a.
URL http://dx.doi.org/10.1007/978-3-540-75757-3_85
- [Taylor et al., 2009] Z. Taylor, O. Comas, M. Cheng, J. Passenger, D. Hawkes, D. Atkinson and S. Ourselin. *On modelling of anisotropic viscoelasticity for soft tissue simulation: Numerical solution and GPU execution*. Medical Image Analysis, vol. 13, no. 2, pages 234–244, 2009.
- [Taylor, 2006] Z. A. Taylor. *A synopsis of the total Lagrangian explicit dynamic finite element algorithm*. Not published, 2006.
- [Taylor et al., 2007b] Z. A. Taylor, T. B. Kirk and K. Miller. *Confocal arthroscopy-based patient-specific constitutive models of cartilaginous tissues - I: development of a microstructural model*. Computer Methods in Biomechanics and Biomedical Engineering, vol. 10, no. 4, pages 307–316, 2007b.
- [Taylor et al., 2008] Z. A. Taylor, M. Cheng and S. Ourselin. *High-Speed Nonlinear Finite Element Analysis for Surgical Simulation Using Graphics Processing Units*. IEEE Transactions on Medical Imaging, vol. 27, no. 5, pages 650–663, 2008.
- [Terzopoulos et al., 1987] D. Terzopoulos, J. Platt, A. Barr and K. Fleischert. *Elastically Deformable Models*. Computer Graphics, vol. 21, pages 205–214, 1987.
- [Terzopoulos et al., 1991] D. Terzopoulos, J. Platt and K. Fleischer. *Heating and melting deformable models*. The Journal of Visualization and Computer Animation, vol. 2, no. 2, pages 68–73, 1991.

- [Wempner, 1969] G. Wempner. *Finite elements, finite rotations and small strains of flexible shells*. International Journal of Solids and Structures, vol. 5, no. 2, pages 117–153, 1969.
- [Wu and Heng, 2004] W. Wu and P. A. Heng. *A hybrid condensed finite element model with GPU acceleration for interactive 3D soft tissue cutting*. Comput. Animat. Virtual Worlds, vol. 15, no. 3-4, pages 219–227, 2004.
- [Wu et al., 2001] X. Wu, M. S. Downes, T. Goktekin and F. Tendick. *Adaptive Nonlinear Finite Elements for Deformable Body Simulation Using Dynamic Progressive Meshes*. In *Computer Graphics Forum*, pages 349–358, 2001.
- [Yan et al., 2007] Z. Yan, L. Gu, P. Huang, S. Lv, X. Yu and X. Kong. *Soft tissue deformation simulation in virtual surgery using nonlinear finite element method*. In *The 29th Annual International Conference of the IEEE EMBS*, pages 3642–3645, 2007.
- [Zhong et al., 2005] H. Zhong, M. P. Wachowiak and T. Peters. *A real time finite element based tissue simulation method incorporating nonlinear elastic behavior*. Comput Methods Biomech Biomed Engin., vol. 8, no. 3, pages 177–189, 2005.