Effective Analysis: Foundations, Implementations, Certification Luminy, France

Formal verification of numerical analysis programs

Sylvie Boldo

Inria

January 12th, 2016



(Applied) Mathematics

(Applied) Mathematics and Formal proof (Coq)

(Applied) Mathematics and Formal proof (Coq) and Floating-point numbers

PDE (Partial Differential Equation)

- \Rightarrow weather forecast
- \Rightarrow nuclear simulation
- \Rightarrow optimal control

 \Rightarrow ...

- - \Rightarrow weather forecast
 - \Rightarrow nuclear simulation
 - \Rightarrow optimal control
 - \Rightarrow ...

Usually too complex to solve by an exact mathematical formula \Rightarrow approximated by numerical scheme over discrete grids

 \Rightarrow mathematical proofs of the convergence of the numerical scheme (we compute something close to the PDE solution if the grids size decreases)

- **PDE** (Partial Differential Equation) \Rightarrow
 - \Rightarrow weather forecast
 - \Rightarrow nuclear simulation
 - $\Rightarrow \quad \mathsf{optimal} \ \mathsf{control}$
 - \Rightarrow ...

Usually too complex to solve by an exact mathematical formula \Rightarrow approximated by numerical scheme over discrete grids

 \Rightarrow mathematical proofs of the convergence of the numerical scheme (we compute something close to the PDE solution if the grids size decreases)

 \Rightarrow C program implementing the scheme

- PDE (Partial Differential Equation) ⇒
 - \Rightarrow weather forecast
 - \Rightarrow nuclear simulation
 - $\Rightarrow \quad \mathsf{optimal} \ \mathsf{control}$
 - \Rightarrow ...

Usually too complex to solve by an exact mathematical formula \Rightarrow approximated by numerical scheme over discrete grids

 \Rightarrow mathematical proofs of the convergence of the numerical scheme (we compute something close to the PDE solution if the grids size decreases)

 \Rightarrow C program implementing the scheme

Let us machine-check this kind of programs!

- François Clément
- Jean-Christophe Filliâtre
- Vincent Martin
- Micaela Mayero
- Guillaume Melquiond
- Pierre Weis

Outline

Introduction

2 Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

5 Conclusion

Numbers



Babylonian clay tablet (1800–1600 BC)

Numbers



Babylonian clay tablet (1800–1600 BC)

$$\begin{array}{rcl} & \swarrow & \\ & & = & (1, 24, 51, 10) \\ & & = & 1 + \frac{24}{60} + \frac{51}{60^2} + \frac{10}{60^3} = 30547/21600 \\ & & \approx & \mathbf{1,41421296} \approx \sqrt{2} \end{array}$$

Numbers



Babylonian clay tablet (1800–1600 BC)

$$\begin{array}{rcl} & \swarrow & (1,24,51,10) \\ & & = & 1 + \frac{24}{60} + \frac{51}{60^2} + \frac{10}{60^3} = 30547/21600 \\ & & \approx & \mathbf{1,41421296} \approx \sqrt{2} \end{array}$$

 \Rightarrow representation of a real number with a finite precision

Sylvie Boldo (Inria)

Using a finite number of bits (the precision p) and based on scientific notation, computers use floating-point (FP) numbers.

A FP number is only a string of bits.

11100011010010011110000111000000

Using a finite number of bits (the precision p) and based on scientific notation, computers use floating-point (FP) numbers.

A FP number is only a string of bits.

11100011010010011110000111000000

We interpret it depending on the respective values of s (sign), e (exponent) and f (fraction).

1	11000110	10010011110000111000000
1	11000110	10010011110000111000000
s	е	f

We associate a real value:



We associate a real value:



except for the special values of $e: \pm 0, \pm \infty$, NaN, subnormals.









For the $+, -, \times, \div, \sqrt{}$, the result is the same as if the infinitely precise mathematical result was computed and then rounded to the nearest floating-point number.

 \Rightarrow guaranteed by the IEEE-754 standard (1985 & 2008).

For the $+, -, \times, \div, \sqrt{}$, the result is the same as if the infinitely precise mathematical result was computed and then rounded to the nearest floating-point number.

 \Rightarrow guaranteed by the IEEE-754 standard (1985 & 2008).

 \Rightarrow portability & accuracy

For the $+, -, \times, \div, \sqrt{}$, the result is the same as if the infinitely precise mathematical result was computed and then rounded to the nearest floating-point number.

 \Rightarrow guaranteed by the IEEE-754 standard (1985 & 2008).

 \Rightarrow portability & accuracy

 \Rightarrow if $x \in \mathbb{R}$ is not too small, $|x - \circ_{\texttt{double}}(x)| \leq 2^{-53}|x|$

Floating-Point Computations

More than one FP operation may lead to incorrect results.

Floating-Point Computations

More than one FP operation may lead to incorrect results.



Outline

Introduction

Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

Conclusion

• ANSI/ISO C Specification Language

- ANSI/ISO C Specification Language
- behavioral specification language for C programs

- ANSI/ISO C Specification Language
- behavioral specification language for C programs
- pre-conditions and post-conditions to functions (and which variables are modified).

- ANSI/ISO C Specification Language
- behavioral specification language for C programs
- pre-conditions and post-conditions to functions (and which variables are modified).
- variants and invariants of the loops.

- ANSI/ISO C Specification Language
- behavioral specification language for C programs
- pre-conditions and post-conditions to functions (and which variables are modified).
- variants and invariants of the loops.
- assertions

- ANSI/ISO C Specification Language
- behavioral specification language for C programs
- pre-conditions and post-conditions to functions (and which variables are modified).
- variants and invariants of the loops.
- assertions
- In annotations, all computations are exact.

- ANSI/ISO C Specification Language
- behavioral specification language for C programs
- pre-conditions and post-conditions to functions (and which variables are modified).
- variants and invariants of the loops.
- assertions
- In annotations, all computations are exact.

 \Rightarrow For the programmer, the specification is easy to understand.

ACSL and floating-point numbers

A floating-point number is a triple:

• the floating-point number, really computed by the program, $x \rightarrow x_f$ floating-point part
- the floating-point number, really computed by the program, $x \rightarrow x_f$ floating-point part
- the value that would have been obtained with exact computations, $x \rightarrow x_e$ exact part

- the floating-point number, really computed by the program, $x \rightarrow x_f$ floating-point part
- the value that would have been obtained with exact computations, $x \rightarrow x_e$ exact part
- the value that we ideally wanted to compute
 - $x \rightarrow x_m$ model part

- the floating-point number, really computed by the program, $x \rightarrow x_f$ floating-point part 1+x+x*x/2
- the value that would have been obtained with exact computations, $x \to x_e$ exact part $1 + x + \frac{x^2}{2}$
- the value that we ideally wanted to compute
 - $x \to x_m \mod part \qquad \exp(x)$

- the floating-point number, really computed by the program, $x \rightarrow x_f$ floating-point part 1+x+x*x/2
- the value that would have been obtained with exact computations, $x \rightarrow x_e$ exact part $1 + x + \frac{x^2}{2}$
- the value that we ideally wanted to compute
 - $x \to x_m$ model part $\exp(x)$
- \Rightarrow easy to split into method error and rounding error

C Program

























Outline



Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

Conclusion

The wave equation

Looking for $u: \mathbb{R}^2 \to \mathbb{R}$ regular enough such that:

$$\frac{\partial^2 u(x,t)}{\partial t^2} - c^2 \frac{\partial^2 u(x,t)}{\partial x^2} = s(x,t)$$

with given values for the initial position $u_0(x)$ and the initial velocity $u_1(x)$.

 \Rightarrow rope oscillation, sound, radar, oil prospection...

Scheme?

We want
$$u_j^k \approx u(j\Delta x, k\Delta t)$$
.
$$\frac{u_j^k - 2u_j^{k-1} + u_j^{k-2}}{\Delta t^2} - c^2 \frac{u_{j+1}^{k-1} - 2u_j^{k-1} + u_{j-1}^{k-1}}{\Delta x^2} = s_j^{k-1}$$

And other horrible formulas to initialize u_j^0 and u_j^1 .

Scheme?

We want
$$u_j^k \approx u(j\Delta x, k\Delta t)$$
.
$$\frac{u_j^k - 2u_j^{k-1} + u_j^{k-2}}{\Delta t^2} - c^2 \frac{u_{j+1}^{k-1} - 2u_j^{k-1} + u_{j-1}^{k-1}}{\Delta x^2} = s_j^{k-1}$$

And other horrible formulas to initialize u_i^0 and u_i^1 .



Program

```
// initialization of p[i][0] and p[i][1]
for (k=1; k<nk; k++) {
    p[0][k+1] = 0.;
    for (i=1; i<ni; i++) {
        dp = p[i+1][k] - 2.*p[i][k] + p[i-1][k];
        p[i][k+1] = 2.*p[i][k] - p[i][k-1] + a*dp;
        }
    p[ni][k+1] = 0.;
}</pre>
```

Program

```
// initialization of p[i][0] and p[i][1]
for (k=1; k<nk; k++) {
    p[0][k+1] = 0.;
    for (i=1; i<ni; i++) {
        dp = p[i+1][k] - 2.*p[i][k] + p[i-1][k];
        p[i][k+1] = 2.*p[i][k] - p[i][k-1] + a*dp;
        }
    p[ni][k+1] = 0.;
}</pre>
```

Two different errors:

round-off errors

due to floating-point roundings

method errors

the scheme only approximates the exact solution

Outline

Introduction

Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

Conclusion

Remainder:

```
\begin{array}{l} dp \ = \ p[\ i + 1][k] \ - \ 2.*p[\ i \ ][k] \ + \ p[\ i - 1][k]; \\ p[\ i \ ][k + 1] \ = \ 2.*p[\ i \ ][k] \ - \ p[\ i \ ][k - 1] \ + \ a*dp; \end{array}
```

If we use a naive technique to bound the rounding errors, we get

Remainder:

$$\begin{aligned} dp &= p[i+1][k] - 2.*p[i][k] + p[i-1][k]; \\ p[i][k+1] &= 2.*p[i][k] - p[i][k-1] + a*dp; \end{aligned}$$

If we use a naive technique to bound the rounding errors, we get

$$|p_i^k - exact(p_i^k)| \leq O\left(2^k 2^{-53}\right)$$

Remainder:

If we use a naive technique to bound the rounding errors, we get

$$|p_i^k - exact(p_i^k)| \leq O\left(2^k 2^{-53}\right)$$

This is too much because the errors do compensate.

Definition of ε_i^k

Remainder:

$$\begin{array}{l} dp \ = \ p[\ i + 1][\ k] \ - \ 2.*p[\ i \][\ k] \ + \ p[\ i - 1][\ k]; \\ p[\ i \][\ k + 1] \ = \ 2.*p[\ i \][\ k] \ - \ p[\ i \][\ k - 1] \ + \ a*dp; \end{array}$$

Let ε_i^{k+1} be the rounding error made during these two lines of computations.

We assume a, p_{i-1}^k , p_i^k , p_{i+1}^k and p_i^{k-1} are exact and we look into the rounding error of these two lines. It is called ε_i^{k+1} .

Definition of ε_i^k

Remainder:

$$\begin{array}{l} dp \ = \ p[\ i + 1][k] \ - \ 2.*p[\ i \][k] \ + \ p[\ i - 1][k]; \\ p[\ i \][k + 1] \ = \ 2.*p[\ i \][k] \ - \ p[\ i \][k - 1] \ + \ a*dp; \end{array}$$

Let ε_i^{k+1} be the rounding error made during these two lines of computations.

We assume a, p_{i-1}^k , p_i^k , p_{i+1}^k and p_i^{k-1} are exact and we look into the rounding error of these two lines. It is called ε_i^{k+1} .

We know (from initializations) that the model values of the $|p_n^m|$ are bounded by 1. We assume that the floating-point values of the $|p_n^m|$ are bounded by 2.

Definition of ε_i^k

Remainder:

$$\begin{array}{l} dp \ = \ p[\ i + 1][k] \ - \ 2.*p[\ i \][k] \ + \ p[\ i - 1][k]; \\ p[\ i \][k + 1] \ = \ 2.*p[\ i \][k] \ - \ p[\ i \][k - 1] \ + \ a*dp; \end{array}$$

Let ε_i^{k+1} be the rounding error made during these two lines of computations.

We assume a, p_{i-1}^k , p_i^k , p_{i+1}^k and p_i^{k-1} are exact and we look into the rounding error of these two lines. It is called ε_i^{k+1} .

We know (from initializations) that the model values of the $|p_n^m|$ are bounded by 1. We assume that the floating-point values of the $|p_n^m|$ are bounded by 2.

 $|\varepsilon_n^m| \le 78 \times 2^{-52}$

Rounding error

$$p_i^k - exact(p_i^k) = \sum_{l=0}^k \sum_{j=-l}^l \alpha_j^l \varepsilon_{i+j}^{k-l}$$

We have an analytical expression of the rounding error with known constants \(\alpha_i^k\).

Rounding error

$$p_i^k - exact(p_i^k) = \sum_{l=0}^k \sum_{j=-l}^l \alpha_j^l \varepsilon_{i+j}^{k-l}$$

- We have an analytical expression of the rounding error with known constants α^k_i.
- It is not that complicated! (we cannot get rid of the pyramidal double summation)

Rounding error

$$p_i^k - exact(p_i^k) = \sum_{l=0}^k \sum_{j=-l}^l lpha_j^l arepsilon_{i+j}^{k-l}$$

- We have an analytical expression of the rounding error with known constants α^k_i.
- It is not that complicated! (we cannot get rid of the pyramidal double summation)
- The rounding error is bounded by $\bigcirc (k^2 2^{-53})$:

$$\left| p_{i}^{k} - \text{exact}\left(p_{i}^{k} \right) \right| \leq 78 \times 2^{-53} \times (k+1) \times (k+2)$$

Outline

Introduction

Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

Conclusion

Method error

We measure that u and u_i^k are close when $(\Delta x, \Delta t) \rightarrow 0$.

We define $e_j^k \stackrel{\text{def}}{=} \overline{u}_j^k - u_j^k$: convergence error where \overline{u}_j^k is the value of u at the (j, k) point of the grid.

Method error

We measure that u and u_i^k are close when $(\Delta x, \Delta t) \rightarrow 0$.

We define $e_j^k \stackrel{\text{def}}{=} \bar{u}_j^k - u_j^k$: convergence error where \bar{u}_i^k is the value of u at the (j, k) point of the grid.

We want to bound $\left\|e_{h}^{k_{\Delta t}(t)}\right\|_{\Delta x}$: the average of the convergence error on all points of the grid at a given time $k_{\Delta t}(t) = \left\lfloor \frac{t}{\Delta t} \right\rfloor \Delta t$.
Method error

We measure that u and u_i^k are close when $(\Delta x, \Delta t) \rightarrow 0$.

We define $e_j^k \stackrel{\text{def}}{=} \bar{u}_j^k - u_j^k$: convergence error where \bar{u}_i^k is the value of u at the (j, k) point of the grid.

We want to bound $\left\|e_{h}^{k_{\Delta t}(t)}\right\|_{\Delta \times}$: the average of the convergence error on all points of the grid at a given time $k_{\Delta t}(t) = \left\lfloor \frac{t}{\Delta t} \right\rfloor \Delta t$.

We want to prove:

$$\left\|e_{h}^{k_{\Delta t}(t)}\right\|_{\Delta x} = O_{[0,t_{\max}]}(\Delta x^{2} + \Delta t^{2})$$

Usually, the big O uses one variable and $f(\mathbf{x}) = O_{\|\mathbf{x}\| \to 0}(g(\mathbf{x}))$ means

 $\exists \alpha, C > 0, \quad \forall \mathbf{x} \in \mathbb{R}^n, \quad \|\mathbf{x}\| \le \alpha \Rightarrow |f(\mathbf{x})| \le C \cdot |g(\mathbf{x})|.$

Usually, the big O uses one variable and $f(\mathbf{x}) = O_{\|\mathbf{x}\| \to 0}(g(\mathbf{x}))$ means

 $\exists \alpha, C > 0, \quad \forall \mathbf{x} \in \mathbb{R}^n, \quad \|\mathbf{x}\| \le \alpha \Rightarrow |f(\mathbf{x})| \le C \cdot |g(\mathbf{x})|.$

Here 2 variables: Δx (grid sizes, tends to 0), and x (time and space). (Think about Taylor expansions)

Usually, the big O uses one variable and $f(\mathbf{x}) = O_{\|\mathbf{x}\| \to 0}(g(\mathbf{x}))$ means

 $\exists \alpha, C > 0, \quad \forall \mathbf{x} \in \mathbb{R}^n, \quad \|\mathbf{x}\| \le \alpha \Rightarrow |f(\mathbf{x})| \le C \cdot |g(\mathbf{x})|.$

Here 2 variables: Δx (grid sizes, tends to 0), and x (time and space). (Think about Taylor expansions)

$$\forall \mathbf{x}, \exists \alpha, C > \mathbf{0}, \quad \forall \mathbf{\Delta} \mathbf{x} \in \mathbb{R}^2, \quad \|\mathbf{\Delta} \mathbf{x}\| \leq \alpha \Rightarrow |f(\mathbf{x}, \mathbf{\Delta} \mathbf{x})| \leq C \cdot |g(\mathbf{\Delta} \mathbf{x})|$$

does not work.

We used a uniform big O:

 $\exists \alpha, C > 0, \quad \forall \mathbf{x}, \Delta \mathbf{x}, \quad \|\Delta \mathbf{x}\| \le \alpha \Rightarrow |f(\mathbf{x}, \Delta \mathbf{x})| \le C \cdot |g(\Delta \mathbf{x})|.$

where variables **x** and $\Delta \mathbf{x}$ are restricted to subsets of \mathbb{R}^2 . (for example such that $\Delta t > 0$) \Rightarrow Taylor expansions The truncation error is defined as how much the exact solution solves the numerical scheme:

$$\varepsilon_j^{k-1} = \frac{\bar{u}_j^k - 2\bar{u}_j^{k-1} + \bar{u}_j^{k-2}}{\Delta t^2} - c^2 \frac{\bar{u}_{j+1}^{k-1} - 2\bar{u}_j^{k-1} + \bar{u}_{j-1}^{k-1}}{\Delta x^2} - s_j^{k-1}$$

The truncation error is defined as how much the exact solution solves the numerical scheme:

$$\varepsilon_j^{k-1} = \frac{\bar{u}_j^k - 2\bar{u}_j^{k-1} + \bar{u}_j^{k-2}}{\Delta t^2} - c^2 \frac{\bar{u}_{j+1}^{k-1} - 2\bar{u}_j^{k-1} + \bar{u}_{j-1}^{k-1}}{\Delta x^2} - s_j^{k-1}$$

The consistency is the boundedness of the truncation error:

$$\left\|\varepsilon_{h}^{k_{\Delta t}(t)}\right\|_{\Delta x} = O_{[0,t_{\max}]}(\Delta x^{2} + \Delta t^{2})$$

By Taylor series and many computations.

Proof idea 2/3: stability

We define a discrete energy by

$$E_{h}(c)(u_{h})^{k+\frac{1}{2}} \stackrel{\text{def}}{=} \frac{1}{2} \left\| \frac{u_{h}^{k+1} - u_{h}^{k}}{\Delta t} \right\|_{\Delta x}^{2} + \frac{1}{2} \left\langle u_{h}^{k}, u_{h}^{k+1} \right\rangle_{A_{h}(c)}$$

kinetic energy potential energy

$$\langle v_h, w_h \rangle_{A_h(c)} \stackrel{\text{def}}{=} \langle A_h(c) v_h, w_h \rangle_{\Delta x} \text{ and } (A_h(c) v_h)_j \stackrel{\text{def}}{=} - c^2 \frac{v_{j+1} - 2v_j + v_{j-1}}{\Delta x^2}$$

Proof idea 2/3: stability

We define a discrete energy by

$$E_{h}(c)(u_{h})^{k+\frac{1}{2}} \stackrel{\text{def}}{=} \frac{1}{2} \left\| \frac{u_{h}^{k+1} - u_{h}^{k}}{\Delta t} \right\|_{\Delta x}^{2} + \frac{1}{2} \left\langle u_{h}^{k}, u_{h}^{k+1} \right\rangle_{A_{h}(c)}$$

kinetic energy potential energy

$$\langle v_h, w_h \rangle_{A_h(c)} \stackrel{\text{def}}{=} \langle A_h(c) v_h, w_h \rangle_{\Delta x} \text{ and } (A_h(c) v_h)_j \stackrel{\text{def}}{=} - c^2 \frac{v_{j+1} - 2v_j + v_{j-1}}{\Delta x^2}.$$

Note that this energy is constant if f = 0. We prove an overestimation and an underestimation of this energy. $\Rightarrow u_h$ does not diverge. The convergence error is solution of the same discrete scheme with inputs

$$u_{0,j} = 0,$$
 $u_{1,j} = \frac{e_j^1}{\Delta t},$ and $s_j^k = \varepsilon_j^{k+1}.$

+ proofs about the initializations.

The convergence error is solution of the same discrete scheme with inputs

$$u_{0,j} = 0,$$
 $u_{1,j} = \frac{e_j^1}{\Delta t},$ and $s_j^k = \varepsilon_j^{k+1}.$

+ proofs about the initializations.

All these proofs require the existence of ζ and ξ in]0,1[with $\zeta \leq 1-\xi$ and we require that $\zeta \leq \frac{c\Delta t}{\Delta x} \leq 1-\xi$ (CFL conditions).

We proved that:

$$\left\| e_h^{k_{\Delta t}(t)} \right\|_{\Delta x} = O_{\substack{t \in [0, t_{\max}] \\ (\Delta x, \Delta t) \to 0 \\ 0 < \Delta x \land 0 < \Delta t \land \\ \zeta \le c \frac{\Delta t}{\Delta x} \le 1 - \xi}} (\Delta x^2 + \Delta t^2).$$

Extraction of the big O constants

The preceding result is a uniform big O defined by:

 $\exists \alpha, C > 0, \quad \forall \mathbf{x}, \Delta \mathbf{x}, \quad \| \Delta \mathbf{x} \| \le \alpha \Rightarrow |f(\mathbf{x}, \Delta \mathbf{x})| \le C \cdot |g(\Delta \mathbf{x})|.$

Extraction of the big O constants

The preceding result is a uniform big O defined by:

 $\exists \alpha, C > 0, \quad \forall \mathbf{x}, \Delta \mathbf{x}, \quad \| \Delta \mathbf{x} \| \le \alpha \Rightarrow |f(\mathbf{x}, \Delta \mathbf{x})| \le C \cdot |g(\Delta \mathbf{x})|.$

Let (α_3, C_3) be the constants for the order-3 Taylor development of the exact solution and (α_4, C_4) for order-4. The initial support is $[\chi_1; \chi_2]$.

Extraction of the big O constants

The preceding result is a uniform big O defined by:

 $\exists \alpha, C > 0, \quad \forall \mathbf{x}, \Delta \mathbf{x}, \quad \| \Delta \mathbf{x} \| \le \alpha \Rightarrow |f(\mathbf{x}, \Delta \mathbf{x})| \le C \cdot |g(\Delta \mathbf{x})|.$

Let (α_3, C_3) be the constants for the order-3 Taylor development of the exact solution and (α_4, C_4) for order-4. The initial support is $[\chi_1; \chi_2]$.

$$\begin{array}{lll} \alpha & = & \min(\alpha_3, \alpha_4, 1, t_{\max}) \\ s_1 & = & \max(1, 2 \cdot C_4 \cdot (c^2 + 1), C_3 \cdot (1 + c^2/2) + 1) \\ s_2 & = & s_1^2 \left(\lfloor \chi_2 \rfloor - \lfloor \chi_1 \rfloor + 2 \cdot c \cdot t_{\max} \cdot \left(1 + \frac{1}{\zeta} \right) + 3 \right) \\ s_3 & = & \frac{1}{\sqrt{2}} \left(C_3 \cdot (1 + c^2/2) + 1 \right) \cdot (\chi_2 - \chi_1 + 1 + (2 \cdot c + 4)) \\ & & \quad + \frac{\sqrt{2}}{2\sqrt{2\xi - \xi^2}} (2 \cdot t_{\max} \cdot s_2 + 2s_2) \\ C & = & \frac{\sqrt{2}}{\sqrt{2\xi - \xi^2}} \cdot 2 \cdot t_{\max} \cdot s_3 \end{array}$$

Sylvie Boldo (Inria)

Outline

Introduction

Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification
- 4 About the Finite Element Method

D Conclusion

Program verification

- 154 lines of annotations for 32 lines of C
- 150 verification conditions:
 - 44 about the behavior
 - 106 about the safety (runtime errors)

Program verification

- 154 lines of annotations for 32 lines of C
- 150 verification conditions:
 - 44 about the behavior
 - 106 about the safety (runtime errors)

Prover	Behavior VC	Safety VC	Total
Alt-Ergo	18	80	98
CVC3	18	89	107
Gappa	2	20	22
Z3	21	63	84
Automatically proved	23	94	117
Coq	21	12	33
Total	44	106	150

- About 90 % of the safety goals (matrix access, Overflow, and so on) are proved automatically.
- 33 theorems are interactively proved using Coq for a total of about 15,000 lines of Coq and 30 minutes of compilation.

Type of proofs	Nb spec lines	Nb lines	Compilation time
Convergence	991	5 275	42 s
Round-off $+$ runtime errors	7 737	13 175	32 min

Outline

Introduction

Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

Conclusion

http://www.ima.umn.edu/~arnold/disasters/sleipner.html

The sinking of the Sleipner A offshore platform

Excerpted from a report of SINTEF, Civil and Environmental Engineering:

The Sleipner A platform produces oil and gas in the North Sea and is supported on the seabed at a water depth of 82 m. It is a Condeep type platform with a concrete gravity base structure consisting of 24 cells and with a total base area of 16 000 m². Four cells are elongated to shafts supporting the platform deck. The first concrete base structure for Sleipner A sprang a leak and sank under a controlled ballasting operation during preparation for deck mating in Gandsfjorden outside Stavanger, Norway on 23 August 1991.

Immediately after the accident, the owner of the platform, Statoil, a Norwegian oil company appointed an investigation group, and SINTEF was contracted to be the technical advisor for this group.

The investigation into the accident is described in 16 reports...

The conclusion of the investigation was that the loss was caused by a failure in a cell wall, resulting in a serious crack and a leakage that the pumps were not able to cope with. The wall failed as a result of a combination of a serious error in the finite element analysis and insufficient anchorage of the reinforcement in a critical zone.

A better idea of what was involved can be obtained from this photo and sketch of the platform. The top deck weighs 57,000 tons, and provides accommodation for about 200 people and support for drilling equipment weighing about 40,000 tons. When the first model sank in August 1991, the crash Sylvie Boldo (Inria)



Verification of numerical analysis programs

Real life applications need solving PDE (Partial Differential Equation) on complex 3D geometries.



Real life applications need solving PDE (Partial Differential Equation) on complex 3D geometries.





Real life applications need solving PDE (Partial Differential Equation) on complex 3D geometries.





Instead of regular 2D/3D grids, we consider meshes made of triangles/tetrahedra.

Sylvie Boldo (Inria)

Verification of numerical analysis programs

FEM encompasses methods for connecting many simple element equations over many small subdomains, named finite elements, to approximate a more complex equation over a larger domain.

(https://en.wikipedia.org/wiki/Finite_element_method)

FEM encompasses methods for connecting many simple element equations over many small subdomains, named finite elements, to approximate a more complex equation over a larger domain.

(https://en.wikipedia.org/wiki/Finite_element_method)

⇒ mathematical proofs of the FEM ⇒ C++ library (Felisce) implementing the FEM

FEM encompasses methods for connecting many simple element equations over many small subdomains, named finite elements, to approximate a more complex equation over a larger domain.

(https://en.wikipedia.org/wiki/Finite_element_method)

⇒ mathematical proofs of the FEM ⇒ C++ library (Felisce) implementing the FEM

Let us machine-check this program!

FEM encompasses methods for connecting many simple element equations over many small subdomains, named finite elements, to approximate a more complex equation over a larger domain.

(https://en.wikipedia.org/wiki/Finite_element_method)

⇒ mathematical proofs of the FEM ⇒ C++ library (Felisce) implementing the FEM

Let us machine-check this program!

First, let us understand/formally prove the mathematics.

Opening a parenthesis



We wanted to prove that

$$u(x,t) = \frac{1}{2} \left(u_0(x+ct) + u_0(x-ct) \right) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi) \, d\xi + \frac{1}{2c} \int_0^t \int_{x-c(t-\tau)}^{x+c(t-\tau)} f(\xi,\tau) \, d\xi \, d\tau$$

is regular and solution of $\frac{\partial^2 u}{\partial t^2}(x,t) - c^2 \frac{\partial^2 u}{\partial x^2}(x,t) = f(x,t)$.

We wanted to prove that

$$u(x,t) = \frac{1}{2} \left(u_0(x+ct) + u_0(x-ct) \right) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi) \, d\xi + \frac{1}{2c} \int_0^t \int_{x-c(t-\tau)}^{x+c(t-\tau)} f(\xi,\tau) \, d\xi \, d\tau$$

is regular and solution of
$$\frac{\partial^2 u}{\partial t^2}(x,t) - c^2 \frac{\partial^2 u}{\partial x^2}(x,t) = f(x,t).$$

We developed a Coq library of real analysis that is:

compatible with the Coq standard library of real numbers,

We wanted to prove that

$$u(x,t) = \frac{1}{2} \left(u_0(x+ct) + u_0(x-ct) \right) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi) \, d\xi + \frac{1}{2c} \int_0^t \int_{x-c(t-\tau)}^{x+c(t-\tau)} f(\xi,\tau) \, d\xi \, d\tau$$

is regular and solution of
$$\frac{\partial^2 u}{\partial t^2}(x,t) - c^2 \frac{\partial^2 u}{\partial x^2}(x,t) = f(x,t).$$

- compatible with the Coq standard library of real numbers,
- with total functions for limit, derivative, integral,

We wanted to prove that

$$u(x,t) = \frac{1}{2} \left(u_0(x+ct) + u_0(x-ct) \right) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi) \, d\xi + \frac{1}{2c} \int_0^t \int_{x-c(t-\tau)}^{x+c(t-\tau)} f(\xi,\tau) \, d\xi \, d\tau$$

is regular and solution of $\frac{\partial^2 u}{\partial t^2}(x,t) - c^2 \frac{\partial^2 u}{\partial x^2}(x,t) = f(x,t)$.

- compatible with the Coq standard library of real numbers,
- with total functions for limit, derivative, integral,
- with parametric integrals, two-dimensional differentiability, asymptotic behaviors,

We wanted to prove that

$$u(x,t) = \frac{1}{2} \left(u_0(x+ct) + u_0(x-ct) \right) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi) \, d\xi + \frac{1}{2c} \int_0^t \int_{x-c(t-\tau)}^{x+c(t-\tau)} f(\xi,\tau) \, d\xi \, d\tau$$

is regular and solution of $\frac{\partial^2 u}{\partial t^2}(x,t) - c^2 \frac{\partial^2 u}{\partial x^2}(x,t) = f(x,t)$.

- compatible with the Coq standard library of real numbers,
- with total functions for limit, derivative, integral,
- with parametric integrals, two-dimensional differentiability, asymptotic behaviors,
- with a tactic dedicated to derivative proofs.

We wanted to prove that

$$u(x,t) = \frac{1}{2} \left(u_0(x+ct) + u_0(x-ct) \right) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi) \, d\xi + \frac{1}{2c} \int_0^t \int_{x-c(t-\tau)}^{x+c(t-\tau)} f(\xi,\tau) \, d\xi \, d\tau$$

is regular and solution of $\frac{\partial^2 u}{\partial t^2}(x,t) - c^2 \frac{\partial^2 u}{\partial x^2}(x,t) = f(x,t).$

- compatible with the Coq standard library of real numbers,
- with total functions for limit, derivative, integral,
- with parametric integrals, two-dimensional differentiability, asymptotic behaviors,
- with a tactic dedicated to derivative proofs.
- $\Rightarrow\,$ Then, it was extended to more than real analysis.

Hierarchy (Lelay, Melquiond)



→: "used to define" --→: "parameter of"
Hierarchy (Lelay, Melquiond)



Hierarchy (Lelay, Melquiond)



Hierarchy (Lelay, Melquiond)



Closing a parenthesis



• more 50 pages of mathematical proofs

- more 50 pages of mathematical proofs
- very detailed!

- more 50 pages of mathematical proofs
- very detailed!
- more than 7,000 lines and 220,000 characters

- more 50 pages of mathematical proofs
- very detailed!
- more than 7,000 lines and 220,000 characters
- with dependencies!

- more 50 pages of mathematical proofs
- very detailed!
- more than 7,000 lines and 220,000 characters
- with dependencies!

- more 50 pages of mathematical proofs
- very detailed!
- more than 7,000 lines and 220,000 characters
- with dependencies!



+ general spaces

- Let us build upon Coquelicot
 - + general spaces
 - + many existing theorems

- + general spaces
- + many existing theorems
 - not always the space we need

- + general spaces
- + many existing theorems
 - not always the space we need

- + general spaces
- + many existing theorems
 - not always the space we need

Please note this is still work in progress.

Summary of the work done

• results about functional spaces, linear and bilinear mappings

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero
- \bullet norm on functions (operator_norm) on $\mathbb{R} \cup \{+\infty\}$

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero
- \bullet norm on functions (operator_norm) on $\mathbb{R} \cup \{+\infty\}$
- 8 equivalences of continuity of linear mappings

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero
- \bullet norm on functions (operator_norm) on $\mathbb{R} \cup \{+\infty\}$
- 8 equivalences of continuity of linear mappings
- definitions of pre-Hilbert and Hilbert spaces inside the Coquelicot hierarchy (and first lemmas)

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero
- \bullet norm on functions (operator_norm) on $\mathbb{R} \cup \{+\infty\}$
- 8 equivalences of continuity of linear mappings
- definitions of pre-Hilbert and Hilbert spaces inside the Coquelicot hierarchy (and first lemmas)
- define clm: the set of the continuous linear mappings

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero
- \bullet norm on functions (operator_norm) on $\mathbb{R} \cup \{+\infty\}$
- 8 equivalences of continuity of linear mappings
- definitions of pre-Hilbert and Hilbert spaces inside the Coquelicot hierarchy (and first lemmas)
- define clm: the set of the continuous linear mappings
- prove it is a NormedModule, to consider clm E (clm $E \mathbb{R}$)

- results about functional spaces, linear and bilinear mappings
- fixed-point theorem in a sub-complete normed module
- decide if a space (\geq NormedModule) is only zero
- \bullet norm on functions (operator_norm) on $\mathbb{R} \cup \{+\infty\}$
- 8 equivalences of continuity of linear mappings
- definitions of pre-Hilbert and Hilbert spaces inside the Coquelicot hierarchy (and first lemmas)
- define clm: the set of the continuous linear mappings
- prove it is a NormedModule, to consider clm E (clm $E \mathbb{R}$)
- state Lax-Milgram theorem

Hierarchy II



Hierarchy II



Hierarchy II



Sylvie Boldo (Inria)

Problems:

• what is a subgroup?

- what is a subgroup?
 - dependent type: consider $\{g \in G \,|\, \chi(g)\}$

- what is a subgroup?
 - dependent type: consider $\{g \in G \mid \chi(g)\} = \{g : G; H : \chi(g)\}$

- what is a subgroup?
 - dependent type: consider {g ∈ G | χ(g)} = {g : G; H : χ(g)} and prove it is a group

- what is a subgroup?
 - dependent type: consider {g ∈ G | χ(g)} = {g : G; H : χ(g)} and prove it is a group
 - $\bullet \ {\rm overgroup} + \chi + {\rm a} \ {\rm property} \ {\rm of} \ \chi$
 - . . .

- what is a subgroup?
 - dependent type: consider {g ∈ G | χ(g)} = {g : G; H : χ(g)} and prove it is a group
 - $\bullet \ {\rm overgroup} + \chi + {\rm a} \ {\rm property} \ {\rm of} \ \chi$
 - . . .
- canonical structures (both a help and a pain)

Problems:

- what is a subgroup?
 - dependent type: consider {g ∈ G | χ(g)} = {g : G; H : χ(g)} and prove it is a group
 - $\bullet \ {\rm overgroup} + \chi + {\rm a} \ {\rm property} \ {\rm of} \ \chi$
 - . . .
- canonical structures (both a help and a pain)

Size:

• maths: more than 7,000 lines and 220,000 characters (50 pages)

Problems:

- what is a subgroup?
 - dependent type: consider {g ∈ G | χ(g)} = {g : G; H : χ(g)} and prove it is a group
 - $\bullet \ {\rm overgroup} + \chi + {\rm a} \ {\rm property} \ {\rm of} \ \chi$
 - . . .
- canonical structures (both a help and a pain)

Size:

- maths: more than 7,000 lines and 220,000 characters (50 pages)
- Coq: more than 2,300 lines and 56,000 characters (117 lemmas)

Problems:

- what is a subgroup?
 - dependent type: consider {g ∈ G | χ(g)} = {g : G; H : χ(g)} and prove it is a group
 - $\bullet \ {\rm overgroup} + \chi + {\rm a} \ {\rm property} \ {\rm of} \ \chi$
 - . . .
- canonical structures (both a help and a pain)

Size:

- maths: more than 7,000 lines and 220,000 characters (50 pages)
- Coq: more than 2,300 lines and 56,000 characters (117 lemmas)

still a lot to do!
Outline



Prerequisite

- Floating-Point Arithmetic
- Deductive Program Verification

3 1-D Wave equation discretization

- Presentation
- Rounding Error
- Method Error
- Program Verification

4 About the Finite Element Method

5 Conclusion

- Very high guarantee
- interdisciplinary (formal methods / numerical analysis)

- Very high guarantee
- interdisciplinary (formal methods / numerical analysis)
- not only rounding errors:

- Very high guarantee
- interdisciplinary (formal methods / numerical analysis)
- not only rounding errors:
 - all other errors such as pointer dereferencing or division by zero

- Very high guarantee
- interdisciplinary (formal methods / numerical analysis)
- not only rounding errors:
 - all other errors such as pointer dereferencing or division by zero
 - link with mathematical properties

- interdisciplinary (formal methods / numerical analysis)
- not only rounding errors:
 - all other errors such as pointer dereferencing or division by zero
 - link with mathematical properties
 - any property can be checked

- interdisciplinary (formal methods / numerical analysis)
- not only rounding errors:
 - all other errors such as pointer dereferencing or division by zero
 - link with mathematical properties
 - any property can be checked
- expressive annotation language (as expressive as Coq)
 ⇒ exactly the specification you want

- interdisciplinary (formal methods / numerical analysis)
- not only rounding errors:
 - all other errors such as pointer dereferencing or division by zero
 - link with mathematical properties
 - any property can be checked
- expressive annotation language (as expressive as Coq)
 ⇒ exactly the specification you want
- an annotated C program to convince numerical analysts

• go deeper into numerical analysis

- go deeper into numerical analysis
- $\Rightarrow\,$ proof of the finite element method

- go deeper into numerical analysis
- \Rightarrow proof of the finite element method
- $\Rightarrow\,$ proof of the finite element method library

- go deeper into numerical analysis
- \Rightarrow proof of the finite element method
- $\Rightarrow\,$ proof of the finite element method library
- \Rightarrow stability (floating-point stability / numerical analysis stability)

- go deeper into numerical analysis
- \Rightarrow proof of the finite element method
- $\Rightarrow\,$ proof of the finite element method library
- \Rightarrow stability (floating-point stability / numerical analysis stability)

• prove and generalize well-known facts/algorithms/programs from the computer arithmetic community

- go deeper into numerical analysis
- \Rightarrow proof of the finite element method
- \Rightarrow proof of the finite element method library
- \Rightarrow stability (floating-point stability / numerical analysis stability)
 - prove and generalize well-known facts/algorithms/programs from the computer arithmetic community
- \Rightarrow basic blocks to build upon