

## FORMAL SOFTWARE ENGINEERING FOR COMPUTATIONAL MODELLING

MAGNE HAVERAAEN

*Institutt for Informatikk, Universitetet i Bergen*

*N-5020 Bergen, Norway*

`Magne.Haveraaen@ii.uib.no`

HELMER ANDRÉ FRIIS\*

*Institutt for Informatikk, Universitetet i Bergen*

*N-5020 Bergen, Norway*

`HelmerAndre.Friis@rf.no`

TOR ARNE JOHANSEN

*Institutt for den faste jords fysikk, Universitetet i Bergen*

*N-5020 Bergen, Norway*

`TorArne.Johansen@ifjf.uib.no`

**Abstract.** Software itself may be considered a formal structure and may be subject to mathematical analysis. This leads to a discipline of formal software engineering (which is not necessarily the same as the use of formal methods in software engineering), where a formal understanding of what software components are and how they may interact is used to engineer both the components themselves and their organisation. A strategy is using the concepts that are suited for organising the problem domain itself to organise the software as well. In this paper we apply this idea in the development of computational modelling software, in particular in the development of a family of related programs for simulation of elastic wave propagation in earth materials. We also discuss some data on the technique's effectiveness.

**CR Classification:** C.0, D.1.5, D.2.1, D.2.11, D.2.13, G.4, I.0

**Key words:** algebraic software methodologies, coordinate free numerics, numerical software, seismic simulation, domain specific languages, software architecture, software process model, software life cycle model

### 1. Introduction

Formal methods and reuse of software components has for a long time been acknowledged as important for improvement of software quality and reduction of software cost. Investigating the formal aspects of software may also lead to significant improvements in these areas. Such a focus may be termed *formal software engineering*. Formal software engineering is not the same as

---

\*Presently at Rogalandforskning, P.O.Box 2503, Ullandhaug, N-4004 Stavanger, Norway.

the use of formal methods in software engineering. Rather it is a focus on the formal, mathematical side of software artifacts, and using this insight in the engineering of software products. Following the tradition of mathematics, applying such an insight does not require the use of formal methods. Informal arguments suffice, as long as they are rooted in a precise semantical understanding. The approach to formal software engineering used in this paper is based on the algebraic “toolbox”, a brief overview is given by Ehrig *et al.* [1998]. This ranges from universal algebra, suited to investigate the language of the problem domain, to category theory, suited to discuss design principles for software architecture.

When designing computer programs for some problem domain, one is at least faced with three problems: (1) what are the concepts that have to be used for the construction of software in this domain, (2) what is a good programming notation for these concepts, and (3) how can they be implemented as software code. This has spawned work in areas like *domain specific embedded languages* (DSEL) [Hudak 1996] and *software architectures* [Bidoit *et al.* 1999]. A *domain specific language* (DSL) provides syntax for the concepts of a problem domain, supporting the expression of problems and solutions for that domain. A DSEL is the merging of a DSL with a general purpose programming language, ensuring full algorithmic and data structure declaration capabilities, but at the cost of constraints on the linguistic elements we may use in the DSL. Software architecture is the organisation of software into packages and modules and how to combine these in order to build specific software. Ideally software should be decomposed in such a way that different considerations may be confined to separate modules. Also, modules that provide alternative implementations of a concept should be easily interchangeable. This requires an extensive analysis of both the problem domain and of the software structuring methodologies available. Doing this well gives a very flexible software structure, easy to tailor for specific needs, as well as adapt to changing requirements – a characteristic of mature software design such as described by Racko [1995].

Domain specific languages are abundant in all application development tasks. Very often the DSL is “invisible”, in the sense that it was embedded already in the design of the programming language we are using. Then we often do not relate consciously towards which DSL is involved. In the 1960’s and 1970’s building domain concepts into a programming language was the common way of providing DSLs. With the advent of advanced modularisation concepts such as abstract data types and object orientation, DSELs may now be provided as separate libraries. Thus a programming language may be adapted to any domain by including the appropriate software libraries for that domain. Often a domain may have several distinct, but coexisting, domain specific languages. Then it becomes very important to be aware of this, and be able to make conscious decisions about which one to use. In the worst case, it may be beneficial to develop an alternative DSL – a DSL that may have better properties for expressing simple solutions to complex problems.

Our idea for a software process model for developing and implementing a DSEL is based on the following steps. As for any proper life cycle model it needs to be amended with additional feedback and control loops, similar to the way the original waterfall software process model by Royce [1970] has been amended, e.g., by Boehm [1976].

- (1) *Establish an appropriate DSL* by defining important properties of the problem domain concepts. Using algebraic specifications for this we identify sorts and functions to express application data and algorithms. It also nails down the semantics of the DSL.
- (2) *Validate the suggested DSL* by checking its
  - *usefulness*: formulate problems using the DSL, embed the DSL into a programming language and design solutions to the problems using the resulting DSEL. (This is mostly a test of the concepts themselves, but may also form part of an application development process, namely the functional requirements and high level design steps, that may be continued independently of the DSEL life cycle process.)
  - *implementability*: design data structures for the sorts and algorithms for the functions. The design should be verified by showing the relationship between the code and the specified concepts. (Note that this only will provide one of many possible designs, each design with different accuracy and resource usage characteristics – but then one such design is sufficient to show the implementability of the concepts, though it may not be the implementation we end up using.)
- (3) *Decide the architecture of the software library* that provides the DSL concepts. This requires grouping the sorts and functions for software components, studying their interplay and trying to parameterise the components as much as possible. The aim is to reduce the number of components and the software complexity of each by focusing on reusability. (*Software complexity* is a measure of the difficulty and cost of writing the software. It does not relate to the resource usage at runtime which is the subject of traditional complexity theory.) This generally increases the versatility of the components, and simplifies the library itself. It may also greatly reduce the work needed to implement the full library. Categorical reasoning is useful for this purpose.
- (4) *Design and implement a prototype of the software library* by coding the data structures and algorithms in a suitable programming language.
- (5) *Maintain the library* by adding variants of the components, as well as correcting (unavoidable) errors. Variants may be demanded by application program development in order to meet efficiency requirements.

The first of these steps identifies the concepts and finds a notation for these. Once the specifications are in place, it is possible to experiment with the notation and concepts to design programs solving actual problems. This allows early checking of the usefulness of the concepts. An adjustment of the

concepts and notation at this stage is comparatively cheap. We also have to check that the concepts are implementable before they are meaningful. Thus we have both outward (usability) and inward (implementability) requirements on a DSL in the validation of its specification. When specifying the domain concepts we develop a DSL, while we transform it to a DSEL when designing solutions.

We firmly believe that the development of a DSL and its implementation as a software library is a separate activity from application program development, and that they have different software life cycle models. An application software process model is typically based on the waterfall life cycle model of Royce [1970] with steps: requirements specification, design, validation of the design with respect to requirements, coding and maintenance. We emphasise the use of a DSL in the problem formulation (requirements specification) and the use of the DSEL in formulating and implementing the solution. While an application program is developed to solve a specific problem for one customer, a DSEL is to serve many application developers. Thus finding the right domain concepts is crucial, and devising a flexible library architecture is important, for the success of both the library development and the application development. Of course one will never in practice proceed along the sketched steps of the life cycle models in linear order. A more reflected understanding of a life cycle model will introduce all kinds of feedback loops, quality assurance steps, detailed guidelines on how to approach standard problems, etc. Such issues are investigated by the software process community, and is captured in models like that of Paulk *et al.* [1993]. The basic steps we have identified are still central in the enhanced models.

This paper is devoted to a case study developed according to the ideas above. The case study is taken from the area of computational modelling, in the typical form of a physical phenomenon described by a partial differential equation (PDE). Computational modelling of real world phenomena is becoming an important research tool in the sciences. Currently this is hampered by the time and effort needed to develop good computational models, and the time and cost needed to port such models onto a high performance computer. We ascribe much of these costs to the use of a less than optimal DSL for conventional numerical software, namely that of indexed array structures, as exemplified by Fortran-66 [1966]. A conventional solver for a PDE embodies the discretisation method, coordinate system, as well as the actual solver algorithm. Typically such software only handles the limited set of problems it was developed for, and can not easily be adapted to related phenomena or account for change of numerical discretisation methods. In many ways such a solver becomes a legacy code for the group that developed it: being their main tool for success, but also limiting what problems they may tackle.

Instead we want to develop an alternative DSL, namely that of *coordinate free numerics*, for this domain. This DSL captures more of the abstract concepts of the underlying mathematics, and may therefore blend more naturally with modern notions of software structure. This work started with the algebraic specification of concepts from partial differential equations by

Haveraaen *et al.* [1992]. Here we continue through the remaining steps of the the development process and we sketch the algebraic tools we have chosen to aid us through this. We also measure the effectiveness of this approach to software development by implementing a collection of application programs for the acoustic imaging of small and large scale geological objects, problems important in oil exploration. This problem is challenging, in that the same set of fundamental PDEs has to be formulated for various kinds of experimental setups, e.g., in cartesian or cylindrical (borehole) coordinates, and with various kinds of geophysical models and boundary conditions.

The choice of computational modelling as our domain was motivated in part by its general importance, but also since this is an area to a large extent neglected by the software methodology community. Showing the relevance of software engineering methodology to this area may then both open it up for more research from the software community, and may also benefit the practitioners in the computational modelling domain. The lack of familiarity with this area for most software engineers has made us include material describing more of the domain background than normal for case studies. We feel this may be needed, in order for the reader to appreciate the complexity of this kind of software and to gain an insight in why the proposed methodology is beneficial for an area which has been self-sufficient for much of the time since the early days of computing.

This paper is organised as follows: Section 2 sketches the use of algebraic techniques for investigating domain specific concepts and defining the coordinate free numerics DSL. Then we sketch the problems we will use as examples, both in order to validate the DSL, and in order to indicate the potential of this approach. In Section 4 we discuss the software architecture for the DSL and describe the Sophus software library which implements the DSL. Section 5 presents how the sample problems were solved, how a version of the DSL library adapted for the problems at hand were developed, and also presents additional variations of the problem and solutions to these. Then, in Section 6, we discuss the results achieved by following this approach. Finally we summarise our findings in Section 7.

## 2. Developing the domain specific language

There is a long tradition in looking at software as a formal entity. This spans from work in programming language semantics [Floyd 1967], via axiomatic formalisms [Hoare 1972] and systematic development techniques [Dijkstra 1976] to development by program refinement [Back 1981] and proof tools [Manna and Waldinger 1980]. One crucial observation is that a program text can be moved between computers (and compilers), but that the results computed may depend on the computer (a problem which is gradually reduced through standardisation). This leads to an acknowledgement of the distinction between syntax and semantics. Such a distinction also exists in the mathematical discipline of universal algebra, which originated with

Whitehead [1898] and had matured by the 1960's [Cohn 1965, Grätzer 1968]. In universal algebra the syntactic entities are called a *signature* and the semantics a *model* for that signature. It took till the late 1970's till universal algebra, in the form of many-sorted universal algebra, was employed in computer science [Goguen *et al.* 1975, Guttag and Horning 1978]. By that time category theory [MacLane 1971] was also making an influence in how to structure mathematical concepts. Now universal algebra and category theory is used together in the form of algebraic development methodologies [Ehrig *et al.* 1998].

### 2.1 Domain investigation tools: algebraic specifications

A *signature*  $\Sigma$  declares a set of *sort names*  $s_1, \dots, s_n$ , and function symbols  $f : s_{i_1}, \dots, s_{i_m} \rightarrow s_{i_{m+1}}$ , where  $s_{i_1}, s_{i_2}, \dots, s_{i_m}$  for  $m \geq 0$  are the argument sorts and  $s_{i_{m+1}}$  is the result sort (including the case  $m = 0$  for a constant). Obviously we may think of a sort as a type name or a class name in programming language terms. A function symbol corresponds to a side-effect free function or typed method. A procedure or method that changes its environment (has side-effects) can be decomposed into one or more side-effect free functions and explicit assignments to program variables. So we may treat a signature as idealised declarations in a program.

A *model*  $A$  for a signature  $\Sigma$  defines for each sort  $s$  of  $\Sigma$  a mathematical set  $A(s)$ , called the carrier, and for each function symbol  $f : s_{i_1}, \dots, s_{i_m} \rightarrow s_{i_{m+1}}$ , a mathematical function  $A(f) : A(s_{i_1}) \times \dots \times A(s_{i_m}) \rightarrow A(s_{i_{m+1}})$ . In a programming context we may let  $A(s)$  be a data structure, perhaps coupled with a data invariant, which defines the set of values that may be stored in the data structure, or the subset thereof that satisfies the data invariant. Likewise,  $A(f)$  may denote an algorithm, which defines a computable function from its argument data values to its result data values.

A specification restricts the class of allowable models for a signature. Algebraic specifications only focuses on the properties that we want satisfied, rather than devising specific constructions of models. Thus it is a rather abstract approach, but permits both mathematical models and programming language oriented models. The CASL specification language [Mosses 1997] is an attempt to bring together and standardise various approaches to algebraic specifications. Given a signature  $\Sigma$  and a specification, we may ask whether an implementation satisfies the specification. Morris [1973] has addressed this problem in a clear way and proposed it as a software development technique. Meyer [1991] presents this technique as *programming by contract*. Object-oriented programming languages have the modularisation mechanisms needed for programming by contract, but not all provide direct support for it.

### 2.2 Domain concepts for partial differential equations

The computational modelling domain is that of mathematics, typically the field of partial differential equations (PDEs).

If we start investigating the problem domain concepts using algebraic methods, we will of course rediscover the basic structures of algebra (the investigation and generalisation of which led Whitehead [1898] to the discovery of universal algebra): monoid, group, ring, field, vector space, linear mappings (matrices), tensors (which generalise rings, fields, vectors, matrices and (multi)linear mappings), etc. Haveraaen *et al.* [1992] specify many of these concepts in a start at analysing this problem domain using algebraic software methodologies. As an example, a ring  $R$  has binary operations  $+$  (addition),  $-$  (subtraction) and  $*$  (multiplication), and constants 0 (zero) and 1 (one). These form the ring signature,

$$\begin{aligned} + & : R, R \rightarrow R, \\ - & : R, R \rightarrow R, \\ * & : R, R \rightarrow R, \\ 0 & : \rightarrow R, \\ 1 & : \rightarrow R. \end{aligned}$$

This is a slight deviation from the presentation in mathematics, where 0 and 1 are perceived as elements of the carrier, rather than symbols of the interface, and the additive inverse is treated the same way. An algebraic specification of a ring could be the following, where  $a, b, c$  range over all ring elements  $R$ .

$$(a + b) + c = a + (b + c), \quad (1)$$

$$a + b = b + a, \quad (2)$$

$$(a * b) * c = a * (b * c), \quad (3)$$

$$(a + b) - b = a, \quad (4)$$

$$0 + a = a, \quad (5)$$

$$1 * a = a, \quad (6)$$

$$a * 1 = a, \quad (7)$$

$$(a + b) * c = (a * c) + (b * c), \quad (8)$$

$$a * (b + c) = (a * b) + (a * c). \quad (9)$$

Here we see that any model for a ring must obey the laws that addition is associative (1) and commutative (2), subtraction is the inverse of addition (4), and multiplication is associative (3) and distributes over addition in the familiar way (8–9). Further, the neutral element with respect to addition is denoted by 0 (5) and the neutral element with respect to multiplication by 1 (6–7). If a neutral element exists it is unique. Likewise the inverse element is unique if it exists. We may then ignore the operations 0, 1 and  $-$ , right up until the moment they, and their properties, are needed again. This simplifies the specification of, e.g., linearity of a function  $L : R \rightarrow R$  to

$$L(a + b) = L(a) + L(b), \quad (10)$$

$$L(a * b) = a * L(b), \quad (11)$$

because equations such as  $L(0) = 0$  and  $L(a - b) = L(a) - L(b)$  can easily be derived. Reducing the textual size of a specification by removing redundant details makes the important points clearer.

Continuing the analysis process for the realm of PDEs we note that a basic assumption is that every spatial point in the physical world, such as a 3-dimensional section of the earth, can be represented by an element of a set  $\mathcal{M}$  called a *manifold*. The physical properties are then ascribed to each point in the form of a *value field*, akin to a function from the manifold to some value domain. A value field has the same algebraic properties as the value domain itself. If the values at each point are reals, such as those for pressure or density, they are said to form a scalar field. A scalar field has ring properties. If they are vectors, such as those for particle displacement, they are said to be vector fields. A vector field likewise has vector properties, with the corresponding scalar being the scalar field. The values may also be matrices, linear or multi-linear mappings, or some other form of data, such as tensors, which generalise scalars, vectors, matrices and multilinear mappings. Tensor fields correspondingly are tensors over vector and scalar fields. If the manifold has sufficient structure, at least a notion of proximity and direction, we may define integration (interior and surface integrals) and differentiation operators on the value fields, such as Lie derivatives, gradients, and divergence, provided the value fields are smooth enough. A time dependent partial differential equation provides a relationship between spatial derivatives of tensor fields representing physical quantities and their time derivatives. Given constraints in the form of the values of the tensor fields at a specific instance in time together with boundary conditions, the aim of a PDE solver is to show how the physical system will evolve over time.

Haveraaen *et al.* [1992] sketched the algebraic specification of many of these concepts. Although they are standard mathematical concepts, they are rarely presented strictly as a signature with axioms. Doing a proper algebraic specification required work which became more involved as we approached the more advanced concepts, but no really hard problems occurred. However, a noticeable effect was that we were moving from the normal indexed-based presentation of the concepts, and in the direction of coordinate free mathematics, as we pursued identifying the concepts involved. This may be because our analysis tool, algebraic specifications, favours the use of high-level, abstract concepts. A coordinate free formulation is valid independent of the choice of coordinate system, i.e., it will be valid whether we use cartesian coordinates, cylindrical coordinates, or some other curvilinear coordinate system. This readily gives a much greater flexibility in choosing implementation strategies: some coordinate systems simplify the implementation of the operators and boundary conditions by introducing symmetries and vanishing terms (terms that will be 0 and thus may be eliminated from the computation), while other coordinate systems may improve numerical accuracy or reduce the amount of computation needed. From this it appears that coordinate free mathematics is a much more versatile framework than the



indexed-based one. An introduction to coordinate free mathematical physics can be found in [Schutz 1980].

### 2.3 Domain notation versus programming notation

Blending the coordinate free DSL into a programming language raises an observation about the relationship between DSLs and DSELs. Mathematics and the specifications' notation favour use of functional and operator style expressions. An object-oriented programming language like C++ [Stroustrup 1997] favours implementations using an imperative style, where argument variables are modified or mutated. As an example, consider an infix binary operator like `+`. C++ allows the declaration of an operator

```
template<class T> T + (const T & a, const T & b) const;
```

which is used in infix notation as `a+b`, returning the sum as a value. The operator preferred by the object-oriented style is

```
template<class T> void += (T & a, const T & b);
```

which is used infix as `a+=b`, but where the left argument is modified to contain the sum of the two variables. C++ itself does not have any intrinsic semantical restrictions on these operators, so there is no formal relationship between the two user-defined operations. (In essence, neither operator may have anything to do with summation at all.)

Since we are working with mathematical concepts it is natural to allow the user to write expressions involving `+`. The template class parameter `T` may represent several megabytes of data, so an implementation using `+=` will definitely be more efficient. We solved this by imposing requirements on the use of symbols such that there would be a clear semantical connection between functions/operators and their mutating counterparts, not just for the built-in operators in C++. A few rules expressing this connection is given in the following table.

<code>a = a + c;</code>	$\longleftrightarrow$	<code>a += c;</code>
<code>a = b + c;</code>	$\longleftrightarrow$	<code>a = b; a += c;</code>
<code>a = a * c;</code>	$\longleftrightarrow$	<code>a *= c;</code>
<code>a = b * c;</code>	$\longleftrightarrow$	<code>a = b; a *= c;</code>
<code>a = b * c + a;</code>	$\longleftrightarrow$	<code>{t = a; a = b * c; a += t;}</code>

These may be used to (1) generate an operator declaration for `+`, `*`, etc., whenever a declaration for `+=`, `*=`, etc. is seen, and (2) rewrite expressions using `+`, `*`, etc., to an equivalent code segment using only `=`, `+=`, `*=`, etc. This is described in more detail by Dinesh *et al.* [1998], who also compare the run-time efficiency of the two styles.

## 3. Validating the domain specific language

We need to check the usefulness of the DSL by making sure that we can use it to express (and solve) problems in the domain and that its concepts are

implementable. Only then will the DSL be valuable as a tool for developing programs. But implementability does not necessarily require exact compliance with the specification, in which case numerical programs would not exist since it is impossible to represent exact real arithmetic on any known computer architecture. Rather, we need a pragmatic notion of compliance, a notion that will accept programs that deliver useful results.

### 3.1 Usability: elastic wave simulation problems

The problems we will use to illustrate requirements of computational modelling are taken from the oil industry. Here the application of elastic wave modelling for the interpretation of various acoustic data from potential hydrocarbon reservoirs, is important. Repeatedly solving the corresponding PDEs for different data sets requires the use of high performance computers.

The recovery of the geological subsurface structure, i.e. the origin and the geometrical picture of the geological layers, is of vital importance in developing prospects of hydrocarbon reservoirs. Also, in order to obtain optimum production strategies for existing oil and gas fields, detailed acoustic illumination of the target zones is required. The target zones are generally quite heterogeneous where the geological and reservoir properties vary within a few metres. Thus, any modelling tool has to be capable of handling heterogeneous models in 2 and 3 dimensions with quite good resolution.

The seismic method is an active remote sensing technique, where the acoustic wave field, generated from a man made elastic impulse, is recorded on so called geophones at different spatial positions in the earth. The most common experimental setups are: marine seismics (source and geophones (hydrophones) are at the sea surface), well-to-well or surface-to-well seismics where the source is in a borehole or at the sea surface and the geophones are within a borehole, and ocean bottom seismics where the geophones are sited at the seabed. In all these experiments, the elastic impulses are of relatively low frequency, e.g. 5 – 100Hz. This energy illuminates geological structures. The picture arises from reflection and conversion of pressure (P) and shear waves (S) at layer interfaces, where the physical properties of the adjacent rocks are discontinuous. By combining the reflection data with estimates of the P or S wave velocity, a geometrical picture of the interfaces is produced. Normally a whole series of recordings are taken with small variations in the position of the source and the geophones.

For all these studies, the assistance of an acoustic wave modelling tool is of vital importance for the confidence of the information retained from the data. The aim is that acoustic wave simulations on the *estimated model* coincide with the acoustic recordings of the real world. This requires the acoustic simulation software to be able to place virtual geophones at any position within the simulated area. The simulation itself extends from sealevel and several thousands of metres downwards into the earth. A simulation which is repeated with small variations in the positioning of the virtual geophones and virtual source.

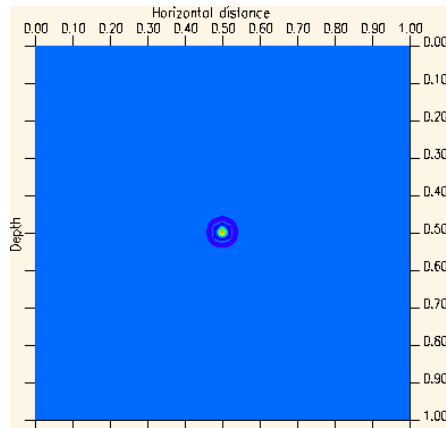
The mathematical equation describing the elastic wave simulation problem is taken to be the elastic wave equation,

$$\begin{aligned}\rho \frac{\partial^2 \vec{u}}{\partial t^2} &= \nabla \cdot \sigma + \vec{f}(t), \\ \sigma &= \Lambda(e), \\ e &= \mathcal{L}_{\vec{u}}(g).\end{aligned}\tag{12}$$

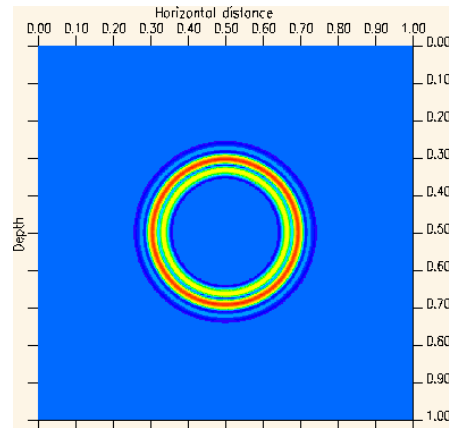
The scalar field density  $\rho$  and the stiffness tensor field  $\Lambda$  are given data that vary within the physical domain, in accordance with the varying geophysical properties of the rocks. The particle displacement vector field  $\vec{u}$  represents the propagation of the seismic wave and will be recomputed at every iteration of the solver algorithm. The tensor field  $g$  defines the coordinate system used, the tensor fields  $\sigma$  and  $e$  are computed intermediate values, and  $\vec{f}(t)$  is a time-varying vector field representing the forces from the elastic impulse. The  $\nabla \cdot$  and  $\mathcal{L}_{\vec{u}}$  are derivation operators, the latter dependent on the displacement  $\vec{u}$ . A fundamental assumption here is that the materials are fully elastic. The elastic wave equation is a standard equation from mathematical physics, and may be found in any textbook on the subject. Its application to seismics is discussed in much detail by Aki and Richards [1980] and Marsden and Hughes [1983]. Eq. (12) is in coordinate free form, i.e., all the entities and operators belong to the coordinate free DSL, validating its usability.

Elastic wave simulation is a very compute-intensive task, where one simulation easily may take several hours. Important factors here are the spatial resolution and frequency of the source. Increased accuracy requires a larger data set, more simulation steps, and consequently increased computation time. Another factor is the complexity of the physical properties of the geological models. In the simplest case, the model is denoted as isotropic inferring that the P and S wave speeds are independent of the wave propagation direction. The isotropy introduces symmetries in the stiffness tensor  $\Lambda$  of Eq. (12), so that the amount of computation can be greatly reduced. An earth model of more general elastic properties, denoted as anisotropic, implies that the wave velocities do generally depend on the wave propagation direction. However, most earth materials have a rotational symmetry in the stiffness properties. In general, this is for an axis perpendicular to the internal layering of the material, but often restricted to be perpendicular to the layer surface. The latter materials are denoted as transverse isotropic, and more symmetries are introduced in the stiffness tensor, but not as much as in the isotropic case.

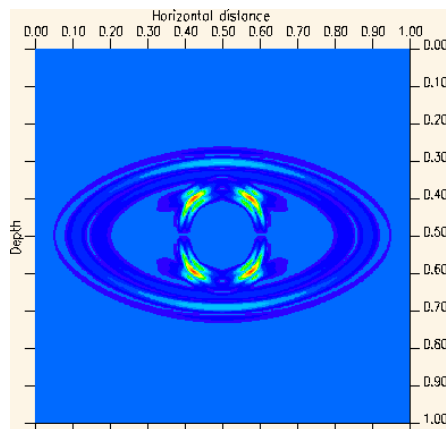
Fig. 1–4, produced by SeisMod (see Section 5), illustrate some examples of elastic waves propagating outward from a point source (producing an elastic impulse) located within various types of materials, taken at different times after the impulse was initiated. The propagation medium is here 2-dimensional, 1km by 1km and homogeneous (i.e. no alteration in the stiffness tensor within the grid). The elastic impulse is described by a Ricker pulse of 30Hz centre-frequency. For the isotropic case, we consider a source in water where no S wave exists. The other simulations are based on media



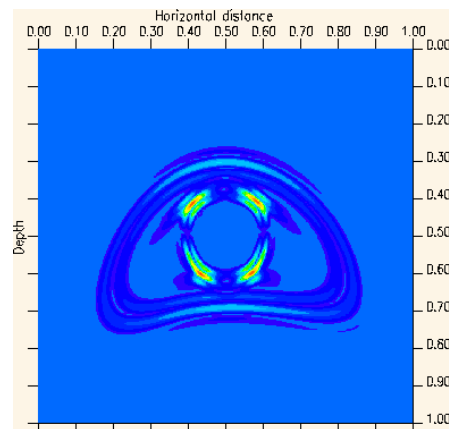
**Fig. 1:** Isotropic case, seismic waves 25ms after the pulse started.



**Fig. 2:** Isotropic case, seismic waves 100ms after the pulse started.



**Fig. 3:** Transverse isotropic case with vertical axis of symmetry, seismic waves 100ms after the pulse started.



**Fig. 4:** Transverse isotropic case with the axis of symmetry curved in the rock, seismic waves 100ms after the pulse started.

with the same density and similar wave propagation velocities. In case of a transverse isotropic material, we clearly see the effect of the directional variation in wave speeds (in particular for the P wave), and a separation between the faster moving P wave and the slower moving S wave. The final example shows a wave field occurring within a strongly anisotropic medium. Here the wavefronts are strongly deformed from the circular shape associated with isotropic media, or the more ellipse-like shape associated with transverse isotropic media.

As the computation time, which may be on the order of several hours, increases by a factor of more than two from the isotropic to the anisotropic case, there is a clear demand for program versions being as specific as possible for the different cases. Also, there is a need for each program version to

run sequentially on single workstations, in parallel on a network of such workstations, and on supercomputers with multiple processors.

### 3.2 Implementability

There are a host of numerical methods that provide implementation strategies for scalar fields, but neither provides an exact representation of a scalar field. A discretisation method will only provide a more or less inaccurate representation of a scalar field, but an approximation that given some assumptions is close enough to provide useful results. The more well known approaches are finite difference methods, finite element methods, finite volume methods, and spectral methods. The discretisation methods vary in software complexity both at the level of scalar fields and the level of equation solvers (see Section 4.3). The choice of numerical discretisation method depends to a large extent on the properties of the PDE to be solved and thus on the applications that are to be developed. Generally, finite difference methods are among the simplest to implement, with finite element methods being among the most difficult with respect to implementation and use.

The elastic wave equation is a good equation to start with from this point of view, as it works well with the finite difference discretisations. This means that we may initially avoid unnecessary complications both when formulating the solver and when implementing the scalar fields – nice properties for a prototype implementation.

## 4. Software architecture and DSL structure

### 4.1 Structuring concepts

A collection of related mathematical structures, such as the data structures of a programming language, typically form a category [Goguen 1991]. A *category*  $\mathbf{C}$  is a collection of objects  $A, B, \dots$ , and morphisms  $f : X \rightarrow Y$ , with an associated associative composition rule  $\circ$  on morphisms and a neutral morphism (with respect to  $\circ$ ) for each object. We will use the category **Prog** as our prime example. The objects of **Prog** are data structures, and the morphisms are all side-effect free algorithms from a data structure to a data structure. The identity morphism and composition rule for morphisms should be obvious. Functions of more than one argument are defined from special data structure objects called *product objects* in the category. The category **Set** is a standard example from mathematics of a category. It has sets as objects and total functions between sets as morphisms.

Categories are related by *functors*, functions between categories. A functor  $F : \mathbf{C} \rightarrow \mathbf{D}$ , from category  $\mathbf{C}$  to category  $\mathbf{D}$ , maps objects to objects and morphisms to morphisms such that identities and compositions are preserved. It is not hard to construct a functor from **Prog** to **Set** that relates the data structures and algorithms (for one specific computer) with the sets of values and mathematical functions being computed.

Functors are in many ways like C++ template classes [Stroustrup 1997] or Ada generic packages [Barstow 1983]. These mechanisms will take a data type as argument and define a new data type based on it. We may for instance define a generic `list` package with a type parameter, such that whenever we instantiate the package with a data structure  $D$ , we get a data structure `list of  $D$` . The idealised functor version of data type constructors have some additional properties. A list data constructing functor  $L : \mathbf{Prog} \rightarrow \mathbf{Prog}$  takes a data structure  $D$  and returns a `list of  $D$`  data structure  $L(D)$ . But in addition to defining the list data type, it will take any function  $f : D \rightarrow E$  and define an iterated function  $L(f) : L(D) \rightarrow L(E)$ . When  $L(f)$  is given a list of  $D$  as argument it will perform  $f$  on every element of the list, returning a `list of  $E$`  with the results. Likewise we may treat array data structure constructors as functors. For every index type  $I$  we have an array constructing functor  $A_I : \mathbf{Prog} \rightarrow \mathbf{Prog}$  which takes an element object  $E$  and defines an `array [ $I$ ] of  $E$` , the array structure with elements of type  $E$ . But we also get the iterated functions, so given for instance a binary operation  $+: E \times E \rightarrow E$  we have  $A_I(+): A_I(E) \times A_I(E) \rightarrow A_I(E)$  which adds, componentwise, i.e., for each index  $i \in I$ , the elements of the two argument arrays, yielding a new array with the summed values. This is very convenient, and may only be simulated by explicit programming of these functions in current programming languages. Unfortunately, this is not fully sufficient, as the generic package mechanisms do not have enough power to let us do this once and for all. (We omit the technical discussion of these deficiencies.) A nice observation is that the array constructing functor can be used to generate the value fields for a manifold  $\mathcal{M}$  by simply applying  $A_{\mathcal{M}}$  to the appropriate value domain, such as the reals, vectors or matrices. We can also use these functors to define finite dimensional vector spaces by the expression  $A_{\{1, \dots, n\}}(\mathcal{R})$ , for appropriate natural numbers  $n$  and the ring  $\mathcal{R}$  of real numbers.

A good modularisation of software is achieved if we minimise the number of distinct functors, and the software complexity of each, needed to build the application software. How we combine the functors to achieve these entities will be a blue-print for the software architecture. Good choices here have large potential for greatly reducing the software development effort. Both by directly reducing our coding effort, and, more importantly, by identifying reusable components for other applications in related problem domains. Carefully structuring the modules reduces our coding effort and reusable software components are identified. This is work at the software architecture level. In the algebraic specification language CASL, for instance, this kind of software architecture can be explicitly defined [Bidoit *et al.* 1999].

#### 4.2 Structuring the problem domain

Now we need to analyse the problem domain in order to find the concepts that we can use in structuring our software. Our approach to formal software engineering is to use algebraic techniques, specifically category theory, for this.

When developing the DSL, Section 2.2, we pointed out that a scalar field has ring properties. A vector field is a value field with vectors, such that the vector fields form a vector space with the scalar field as the ring. As a consequence,  $n$ -dimensional vector fields over a manifold  $\mathcal{M}$  may be constructed by either of the two approaches:

- (1) applying the value domain construction to vectors,  $A_{\mathcal{M}}(A_{\{1,\dots,n\}}(\mathcal{R}))$ ,  
or
- (2) applying the vector construction to scalar fields,  $A_{\{1,\dots,n\}}(A_{\mathcal{M}}(\mathcal{R}))$ ,

and similarly for tensor fields. There does not seem to be any immediate reason to prefer one over the other, and conventional numerical software uses the first construction. However, a closer scrutiny of the problem domain reveals that a tensor field contains advanced integration and derivation operations which are not definable from the tensor abstraction, but requires access to the value field properties, i.e., to the discretisation.

- Applying the value domain construction  $A_{\mathcal{M}}$  to vectors ( $A_{\{1,\dots,n\}}(\mathcal{R})$ ) as in construction (1) makes us rebuild the discretisation for every level of construction, i.e., one for scalar fields, another for tensor fields etc. This was observed in the tensor oriented implementation reported by Verner *et al.* [1993].
- But the integration and differentiation operators of the vector field  $A_{\{1,\dots,n\}}(R)$  for arbitrary scalar field  $R$  in construction (2) may be expressed from operations like integration and partial derivatives on the scalar fields  $R = A_{\mathcal{M}}(\mathcal{R})$ .

This reveals that apparently equivalent constructions from a data structure and algorithmic complexity viewpoint, may have dramatically different software complexity. Based on these observations it is clear that a more fruitful approach is to use the second construction above as starting point. Instead of building many different constructors for the value domains (vectors, matrices, linear mappings, etc.), we note that it suffices to build a tensor constructor, which, given certain assumptions, encompasses all these. Tensors also give us the building blocks needed to define coordinate free operators. The implementation may then be reduced to build a constructor for scalar fields and a constructor for tensor fields.

The functor  $S_{\mathcal{M}} : \mathbf{Prog} \rightarrow \mathbf{Prog}$  for the construction of scalar fields may be implemented by amending the construction  $A_{\mathcal{M}}$  such that it also includes the definition of integration and partial differential operators. The tensor constructor  $T_{\{1,\dots,n\}} : \mathbf{Prog} \rightarrow \mathbf{Prog}$  amends  $A_{\{1,\dots,n\}}$  with the integration, general derivation operators and other tensor operators, assuming that the template parameter has an appropriate interface. The tensor field construction for a manifold  $\mathcal{M}$  then becomes  $T_{\{1,\dots,n\}}(S_{\mathcal{M}}(\mathcal{R}))$ . If the tensor exhibits symmetries, we may be able to use a constructor  $A_{\{1,\dots,m\}}$ , where  $m \leq n$ , as data structure.

### 4.3 Software architecture for PDE problems

Our analysis of the PDE domain has given us the components we need for the software architecture. It has also provided us with a problem domain specific language and specifications of the types and operations that we need. But we also have to make sure that this really can work as a framework for implementing numerical methods. In this analysis we need to consider the issue of developing both sequential and parallel versions of the software.

The different components involved in a numerical solution of a PDE can be factored in three layers:

- (1) The numerical discretisation methods which makes it possible to represent the value fields  $S_{\mathcal{M}}$  for the infinite set  $\mathcal{M}$  by a finite approximation. The discretisation will need to provide the ring operations and the partial differentiation operations.  $S_{\mathcal{M}}$  will need to contain a large set of values, often up to a million or more, in order to provide a good approximation. To represent the data values on  $\mathcal{M}$  it is convenient to use the functor  $A_I$ , for some suitable index set  $I$ , as data structure.
- (2) The tensor construction  $T_{\{1,\dots,n\}}$  is where coordinate systems are handled and the advanced differentiation operations are implemented. The construction should work with any scalar field, i.e., with any discretisation of a scalar field as well. Conceptually  $T_{\{1,\dots,n\}}$  is an extension of  $A_{\{1,\dots,m\}}$ , for some  $m \leq n$ , and using the latter as the data structure seems natural.
- (3) The uppermost numerical layer is the solver algorithm itself. Here the time discretisation is decided, and the iterator that will generate the transient behaviour (such as for a seismic simulator where we are interested to know how the seismic wave propagates) or the steady state solution (for instance if we want to find a steady state flow pattern) is implemented. These algorithms are often normalised, i.e., the numbers they work with are scaled to be around 1.0, where the numerical resolution of the machine is best. The numbers are then scaled back for input/output purposes and as needed by the solver algorithm.

Using the array constructor to implement both the numerical discretisation and the tensor construction allows for a reuse of the array module. But more importantly it allows a separation of concerns when implementing these modules: the array constructor may focus on the data layout pattern, while the numerical modules may focus on the numerical aspects, using the array construction for the storage aspects. The architecture also implies that we only need to relate to, and thus implement, the discretisation method when we implement the scalar field, and that the vector and tensor field implementations are independent of this choice. If we need to change discretisation method, this will be localised to one module, and not being spread out all over the code, which is the normal case with traditional numerical software.

This also provides a route to parallelisation. We will, at the scalar field level at least, have a large collection of data values that may be distributed in



a dataparallel fashion [Bougé 1996]. Actually, it suffices to provide a parallel implementation of the array constructor to get a parallel version of the whole program. Haveraaen [1998] discusses this.

#### 4.4 The Sophus library

The software architecture developed in the previous section is implemented by the Sophus software library. It provides the abstract mathematical concepts from PDE theory as programming entities. Its concepts are based on the notions of manifold, scalar field and tensor field, while the implementations are based on the conventional numerical algorithms and discretisations. Sophus is structured around the following concepts:

- Basic  $n$ -dimensional mesh structures  $M_n : \mathbf{Prog} \rightarrow \mathbf{Prog}$  taking a ring  $R$  as argument. A mesh is an array constructor  $A_{\{1, \dots, k_1\} \times \dots \times \{1, \dots, k_n\}}$ , and includes the definition of the general iterated operations. Specifically, operations like  $+$ ,  $-$  and  $*$  are iterated over all elements (much like Fortran-90 array operators [Adams *et al.* 1992]), and operations to add, subtract and multiply all elements of the mesh by a scalar are included. There are also operations for shifting meshes in one or more dimensions. Operations like multidimensional matrix multiplication and equation solvers may easily be implemented for the meshes. Sparse meshes, i.e., meshes where most of the elements are 0 or have some other fixed value, may also be provided. Parallel and sequential implementations of mesh structures can be used interchangeably, allowing easy porting between computer architectures of any program built on top of the mesh abstraction.
- Manifolds  $\mathcal{M}$ . These define sets with a notion of proximity and direction which represent the physical space where the problem to be solved takes place.
- Scalar fields  $S_{\mathcal{M}}$ . They describe the measurable quantities of the physical problem to be solved. As the basic layer of “continuous mathematics” in the library, they provide the partial derivation and integration operations. Also, two scalar fields on the same manifold may be pointwise added, subtracted and multiplied. The different discretisation methods, such as finite difference, finite element and finite volume methods, provide different designs for the implementation of scalar fields. Scalar fields are typically implemented using the basic mesh structures for data.
- Tensors  $T_{\{1, \dots, n\}}$ . These provide coordinate free mathematics based on the knowledge of the coordinate system, whether it is cartesian, axisymmetric or general curvilinear. The tensor module provides the advanced differentiation and integration operations, based on the partial derivatives and integrals of the scalar fields. Tensors also provide operations such as componentwise addition, subtraction and multiplication, as well as tensor product, composition and application. The

implementation uses the basic mesh structures, with scalar fields as the ring parameter.

Using Sophus, the numerical equation solvers are formulated on top of the coordinate-free layer, forming an abstract, high level program for the solution of the problem.

## 5. Implementations: DSL library and elastic wave simulators

In Section 4 we have investigated the structure of the DSL and decided what basic components it consists of and how the DSL can be built from these components. We also studied, in Section 3.2, the implementability of the DSL, but postponed the decision about which discretisation method to use till we had an application that would require the use of a specific method. The usefulness of the DSL was validated in Section 3.1 where we formulated the elastic wave problem.

Here we will start the design of SeisMod, a collection of elastic wave simulators, and use this to decide which variant of the library modules to implement. The piecewise implementation of the library this implies should be seen as a normal way of maintaining the library, making it more complete and versatile as time goes. But such an incremental building will only work out if the library architecture is well thought out, i.e., mature, so that a reorganisation of the library will not be required as the library is gradually built.

The elastic wave simulators we are using as our example provide some additional complexity for the solver design. This is due to the boundary conditions of the problem. Since we must do our simulations in a finite computational domain, the seismic waves should ideally leave the domain when reaching the boundaries. This can not be achieved without special numerical treatment. The upper boundary (for instance the sea/air border) must also be handled specially in order to obtain the proper physical behaviour. These problems are solved by implementing a boundary-aware scalar field on top of the plain scalar field discretisation. This boundary aware scalar field is then used when instantiating the tensors.

This leaves us with the following main software modules for the seismic simulators:

- **Mesh**: domain information, index set, and the basic  $n$ -dimensional meshes in sequential and parallel version.
- **Tn**: domain information, manifold, and  $n$ -dimensional toroidal scalar field (no boundaries).
- **Bn**: domain information, manifold, and  $n$ -dimensional scalar field with boundaries.
- **Tensor**: tensors with differentiation operators.
- **Seismod**: the solver algorithm in coordinate free form.

The following subsection describes a basic collection of solver programs for the seismic problem. Then we describe variations of the seismic problem which exercise the flexibility of the software architecture developed above. All the variations we develop are captured in Table I, which identifies 32 versions of elastic wave simulators.

TABLE I: Some configurations for SeisMod.

Configuration	seismic	ultrasonic in borehole
<b>Mesh</b>	D S or P	D S or P
<b>Tn</b>	D S	D S
<b>Bn</b>	D S or U	D U
<b>Tensor</b>	SI or TI or TA	UI or UTI
<b>Seismod</b>	SE or PE	SE or PE

Legend:

- D is domain information, such as shape and indices or elements of a manifold.
- **Mesh** S is sequential implementation.
- **Mesh** P is parallel implementation.
- **Tn** S is staggered finite difference implementation.
- **Bn** S is sea surface boundary handler, i.e., the scalar field extends all the way up to the sea surface.
- **Bn** U is underground boundary handler, i.e., the scalar field is surrounded by rock.
- **Tensor** SI is standard isotropic.
- **Tensor** TI is transverse isotropic with vertical axis of symmetry.
- **Tensor** TA is transverse isotropic with arbitrary axis of symmetry (very close to full anisotropic).
- **Tensor** UI is ultrasonic for borehole with isotropy (cylindrical coordinates).
- **Tensor** UTI is ultrasonic for borehole with transverse isotropy (cylindrical coordinates).
- **Seismod** SE is standard elastic.
- **Seismod** PE is poro-elastic.

### 5.1 A basic collection of simulators

The original elastic wave simulator problem, see Section 3.1, was described for a general 3-dimensional earth with any complexity and heterogeneity in the physical properties (density and stiffness) within this model. In full generality, this should infer a stiffness tensor  $\Lambda$  of Eq. (12) at each grid point to have  $3^4 = 81$  components. However, the physical constraints implied by considering a linear and infinitesimal deformation theory, cause the stiffness tensor to contain at most 21 different components. The computational effort may in some cases be reduced by using 2D models. Such models are good approximations when the earth models vary only slowly perpendicular to the considered 2D plane. Implementing this, using cartesian coordinates and for the isotropic simplification, reduces the number of distinct components in

any tensor involved in the computation to at most 3. Even if we assume the more general anisotropic case in 3D, we never need more than 7 distinct components. But this reduction of storage and computations comes at a price: all the algorithms in `Tensor` must be specially adapted to the combined choices. So we need three distinct `Tensor` implementations, one for each of isotropic (SI), transverse isotropic (TI) (assuming vertical axis of symmetry), and anisotropic (TA) cases.

Using a sequential `Mesh` in the implementation of `Tn` gives a sequential implementation of the seismic simulators which will run on any machine with a decent compiler. We may then exchange this with a parallel `Mesh` to achieve a parallel version of the simulators. Implementing `Mesh` using MPI [Snir *et al.* 1996] gives us a code that can run on any machine configuration supporting MPI, such as a network of workstations or a multiprocessor.

The seismic simulator has been implemented for all the three rock model complexities both for sequential execution and parallel execution on a network of workstations and on a multiprocessor supercomputer. This provides a total of 6 source versions for these different machine architectures: `Mesh S` or `P`, `Bn S`, `Tensor SI`, `TI` or `TA`, and `Seismod SE`.

### 5.2 Modelling of ultrasonic measurements in boreholes

While the typical seismic experiments reveal earth models several kilometres wide and deep, acoustic measurements in boreholes are performed to retrieve details of the material properties in the close vicinity of the well. Here both high-frequency P, S and surface waves are studied in order to depict the zones to be carefully handled during the production and injection phase. The near borehole reservoir properties are important both for the estimate and production of the hydrocarbon fluids. In an analogous way to seismics, these properties are sought for by ultrasonic acoustic experiments. In this, a source and an array of transducers are placed within the borehole. Here acoustic waves within the range of 5 – 20kHz are considered. For these frequencies the acoustic waves illuminate only some few decimetres of the surrounding formation. In these experiments P and S waves and boundary waves occurring in the borehole wall are considered.

From a modelling viewpoint, the problem is simplified if we consider the earth model to be symmetric around the borehole, i.e., it is axisymmetric with the well bore itself as the axis of symmetry. This is very useful since then the 3D problem becomes a 2D problem using cylindrical coordinates, resulting in a corresponding reduction of CPU-time. This implies the need to create a new implementation of `Tensor`. Unfortunately, since the geological model complexity also influences the implementation of `Tensor`, we need to create different versions of it, if we are to cater for the various geological models (isotropy, transverse isotropy, anisotropy). However, two versions is enough, since we in the case of transverse isotropy, only need to consider an axis of symmetry parallel with the borehole.

Furthermore, we here also have to take into account that the boundary conditions in the borehole are different from those on the surface. So we need a new `Bn` with open boundaries in all directions.

Thus we need 3 modified implementations to cater for the requirements of the borehole problem (`Bn U` and `Tensor UI` and `UTI`). But the changes are not radical from the original versions, and more than 50% of the code of the modules we replace can be reused in the creation of these modules. A fact which in some cases may be made explicit by the use of inheritance. The changes do not affect any other modules, nor how they should be fitted together. We also get similar parallel and sequential versions without any further implementation effort, giving us, all in all, 4 new seismic simulators.

### *5.3 Modelling of seismic tomography experiments*

The variants of modules we have defined allow us to configure yet another version of `SeisMod`, a system for seismic tomography simulation. This is when an elastic impulse is initiated at the surface or within a borehole, and the geophones are placed in surrounding boreholes. The modelling is here to compute the purely transmitted wavefield. We already covered the case when the tomography extends to the surface in our original simulators. But if the target area is entirely subsurface, we only have to apply the underground boundary handler `Bn U`. This may be combined with any of the `Tensor` versions `SI`, `TI` or `TA`, and of course with a sequential or parallel mesh implementation, increasing the number of specific simulators from 6 seismic and 4 ultrasonic to 12 seismic and 4 ultrasonic.

### *5.4 Modelling of wave propagation in poro-elastic materials*

The hydrocarbons we want to investigate are stored in porous rocks. So the elastic waves will propagate in two-phase materials, or poro-elastic materials. In particular, the use of a poro-elastic modelling may describe several wave phenomena observed in ultrasonic borehole data. Typical poro-elastic materials are high porosity rocks containing viscous fluids. In such materials the particle displacements, associated to an induced stress field, of the solid skeleton and the fluid may differ. This physical system is described by two coupled equations of motion, which, accordingly, have to be solved simultaneously. The pioneering, mostly analytical work of Biot [1955, 1956a, 1956b] reveals the existence of two P waves and one S wave. The two P waves are usually denoted the fast and the slow P wave, where the fast P wave propagates mainly within the skeleton, but modified by the presence of the fluid. The slow P wave, accordingly, propagates mainly along the fluid phase and is modified by the skeleton. The fast P wave carries the main P wave energy, while the slow wave is mainly of phenomenological character.

The modelling of poro-elastic materials is receiving increased attention with the demand for a more precise elastic wave simulation in reservoir zones. The coupled set of equations to be considered are

$$\begin{aligned}
\rho_{11} \frac{\partial^2 \vec{u}}{\partial t^2} + \rho_{12} \frac{\partial^2 \vec{U}}{\partial t^2} &= \nabla \cdot \sigma + (1 - q) \vec{f}(t) + b \left( \frac{\partial \vec{U}}{\partial t} - \frac{\partial \vec{u}}{\partial t} \right), \\
\rho_{12} \frac{\partial^2 \vec{u}}{\partial t^2} + \rho_{22} \frac{\partial^2 \vec{U}}{\partial t^2} &= \nabla s + q \vec{f}(t) - b \left( \frac{\partial \vec{U}}{\partial t} - \frac{\partial \vec{u}}{\partial t} \right), \\
\sigma &= \Lambda(e) + Q(\epsilon), \\
s &= M(e) + R(\epsilon), \\
e &= \mathcal{L}_{\vec{u}}(g), \\
\epsilon &= \nabla \cdot \vec{U}.
\end{aligned} \tag{13}$$

Here the displacement vector fields are  $\vec{u}$  for the solid and  $\vec{U}$  for the fluid. Likewise the tensor fields  $\sigma$  and  $e$  of Eq. (12) get companions  $s$  and  $\epsilon$ , respectively. Further the stiffness tensor field  $\Lambda$  is split into four tensor fields  $\Lambda$ ,  $Q$ ,  $M$  and  $R$ , representing the stiffness of the two phases and their interaction. We also get three density coefficients  $\rho_{ij}$  and a friction tensor field  $b$ , as well as a porosity  $q$ . The operators remain basically the same, but we get an additional derivation operation  $\nabla$ . Still, the formulation in Eq. (13) is in coordinate free form, as required for our DSL.

Even though these equations are very different from Eq. (12), we see that all the concepts involved are tensor fields and tensor field operations. The new equations initiate a new implementation of `Seismod`, the PE version. In principle *no other module* needs modification, and this would be the case if we had used a fully general tensor field class implementation. However, as we decided to implement specialised tensor class versions for SI, TI, TA, UI and UTI, our incremental implementation policy means that these may not cater for the requirements on data formats and algorithms from the new tensor fields and tensor field operations. These additions will mostly be the same across versions, but not quite, implying that we should analyse the situation with specialised implementation versions further. With these amendments, poro-elastic materials may now be embedded in geological models for either large-scale seismics or ultrasonic experiments in boreholes, doubling the number of elastic wave simulators from 16 to 32.

## 6. Results achieved

### 6.1 Program development productivity

In the previous section we have shown that the Sophus PDE software architecture has made it easy to pin-point exactly which module needs to be modified when requirements are changed. Besides the obvious intuitive benefits from this, there are definite quantitative benefits as well. This can be seen when considering various cost estimating techniques, of which the COCOMO model by Boehm [1981] is the most well known. Essentially, a model like COCOMO correlates the cost of developing a piece of software with the size of the software. More advanced estimation techniques take into account the software complexity as well. Experience, as captured by the cost models, shows that cost grows more than linearly in program size, both

for development and maintenance. Thus halving the code that needs to be modified in order to meet new requirements, more than halves the cost of the modification.

Table II relates different modules in the Sophus library to the size of the code that implements them. The versions that can be generated were discussed in the previous section, and summed up in Table I. The total number of lines has been tallied under four sets of columns. The first tally represents the lines of code needed for the standard seismic simulator (cartesian coordinates, isotropic rock model, sequential), disregarding other support modules and configuration files. This tally will be used as the reference cost of developing a seismic simulator. The second tally gives the total number of lines needed to create the 5 additional simulator versions described by the original requirements specification in Section 3.1. The third tally gives the total number of lines for the 4 versions of ultrasonic, axisymmetric bore-hole acoustics. This also gives us 6 versions of underground tomography, giving a total of 16 versions of elastic wave simulators. The last tally, for the PE column, gives the cost of the poro-elastic version. In the tensor row for this column we have included the total number of lines to be amended to the tensor classes due to the new tensor operations and data required by poro-elasticity. Note that the size of the variants in general gives a too high estimate of the cost for creating a version, since much of the code from the original version of a module can be reused.

TABLE II: Indicative sizes of the modules in lines of code. The legend is given in Table I.

Module	D	S/SI/SE	P	TI	TA	U/UI	UTI	PE
Mesh	2000	2000	2700					
Tn	1700	1600						
Bn	1500	1900				2100		
Tensor		1000		1000	1800	1000	1100	1200
Seismod		600						700
Total:	12300		5500			4200		1900

Investigating some specific example configurations, we see that creating a parallel version of `Mesh` only costs two thirds of that of implementing a full `Mesh`, since we do not need to reimplement the domain and index types (column D). The cost of the parallel code represents redeveloping less than 25% of the full program. Once the parallel module has been developed, it can be freely mixed with any version of the program, yielding a parallel version of that program for no additional cost.

Shifting from the isotropic version to one of the anisotropic version represents a redevelopment cost of 8%–15% of the software cost, as can be seen by the changes needed, which only affects `Tensor`. Achieving parallel versions of the anisotropic codes now comes for free. So the development of all the 6

versions asked for in Section 3.1 adds less than 50% to the cost of developing the first version.

Creating the borehole version of the seismic simulator means we have to upgrade both the `Bn` and the `Tensor` modules. The former due to the change of boundary conditions, the latter due to the use of cylindrical coordinates to improve efficiency. These changes represent a total cost of less than 35% of the reference program code. In fact the modifications involve much less than 50% of the 4200 lines listed, bringing the cost of the borehole version down to about 15% of the development costs. We get a parallel borehole version for free. More interestingly, the revised `Bn` implementation may be combined with the other tensor classes to yield a underground tomography simulator. This adds 6 new seismic simulator versions for free, in addition to the 4 borehole versions we consciously developed.

This gives us a tally of 16 versions of the seismic simulator, counting sequential and parallel versions. The total cost of the 15 extra versions developed represent less than 80% of the development cost of the initial simulator. If we amortise the total cost of developing all the 16 versions, we end up with a development cost of each in the range of 10%–15% of the development of the initial program. Recognising that the number of versions will double when coding the poro-elastic version (the last column), we see that the average development cost of a version drops to less than 10% of the initial development costs. Assuming the existence of a fully developed library, we see that the marginal cost of writing a new program may easily become less than 10% of the development cost of a full program. The most notable case is when we may obtain 16 poro-elastic versions at a total cost about 15% of the cost of the reference program, even though the change in the mathematical equations are large. This is a tremendous gain in software development productivity over traditional development. A clear demonstration of the benefits of software reuse, contingent on a flexible and robust software architecture, i.e., a mature software design.

## 6.2 Discussion

It can be argued that all of these examples are taken from the same small area within PDE. Thus we have only shown that these gains are achievable when many closely related versions of a program are to be developed. This may be the case, but Grant *et al.* [1998] applied the Sophus software architecture to solve computational fluid dynamics problems. In their case study, many new implementations were needed. Almost all were instances of the Sophus modules already identified here, and none required a modification of the library architecture. This supports our view that the Sophus software architecture is mature. In general we see that a new discretisation method will obviously require a new scalar field implementation, new equations to be solved require new top level solvers, etc, but all within the same framework.

Comparing this with experience from other numerical libraries for solving PDEs, e.g., reports by Budge *et al.* [1992] and Verner *et al.* [1993], we note



that Sophus seems to avoid many of the difficulties and added costs that often appear when extending the domain of a library. Typical problems that have been identified by library designers include inappropriate abstractions and inappropriate architectures. In the problem domain we have presented here, an inadequate abstraction could be a coordinate system dependent definition of the operators. Then Eq. (12) and Eq. (13) would need to be changed according to what coordinate system was to be used (this is the normal case in current computational models). The choice of an inappropriate architecture shows when using the construction  $A_{\mathcal{M}}(A_{\{1,\dots,n\}}(\mathcal{R}))$  directly as data structure for vector fields. Then both scalar fields and tensor fields embody the discretisation method, meaning that discretisation code is duplicated in the library. When expanding the library with another discretisation method most of the library has to be reprogrammed for the new method. Budge *et al.* [1992] report that even user programs may have to be rewritten for an improved library.

Our case study of the Sophus software architecture for PDEs does not show these symptoms. The Sophus library is designed to mimic the abstract structure of the mathematics of partial differential equations, as used in the description of many natural and industrial phenomena. Its apparatus supports coordinate free numerics in the formulation and development of solvers to the problems. By requiring strict adherence to specified interfaces, we have been able to achieve that different implementations of the same mathematical concepts are basically interchangeable. The Sophus library components can be organised in different layers of abstractions. The interchangeability of modules within a layer allows software developers to experiment with different solution strategies and high performance computer architectures without the need for extensive reprogramming. For example, the change from a sequential to a parallel version of a seismic simulator does not involve any reprogramming of the rest of the solver application, just a small change in the configuration. Currently the Sophus library and application software is implemented using C++.

## 7. Summary

We have proposed a methodology that focuses on developing domain specific languages (DSL) as an important basis for software library development. The software library is seen as providing a domain specific embedded language (DSEL) [Hudak 1996]. The DSL will aid when formulating problem descriptions, while the DSEL is useful for the design and implementation of application programs. In our opinion development of a DSL and its implementation as a software library has a software life cycle model, see Section 1, which is distinct from that of application software development [Boehm 1976].

Software itself is a formal entity, with strict rules for syntax and a precise interpretation (ultimately this is the semantics given by the formal system that a computer represents, but hopefully it is given as a formal semantics

that can be investigated by mathematical tools). This leads to a discipline of formal software engineering, i.e., engineering of the formal side of software. Formal software engineering does not imply the use of formal methods throughout. We see formal methods as especially useful when developing a DSL and when investigating the structure for the software library that will provide the resulting DSEL. We strongly believe that algebraic methods are appropriate for this, and have used them extensively in our case study.

The major part of this paper is a case study from the field of computational modelling, specifically that of partial differential equations (PDE), an area which is of significance for the industrial sector. In [Haveraaen *et al.* 1992] we have used algebraic specification techniques to identify types and operations from the problem domain. This gave us a coordinate free language for PDEs. This DSL is distinct from the one commonly used for implementing PDE software, but it is a more abstract language with much better software structuring properties. The fruitfulness of the DSL was checked with the development of a collection of application programs. The DSL was also studied with the purpose of structuring the concepts in order to find an optimal architecture for the supporting software library. The library architecture was designed to eliminate a large part of the implementation work by focusing on reusable template classes.

In summary the concepts of universal algebra allowed us to identify the basic entities – the types and functions – of the problem domain, giving them a syntax and a meaning, yielding a domain specific language. This language was used to formulate problems, and to check its usability for writing software – programming in the small. Then we used categories and functor concepts to pinpoint how the language elements should be organised into a collection of modules. By carefully structuring the modules our coding effort can be reduced, and reusable software components be defined. This is work at the software architecture level – programming in the large.

Later, after the initial development of the DSL library and the application programs, additional applications were specified, and the impact of the new requirements on the library were studied. The case study showed that the library architecture remained stable throughout these changing requirements. It was possible to incrementally add interchangeable versions of new implementations as the need arose, and then use these with the other constructors to build the needed software. This seems to yield a radical improvement in software development and maintenance productivity. After developing a family of related problem solving software, we developed a family of 16 new application programs. The new family of solvers were so different from the existing applications that they traditionally would have required a full development process to be initiated (but some reuse of code should be expected even in this case). With our approach a total cost of 15% of the development cost of a single application was achieved for the new family of applications. Moreover, the library architecture makes it clear exactly which components need to be reconsidered when requirements change. In most cases a component would relate to a specific requirement, but in some

cases a component was influenced by several independent requirements and required reimplementations whenever one of them changed. Further study of structuring techniques is needed to see if these influences can be factored out in a more satisfactory way.

Our software methodology relies on concepts packaged under the term object-orientation, and thus the elastic wave simulator is *object-oriented numerics* (OON) in the sense of Wong *et al.* [1993] and Arge *et al.* [1997]. However, we end up with a software structure which deviates much from most packages claiming to be OON, such as those of Verner *et al.* [1993] and Bruaset and Langtangen [1997]. The OON approaches represent a clear improvement over the traditional numerical software strategy, but many of the OON approaches still remain within that tradition when it comes to the conceptual decomposition of the software architecture.

The development of software libraries for important industrial problems is no simple task. One of the most important aspects is to find a mature library structure [Racko 1995]. Maturity means the library structure is robust faced with new demands and further development. We approached library development using formal software engineering, specifically algebraic methods. The software library produced this way has shown a high degree of maturity. This manifests itself by the fact that the structure of the library, and a large part of its code, is unaltered when the basic discretisation method is changed, and that a radical change in the assumptions of the problem (moving from elasticity to poro-elasticity) requires modifications in one of the components. Moreover, only one component needs to be changed when the software is moved onto a parallel machine architecture.

Formal software engineering is just a small part of the field of software engineering, but a sharper focus on this small area seems to be more than worthwhile in the benefits that may be gained in reduced software cost and improved reusability.

It may be argued that our example is taken from an area, partial differential equations, with a highly developed formal theory, and that this example does not indicate any general applicability of formal software engineering and algebraic methods. Our reply is that we promote algebraic technology as a language discovery device, and that any domain that is to be analysed for the purpose of developing software needs its language concepts to be discovered and formalised. Secondly, we promote reasoning about software as a formal entity. This again is a property of software as such, and is independent of the problem domain the software is developed for. Thus both of these arguments should motivate the investigation of the potential of formal software engineering in other application areas.

### Acknowledgements

This research was supported in part by The Research Council of Norway (NFR), the European Union under ESPRIT Project 21871 (Scientific Com-

puting and Algebraic Abstractions), Statoil, Saga Petroleum, and by computing resources grants from NFR's Programme for Supercomputing.

SeisMod was developed in close cooperation with Åsmund Drottning (contact address UNIGEO, Thormøhlensgt. 55, N-5008 BERGEN, Norway, e-mail aasmund@unigeo.no), who also provided data for the example illustrations in cooperation with Roger Bakke. Numerous people have contributed to the design and implementation of Sophus and SeisMod, of which we will especially mention Hans Munthe-Kaas who was among the initiators of this work, Victor Aarre, Kristin G. Frøysa, Helge Gunnarsli, Kristian Stewart and Steinar Søreide.

## References

- ADAMS, J.C., BRAINERD, W.S., AND MARTIN, J.T. 1992. *Fortran 90 Handbook: Complete ANSI/ISO Reference*. Intertext Publications.
- AKI, KEIITI AND RICHARDS, PAUL G. 1980. *Quantitative Seismology Theory and Methods*. Volume 1. W.H.Freeman and Company.
- ARGE, E., BRUASET, A.M., AND LANGTANGEN, H.P. 1997. Object-Oriented Numerics. In *Numerical Methods and Software Tools in Industrial Mathematics*, Dæhlen, Morten and Tveito, Aslak, Editors. Birkhäuser, Boston, 7–26.
- BACK, R.J.R. 1981. On Correct Refinement of Programs. *Journal of Computer and System Sciences* 23, 49–68.
- BARSTOW, D., EDITOR. 1983. *The Programming Language Ada - Reference Manual*. Springer LNCS 155.
- BIDOIT, M., SANNELLA, D., AND TARLECKI, A. 1999. Architectural Specifications in CASL. In *Algebraic Methodology and Software Technology*, Haeberer, Armando Martín, Editor. Volume 1548 of *Lecture Notes in Computer Science*. Springer, 341–357.
- BIOT, M.A. 1955. Theory of Elasticity and Consolidation for a Porous Anisotropic Solid. *Journal of Applied Physics* 26, 2, 182–185.
- BIOT, M. A. 1956a. Theory of Propagation of Elastic Waves in a Fluid-Saturated Porous Solid. I. Low-Frequency Range. *J. Acoust. Soc. Am.* 28, 1, 168–178.
- BIOT, M. A. 1956b. Theory of Propagation of Elastic Waves in a Fluid-Saturated Porous Solid. II. Higher Frequency Range. *J. Acoust. Soc. Am.* 28, 2, 179–191.
- BOEHM, B.W. 1976. Software Engineering. *IEEE Trans. Comput. C-25*, 1226–1241.
- BOEHM, B.W. 1981. *Software Engineering Economics*. Prentice-Hall.
- BOUGÉ, L. 1996. The Data Parallel Model: A semantic Perspective. In *The Data Parallel Programming Model*, Perrin, Guy-René and Darté, Alain, Editors. Volume 1132 of *Lecture Notes in Computer Science*. Springer, 4–26.
- BRUASET, A.M. AND LANGTANGEN, H.P. 1997. A Comprehensive Set of Tools for Solving Partial Differential Equations; Diffpack. In *Numerical Methods and Software Tools in Industrial Mathematics*, Dæhlen, Morten and Tveito, Aslak, Editors. Birkhäuser, Boston, 61–90.
- BUDGE, K.G., PEERY, J.S., AND ROBINSON, A.C. 1992. High-Performance Scientific Computing Using C++. In *Proceedings from C++ Technical Conference, August 10–13 1992*. USENIX Association, 131–150.
- COHN, P.M. 1965. *Universal algebra*. Harper and Row, New York.
- DIJKSTRA, E.W. 1976. *A Discipline of Programming*. Prentice-Hall.
- DINESH, T.B., HAVERAAEN, M., AND HEERING, J. 1998. An Algebraic Programming Style for Numerical Software and its Optimisation. Technical Report SEN-R9844, CWI, Amsterdam, NL.
- EHRIG, H., GAIEWSKY, M., AND WOLTER, U. 1998. From Abstract Data Types to Algebraic Development Techniques: A Shift of Paradigms. In *Recent Trends in*

- Algebraic Development Techniques*, Presicce, Francesco Parisi, Editor. Volume 1376 of *Lecture Notes in Computer Science*. Springer Verlag, 1–17.
- FLOYD, R.W. 1967. Assigning Meanings to Programs. In *Mathematical Aspects of Computer Science*, Schwartz, J.T., Editor. Volume XIX of *Proceedings of symposia in applied mathematics*. American Mathematical Society, Providence, Rhode Island, 19–31.
- FORTRAN-66 1966. *USA standard FORTRAN: approved March 7, 1966*.
- GOGUEN, J.A., THATCHER, J.W., WAGNER, E.G., AND WRIGHT, J. B. 1975. Abstract data types as initial algebras and the correctness of data representation. In *Proceedings of the Conference on Computer Graphics, Pattern Recognition and Data Representations*.
- GOGUEN, J.A. 1991. A categorical manifesto. *Mathematical Structures in Computer Science 1*, 49–67.
- GRANT, P.W., HAVERAAEN, M., AND WEBSTER, M.F. 1998. Tensor Abstraction Programming of Computational Fluid Dynamics Problems. Technical Report CSR3-98, Department of Computer Science, University of Wales Swansea, Swansea SA2 8PP, United Kingdom.
- GRÄTZER, G. 1968. *Universal algebra*. The University series in higher mathematics. Van Nostrand, Princeton, N.J.
- GUTTAG, J.V. AND HORNING, J.J. 1978. The algebraic specification of abstract data types. *Acta Informatica 10*, 27–52.
- HAVERAAEN, M. 1998. Abstractions for Programming Parallel Machines. Technical Report 162, Department of Informatics, University of Bergen, P.O.Box 7800, N-5020 Bergen, Norway.
- HAVERAAEN, M., MADSEN, V., AND MUNTHE-KAAS, H. 1992. Algebraic Programming Technology for Partial Differential Equations. In *Norsk Informatikk Konferanse – NIK'92*. Tapir, Norway, 55–68.
- HOARE, C.A.R. 1972. Proofs of correctness of data representations. *Acta Informatica 1*, 4, 271–281.
- HUDAK, P. 1996. Building Domain-Specific Embedded Languages. *Computing Surveys 28*, 4es (December), electronic supplement, only available from ACM on the WWW via <http://www.acm.org/pubs/articles/journals/surveys/1996-28-4es/a196-hudak/a196-hudak.html>.
- MACLANE, S. 1971. *Categories for the Working Mathematician*. Springer-Verlag, New York.
- MANNA, Z. AND WALDINGER, R. 1980. A deductive approach to program synthesis. *ACM Transactions on Programming Languages and Systems 2*, 1, 90–121.
- MARSDEN, J.E. AND HUGHES, T.J.R. 1983. *Mathematical Foundations of Elasticity*. Prentice-Hall, Englewood Cliffs, New Jersey.
- MEYER, B. 1991. Design by Contract. In *Advances in object-oriented software engineering*, Mandrioli, Dino and Meyer, Bertrand, Editors. Prentice Hall, Englewood Cliff N.J., 1–50.
- MORRIS, J.H. 1973. Types are not sets. In *Proceedings of the ACM Symposium on the Principles of Programming Languages*, 120–124.
- MOSSES, P.D. 1997. CoFI: The Common Framework Initiative for Algebraic Specification and Development. In *TAPSOFT'97: Theory and Practice of Software Development*, Volume 1214 of *Lecture Notes in Computer Science*. Springer-Verlag, 115–137.
- PAULK, M.C., CURTIS, B., CHRISSIS, M.B., AND WEBER, C.V. 1993. The Capability Maturity Model For Software, Version 1.1. Report CMU/SEI-93-TR-24 and ESC-TR-93-177, Software Engineering Institute, Carnegie Mellon University, Pittsburgh, Pennsylvania 15213.
- RACKO, J.T. 1995. Choose Me. *Software Development*, (April), 81–84.
- ROYCE, W.W. 1970. Managing the Development of Large Software Systems: Concepts and Techniques. In *Proc. WESCON*. IEEE Press, New York, 1–9.
- SCHUTZ, B. 1980. *Geometrical Methods of Mathematical Physics*. Cambridge University Press.

- SNIR, M., OTTO, S.W., HUSS-LEDERMAN, S., WALKER, D.W., AND DONGARRA, J. 1996. *MPI – The complete reference*. MIT Press, Cambridge, Mass., USA.
- STROUSTRUP, B. 1997. *The C++ Programming Language*, 3rd edition. Addison-Wesley.
- VERNER, D.A., HEILEMAN, G.L., BUDGE, K.G., AND ROBINSON, ALLEN C. 1993. Development of Generic Field Classes for Finite Element and Finite Difference Problems. In *OON-SKI'93 Conference Proceedings*. Rogue Wave Software/SIAM, 354–363.
- WHITEHEAD, A.N. 1898. *A Treatise on Universal Algebra, with Applications I*. Cambridge University Press.
- WONG, M. K.W., BUDGE, K.G., PEERY, J.S., AND ROBINSON, A.C. 1993. Object-Oriented Numerics: A Paradigm for Numerical Object-Oriented Programming. *Computers in Physics* 7, 6 (Nov./Dec.), 655–663.