19th IEEE Symposium on Computer Arithmetic (ARITH'19) Portland, Oregon, USA, June 8-10, 2009

# A new binary floating-point division algorithm and its software implementation on the ST231 processor

Claude-Pierre Jeannerod<sup>1,2</sup> Hervé Knochel<sup>4</sup> Christophe Monat<sup>4</sup> **Guillaume Revy**<sup>2,1</sup> Gilles Villard<sup>3,2,1</sup>

Université de Lvon<sup>2</sup> Arénaire Inria project-team (LIP, ENS Lvon)1 CNRS<sup>3</sup>













Compilation Expertise Centre (STMicroelectronics Grenoble)<sup>4</sup>



## Context and objectives

#### Context

- FLIP software library
  - → http://flip.gforge.inria.fr/
  - $\rightarrow$  support for floating-point arithmetic on integer processors
- low latency implementation of binary floating-point division
  - $\rightarrow$  targets a VLIW integer processor of the ST200 family
- no support of subnormal numbers
  - $\rightarrow$  input/output:  $\pm 0, \pm \infty$ , NaN or *normal* number

#### Objectives

- faster software implementation (compared to FLIP 0.3)
  - $\rightarrow$  expose instruction-level parallelism via bivariate polynomial evaluation
- correctly rounded
  - $\rightarrow$  rounding-to-nearest even

2/25

#### Notation and assumptions

• Input (x, y): two positive normal numbers

 $\rightarrow$  precision p, extremal exponents ( $e_{\min}, e_{\max}$ )

$$x = (-1)^{s_x} \cdot m_x \cdot 2^{e_x} \text{ with } \begin{cases} s_x \in \{0, 1\} \\ m_x = 1.m_{x,1} \dots m_{x,p-1} \in [1, 2) \\ e_x \in \{e_{\min}, \dots, e_{\max}\} \end{cases}$$

Computation: k-bit unsigned integers

 $\rightarrow$  register size k

• Example for binary32 format:  $(k, p, e_{max}) = (32, 24, 127)$ 

### Outline of the talk

Division via polynomial evaluation

Generation of an efficient evaluation program

Validation of the generated evaluation program

**Experimental results** 

Concluding remarks

## Outline of the talk

#### Division via polynomial evaluation

Generation of an efficient evaluation program

Validation of the generated evaluation program

**Experimental results** 

**Concluding remarks** 

## Division algorithm flowchart

Definition

$$c = \begin{cases} 1 & \text{if } m_x \ge m_y, \\ 0 & \text{if } m_x < m_y. \end{cases}$$

6/25

## Division algorithm flowchart

- ▶ Definition  $c = \begin{cases} 1 & \text{if } m_x \ge m_y, \\ 0 & \text{if } m_x < m_y. \end{cases}$
- Range reduction



## Division algorithm flowchart

- ▶ Definition  $c = \begin{cases} 1 & \text{if } m_x \geq m_y, \\ 0 & \text{if } m_x < m_y. \end{cases}$
- Range reduction



#### How to compute the correctly rounded significand $RN_p(\ell)$ ?

Guillaume Revy

## How to compute a correctly rounded significand ?

- Iterative methods (restoring, non-restoring, ...)
  - Oberman and Flynn (1997)
  - minimal instruction-level parallelism exposure, sequential algorithm

### How to compute a correctly rounded significand ?

- Iterative methods (restoring, non-restoring, ...)
  - Oberman and Flynn (1997)
  - minimal instruction-level parallelism exposure, sequential algorithm
- Multiplicative methods (Newton-Raphson, Goldschmidt)
  - Piñeiro and Bruguera (2002) Raina's Ph.D/FLIP (2006)
  - more instruction-level parallelism exposure
  - previous implementation of division (FLIP 0.3)

## How to compute a correctly rounded significand ?

- Iterative methods (restoring, non-restoring, ...)
  - Oberman and Flynn (1997)
  - minimal instruction-level parallelism exposure, sequential algorithm
- Multiplicative methods (Newton-Raphson, Goldschmidt)
  - Piñeiro and Bruguera (2002) Raina's Ph.D/FLIP (2006)
  - more instruction-level parallelism exposure
  - previous implementation of division (FLIP 0.3)
- Polynomial-based methods
  - Agarwal, Gustavson and Schmookler (1999)

     — univariate polynomial evaluation
  - Our approach
    - $\rightarrow$  single bivariate polynomial evaluation

## Truncated one-sided approximation

- See for example, Ercegovac and Lang (2004)
- 3 steps
  - 1. compute  $v = (01.v_1 \dots v_{k-2})$  such that

 $-2^{-p} \le \ell - v < 0$  that is implied by  $|(\ell + 2^{-p-1}) - v| < 2^{-p-1}$ 

- 2. truncate v after p fraction bits
- 3. obtain  $\mathsf{RN}_p(\ell)$  after possibly adding  $2^{-p}$

## Truncated one-sided approximation

- See for example, Ercegovac and Lang (2004)
- 3 steps
  - 1. compute  $v = (01.v_1 \dots v_{k-2})$  such that

 $-2^{-p} \le \ell - v < 0$  that is implied by  $|(\ell + 2^{-p-1}) - v| < 2^{-p-1}$ 

- 2. truncate v after p fraction bits
- 3. obtain  $\mathsf{RN}_p(\ell)$  after possibly adding  $2^{-p}$

#### How to compute the one-sided approximation v?

1. Consider  $\ell + 2^{-p-1}$  as the exact result of the function

 $F(s,t) = s/(1+t) + 2^{-p-1},$  at the points  $s^* = 2^{1-c}m_x$  and  $t^* = m_y - 1$ :  $\ell + 2^{-p-1} = F(s^*,t^*).$ 

1. Consider  $\ell + 2^{-p-1}$  as the exact result of the function

 $F(s,t) = s/(1+t) + 2^{-p-1},$  at the points  $s^* = 2^{1-c}m_x$  and  $t^* = m_y - 1$ :  $\ell + 2^{-p-1} = F(s^*,t^*).$ 

2. Approximate F(s,t) by a bivariate polynomial P(s,t)

 $P(s,t) = s \cdot a(t) + 2^{-p-1}.$ 

→ a(t): univariate polynomial approximant of 1/(1+t)→ approximation entails an error  $\epsilon_{approx}$ 

1. Consider  $\ell + 2^{-p-1}$  as the exact result of the function

 $F(s,t) = s/(1+t) + 2^{-p-1},$  at the points  $s^* = 2^{1-c}m_x$  and  $t^* = m_y - 1$ :  $\ell + 2^{-p-1} = F(s^*,t^*).$ 

2. Approximate F(s,t) by a bivariate polynomial P(s,t)

 $P(s,t) = s \cdot a(t) + 2^{-p-1}.$ 

- $\rightarrow a(t)$ : univariate polynomial approximant of 1/(1+t)
- ightarrow approximation entails an error  $\epsilon_{
  m approx}$
- 3. Evaluate P(s,t) by a well-chosen efficient evaluation program  $\mathcal{P}$

 $v = \mathcal{P}(s^*, t^*).$ 

 $\rightarrow$  evaluation entails an error  $\epsilon_{\text{eval}}$ 

1. Consider  $\ell + 2^{-p-1}$  as the exact result of the function

 $F(s,t) = s/(1+t) + 2^{-p-1},$  at the points  $s^* = 2^{1-c}m_x$  and  $t^* = m_y - 1$ :  $\ell + 2^{-p-1} = F(s^*,t^*).$ 

2. Approximate F(s,t) by a bivariate polynomial P(s,t)

 $P(s,t) = s \cdot a(t) + 2^{-p-1}.$ 

- $\rightarrow a(t)$ : univariate polynomial approximant of 1/(1+t)
- ightarrow approximation entails an error  $\epsilon_{
  m approx}$
- 3. Evaluate P(s,t) by a well-chosen efficient evaluation program  $\mathcal{P}$

 $v = \mathcal{P}(s^*, t^*).$ 

 $\rightarrow$  evaluation entails an error  $\epsilon_{\text{eval}}$ 

How to ensure that 
$$|(\ell + 2^{-p-1}) - v| < 2^{-p-1}$$
?

## Sufficient error bounds

#### Since by triangular inequality

$$|(\ell + 2^{-p-1}) - v| \le \mu \cdot \epsilon_{\text{approx}} + \epsilon_{\text{eva}}$$

with

$$\mu = \max\{s^*\} = \max\{2^{1-c}m_x\} = (4 - 2^{3-p})$$

### Sufficient error bounds

#### Since by triangular inequality

$$|(\ell + 2^{-p-1}) - v| \le \mu \cdot \epsilon_{\text{approx}} + \epsilon_{\text{eval}}$$

with

$$\mu = \max\{s^*\} = \max\{2^{1-c}m_x\} = (4 - 2^{3-p})$$

$$\mu \cdot \epsilon_{\text{approx}} + \epsilon_{\text{eval}} < 2^{-p-1}$$

Sufficient conditions can be obtained

$$\epsilon_{\mathsf{approx}} < 2^{-p-1}/\mu \quad \mathsf{and} \quad \epsilon_{\mathsf{eval}} < 2^{-p-1} - \mu \cdot \epsilon_{\mathsf{approx}}$$

## Outline of the talk

#### Division via polynomial evaluation

#### Generation of an efficient evaluation program

Validation of the generated evaluation program

**Experimental results** 

**Concluding remarks** 

#### Automatic generation of an efficient evaluation program

- $\blacktriangleright$  Evaluation program  $\mathcal P$  = main part of the full software implementation
  - $\rightarrow~$  dominates the cost
- By efficient, one means an evaluation program that
  - $\rightarrow\,$  reduces the evaluation latency
  - $\rightarrow$  reduces the number of multiplications
  - $\rightarrow$  is accurate enough

#### Automatic generation of an efficient evaluation program

 $\blacktriangleright$  Evaluation program  $\mathcal P$  = main part of the full software implementation

- $\rightarrow~$  dominates the cost
- By efficient, one means an evaluation program that
  - $\rightarrow\,$  reduces the evaluation latency
  - $\rightarrow$  reduces the number of multiplications
  - $\rightarrow$  is accurate enough
- Target architecture : ST231
  - $\rightarrow$  4-issue VLIW integer processor with at most 2 mul. per cycle
  - $\rightarrow$  latencies: addition = 1 cycle, multiplication = 3 cycles

#### Automatic generation of an efficient evaluation program

 $\blacktriangleright$  Evaluation program  $\mathcal P$  = main part of the full software implementation

- $\rightarrow~$  dominates the cost
- By efficient, one means an evaluation program that
  - $\rightarrow\,$  reduces the evaluation latency
  - $\rightarrow$  reduces the number of multiplications
  - $\rightarrow$  is accurate enough
- Target architecture : ST231
  - $\rightarrow$  4-issue VLIW integer processor with at most 2 mul. per cycle
  - $\rightarrow$  latencies: addition = 1 cycle, multiplication = 3 cycles

#### Which evaluation program to evaluate the polynomial P(s,t)?

$$P(s,t) = 2^{-p-1} + s \cdot \sum_{i=0}^{10} a_i t^i$$

• Horner's scheme:  $(3 + 1) \times 11 = 44$  cycles

 $\rightarrow$  sequential scheme, no instruction-level parallelism exposure

$$P(s,t) = 2^{-p-1} + s \cdot \sum_{i=0}^{10} a_i t^i$$

- Horner's scheme:  $(3 + 1) \times 11 = 44$  cycles
  - $\rightarrow$  sequential scheme, no instruction-level parallelism exposure
- Estrin's scheme: 20 cycles
  - $\rightarrow$  more instruction-level parallelism
  - $\rightarrow$  a last multiplication by s
  - $\rightarrow$  2 cycles save by distributing the multiplication by s in the evaluation of the univariate polynomial a(t)

$$P(s,t) = 2^{-p-1} + s \cdot \sum_{i=0}^{10} a_i t^i$$

- Horner's scheme:  $(3 + 1) \times 11 = 44$  cycles
  - $\rightarrow$  sequential scheme, no instruction-level parallelism exposure
- Estrin's scheme: 20 cycles
  - $\rightarrow$  more instruction-level parallelism
  - $\rightarrow$  a last multiplication by s
  - $\rightarrow\,$  2 cycles save by distributing the multiplication by s in the evaluation of the univariate polynomial a(t)

#### We can do much better.

- But how to explore the solution space and choose an efficient evaluation program ?
  - $\rightarrow$  interest of automatic generation

...

Similar to Harrison, Kubaska, Story and Tang (1999)

- Assumption
  - $\rightarrow$  unbounded parallelism
  - $\rightarrow\,$  latencies of arithmetic operators: + and  $\times\,$

Similar to Harrison, Kubaska, Story and Tang (1999)

- Assumption
  - $\rightarrow$  unbounded parallelism
  - $\rightarrow\,$  latencies of arithmetic operators: + and  $\times\,$
- Two sub-steps
  - 1. determine a target latency  $\tau$

ie. 
$$\tau = 3 \times \lceil \log_2(\deg(P)) \rceil + 1$$

2. generate automatically a set of evaluation trees, with height  $\leq \tau$ 

Similar to Harrison, Kubaska, Story and Tang (1999)

- Assumption
  - $\rightarrow$  unbounded parallelism
  - $\rightarrow\,$  latencies of arithmetic operators: + and  $\times\,$
- Two sub-steps
  - 1. determine a target latency  $\tau$

ie. 
$$\tau = 3 \times \lceil \log_2(\deg(P)) \rceil + 1$$

- 2. generate automatically a set of evaluation trees, with height  $\leq \tau$
- $\Rightarrow$  if no tree satisfies au then increase au and restart

Similar to Harrison, Kubaska, Story and Tang (1999)

- Assumption
  - $\rightarrow$  unbounded parallelism
  - $\rightarrow\,$  latencies of arithmetic operators: + and  $\times\,$
- Two sub-steps
  - 1. determine a target latency  $\tau$

ie. 
$$\tau = 3 \times \lceil \log_2(\deg(P)) \rceil + 1$$

- 2. generate automatically a set of evaluation trees, with height  $\leq \tau$
- $\Rightarrow$  if no tree satisfies  $\tau$  then increase  $\tau$  and restart
- ▶ Number of evaluation trees = extremely large → several filters

$$P(s,t) = 2^{-p-1} + s \cdot \sum_{i=0}^{10} a_i t^i$$

• first target latency 
$$\tau = 13$$

 $\rightarrow$  no tree found

$$P(s,t) = 2^{-p-1} + s \cdot \sum_{i=0}^{10} a_i t^i$$

- first target latency  $\tau = 13$  $\rightarrow$  no tree found
- second target latency  $\tau = 14$ 
  - $\rightarrow$  obtained in about 10 sec.



$$P(s,t) = 2^{-p-1} + s \cdot \sum_{i=0}^{10} a_i t^i$$

- first target latency  $\tau = 13$  $\rightarrow$  no tree found
- second target latency τ = 14 → obtained in about 10 sec.
- distribute the multiplication by s
  - $\rightarrow$  otherwise: 18 cycles
- too difficult to find such tree by hand



### Arithmetic operator choice

- Polynomial coefficients implemented in absolute value
- All intermediate values have constant sign
  - $\Rightarrow$  not store the sign: more accuracy

## Arithmetic operator choice

- Polynomial coefficients implemented in absolute value
- All intermediate values have constant sign
  - $\Rightarrow$  not store the sign: more accuracy
- Label evaluation trees by appropriate arithmetic operator: + or -

## Arithmetic operator choice

- Polynomial coefficients implemented in absolute value
- All intermediate values have constant sign
  - $\Rightarrow$  not store the sign: more accuracy
- Label evaluation trees by appropriate arithmetic operator: + or -
- If the sign of an intermediate value changes when the input varies then the evaluation tree is rejected
  - ⇒ implementation with certified interval arithmetic (MPFI)

## Practical scheduling checking

- Schedule the evaluation trees on a simplified model of a real target architecture
  - $\rightarrow\,$  operator costs, nb. issues, constraints on operators
  - $\rightarrow$  no syllables constraint

## Practical scheduling checking

- Schedule the evaluation trees on a simplified model of a real target architecture
  - $\rightarrow$  operator costs, nb. issues, constraints on operators
  - $\rightarrow$  no syllables constraint
- Check if no increase of latency in practice compared to the latency on unbounded parallelism
  - $\Rightarrow$  if practical latency > theoretical latency then the evaluation tree is rejected
  - $\Rightarrow$  implementation using naive list scheduling algorithm is enough

### Outline of the talk

Division via polynomial evaluation

Generation of an efficient evaluation program

Validation of the generated evaluation program

Experimental results

**Concluding remarks** 

► Approximation of 1/(1 + t) by truncated Remez' polynomial of degree 10



► Approximation of 1/(1 + t) by truncated Remez' polynomial of degree 10



Deduction of the evaluation error bound from 
elsiphered elsiph

$$\epsilon_{\rm eval} < 2^{-25} - (4 - 2^{-21}) \cdot 2^{-27.41...} \approx 2^{-26.9999...} \approx 7.4 \text{ e-}9.$$

- Case 1:  $m_x \ge m_y \rightarrow$  condition satisfied
- Case 2:  $m_x < m_y \rightarrow$  condition not satisfied

ie.  $s^* = 3.935581684112548828125$  and  $t^* = 0.97490441799163818359375$ 



- Case 1:  $m_x \ge m_y \rightarrow$  condition satisfied
- Case 2:  $m_x < m_y \rightarrow$  condition not satisfied

ie.  $s^* = 3.935581684112548828125$  and  $t^* = 0.97490441799163818359375$ 





- Case 1:  $m_x \ge m_y \rightarrow$  condition satisfied
- Case 2:  $m_x < m_y \rightarrow$  condition not satisfied

ie.  $s^* = 3.935581684112548828125$  and  $t^* = 0.97490441799163818359375$ 



- 1. determine an interval  $\mathcal{I}$  around this point
- 2. compute  $\epsilon_{approx}$  over  $\mathcal{I}$
- 3. determine an evaluation error bound  $\eta$

<sup>4.</sup> check if  $\epsilon_{\text{eval}} < \eta$  ?

## Evaluation program validation strategy

Find a splitting of the input interval into n subinterval(s) T<sup>(i)</sup>, and check that

$$\mu \cdot \epsilon_{\rm approx}^{(i)} + \epsilon_{\rm eval}^{(i)} < 2^{-p-1}$$

on each subinterval.

### Evaluation program validation strategy

Find a splitting of the input interval into n subinterval(s) T<sup>(i)</sup>, and check that

$$\mu \cdot \epsilon_{\text{approx}}^{(i)} + \epsilon_{\text{eval}}^{(i)} < 2^{-p-1}$$

on each subinterval.

- Implementation of the splitting by dichotomy
  - for each  $T^{(i)}$ 
    - 1. compute a certified approximation error bound  $\epsilon_{approx}^{(i)}$
    - 2. determine an evaluation error bound  $\epsilon_{\text{eval}}^{(i)}$
    - 3. check this bound
  - $\Rightarrow$  if this bound is not satisfied,  $\mathcal{T}^{(i)}$  is split up into 2 subintervals
    - implemented using Sollya (steps 1 and 2) and Gappa (step 3)

## Evaluation program validation strategy

► Find a splitting of the input interval into n subinterval(s) T<sup>(i)</sup>, and check that

$$\mu \cdot \epsilon_{\text{approx}}^{(i)} + \epsilon_{\text{eval}}^{(i)} < 2^{-p-1}$$

on each subinterval.

- Implementation of the splitting by dichotomy
  - for each  $T^{(i)}$ 
    - 1. compute a certified approximation error bound  $\epsilon_{approx}^{(i)}$
    - 2. determine an evaluation error bound  $\epsilon_{\text{eval}}^{(i)}$
    - 3. check this bound
  - $\Rightarrow$  if this bound is not satisfied,  $\mathcal{T}^{(i)}$  is split up into 2 subintervals
    - implemented using Sollya (steps 1 and 2) and Gappa (step 3)
- Example of binary32 implementation
  - $\rightarrow$  launched on a 64 processor grid
  - $\rightarrow$  36127 subintervals found in several hours ( $\approx$  5h.)

## Outline of the talk

Division via polynomial evaluation

Generation of an efficient evaluation program

Validation of the generated evaluation program

Experimental results

**Concluding remarks** 

## **Experimental results**

#### Performances on ST231

	Nb. of instructions	Latency (# cycles)	IPC	Code size (bytes)
rounding to nearest	86	27	3.18	416

- speed-up by a factor of about 1.78 in rounding to nearest compared to the previous implementation (48 cycles)
  - optimized implementation
  - efficient ST200 compiler (st200cc)
- high IPC value: confirms the parallel nature of our approach

## Outline of the talk

Division via polynomial evaluation

Generation of an efficient evaluation program

Validation of the generated evaluation program

Experimental results

Concluding remarks

25/25

## Concluding remarks

#### Contributions

- New approach for the implementation of binary floating-point division
  - $\rightarrow$  bivariate polynomial-based algorithm
  - $\rightarrow\,$  automatic generation and validation of efficient evaluation program
  - $\rightarrow$  implementation targeted ST231 VLIW integer processor
- Speed-up by a factor of about 1.78 in rounding to nearest compared to the previous implementation

#### Since then

- Extension to subnormal numbers support
  - $\rightarrow$  implementation in 31 cycles: 4 extra cycles
- Implementation of other functions

	Latency (# cycles)	IPC	Code size (bytes)	Speed-up
square root	21	2.47	276	2.38
reciprocal	22	2.59	336	1.75
reciprocal square root	29	2.24	368	2.27