Optimizing correctly-rounded reciprocal square roots for embedded VLIW cores

Claude-Pierre Jeannerod^{1, 2} and Gu

Guillaume Revy^{2, 1}

¹INRIA Grenoble - Rhône-Alpes (Arénaire project-team, LIP, ENS Lyon) ²Université de Lyon

INRIA LO UNIVERSITE DE LYON

1. Introduction

Context & Motivation

- Implementation of an efficient software support for IEEE 754 floating-point arithmetic on integer processors
- set of correctly-rounded mathematical operators, handling of subnormal numbers, and handling of special inputs
- development of the FLIP library [1] for the *binary32* floating-point format
- Optimized for the ST231 processor
- 4-issue VLIW integer processor from the ST200 processor family (STMicroelectronics) → no FPU

Purpose of our work

- Software implementation of correctly-rounded reciprocal square root $(x^{-1/2}) \rightarrow 29$ cycles
- frequently used in digital signal processing [6]

parallelism (ILP) of the ST231

- correctly-rounded implementation recommended by the latest revision of the IEEE 754 standard [2, §9.2]
- Optimized for the *binary32* format and the ST231 core
- Correctly-rounded RoundTiesToEven (rounding to nearest)
- Extension of our bivariate polynomial evaluation-based method introduced in [4] for square root

Efficiency achieved by exploiting at best the instruction-level

Example of application

Typical use of reciprocal square root = 3D vector normalization

$$[x, y, z] \quad \mapsto [x/w, y/w, z/w]$$

with

$$w = \sqrt{x^2 + y^2 + z^2}$$

- Without reciprocal square root:
- \rightarrow 1 sqrt (23 cycles) + 3 div (3× 32 cycles) = 119 cycles With reciprocal square root:

integer processor for embedded media systems
 → highly used in audio and video domains (HD-IPTV, cell phones, wireless terminals, PDAs)

2. ST231 architecture and compiler

ST231, a 4-issue VLIW 32-bit embedded integer architecture

- \blacktriangleright 4 parallel ALU's / 2 parallel pipelined $32\times32\rightarrow32\text{-bit}$ multipliers
- ► 1 leading zero counter
- \blacktriangleright Predicate execution \rightarrow select instruction to remove branch penalty
- ► 64 general purpose 32-bit registers / 8 1-bit *branch* (condition) registers
- Efficient 32-bit immediate operand encoding

ST231 Compiler

- Open64 compiler technology
- Instruction level parallelism extractor and scheduler
- Select-based *if-conversion* \rightarrow straightline assembly code
- \rightarrow sequences of *select* instructions instead of costly control flow
- ► Linear Assembly Optimizer (LAO): generates schedule very close to the optimal



- \rightarrow 1 rsqrt (29 cycles) + 3 mul (3× 21 cycles) = 92 cycles
- ► Latency reduction by over 20 %

3. Some properties of reciprocal square root

Handling of special operands $x \in \{x < 0, \pm 0, \pm \infty, NaN\}$

- Filter out special operands using the standard binary interchange encoding format
- ► Compute special results required by [2] in parallel with the generic case
- **Positive finite operand** x (precision $p \ge 2$)
- **Input:** binary32 floating-point number: $x = m \cdot 2^e = m' \cdot 2^{e'}$, with $m' = m \cdot 2^{\lambda}$, $e' = e \lambda$, and

$$e' \in \mathbb{N} \cap [e_{\min} - p + 1, e_{\max}] \quad \text{and} \quad m' = 1.m_{\lambda+1} \cdots m_{p-1}$$

Output: $RN(x^{-1/2}) = correct rounding-to-nearest of <math>x^{-1/2}$

$$x^{-1/2} = \ell \cdot 2^d \quad \text{and} \quad \mathsf{RN}(x^{-1/2}) = \mathsf{RN}(\ell) \cdot 2^d, \quad \text{and} \quad \ell \in [1, 2] \quad \text{and} \quad e_{\min} \le d \le e_{\max}$$

and where

$$d = -(e+1-c)/2, \quad \ell = s\sqrt{2/(1+t)}, \quad \text{and} \quad c = [e' \text{ is odd}]$$

Two useful properties

- $ightarrow x^{-1/2}$ falls in the range of normal floating-point numbers
- $ightarrow x^{-1/2}$ cannot be halfway between two consecutive floating-point numbers

4. How to approximate the exact value ℓ ?

Existing method in FLIP 0.3: multiplicative method

Initial approximation: degree-3 univariate polynomial

Refinement by Goldschmidt's iteration [3]

Our approach: one-sided truncated approximation [4]

 \blacktriangleright Approximation of ℓ from above by v

 $2^{-p} < \ell - v \le 0$ and $|\ell + 2^{-p-1} - v| < 2^{-p-1}$

Computation of u =truncation of v after p fraction bits

 $0 \le v - u < 2^{-p}$ and $|\ell - u| < 2^{-p}$

Approximation v = result of the evaluation of a single bivariate polynomial $P(s,t) = 2^{-p-1} + s \cdot a(t)$

with a(t) a degree-9 truncated Remez approximant computed with $\ensuremath{\mathsf{Sollya}}$

How to evaluate P(s,t) efficiently?

Horner's rule: 38 cycles, no ILP exposure

Efficient and certified parenthesization automatically generated using CGPE [5]

• Reduction of evaluation latency

 \rightarrow 13 cycles on unbounded parallelism, 14 cycles on ST231

• Evaluation error checked with Gappa

 \rightarrow ensure correct rounding

CGPE generation flowchart



uint32_tr	<pre>sqrt_eval(uint32_t T, uint32_</pre>	t S)
<pre>uint32_t r uint32_t r uint32</pre>	<pre>0 = mul(T, 0x5a82685d); 1 = 0xb504f31f - r0; 2 = mul(S, r1); 3 = 0x00000020 + r2; 4 = mul(T, T); 5 = mul(S, r4); 6 = mul(T, 0x386fd5f4); 7 = 0x43df72f7 - r6; 8 = mul(r5, r7); 9 = r3 + r8; 10 = mul(T, 0x28724100); 11 = 0x308b1798 - r10; 12 = mul(r4, r11); 13 = mul(r5, r12); 14 = r9 + r13; 15 = mul(r4, r4); 16 = mul(r5, r15); 17 = mul(T, 0x106c5cd9); 18 = 0x1d7bf968 - r17; 19 = mul(T, 0x00fa9aa4); 20 = 0x05dfffa4 - r19; 21 = mul(r4, r20); 22 = r18 + r21; 23 = mul(r16, r22); 24 = r14 + r23; 4; scheduling</pre>	<pre>// 1.31 // 1.31 // 2.30 // 2.30 // 0.32 // 1.31 // 2.30 // 2.30</pre>
 	Is1 Is2 Is3 Is4	/
Cycle 0 Cycle 1 Cycle 2 Cycle 3 Cycle 4 Cycle 5 Cycle 6 Cycle 7 Cycle 8 Cycle 9 Cycle 10 Cycle 11 Cycle 12 Cycle 13	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	

6. Validation and performances

5. How to deduce $RN(\ell)$ from the approximation v?

Deduce RN(ℓ): decide whether $u \ge \ell$, which is equivalent to $u^2 \cdot (1+t) \ge 2 \cdot s^2$, with $2 \cdot s^2 \in \{2,4\}$

▶ u, t, and $2 \cdot s^2$ exactly representable with 32 bits



Computing the first 64 bits of the exact product are enough

Test done on the first 32 bits of the exact product

Exhaustive comparison with Glibc and MPFR

Performances on ST231	FLIP 1.0		FLIP 0.3	Speedup
	(with subnormal)	(without subnormal)	(without subnormal)	
	29 cycles	28 cycles	67 cycles	2.3

Interest of the specialization of	Code se
the reciprocal square root operator	for com
	div(1.0

Code sequence used		Number N of	Latency L	N/L
for computing $x^{-1/2}$		instructions	(cycles)	
	<pre>div(1.0f,sqrt(x))</pre>	147 [124]	53 [47]	2.7 [2.6]
	inv(sqrt(x))	121 [115]	49 [47]	2.5 [2.4]
	rsqrt(x)	68 [63]	29 [28]	2.3 [2.2]

Some references

[1] FLIP - Floating-point Library for Intger Processors. Available at http://flip.gforge.inria.fr/.

[2] IEEE standard for floating-point arithmetic. IEEE Std. 754-2008, pp.1-58, Aug. 29 2008.

[3] M. D. Ercegovac and T. Lang. *Digital Arithmetic*. Morgan Kaufmann, 2004.

[4] Claude-Pierre Jeannerod, Hervé Knochel, Christophe Monat, and Guillaume Revy. Computing floating-point square roots via bivariate polynomial evaluation. Technical Report RR2008-38, LIP, October 2008.

[5] Guillaume Revy. CGPE - Code Generation for Polynomial Evaluation. Available at http://cgpe.gforge.inria.fr/.

[6] Michael J. Schulte and Kent E. Wires. High-speed inverse square roots. In Proceedings of the 14th IEEE Symposium on Computer Arithmetic (ARITH-14), pages 124–131. IEEE Computer Society, 1999.

43rd Asilomar Conference on Signals, Systems and Computers – Pacific Grove, CA, USA, 1-4 November, 2009.