Automatic Generation of Fast and Certified Code for Polynomial Evaluation

Guillaume Revy

DALI project-team

Université de Perpignan Via Domitia LIRMM (CNRS: UMR 5506 - UM2)







Joint work with Christophe Mouilleron (LIP, ENS Lyon)

- Embedded systems are ubiquitous
 - microprocessors dedicated to one or a few specific tasks
 - satisfy constraints: area, energy consumption, conception cost

- Embedded systems are ubiquitous
 - microprocessors dedicated to one or a few specific tasks
 - satisfy constraints: area, energy consumption, conception cost
- Some embedded systems do not have any FPU (floating-point unit)



- Embedded systems are ubiquitous
 - microprocessors dedicated to one or a few specific tasks
 - satisfy constraints: area, energy consumption, conception cost
- Some embedded systems do not have any FPU (floating-point unit)





- Highly used in audio and video applications
 - demanding on floating-point computations

- Embedded systems are ubiquitous
 - microprocessors dedicated to one or a few specific tasks
 - satisfy constraints: area, energy consumption, conception cost
- Some embedded systems do not have any FPU (floating-point unit)



- Highly used in audio and video applications
 - demanding on floating-point computations

On the one side: the IEEE 754-2008 standard, ...

Definition of IEEE floating-point arithmetic

- floating-point formats: single precision, double precision, ...
- special values: ±0, ±∞, NaN
- 4 rounding modes: to nearest even, upward, downward, and toward zero
- mathematical function behavior
 - \rightarrow special input (ex: $\sqrt{-0} = -0$)
 - \rightarrow requires / recommends correct rounding

- make computations reproducible
- and make results architecture-independent

... on the other side: the ST231 processor



- 4-issue VLIW 32-bit integer processor
 → no FPU
- Parallel execution unit
 - 4 integer ALUs
 - 2 pipelined multipliers $32 \times 32 \rightarrow 32$
- Latencies: ALU = 1 cycle / Mul = 3 cycles

... on the other side: the ST231 processor



- 4-issue VLIW 32-bit integer processor
 → no FPU
- Parallel execution unit
 - 4 integer ALUs
 - 2 pipelined multipliers $32 \times 32 \rightarrow 32$
- Latencies: ALU = 1 cycle / Mul = 3 cycles

VLIW (Very Long Instruction Word)

- \rightarrow instructions grouped into bundles
- \rightarrow Instruction-Level Parallelism (ILP) explicitly exposed by the compiler

- This work takes mainly part in the context of the development of FLIP
 - \hookrightarrow software support for binary32 floating-point arithmetic on integer processors
- Underlying problem: development "by hand"
 - long and tedious, error prone
 - new target ? new floating-point format ?

- This work takes mainly part in the context of the development of FLIP
 - → software support for binary32 floating-point arithmetic on integer processors
- Underlying problem: development "by hand"

 - long and tedious, error prone
 new target ? new floating-point format ?

- This work takes mainly part in the context of the development of FLIP
 - → software support for binary32 floating-point arithmetic on integer processors
- Underlying problem: development "by hand"

 - long and tedious, error prone
 new target ? new floating-point format ?

- Current challenge: tools and methodologies for the automatic generation of efficient and certified programs
 - optimized for a given format, for the target architecture

- Some works on code generation and transformation:
 - \hookrightarrow code generators: hardware (FloPoCo) and software (Sollya, Metalibm)
 - \hookrightarrow code transformation for increasing numerical accuracy [Martel, 2009]
- LEMA project [Lefèvre et al., 2010]: language and library
 - \hookrightarrow design easily a generation toolchain

- Some works on code generation and transformation:
 - \hookrightarrow code generators: hardware (FloPoCo) and software (Sollya, Metalibm)
 - \hookrightarrow code transformation for increasing numerical accuracy [Martel, 2009]
- LEMA project [Lefèvre et al., 2010]: language and library
 - \hookrightarrow design easily a generation toolchain
- Spiral project: hardware and software code generation for DSP algorithms

Can we teach computers to write fast libraries?

- Some works on code generation and transformation:
 - \hookrightarrow code generators: hardware (FloPoCo) and software (Sollya, Metalibm)
 - \hookrightarrow code transformation for increasing numerical accuracy [Martel, 2009]
- LEMA project [Lefèvre et al., 2010]: language and library
 - \hookrightarrow design easily a generation toolchain
- Spiral project: hardware and software code generation for DSP algorithms Can we teach computers to write fast libraries?

Our tool: CGPE (Code Generation for Polynomial Evaluation)

In the particular case of **polynomial evaluation**, can we teach computers to write **fast and certified** codes, for a given target and optimized for a given format?

 \hookrightarrow adding a systematic certification step

Basic blocks for implementing correctly-rounded operators



Objectives

- → Low latency, correctly-rounded implementations
- → ILP exposure

Basic blocks for implementing correctly-rounded operators















Outline of the talk

- 1. Background on polynomial evaluation
- 2. The CGPE tool
- 3. Experimental results
- 4. Conclusion

Outline of the talk

1. Background on polynomial evaluation

2. The CGPE tool

- 3. Experimental results
- 4. Conclusion

- Compute fast and certified schemes for evaluating a polynomial $P(x, y) = \alpha + y \cdot a(x)$
 - $\rightarrow\,$ using only additions and multiplications
 - $\rightarrow\,$ reducing the evaluation latency on unbounded parallelism

- Compute fast and certified schemes for evaluating a polynomial $P(x,y) = \alpha + y \cdot a(x)$
 - \rightarrow using only additions and multiplications
 - $\rightarrow\,$ reducing the evaluation latency on unbounded parallelism
- Evaluation program = main part of the full software implementation
 - \rightarrow dominates the cost
 - \hookrightarrow make it as efficient as possible

- Compute fast and certified schemes for evaluating a polynomial $P(x,y) = \alpha + y \cdot a(x)$
 - $\rightarrow\,$ using only additions and multiplications
 - $\rightarrow\,$ reducing the evaluation latency on unbounded parallelism
- Evaluation program = main part of the full software implementation
 - \rightarrow dominates the cost
 - \hookrightarrow make it as efficient as possible

Two families of algorithms

- algorithms with coefficient adaptation: Knuth and Eve (60's), Paterson and Stockmeyer (1973), ...
 - \rightarrow ill-suited in the context of fixed-point arithmetic
- algorithms without coefficient adaptation

- Compute fast and certified schemes for evaluating a polynomial $P(x,y) = \alpha + y \cdot a(x)$
 - $\rightarrow\,$ using only additions and multiplications
 - $\rightarrow\,$ reducing the evaluation latency on unbounded parallelism
- Evaluation program = main part of the full software implementation
 - \rightarrow dominates the cost
 - \hookrightarrow make it as efficient as possible

Two families of algorithms

- algorithms with coefficient adaptation: Knuth and Eve (60's), Paterson and Stockmeyer (1973), ...
 - \rightarrow ill-suited in the context of fixed-point arithmetic
- algorithms without coefficient adaptation



Horner's rule

- → 3 additions
 3 multiplications
- \hookrightarrow latency: 12 cycles
 - optimal in terms of multiplication number [Pan, 1966], [Borodin, 1971],
 - ⊖ fully sequential



Second-order Horner's rule

- → 3 additions
 4 multiplications
- \hookrightarrow latency: 11 cycles
 - some ILP exposure
- ⊖ subparts evaluated in a fully sequential way: at most 2 ways used



Estrin's rule

- → 3 additions
 4 multiplications
- → latency: 8 cycles
 - ⊕ "divide and conquer" strategy
- ⊕ more ILP exposure



Definition of evaluation schemes



1

1

Definition of evaluation schemes



Definition of evaluation schemes



Remarks on polynomial evaluation

- There are several other schemes for evaluating a polynomial a(x)
 - ► can be adapted for bivariate polynomial $P(x, y) = \alpha + y \cdot a(x)$
- Constant number of +, while number of × is non-constant
 - ► reducing the latency ⇔ increasing the number of × to expose ILP
 - trade-off latency / number of multiplications

Remarks on polynomial evaluation

- There are several other schemes for evaluating a polynomial a(x)
 - ► can be adapted for bivariate polynomial $P(x, y) = \alpha + y \cdot a(x)$
- Constant number of +, while number of × is non-constant
 - ► reducing the latency ⇔ increasing the number of × to expose ILP
 - trade-off latency / number of multiplications

Evaluation error

- different theoretical error bounds
- difference between numerical quality in practice [Revy, 2006]

Remarks on polynomial evaluation

- There are several other schemes for evaluating a polynomial a(x)
 - ► can be adapted for bivariate polynomial $P(x, y) = \alpha + y \cdot a(x)$
- Constant number of +, while number of × is non-constant
 - ► reducing the latency ⇔ increasing the number of × to expose ILP
 - trade-off latency / number of multiplications

Evaluation error

- different theoretical error bounds
- difference between numerical quality in practice [Revy, 2006]

\rightsquigarrow We need a tool for exploring the space of evaluation schemes.

How many schemes for evaluating a polynomial?

n	$\mu_n ightarrow a(x)$	$\mu'_n ightarrow lpha + y \cdot a(x)$	
1	1	10	
2	7	481	
3	163	88384	
4	11602	57363910	
5	2334244	122657263474	
6	1304066578	829129658616013	
7	1972869433837	17125741272619781635	
8	8012682343669366	1055157310305502607244946	
9	86298937651093314877	190070917121184028045719056344	
10	2449381767217281163362301	98543690848554380947490522591191672	

How many schemes for evaluating a polynomial?

n	$\mu_n ightarrow a(x)$	$\mu'_n ightarrow lpha + y \cdot a(x)$	wn	
1	1	10	1	
2	7	481	1	
3	163	88384	1	
4	11602	57363910	2	
5	2334244	122657263474	3	
6	1304066578	829129658616013	6	
7	1972869433837	17125741272619781635	11	
8	8012682343669366	1055157310305502607244946	23	
9	86298937651093314877	190070917121184028045719056344	46	
10	2449381767217281163362301	98543690848554380947490522591191672	98	

Two well-known special cases

the number of evaluation schemes for xⁿ [Wedderburn, Etherington]

$$w_n \sim rac{\eta \xi^n}{n^{3/2}}$$
 or $\left\{ egin{array}{c} \xi pprox 2.48325 \ \eta pprox 0.31877 \end{array}
ight.$ [Otter, 1948],

How many schemes for evaluating a polynomial?

п	$\mu_n ightarrow a(x)$	$\mu'_n ightarrow lpha + y \cdot a(x)$	wn	(2 <i>n</i> – 1)!!
1	1	10	1	1
2	7	481	1	3
3	163	88384	1	15
4	11602	57363910	2	105
5	2334244	122657263474	3	945
6	1304066578	829129658616013	6	10395
7	1972869433837	17125741272619781635	11	135135
8	8012682343669366	1055157310305502607244946	23	2027025
9	86298937651093314877	190070917121184028045719056344	46	34459425
10	2449381767217281163362301	98543690848554380947490522591191672	98	654729075

Two well-known special cases

the number of evaluation schemes for xⁿ [Wedderburn, Etherington]

$$w_n \sim rac{\eta \xi^n}{n^{3/2}}$$
 or $\left\{ egin{array}{c} \xi pprox 2.48325 \ \eta pprox 0.31877 \end{array}
ight.$ [Otter, 1948],

• the number of evaluation schemes for $\sum_{i=1}^{n} a_i \operatorname{est} (2n-1)!! \sim \sqrt{2} \left(\frac{2n}{e}\right)^n$.

Schemes of low evaluation latency





Total number of schemes: 2334244

- → minimal latency for degree-5 univariate polynomial: 10 cycles
- → number of schemes of minimal latency: 36

Outline of the talk

1. Background on polynomial evaluation

- 2. The CGPE tool
- 3. Experimental results
- 4. Conclusion

Overview of CGPE and related works

- Goal of CGPE [Mouilleron and Revy, 2011]: automate the design of fast and certified C codes for evaluating univariate/bivariate polynomials
 - in fixed-point arithmetic
 - by using the target architecture features (as much as possible)

Remarks:

- \rightsquigarrow fast = that reduces the evaluation latency on a given target
- certified = we can bound the error entailed by the evaluation within the given target's arithmetic

Overview of CGPE and related works

Some related works

- [Cheung et al., 2005] and [Lee and Villasenor, 2009]: methodology for implementing automatically mathematical function in a given precision
 - $\ominus\;$ based on small degree polynomial evaluation using Horner's \rightarrow no ILP
- [Harrison et al., 1999]: method for generating optimal evaluation scheme to evaluate univariate polynomials on Itanium[®] using fma
 - \ominus ST231 has only addition and multiplication, but no fma
- [Green, 2002]: brute force method for generating polynomial evaluation schemes using at best SIMD instructions of the processor of PlayStation[®] 2
 - \ominus objective: generation at compile-time \rightarrow *brute force* method is unfeasible

Global architecture of CGPE

Input of CGPE

- 1. polynomial coefficients and variables: value intervals, fixed-point format, ...
- 2. set of criteria: maximum error bound and bound on latency (or the lowest)
- 3. delay of one of the variable
- 4. some architectural constraints: operator cost, parallelism, ...

Global architecture of CGPE

Input of CGPE

- 1. polynomial coefficients and variables: value intervals, fixed-point format, ...
- 2. set of criteria: maximum error bound and bound on latency (or the lowest)
- 3. delay of one of the variable
- 4. some architectural constraints: operator cost, parallelism, ...

CGPE works in two steps:

- 1. computation of evaluation schemes: reducing evaluation latency on unbounded parallelism and exposing as much ILP as possible
- 2. selection among the generated schemes, according to different criteria:
 - evaluation using only unsigned fixed-point arithmetic
 - scheduling feasible on ST231
 - evaluation error bound satisfying the required error bound

At the end: CGPE automatically writes C codes with accuracy certificates

Heuristics in DAG set computation

Determination of the minimal target latency on unbounded parallelism

- gives a good estimation of the best evaluation latency of the polynomial on the target architecture
- takes some problem parameters (operator costs, delay, ...) into account

Heuristics in DAG set computation

Determination of the minimal target latency on unbounded parallelism

- gives a good estimation of the best evaluation latency of the polynomial on the target architecture
- takes some problem parameters (operator costs, delay, ...) into account
- Non exhaustive computation of evaluation schemes
 - elimination of the schemes that do not satisfy latency constraint
 - Imitation to some splittings: evaluation of high and low parts separately
 - restriction to N schemes at each step of the computation

Heuristics in DAG set computation

Determination of the minimal target latency on unbounded parallelism

- gives a good estimation of the best evaluation latency of the polynomial on the target architecture
- takes some problem parameters (operator costs, delay, ...) into account
- Non exhaustive computation of evaluation schemes
 - elimination of the schemes that do not satisfy latency constraint
 - Imitation to some splittings: evaluation of high and low parts separately
 - restriction to N schemes at each step of the computation
- At the end of the computation: set of DAG evaluating the input polynomial and satisfying the latency constraint on unbounded parallelism.

Filters for adding numerical constraints

- 1. Arithmetic operator choice
 - ensure that all intermediate variables are of constant sign
 - avoid an extra cost due to sign handling / gain 1 bit of accuracy

Filters for adding numerical constraints

- 1. Arithmetic operator choice
 - ensure that all intermediate variables are of constant sign
 - avoid an extra cost due to sign handling / gain 1 bit of accuracy
- 2. Scheduling on a simplified model of the target (like the ST231)
 - constraints of architecture: cost of operators, instruction bundling, ...
 - delays on variables

Filters for adding numerical constraints

- 1. Arithmetic operator choice
 - ensure that all intermediate variables are of constant sign
 - avoid an extra cost due to sign handling / gain 1 bit of accuracy
- 2. Scheduling on a simplified model of the target (like the ST231)
 - constraints of architecture: cost of operators, instruction bundling, ...
 - delays on variables
- 3. Evaluation error bound checking
 - straightline polynomial evaluation program
 - "C code certification" using Gappa
 - \rightsquigarrow we can bound the evaluation error in integer arithmetic

Outline of the talk

1. Background on polynomial evaluation

2. The CGPE tool

- 3. Experimental results
- 4. Conclusion

	x ^{1/2}	x ^{-1/2}	x ^{1/3}	x ^{-1/3}	$\log_2(1+x)$	$\frac{1}{\sqrt{1+x^2}}$	$\frac{\exp(1+x)}{1+x}$
Degree (d_x, d_y)	(8,1)	(9,1)	(8,1)	(9,1)	(6,0)	(7,0)	(10,0)
Target / Minimal latency	13/13	13/?	16/16	16 / 16	10 / 11	10 / 11	13 / 13
Achieved latency	13	14	16	16	11	11	13
Scheme computation	195ms	73ms	26s	25s	17ms	10ms	40ms
	[50]	[50]	[50]	[50]	[50]	[50]	[50]
Arithmetic operator choice	3ms	3ms	7ms	11ms	1ms	2ms	3ms
	[35]	[29]	[30]	[26]	[2]	[12]	[27]
Scheduling checking	16s	1m33s	43ms	439ms	2ms	64ms	49s
	[11]	[1]	[30]	[24]	[1]	[5]	[5]
Certification (Gappa)	10s	1s	27s	27s	230ms	1s	7s
	[11]	[1]	[30]	[24]	[1]	[5]	[4]
Total time ($pprox$)	27s	1m35s	55s	53s	1s	2s	57s

	x ^{1/2}	x ^{-1/2}	x ^{1/3}	x ^{-1/3}	$\log_2(1+x)$	$\frac{1}{\sqrt{1+x^2}}$	$\frac{\exp(1+x)}{1+x}$
Degree (d_x, d_y)	(8,1)	(9,1)	(8,1)	(9,1)	(6,0)	(7,0)	(10,0)
Target / Minimal latency	13/13	13/?	16 / 16	16 / 16	10 / 11	10 / 11	13 / 13
Achieved latency	13	14	16	16	11	11	13
Scheme computation	195ms	73ms	26s	25s	17ms	10ms	40ms
	[50]	[50]	[50]	[50]	[50]	[50]	[50]
Arithmetic operator choice	3ms	3ms	7ms	11ms	1ms	2ms	3ms
	[35]	[29]	[30]	[26]	[2]	[12]	[27]
Scheduling checking	16s	1m33s	43ms	439ms	2ms	64ms	49s
	[11]	[1]	[30]	[24]	[1]	[5]	[5]
Certification (Gappa)	10s	1s	27s	27s	230ms	1s	7s
	[11]	[1]	[30]	[24]	[1]	[5]	[4]
Total time (\approx)	27s	1m35s	55s	53s	1s	2s	57s

1. Impact of the target latency on the first step of the generation

	x ^{1/2}	x ^{-1/2}	x ^{1/3}	x ^{-1/3}	$\log_2(1+x)$	$\frac{1}{\sqrt{1+x^2}}$	$\frac{\exp(1+x)}{1+x}$
Degree (d_x, d_y)	(8,1)	(9,1)	(8,1)	(9,1)	(6,0)	(7,0)	(10,0)
Target / Minimal latency	13 / 13	13/?	16/16	16 / 16	10 / 11	10 / 11	13 / 13
Achieved latency	13	14	16	16	11	11	13
Scheme computation	195ms	73ms	26s	25s	17ms	10ms	40ms
	[50]	[50]	[50]	[50]	[50]	[50]	[50]
Arithmetic operator choice	3ms	3ms	7ms	11ms	1ms	2ms	3ms
	[35]	[29]	[30]	[26]	[2]	[12]	[27]
Scheduling checking	16s	1m33s	43ms	439ms	2ms	64ms	49s
	[11]	[1]	[30]	[24]	[1]	[5]	[5]
Certification (Gappa)	10s	1s	27s	27s	230ms	1s	7s
	[11]	[1]	[30]	[24]	[1]	[5]	[4]
Total time (\approx)	27s	1m35s	55s	53s	1s	2s	57s

- 1. Impact of the target latency on the first step of the generation
- 2. What may dominate the cost: scheduling and certification using Gappa

	x ^{1/2}	x ^{-1/2}	x ^{1/3}	x ^{-1/3}	$\log_2(1+x)$	$\frac{1}{\sqrt{1+x^2}}$	$\frac{\exp(1+x)}{1+x}$
Degree (d_x, d_y)	(8,1)	(9,1)	(8,1)	(9,1)	(6,0)	(7,0)	(10,0)
Target / Minimal latency	13 / 13	13 / ?	16 / 16	16 / 16	10 / <mark>11</mark>	10 / <mark>11</mark>	13 / 13
Achieved latency	13	14	16	16	11	11	13
Scheme computation	195ms	73ms	26s	25s	17ms	10ms	40ms
	[50]	[50]	[50]	[50]	[50]	[50]	[50]
Arithmetic operator choice	3ms	3ms	7ms	11ms	1ms	2ms	3ms
	[35]	[29]	[30]	[26]	[2]	[12]	[27]
Scheduling checking	16s	1m33s	43ms	439ms	2ms	64ms	49s
	[11]	[1]	[30]	[24]	[1]	[5]	[5]
Certification (Gappa)	10s	1s	27s	27s	230ms	1s	7s
	[11]	[1]	[30]	[24]	[1]	[5]	[4]
Total time (\approx)	27s	1m35s	55s	53s	1s	2s	57s

- 1. Impact of the target latency on the first step of the generation
- 2. What may dominate the cost: scheduling and certification using Gappa
- 3. Optimality of some generated codes, in terms of evaluation latency

Outline of the talk

1. Background on polynomial evaluation

2. The CGPE tool

3. Experimental results

4. Conclusion

Conclusions

Code generation for fast and certified polynomial evaluation

- in fixed-point arithmetic
- methodologies and tools to automate polynomial evaluation implementation
- heuristics and techniques for generating quickly fast and certified C codes
- implemented in the tool CGPE (Code Generation for Polynomial Evaluation)

```
http://cgpe.gforge.inria.fr/
```

Speed-up significantly the development time of mathematical library

• CGPE: allows to write and certify automatically \approx 50 % of the codes of FLIP

Current work and perspectives

Current work

- precomputation in order to help the DAG set computation in choosing the appropriate splittings: the ones leading to DAGs with optimal latency on unbounded parallelism
- earlier DAG elimination, by checking accuracy during generation step

Current work and perspectives

Current work

- precomputation in order to help the DAG set computation in choosing the appropriate splittings: the ones leading to DAGs with optimal latency on unbounded parallelism
- earlier DAG elimination, by checking accuracy during generation step

Perspectives

- extend CGPE to handle floating-point arithmetic,
- make CGPE more general to tackle other problems, like evaluation of a polynomial at a matrix point.



Borodin, A. (1971).

Horner's rule is uniquely optimal.

Theory of machines and computations, pages 45–58.



Automating custom-precision function evaluation for embedded processors.

In CASES'05: Proceedings of the 2005 international conference on Compilers, architectures and synthesis for embedded systems, pages 22-31, New York, NY, USA. ACM.

Green, R. (2002).

Faster Math Functions.

Tutorial at Game Developers Conference.

Harrison, J., Kubaska, T., Story, S., and Tang, P. (1999).

The computation of transcendental functions on the IA-64 architecture.

Intel Technology Journal, 1999-Q4:1-7.



Lee, D.-U. and Villasenor, J. D. (2009).

Optimized Custom Precision Function Evaluation for Embedded Processors.

IEEE Transactions on Computers, 58(1):46-59.

Lefèvre, V., Théveny, P., de Dinechin, F., Jeannerod, C.-P., Mouilleron, C., Pfannholzer, D., and Revol, N. (2010).

LEMA: towards a language for reliable arithmetic.

ACM Communications in Computer Algebra, 44:41–52.

Martel, M. (2009).

Enhancing the Implementation of Mathematical Formulas for Fixed-Point and Floating-Point Arithmetics.

In Journal of Formal Methods in System Design, volume 35, pages 265–278. Springer.



Mouilleron, C. and Revy, G. (2011).

Automatic Generation of Fast and Certified Code for Polynomial Evaluation.

In <u>Proc. of the 20th IEEE Symposium on Computer Arithmetic (ARITH'20)</u>, Tuebingen, Germany.

Otter, R. (1948).

The number of trees.

The Annals of Mathematics, 49(3):pp. 583–599.

Guillaume Revy (Groupe de travail PEQUAN – July 7, 2011)



Pan, V. Y. (1966).

Methods of Computing Values of Polynomials.

Russian Mathematical Surveys, 21(1):105–136.

Revy, G. (2006).

Analyse et implantation d'algorithmes rapides pour l'évaluation polynomiale sur les nombres flottants.

Master's thesis, École normale supérieure de Lyon, 46 allée d'Italie, F-69364 Lyon cedex 07, France.