*Ph.D. defense* Perpignan, December 10<sup>th</sup>, 2014

# Synthesis of certified programs in fixed-point arithmetic, and its application to linear algebra basic blocks

#### Mohamed Amine Najahi

Advisors: Matthieu Martel and Guillaume Revy

Univ. Perpignan Via Domitia, DALI project-team Univ. Montpellier 2, LIRMM, UMR 5506 CNRS, LIRMM, UMR 5506









M. A. Najahi (DALI UPVD/LIRMM, UM2, CNR

The embedded systems market is growing

### Around 30% of an airplane's cost, around 40% of a car's cost



The embedded systems market is growing

### Around 30% of an airplane's cost, around 40% of a car's cost



Embedded systems are often dedicated to computational tasks

The embedded systems market is growing

### Around 30% of an airplane's cost, around 40% of a car's cost



- Embedded systems are often dedicated to computational tasks
- Embedded systems face multiple constraints
  - efficiency
  - cost

- limits on hardware resources
- limits on energy consumption

The embedded systems market is growing

### Around 30% of an airplane's cost, around 40% of a car's cost



- Embedded systems are often dedicated to computational tasks
- Embedded systems face multiple constraints
  - efficiency cost

►

- limits on hardware resources
- limits on energy consumption

#### How to implement these computational tasks in embedded systems?

Floating-point computations	Fixed-point computations

### Floating-point computations

- © Easy and fast to implement
- Easily portable [?]

### **Fixed-point computations**

- © Tedious and time consuming to implement
  - > 50% of design time [?]

### Floating-point computations

- © Easy and fast to implement
- Easily portable [?]
- C Requires dedicated hardware
- Slow if emulated in software

### **Fixed-point computations**

- © Tedious and time consuming to implement
  - >50% of design time [?]
- © Relies only on integer instructions
- Sefficient

# Floating-point computations

- © Easy and fast to implement
- Easily portable [?]
- C Requires dedicated hardware
- Slow if emulated in software

### Fixed-point computations

- © Tedious and time consuming to implement
  - > 50% of design time [?]
- © Relies only on integer instructions
- © Efficient

#### Embedded systems targets





 $\mu$ -controllers

DSF



→ have efficient integer instructions

Fixed-point arithmetic is well suited for embedded systems

# Floating-point computations

- Easy and fast to implement  $\odot$
- Easily portable [?]
- Requires dedicated hardware  $\odot$
- Slow if emulated in software

### Fixed-point computations

- C Tedious and time consuming to implement
  - >50% of design time [?]
- Relies only on integer instructions
- Efficient

#### Embedded systems targets





u-controllers



have efficient integer instructions

Fixed-point arithmetic is well suited for embedded systems

#### But, how to make it easy, fast, and numerically safe to use by non-expert programmers?

- To make fixed-point programming easy and fast
  - develop automated code synthesis tools
- To make fixed-point programming safe numerically
  - the tools must generate bounds on rounding errors

- To make fixed-point programming easy and fast
  - develop automated code synthesis tools
- To make fixed-point programming safe numerically
  - the tools must generate bounds on rounding errors

#### Float-to-fix conversion



- Works for a generic input
- IDFix [?], GUSTO [?], ...
- Based on floating-point simulations
  - no numeric certification
  - do not scale

- To make fixed-point programming easy and fast
  - develop automated code synthesis tools
- To make fixed-point programming safe numerically
  - the tools must generate bounds on rounding errors

# Float-to-fix conversion



- Works for a generic input
- IDFix [?], GUSTO [?], ...
- Based on floating-point simulations
  - no numeric certification
  - do not scale

#### Fixed-point code synthesis for IP blocks



- IP blocks: polynomials [?], filters [?], ...
- Based on analytic techniques
  - interval and affine arithmetics [?], [?], differentiation [?], ...
  - fast and scalable

- To make fixed-point programming easy and fast
  - develop automated code synthesis tools
- To make fixed-point programming safe numerically
  - the tools must generate bounds on rounding errors





- Works for a generic input
- IDFix [?], GUSTO [?], ...
- Based on floating-point simulations
  - no numeric certification
  - do not scale

#### Fixed-point code synthesis for IP blocks

$$P(x) = x^2 - x + 3$$

$$x \in [-5, 1.5]$$

- IP blocks: polynomials [?], filters [?], ...
- Based on analytic techniques
  - interval and affine arithmetics [?], [?], differentiation [?], ...
  - fast and scalable

#### DEFIS (ANR, 2011-)

Goal: develop techniques and tools to automate fixed-point programming

DEFIS (ANR, 2011-)

Goal: develop techniques and tools to automate fixed-point programming

Combines conversion and IP block synthesis

- Ménard et al. (CAIRN, Univ. Rennes) [?]:
  - automatic float-to-fix conversion
- Didier et al. (PEQUAN, Univ. Paris) [?]:
  - code generation for the linear filter IP block



DEFIS (ANR, 2011-)

Goal: develop techniques and tools to automate fixed-point programming

Combines conversion and IP block synthesis

- Ménard et al. (CAIRN, Univ. Rennes) [?]:
  - automatic float-to-fix conversion
- Didier et al. (PEQUAN, Univ. Paris) [?]:
  - code generation for the linear filter IP block
- Our approach (DALI, Univ. Perpignan):
  - certified fixed-point synthesis for:
    - Fine grained IP blocks: dot-products, polynomials, ...
    - High level IP blocks: matrix multiplication, triangular matrix inversion, Cholesky decomposition



DEFIS (ANR, 2011-)

Goal: develop techniques and tools to automate fixed-point programming

Combines conversion and IP block synthesis

- Ménard et al. (CAIRN, Univ. Rennes) [?]:
  - automatic float-to-fix conversion
- Didier et al. (PEQUAN, Univ. Paris) [?]:
  - code generation for the linear filter IP block
- Our approach (DALI, Univ. Perpignan):
  - · certified fixed-point synthesis for:
    - Fine grained IP blocks: dot-products, polynomials, ...
    - High level IP blocks: matrix multiplication, triangular matrix inversion, Cholesky decomposition





- 1. Specify an arithmetic model
  - Contributions:
    - formalization of  $\surd$  and /



- 1. Specify an arithmetic model
  - Contributions:
    - formalization of  $\surd$  and /
- 2. Build a synthesis tool, CGPE, for fine grained IP blocks:
  - it adheres to the arithmetic model
  - Contributions:
    - implementation of the arithmetic model
    - instruction selection



- 1. Specify an arithmetic model
  - Contributions:
    - formalization of  $\surd$  and /
- 2. Build a synthesis tool, CGPE, for fine grained IP blocks:
  - it adheres to the arithmetic model
  - Contributions:
    - implementation of the arithmetic model
    - instruction selection
- 3. Build a second synthesis tool, FPLA, for algorithmic IP blocks:
  - it generates code using CGPE
  - Contributions:
    - trade-off implementations for matrix multiplication
    - code synthesis for Cholesky decomposition and triangular matrix inversion



# Outline of the talk

# Outline of the talk

A fixed-point number *x* is defined by two integers:

- ▷ X the k-bit integer representation of x
- $\triangleright$  f the implicit scaling factor of x



$$\rightarrow$$
 The value of x is given by  $x = \frac{X}{2^f} = \sum_{\ell=-f}^{k-1-f} X_{\ell+f} \cdot 2^\ell$ 

A fixed-point number *x* is defined by two integers:

- ▷ X the k-bit integer representation of x
- $\triangleright$  f the implicit scaling factor of x



$$\rightarrow$$
 The value of x is given by  $x = \frac{X}{2^f} = \sum_{\ell=-f}^{k-1-f} X_{\ell+f} \cdot 2^\ell$ 

#### Notation

A fixed-point number with *i* bits of integer part and *f* bits of fraction part is in the **Q**<sub>*i*,*f*</sub> format

A fixed-point number *x* is defined by two integers:

- X the k-bit integer representation of x
- $\triangleright$  f the implicit scaling factor of x



→ The value of x is given by 
$$x = \frac{X}{2^f} = \sum_{\ell=-f}^{k-1-f} X_{\ell+f} \cdot 2^\ell$$

#### Notation

A fixed-point number with *i* bits of integer part and *f* bits of fraction part is in the **Q**<sub>*i*,*f*</sub> format

#### Example:

▶ x in  $\mathbf{Q}_{3.5}$  and X = (10011000)<sub>2</sub> = (152)<sub>10</sub>  $\longrightarrow$  x = (100.11000)<sub>2</sub> = (4.75)<sub>10</sub>

A fixed-point number *x* is defined by two integers:

- X the k-bit integer representation of x
- $\triangleright$  f the implicit scaling factor of x



→ The value of x is given by 
$$x = \frac{X}{2^f} = \sum_{\ell=-f}^{k-1-f} X_{\ell+f} \cdot 2^\ell$$

#### Notation

A fixed-point number with *i* bits of integer part and *f* bits of fraction part is in the **Q**<sub>*i*,*f*</sub> format

#### Example:

▶ x in  $\mathbf{Q}_{3.5}$  and X = (1001 1000)  $_2$  = (152)  $_{10}$   $\longrightarrow$  x = (100.11000)  $_2$  = (4.75)  $_{10}$ 

#### How to compute with fixed-point numbers?

$$\square P(x) = a_0 + (x \cdot a_1)$$

	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	
x	<b>Q</b> <sub>3.5</sub>	[16,208]	

$$\square P(x) = a_0 + (x \cdot a_1)$$

	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	3.5
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	1.5
x	<b>Q</b> <sub>3.5</sub>	[16,208]	[0.5, 6.5]

Toy example: degree-1 polynomial evaluation with 8 bit numbers



1. Pick an evaluation scheme



- 1. Pick an evaluation scheme
- 2. Determine the range and fixed-point formats of intermediate variables



	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	3.5
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	1.5
x	<b>Q</b> <sub>3.5</sub>	[16,208]	[0.5, 6.5]
$u = (x \cdot a_1)$	<b>Q</b> <sub>4.4</sub>	[12, 156]	[0.75, 9.75]
$v = (a_0 \gg 2)$	<b>Q</b> <sub>4.4</sub>	56	3.5

- 1. Pick an evaluation scheme
- 2. Determine the range and fixed-point formats of intermediate variables



	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	3.5
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	1.5
x	<b>Q</b> <sub>3.5</sub>	[16,208]	[0.5, 6.5]
$u = (x \cdot a_1)$	<b>Q</b> <sub>4.4</sub>	[12, 156]	[0.75, 9.75]
$v = (a_0 \gg 2)$	<b>Q</b> <sub>4.4</sub>	56	3.5
w = (v + u)	<b>Q</b> <sub>4.4</sub>	[68,212]	[4.25, 13.25]

- 1. Pick an evaluation scheme
- 2. Determine the range and fixed-point formats of intermediate variables

$\square P(x) = a_0 + $	$(x \cdot a_1)$
-------------------------	-----------------



	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	3.5
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	1.5
x	<b>Q</b> <sub>3.5</sub>	[16,208]	[0.5, 6.5]
$u = (x \cdot a_1)$	<b>Q</b> <sub>4.4</sub>	[12, 156]	[0.75, 9.75]
$v = (a_0 \gg 2)$	<b>Q</b> <sub>4.4</sub>	56	3.5
w = (v + u)	<b>Q</b> <sub>4.4</sub>	[68,212]	[4.25, 13.25]

uint8_t pol(uint8_t x	/*Q3.5*/)
$\{ uint8 + u = mul (x)$	. 224 ):
uint8_t v = 224 >> 2	2;
uint8_t $w = v + u;$	
return w; }	/*Q4.4*/

$$\square P(x) = a_0 + (x \cdot a_1)$$



	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	3.5
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	1.5
x	<b>Q</b> <sub>3.5</sub>	[16,208]	[0.5, 6.5]
$u = (x \cdot a_1)$	<b>Q</b> <sub>4.4</sub>	[12, 156]	[0.75, 9.75]
$v = (a_0 \gg 2)$	<b>Q</b> <sub>4.4</sub>	56	3.5
w = (v + u)	<b>Q</b> <sub>4.4</sub>	[68,212]	[4.25, 13.25]

```
Floating-point version
```

```
float pol(float x)
{
    return 3.5+1.5*x;
}
```
# How to compute with fixed-point numbers?

Toy example: degree-1 polynomial evaluation with 8 bit numbers

$\square P(x) = a_0$	$+(x \cdot a_1)$
----------------------	------------------



	Format	Int. repr.	Decimal value
$a_0$	<b>Q</b> <sub>2.6</sub>	224	3.5
a <sub>1</sub>	<b>Q</b> <sub>1.7</sub>	192	1.5
x	<b>Q</b> <sub>3.5</sub>	[16,208]	[0.5, 6.5]
$u = (x \cdot a_1)$	<b>Q</b> <sub>4.4</sub>	[12, 156]	[0.75, 9.75]
$v = (a_0 \gg 2)$	<b>Q</b> <sub>4.4</sub>	56	3.5
w = (v + u)	<b>Q</b> <sub>4.4</sub>	[68,212]	[4.25, 13.25]

Floating-point version			
float pol	(float x)		
{ return }	3.5+1.5*x;		

Even for small problems, this process is tedious: How to automate it?

Val(v) is the range of v

- For each coefficient or variable v, we keep track of 2 intervals **Val**(v) and **Err**(v)
- Our model assumes a fixed word-length k

**Err**(v) encloses the rounding error of computing v

- For each coefficient or variable v, we keep track of 2 intervals Val(v) and Err(v)
- Our model assumes a fixed word-length k

#### Val(v) is the range of v

• the format  $\mathbf{Q}_{i,f}$  of v is deduced from  $\mathbf{Val}(v) = [\underline{v}, \overline{v}]$ 

$$\bullet \quad i = \left[ \log_2 \left( \max\left( \left| \underline{\mathbf{v}} \right|, \left| \overline{\mathbf{v}} \right| \right) \right) \right] + \alpha \qquad \bullet \quad f = k - i$$

$$\alpha = \begin{cases} 1, & \text{if mod}\left(\log_2(\overline{\mathbf{v}}), 1\right) \neq 0, \\ 2, & \text{otherwise} \end{cases}$$

 $\mathbf{Err}(v)$  encloses the rounding error of computing v

• a bound  $\epsilon$  on rounding errors is deduced from  $\mathbf{Err}(v) = [\underline{\mathbf{e}}, \overline{\mathbf{e}}]$ 

• 
$$\epsilon = \max(|\underline{\mathbf{e}}|, |\overline{\mathbf{e}}|)$$

- For each coefficient or variable v, we keep track of 2 intervals Val(v) and Err(v)
- Our model assumes a fixed word-length k

#### Val(v) is the range of v

• the format  $\mathbf{Q}_{i,f}$  of v is deduced from  $\mathbf{Val}(v) = [\underline{v}, \overline{v}]$ 

$$i = \left[\log_2\left(\max\left(\left|\underline{\mathbf{v}}\right|, \left|\overline{\mathbf{v}}\right|\right)\right)\right] + \alpha \qquad F = k - i$$

$$=\begin{cases} 1, & \text{if mod}\left(\log_2(\overline{\mathbf{v}}), 1\right) \neq 0, \\ 2, & \text{otherwise} \end{cases}$$

2, otherwise

 $\mathbf{Err}(v)$  encloses the rounding error of computing v

 a bound *ε* on rounding errors is deduced from
 Err(v) = [e, ē]

• 
$$\epsilon = \max(|\underline{\mathbf{e}}|, |\overline{\mathbf{e}}|)$$



α

- For each coefficient or variable v, we keep track of 2 intervals Val(v) and Err(v)
- Our model assumes a fixed word-length k

### Val(v) is the range of v

the format Q<sub>if</sub> of v is deduced from  $Val(v) = [v, \overline{v}]$ 

$$i = \left[\log_2\left(\max\left(|\underline{\mathbf{v}}|, |\overline{\mathbf{v}}|\right)\right)\right] + \alpha \qquad f = k - i$$

$$x = \begin{cases} 1, & \text{if } \mod(\log_2(\overline{\mathbf{v}}), 1) \neq 0\\ 2, & \text{otherwise} \end{cases}$$

 $\mathbf{Err}(v)$  encloses the rounding error of computing v

 $\blacksquare$  a bound  $\epsilon$  on rounding errors is deduced from  $\mathbf{Err}(v) = [\mathbf{e}, \overline{\mathbf{e}}]$ 

$$\bullet \quad \epsilon = \max\left(\left|\underline{\mathbf{e}}\right|, \left|\overline{\mathbf{e}}\right|\right)$$



How to propagate **Val**(v) and **Err**(v) for  $\diamond \in \{+, -, \times, \ll, \gg, \sqrt{2}, //\}$ ?

- The two variables  $v_1$  and  $v_2$  must be aligned: in the same fixed-point format  $\mathbf{Q}_{i,f}$
- 1. If  $Val(v_1) + Val(v_2) \subseteq \mathscr{R}ange(\mathbf{Q}_{i,f})$ , overflow cannot occur



- The two variables  $v_1$  and  $v_2$  must be aligned: in the same fixed-point format  $\mathbf{Q}_{i,f}$
- 1. If  $Val(v_1) + Val(v_2) \subseteq \mathscr{R}$ ange  $(\mathbf{Q}_{i,f})$ , overflow cannot occur



- The two variables  $v_1$  and  $v_2$  must be aligned: in the same fixed-point format  $\mathbf{Q}_{i,f}$
- 1. If  $Val(v_1) + Val(v_2) \subseteq \mathscr{R}ange(\mathbf{Q}_{i,f})$ , overflow cannot occur



2. If  $\mathscr{R}$ ange  $(\mathbf{Q}_{i,f}) \subset \operatorname{Val}(v_1) + \operatorname{Val}(v_2) \subset \mathscr{R}$ ange  $(\mathbf{Q}_{i+1,f-1})$ , overflow may occur



- The two variables  $v_1$  and  $v_2$  must be aligned: in the same fixed-point format  $\mathbf{Q}_{i,f}$
- 1. If  $Val(v_1) + Val(v_2) \subseteq \mathscr{R}ange(\mathbf{Q}_{i,f})$ , overflow cannot occur



2. If  $\mathscr{R}$ ange  $(\mathbf{Q}_{i,f}) \subset \operatorname{Val}(v_1) + \operatorname{Val}(v_2) \subset \mathscr{R}$ ange  $(\mathbf{Q}_{i+1,f-1})$ , overflow may occur



- The two variables  $v_1$  and  $v_2$  must be aligned: in the same fixed-point format  $\mathbf{Q}_{i,f}$
- 1. If  $Val(v_1) + Val(v_2) \subseteq \mathscr{R}ange(\mathbf{Q}_{i,f})$ , overflow cannot occur



2. If  $\mathscr{R}$ ange  $(\mathbf{Q}_{i,f}) \subset \operatorname{Val}(v_1) + \operatorname{Val}(v_2) \subset \mathscr{R}$ ange  $(\mathbf{Q}_{i+1,f-1})$ , overflow may occur



• The output format of a  $\mathbf{Q}_{i_1.f_1} \times \mathbf{Q}_{i_2.f_2}$  is  $\mathbf{Q}_{i_1+i_2.f_1+f_2}$ 



• The output format of a  $\mathbf{Q}_{i_1.f_1} \times \mathbf{Q}_{i_2.f_2}$  is  $\mathbf{Q}_{i_1+i_2.f_1+f_2}$ 





- The output format of a  $\mathbf{Q}_{i_1,f_1} \times \mathbf{Q}_{i_2,f_2}$  is  $\mathbf{Q}_{i_1+i_2,f_1+f_2}$
- But, doubling the word-length is costly



**Err**<sub>×</sub> = 
$$\left[0, 2^{-f_r} - 2^{-(f_1 + f_2)}\right]$$



- The output format of a  $\mathbf{Q}_{i_1.f_1} \times \mathbf{Q}_{i_2.f_2}$  is  $\mathbf{Q}_{i_1+i_2.f_1+f_2}$
- But, doubling the word-length is costly



**Err**<sub>×</sub> = 
$$\left[0, 2^{-f_r} - 2^{-(f_1 + f_2)}\right]$$

This multiplication is available on integer processors and DSPs

int32\_t mul (int32\_t v1, int32\_t v2) {
 int64\_t prod = ((int64\_t) v1) \* ((int64\_t) v2);
 return (int32\_t) (prod >> 32);
}

The output integer part of  $\mathbf{Q}_{i_1,f_1}/\mathbf{Q}_{i_2,f_2}$  may be as large as  $i_1 + f_2$ 



The output integer part of  $\mathbf{Q}_{i_1,f_1}/\mathbf{Q}_{i_2,f_2}$  may be as large as  $i_1 + f_2$ 



- The output integer part of  $\mathbf{Q}_{i_1,f_1}/\mathbf{Q}_{i_2,f_2}$  may be as large as  $i_1 + f_2$
- But, doubling the word-length is costly



- The output integer part of  $\mathbf{Q}_{i_1,f_1}/\mathbf{Q}_{i_2,f_2}$  may be as large as  $i_1 + f_2$
- But, doubling the word-length is costly
- How to obtain sharper a error bounds on Err/?



- The output integer part of  $\mathbf{Q}_{i_1,f_1}/\mathbf{Q}_{i_2,f_2}$  may be as large as  $i_1 + f_2$
- But, doubling the word-length is costly
- How to obtain sharper a error bounds on Err/?



#### How to decide of the output format of division?

- A large integer part
  - prevents overflow
  - loose error bounds and loss of precision

- A small integer part
  - X may cause overflow
  - sharp error bounds and more accurate computations

# The propagation rule and implementation of division

Once the output format decided Q<sub>ir.fr</sub>



$$\overbrace{\operatorname{Val}(v_2)}^{\mathsf{Val}(v_1)} = \frac{\operatorname{Val}(v_1)}{\widetilde{\operatorname{Val}(v)} + \operatorname{Err}_{/}} \cap \operatorname{Val}(v_2) \text{ and } \widetilde{\operatorname{Val}(v)} = [-2^{i_r-1}, -2^{-f_r}] \cup [2^{-f_r}, 2^{i_r-1} - 2^{f_r}]$$

# The propagation rule and implementation of division

Once the output format decided Q<sub>ir.fr</sub>



$$\overrightarrow{\text{Val}(v_2)} = \frac{\text{Val}(v_1)}{\overbrace{\text{Val}(v)} + \text{Err}_{/}} \cap \text{Val}(v_2) \text{ and } \widetilde{\text{Val}(v)} = [-2^{i_r-1}, -2^{-f_r}] \cup [2^{-f_r}, 2^{i_r-1} - 2^{f_r}]$$

```
int32_t div (int32_t V1, int32_t V2, uint16_t eta)
{
    int64_t t1 = ((int64_t)V1) << eta;
    int64_t V = t1 / V2;
    return (int32_t) V;
}</pre>
```

# The propagation rule and implementation of division

Once the output format decided Q<sub>ir.fr</sub>



$$\overrightarrow{\text{Val}(v_2)} = \frac{\text{Val}(v_1)}{\overbrace{\text{Val}(v) + \text{Err}_{/}}} \cap \text{Val}(v_2) \text{ and } \overbrace{\text{Val}(v)} = [-2^{i_r-1}, -2^{-f_r}] \cup [2^{-f_r}, 2^{i_r-1} - 2^{f_r}]$$

#### Additional code to check for run-time overflows

Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2.30}$ 

Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2.30}$ 

• Cramer's rule: if 
$$\Delta = ad - bc \neq 0$$
 then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$ 

Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2.30}$ 

Cramer's rule: if 
$$\Delta = ad - bc \neq 0$$
 then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$ 



Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2.30}$ 

• Cramer's rule: if 
$$\Delta = ad - bc \neq 0$$
 then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$ 



Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2,30}$ 

Cramer's rule: if 
$$\Delta = ad - bc \neq 0$$
 then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$ 



Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2.30}$ 

Cramer's rule: if 
$$\Delta = ad - bc \neq 0$$
 then  $A^{-1} = \begin{pmatrix} \frac{a}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$ 



• Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2.30}$   
• Cramer's rule: if  $\Delta = ad - bc \neq 0$  then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & -\frac{b}{\Delta} \\ -\frac{c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$   
•  $\frac{2^{-11}}{2^{-21}}$   
•  $\frac{2^{-11}}{2^{-21}}$   
•  $\frac{2^{-11}}{2^{-21}}$   
•  $\frac{2^{-11}}{2^{-21}}$   
•  $\frac{2^{-21}}{2^{-21}}$   
•  $\frac{2^{-21}}{2^{-21}}$ 



• Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2:30}$   
• Cramer's rule: if  $\Delta = ad - bc \neq 0$  then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$   
•  $\frac{2^{-11}}{(-1, 1)}$  •  $\frac{100\%}{(-1, 1)}$  1)}$  • 00\%

DIVISION OUTPUT FORMAT







• Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2:30}$   
• Cramer's rule: if  $\Delta = ad - bc \neq 0$  then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$   
 $\begin{bmatrix} -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} \mathbf{Q}_{3:29} \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -1, 1 \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -2, 2 \\ -2, 2 \end{bmatrix} \begin{bmatrix} -1, 1 \\ -2, 2 \\ -2, 2 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -2, 2 \\ -2, 3 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -2, 2 \\ -2, 3 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -2, 2 \\ -2, 3 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -2, 2 \\ -2, 3 \end{bmatrix} \xrightarrow{\mathbf{Q}} \xrightarrow{\mathbf$ 

DIVISION OUTPUT FORMAT

• Consider 
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
 with  $a, b, c, d \in [-1, 1]$  in the format  $\mathbf{Q}_{2:30}$   
• Cramer's rule: if  $\Delta = ad - bc \neq 0$  then  $A^{-1} = \begin{pmatrix} \frac{d}{\Delta} & \frac{-b}{\Delta} \\ \frac{-c}{\Delta} & \frac{a}{\Delta} \end{pmatrix}$   
•  $\begin{bmatrix} -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} \mathbf{Q}_{3:29} \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -2, 2 \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \begin{bmatrix} -1, 1 \\ -1, 1 \end{bmatrix} \xrightarrow{\mathbf{Q}} \xrightarrow{\mathbf{Q}} \xrightarrow{\mathbf{Q}}$ 

DIVISION OUTPUT FORMAT
# Outline of the talk

CGPE (Code Generation for Polynomial Evaluation): initiated by Revy [?]

synthesizes fixed-point code for polynomial evaluation

- CGPE (Code Generation for Polynomial Evaluation): initiated by Revy [?]
  - synthesizes fixed-point code for polynomial evaluation
- 1. Computation step  $\rightsquigarrow$  front-end
  - computes evaluation schemes ~> DAGs



- CGPE (Code Generation for Polynomial Evaluation): initiated by Revy [?]
  - synthesizes fixed-point code for polynomial evaluation
- 1. Computation step ~> front-end
  - computes evaluation schemes ~> DAGs
- 2. Filtering step  $\rightsquigarrow$  middle-end
  - applies the arithmetic model
  - prunes the DAGs that do not satisfy different criteria:
    - latency → scheduling filter
    - accuracy → numerical filter
    - ...



- CGPE (Code Generation for Polynomial Evaluation): initiated by Revy [?]
  - synthesizes fixed-point code for polynomial evaluation

### 1. Computation step $\rightsquigarrow$ front-end

computes evaluation schemes ~>> DAGs

### 2. Filtering step $\rightsquigarrow$ middle-end

- applies the arithmetic model
- prunes the DAGs that do not satisfy different criteria:
  - latency ~>> scheduling filter
  - accuracy ~> numerical filter
  - ...

### 3. Generation step $\rightsquigarrow$ back-end

 generates C codes and Gappa accuracy certificates



• Low-pass Butterworth filter with cutoff frequency  $0.3 \cdot \pi$ :

$$y[k] = \sum_{i=0}^{3} b_i \cdot u[k-i] - \sum_{i=1}^{3} a_i \cdot y[k-i]$$

• Low-pass Butterworth filter with cutoff frequency  $0.3 \cdot \pi$ :

$$y[k] = \sum_{i=0}^{3} b_i \cdot u[k-i] - \sum_{i=1}^{3} a_i \cdot y[k-i]$$



M. A. Najahi (DALI UPVD/LIRMM, UM2, CNRS)

• Low-pass Butterworth filter with cutoff frequency  $0.3 \cdot \pi$ :

$$y[k] = \sum_{i=0}^{3} b_i \cdot u[k-i] - \sum_{i=1}^{3} a_i \cdot y[k-i]$$



M. A. Najahi (DALI UPVD/LIRMM, UM2, CNRS)

• Low-pass Butterworth filter with cutoff frequency  $0.3 \cdot \pi$ :

 $y[k] = \sum_{i=0}^{3} b_i \cdot u[k-i] - \sum_{i=1}^{3} a_i \cdot y[k-i]$ 

<pre>int32_t filter(</pre>	int32_t u0 /*	* Q5.27*/	, int32_	_t ul	/*Q5.27*/ ,
	int32_t u2 /*	*Q5.2/*/	, int32_	_t u3	/*Q5.2/*/ ,
	int32_t y1 /*	*Q6.26*/	, int32_	_t y2	/*Q6.26*/ ,
	int32_t y3 /*	*Q6.26*/	)		
{			Formats	Err	
int32_t r0 =	mul(0x4a5cdb26,	y1); //	Q8.24	$[-2^{\{-24\}}]$	,0]
int32 t r1 =	mul(0xa6eb5908,	v2); //	07.25	$[-2^{+}]$	,01
int32 t r2 =	mul(0x4688a637.	v3): //	05.27	[-2^{-27}]	. 01
int 32 t r3 =	mul(0x65718e3b.	10): 1/	02.30	1-201-301	. 01
int 32 t r4 =	mul(0x65718e3b)	113) : //	02 30	1-2-1-301	. 01
int 32 + r5 =	$r_{3} + r_{4}$		02 30	$[-2^{4}]$	. 01
int 32 + r6 =	r5 >> 2:		04 28	[-2^(-27)	67813.01
int 32 t x7 -	mul (0x4e152aad		01 20	[_2^[_2].	01
int 32 t 17 -	mul (0x4c152aad,	u1), //	04.20		,01
111C32_C 10 =	mui (0x4Ci5zaau,	uz); //	Q4.20		,01
int32_t r9 =	r/ + r8;		Q4.28	[-20][-27]	,0]
$int32_t r10 =$	r6 + r9;	11	Q4.28	[-2^{-26.	2996},0]
int32_t <b>r11</b> =	r10 >> 1;		Q5.27	[-2^{-25.	9125},0]
int32_t r12 =	r2 + r11;	/ /	Q5.27	[-2^{ -25.	3561},0]
int32_t r13 =	r12 >> 2;		Q7.25	[-2^{ -24.	3853},0]
int32_t r14 =	r1 + r13;		Q7.25	[-2^{ -23.	6601},0]
int32_t r15 =	r14 >> 1;		Q8.24	[-2^{ -23.	1798},0]
int32 t r16 =	r0 + r15;		08.24	[-2^{-22.	5324},01
int32_t r17 =	r16 << 2;	11	Q6.26	[-2^{-22.	5324},0]
return r17;					
}					

• Low-pass Butterworth filter with cutoff frequency  $0.3 \cdot \pi$ :

 $y[k] = \sum_{i=0}^{3} b_i \cdot u[k-i] - \sum_{i=1}^{3} a_i \cdot y[k-i]$ 

int32_t filter(	int32_t u0 /*Q5.27*	/ , int32_t u1 /*Q5.27*/ ,
	int32_t u2 /*Q5.27*	/ , int32_t u3 /*Q5.27*/ ,
	int32 t v1 /*06.26*	/. int32 t v2 /*06.26*/.
	int32 t v3 /*06.26*	/ )
{		//Formats Err
int32 t r0 =	<pre>mul(0x4a5cdb26, v1);</pre>	//08.24 [-2^{-24}.0]
int 32 t r1 =	mul(0xa6eb5908, v2):	//07.25 [-2^{-25}.0]
$int 32 \pm r^2 =$	mu1(0x4688a637, v3):	//05 27 [-2^{-27}.0]
int32 t x3 =	mu1(0x65718e3b, u0);	//02 30 [-2^/-301 0]
$int 32 \pm rA =$	mul(0x65718e3b, u3);	//02 30 [-2^/_301 0]
int 32 t r5 -	mai(0x03/10030, a3),	//02.30 [_2^[_30],0]
int 32 t x6 -	x5 xx 2.	//04 20 [-2^[-27 6701] 0]
int 32 t 10 -	10 // 2, mul (0.4 cl 52 cod	//04.20 [-2 (-27.0701),0]
int 32_t 17 =	mul(0x4c152aad, u1);	//04.20 [=2 {=20},0]
int32_t r8 =	mul(0x4c152aad, u2);	//Q4.28 [-2^{-28},0]
int32_t r9 =	r7 + r8;	//Q4.28 [-2^{-27},0]
int32_t <u>r10 =</u>	r6 + r9;	//Q4.28 [-2^{-26.2996},0]
int32_t r11 =	r10 >> 1;	//Q5.27 [-2^{-25.9125},0]
int32_t r12 =	r2 + r11;	//Q5.27 [-2^{-25.3561},0]
int32 t r13 =	r12 >> 2;	//07.25 [-2^{-24.3853},0]
int32 t r14 =	r1 + r13;	//07.25 [-2^{-23.6601}.0]
int32 t r15 =	r14 >> 1;	//08.24 [-2^{-23.1798}.0]
int32 t r16 =	r0 + r15;	//08.24 [-2^{-22.5324}.0]
int32 t r17 =	r16 << 2;	//06.26 [-2^{-22.5324},0]
return r17;		
}		

## Instruction selection to handle advanced instructions

- It is a well known problem in compilation ~>> proven to be NP-complete on DAGs
- Usually solved using a tiling algorithm:
  - ▶ input:
    - a DAG representing an arithmetic expression,
    - a set of tiles, with a cost for each,
    - a function that associates a cost to a DAG.
  - output: a set of covering tiles that minimize the cost function.

## Instruction selection to handle advanced instructions

- It is a well known problem in compilation ~>> proven to be NP-complete on DAGs
- Usually solved using a tiling algorithm:
  - ▶ input:
    - a DAG representing an arithmetic expression,
    - a set of tiles, with a cost for each,
    - a function that associates a cost to a DAG.
  - output: a set of covering tiles that minimize the cost function.

#### Implementation in CGPE

- XML architecture description file that contains for each instruction:
  - its name, its type (signed or unsigned), its latency (# cycles),
  - a description of the pattern it matches,
  - a C macro to emulate it in software,
  - and a piece of Gappa script to compute the error entailed by its evaluation in fixed-point arithmetic.

- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$





- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()
  - Example: how to evaluate  $a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$



- 1: BottomUpDP() + TopDownSelect()
- 2: ImproveCSEDecision()
- 3: BottomUpDP() + TopDownSelect()

Example: how to evaluate 
$$a_0 + ((a_1 \cdot x) + ((a_2 \cdot (x \cdot x)) \ll 1))?$$

- In our case, only the first step of NOLTIS is valuable.
- NOLTIS algorithm mainly relies on the evaluation of a cost function. We have implemented three different cost functions:
  - number of operator (regardless common subexpressions)
  - evaluation latency on unbounded parallelism
  - evaluation accuracy, computed using the Gappa script of each instruction

# Impact on the number of instructions



Figure: Average number of instructions in 50 synthesized codes, for the evaluation of polynomials of degree 5 up to 12 for various elementary functions.

Remark 1: average reduction of 8.7 % up to 13.75 %

■ Remark 2: interest of ST231 shift-and-add with left shift for sin(x) implementation ~> reduction of 8.7 %

Remark 3: interest of shift-and-add with right shift for cos(x) and  $log_2(1+x)$ implementation  $\rightarrow$  reduction of 12.8 % and 13.75 %, respectively

### Impact on the accuracy of some functions

f(x)	I	d	$\log_2(\epsilon)$		
		not optimized optimiz			
$\exp(x) - 1$	[-0.25, 0.25]	7	-26.98	-27.34	
exp(x)	[0, 1]	7	-13.94	-14.90	
sin(x)	[-0.5, 0.5]	9	-18.95	-19.91	
$\cos(x)$	[-0.5, 0.25]	5	-27.01	-27.26	
tan(x)	[0.25, 0.5]	9	-18.81	-19.64	
$\log_2(1+x)/x$	[2 <sup>-23</sup> , 1]	7	-13.94 -14.8		
$\sqrt{1+x}$	[2 <sup>-23</sup> , 1]	7	-13.94	-14.90	

Table: Impact of the accuracy based selection step on the certified accuracy of the generated code for various functions.

- Remark 1: with a mulace that computes  $(a * b) + (c \gg n)$  with  $n \in \{1, \dots, 31\}$  with one final rounding
- Remark 2: more accurate results for all the cases, almost up to 1 bit of accuracy for *exp*, *sin*,  $\log_2$  and  $\sqrt{1+x}$

# Outline of the talk

Let *M* be a matrix of fixed-point variables,

to generate certified code that inverts  $M' \in M$  a symmetric positive definite, we need to:

- 1. Generate certified code to compute B a lower triangular s.t.  $M' = B \cdot B^T$
- 2. Generate certified code to compute  $N = B^{-1}$
- 3. Generate certified code to compute  $M'^{-1} = N^T \cdot N$

Let *M* be a matrix of fixed-point variables,

to generate certified code that inverts  $M' \in M$  a symmetric positive definite, we need to:

- 1. Generate certified code to compute B a lower triangular s.t.  $M' = B \cdot B^T$
- 2. Generate certified code to compute  $N = B^{-1}$
- 3. Generate certified code to compute  $M'^{-1} = N^T \cdot N$

### The basic blocks we need to include in our tool-chain





Let *M* be a matrix of fixed-point variables,

to generate certified code that inverts  $M' \in M$  a symmetric positive definite, we need to:

- 1. Generate certified code to compute B a lower triangular s.t.  $M' = B \cdot B^T$
- 2. Generate certified code to compute  $N = B^{-1}$
- 3. Generate certified code to compute  $M'^{-1} = N^T \cdot N$

### The basic blocks we need to include in our tool-chain

- Certified code synthesis for Cholesky decomposition
- Certified code synthesis for triangular matrix inversion



Let *M* be a matrix of fixed-point variables,

to generate certified code that inverts  $M' \in M$  a symmetric positive definite, we need to:

- 1. Generate certified code to compute B a lower triangular s.t.  $M' = B \cdot B^T$
- 2. Generate certified code to compute  $N = B^{-1}$
- 3. Generate certified code to compute  $M'^{-1} = N^T \cdot N$

### The basic blocks we need to include in our tool-chain

- Certified code synthesis for Cholesky decomposition
- Certified code synthesis for triangular matrix inversion
- Certified code synthesis for matrix multiplication



## Linear algebra basic blocks



### Linear algebra basic blocks



# Matrix multiplication

### Inputs

Two matrices A and B of fixed-point variables

 $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ 

- A bound C<sub>1</sub> on the roundoff error
- A bound  $\mathscr{C}_2$  on the code size

# Matrix multiplication

### Inputs

Two matrices A and B of fixed-point variables

 $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ 

- A bound C<sub>1</sub> on the roundoff error
- A bound C2 on the code size

### Output

Fixed-point code (C, VHDL, ...) that evaluates the product

$$C' = A' \cdot B'$$
, where  $A' \in A$  and  $B' \in B$ 

that satisfy both  $\mathscr{C}_1$  and  $\mathscr{C}_2$ 

Accuracy certificate (verifiable by a formal proof checker)

# Straightforward algorithms

### Accurate algorithm

 Main idea: a dot product code for each coefficient of the resulting matrix

#### Accurate algorithm

#### Inputs:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ 

#### Outputs:

C code to compute the product  $A \cdot B$ 

 $n \cdot n$  accuracy certificates

#### Steps:

- 1: for  $1 < i \le n$  do
- 2: **for**  $1 < j \le n$  **do**
- 3:  $DPSynthesis(A_{i,:}, B_{:,j})$
- 4: end for
- 5: end for
- 6: Check  $C_1$  and  $C_2$

# Straightforward algorithms

### Accurate algorithm

 Main idea: a dot product code for each coefficient of the resulting matrix

#### Accurate algorithm

#### Inputs:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ 

#### Outputs:

- C code to compute the product  $A \cdot B$
- $n \cdot n$  accuracy certificates

#### Steps:

- 1: for  $1 < i \le n$  do
- 2: **for**  $1 < j \le n$  **do**
- 3:  $DPSynthesis(A_{i,:}, B_{:,j})$
- 4: end for
- 5: end for
- 6: Check  $C_1$  and  $C_2$

### Compact algorithm

 Main idea: a unique dot product code for all the coefficient of the resulting matrix

#### Compact algorithm

#### Inputs:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ 

#### Outputs:

- C code to compute the product  $A \cdot B$
- 1 accuracy certificate

#### Steps:

- 1:  $\mathcal{U} = A_{1,:} \cup A_{2,:} \cup \cdots \cup A_{n,:}$ , with  $\mathcal{U} \in \mathbb{F}ix^{1 \times n}$
- 2:  $\mathcal{V} = B_{:,1} \cup B_{:,2} \cup \cdots \cup B_{:,n}$ , with  $\mathcal{V} \in \mathbb{F}ix^{n \times 1}$
- 3: DPSynthesis( $\mathcal{U}, \mathcal{V}$ )
- 4: Check C<sub>1</sub> and C<sub>2</sub>

## Example of a multiplication of $2 \times 2$ matrices

We consider code synthesis for the multiplication of matrices  $A' \in A$  and  $B' \in B$ :

$$A = \begin{pmatrix} [-1000, 1000] & [-3000, 3000] \\ [-1, 1] & [-1, 1] \end{pmatrix} \text{ and } B = \begin{pmatrix} [-2000, 2000] & [-2, 2] \\ [-4000, 4000] & [-10, 10] \end{pmatrix}$$

Coefficient	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	B <sub>1,1</sub>	B <sub>1,2</sub>	B <sub>2,1</sub>	B <sub>2,2</sub>
Fixed-point format	<b>Q</b> <sub>11.21</sub>	<b>Q</b> <sub>12.20</sub>	<b>Q</b> <sub>2.30</sub>	<b>Q</b> <sub>2.30</sub>	<b>Q</b> <sub>11.21</sub>	<b>Q</b> <sub>3.29</sub>	<b>Q</b> <sub>2.30</sub>	<b>Q</b> <sub>5.27</sub>

## Example of a multiplication of $2 \times 2$ matrices

We consider code synthesis for the multiplication of matrices  $A' \in A$  and  $B' \in B$ :

$$A = \begin{pmatrix} [-1000, 1000] & [-3000, 3000] \\ [-1, 1] & [-1, 1] \end{pmatrix} \text{ and } B = \begin{pmatrix} [-2000, 2000] & [-2, 2] \\ [-4000, 4000] & [-10, 10] \end{pmatrix}$$

Coefficient	A <sub>1,1</sub>	A <sub>1,2</sub>	A <sub>2,1</sub>	A <sub>2,2</sub>	B <sub>1,1</sub>	B <sub>1,2</sub>	B <sub>2,1</sub>	B <sub>2,2</sub>
Fixed-point format	<b>Q</b> <sub>11.21</sub>	<b>Q</b> <sub>12.20</sub>	<b>Q</b> <sub>2.30</sub>	<b>Q</b> <sub>2.30</sub>	<b>Q</b> <sub>11.21</sub>	<b>Q</b> <sub>3.29</sub>	<b>Q</b> <sub>2.30</sub>	<b>Q</b> <sub>5.27</sub>

1								
	Dot-product $A_{1,:} \cdot B_{:,1}$ $A_{1,:} \cdot B_{:,2}$ $A_{2,:} \cdot B_{:,1}$ $A_{2,:} \cdot B_{:,2}$							
	Evaluated using	DPCode <sub>1,1</sub>	DPCode <sub>1,2</sub>	DPCode <sub>2,1</sub>	DPCode <sub>2,2</sub>			
	Output format	Q <sub>26,6</sub>	Q <sub>18,14</sub>	Q <sub>15,17</sub>	Q <sub>7,25</sub>			
	Certified error	≈ 2 <sup>-5</sup>	≈ 2 <sup>-14</sup>	≈ 2 <sup>-16</sup>	≈ 2 <sup>-24</sup>			
	Maximum error		2-5					
	Average error	≈ 2 <sup>-7</sup>						

### Compact algorithm

Dot-product	A <sub>1,:</sub> · B <sub>:,1</sub>	$A_{1,:} \cdot B_{:,2}$	A <sub>2,:</sub> · B <sub>:,1</sub>	$A_{2,:} \cdot B_{:,2}$			
Evaluated using	DPCode <sub>W,V</sub>						
Output format	Q <sub>26,6</sub>						
Certified error	≈ 2 <sup>-5</sup>						
Maximum error	≈ 2 <sup>-5</sup>						
Average error	≈ 2 <sup>-5</sup>						

Accurate algorithm










### Number of possible trade-off algorithms

Given by  $\mathscr{B}(n)^2$  with  $\mathscr{B}(n)$  the  $n^{\text{th}}$  Bell number

n	5	6	10	16	25	64	
ℬ(n) <sup>2</sup>	2704	41 209	$\approx 2^{34}$	$\approx 2^{66}$	≈ 2 <sup>124</sup>	≈ 2 <sup>433</sup>	



n	5	6	10	16	25	64	
𝔅(n)²	2704	41 209	$\approx 2^{34}$	$\approx 2^{66}$	≈ 2 <sup>124</sup>	≈ 2 <sup>433</sup>	

#### How to find in reasonable time a partition that reduces cost size without harming the accuracy?

### Distances between fixed-point variables

The Hausdorff distance

$$d_{H}: \mathbb{F}ix \times \mathbb{F}ix \to \mathbb{R}^{+}$$
$$d_{H}(I_{1}, I_{2}) = max\left\{ \left| \underline{I_{1}} - \underline{I_{2}} \right|, \left| \overline{I_{1}} - \overline{I_{2}} \right| \right\}$$

#### **Fixed-point distance**

$$d_F : \mathbb{F}ix \times \mathbb{F}ix \to \mathbb{N}$$
$$d_F (l_1, l_2) = \left| IntegerPart(l_1) - IntegerPart(l_2) \right|$$

## Distances between fixed-point variables

The Hausdorff distance

$$d_{H}: \mathbb{F}ix \times \mathbb{F}ix \to \mathbb{R}^{+}$$
$$d_{H}(I_{1}, I_{2}) = max\left\{ \left| \underline{I_{1}} - \underline{I_{2}} \right|, \left| \overline{I_{1}} - \overline{I_{2}} \right| \right\}$$



**Fixed-point distance** 

$$d_F : \mathbb{F}ix \times \mathbb{F}ix \to \mathbb{N}$$
$$d_F(I_1, I_2) = \left| IntegerPart(I_1) - IntegerPart(I_2) \right|$$

#### Input:

Two matrices  $A \in \mathbb{F}|x^{n \times n}$  and  $B \in \mathbb{F}|x^{n \times n}$ An accuracy bound  $\mathscr{C}_1$  (ex. the average error bound is  $< \epsilon$ ) A code size bound  $\mathscr{C}_2$ 

A metric d

#### Output:

Code to compute  $A \cdot B$  s.t.  $\mathscr{C}_1$  and  $\mathscr{C}_2$  are satisfied, or no code otherwise

#### Algorithm:

1:  $\mathscr{S}_{\Delta} \leftarrow \{A_0, \ldots, A_{n-1}\}$ 2:  $\mathscr{S}_{B} \leftarrow \{B_{0}, \ldots, B_{n-1}\}$ 3: while C1 is satisfied do  $(u_{\Delta}, v_{\Delta}), d_{\Delta} \leftarrow findClosestPair(\mathcal{S}_{\Delta}, d)$ 4:  $(u_B, v_B), d_B \leftarrow findClosestPair(\mathcal{S}_B, d)$ 5: 6: if  $d_A \leq d_B$  then 7: remove( $u_{\Delta}, v_{\Delta}, \mathscr{S}_{\Delta}$ ) 8:  $insert(u_{A} \cup v_{A}, \mathscr{S}_{A})$ 9: else 10. remove(up, vp, Sp) 11. insert( $u_B \cup v_B, \mathscr{S}_B$ ) 12: end if 13: for  $(A_i, B_i) \in \mathscr{S}_A \times \mathscr{S}_B$  do 14.  $DPSynthesis(A_i, B_i)$ 15: end for 16: end while 17: /\* Revert the last merging step, and check the bound  $\mathscr{C}_2$ . \*/



#### Input:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ An accuracy bound  $\mathscr{C}_1$  (ex. the average error bound is  $< \epsilon$ ) A code size bound  $\mathscr{C}_2$ 

A metric d

#### Output:

Code to compute  $A \cdot B$  s.t.  $\mathscr{C}_1$  and  $\mathscr{C}_2$  are satisfied, or no code otherwise

#### Algorithm:

1:  $\mathscr{S}_{\Delta} \leftarrow \{A_0, \ldots, A_{n-1}\}$ 2:  $\mathscr{S}_{B} \leftarrow \{B_{0}, \ldots, B_{n-1}\}$ 3: while C1 is satisfied do  $(u_{\Delta}, v_{\Delta}), d_{\Delta} \leftarrow findClosestPair(\mathcal{S}_{\Delta}, d)$ 4:  $(u_B, v_B), d_B \leftarrow findClosestPair(\mathcal{S}_B, d)$ 5: 6: if  $d_A \leq d_B$  then 7: remove( $u_{\Delta}, v_{\Delta}, \mathscr{S}_{\Delta}$ ) 8:  $insert(u_{A} \cup v_{A}, \mathscr{S}_{A})$ 9: else 10. remove(up, vp, Sp) 11: insert( $u_B \cup v_B, \mathscr{S}_B$ ) 12: end if 13: for  $(A_i, B_i) \in \mathscr{S}_A \times \mathscr{S}_B$  do 14.  $DPSynthesis(A_i, B_i)$ 15: end for 16: end while 17: /\* Revert the last merging step, and check the bound  $\mathscr{C}_2$ . \*/



#### Input:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ An accuracy bound  $\mathscr{C}_1$  (ex. the average error bound is  $< \epsilon$ ) A code size bound  $\mathscr{C}_2$ 

A metric d

#### Output:

Code to compute  $A \cdot B$  s.t.  $\mathscr{C}_1$  and  $\mathscr{C}_2$  are satisfied, or no code otherwise

#### Algorithm:

1:  $\mathscr{S}_{\Delta} \leftarrow \{A_0, \ldots, A_{n-1}\}$ 2:  $\mathscr{S}_{B} \leftarrow \{B_{0}, \ldots, B_{n-1}\}$ 3: while C1 is satisfied do  $(u_{\Delta}, v_{\Delta}), d_{\Delta} \leftarrow findClosestPair(\mathcal{S}_{\Delta}, d)$ 4:  $(u_B, v_B), d_B \leftarrow findClosestPair(\mathcal{S}_B, d)$ 5: 6: if  $d_A \leq d_B$  then 7: remove( $u_{\Delta}, v_{\Delta}, \mathscr{S}_{\Delta}$ ) 8:  $insert(u_{A} \cup v_{A}, \mathscr{S}_{A})$ 9: else 10. remove(up, vp, Sp) 11: insert( $u_B \cup v_B, \mathscr{S}_B$ ) 12: end if 13: for  $(A_i, B_i) \in \mathscr{S}_A \times \mathscr{S}_B$  do 14.  $DPSynthesis(A_i, B_i)$ 15: end for 16: end while





 $\mathscr{C}_1$  is satisfied

17: /\* Revert the last merging step, and check the bound  $\mathscr{C}_2$ . \*/

#### Input:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ An accuracy bound  $\mathscr{C}_1$  (ex. the average error bound is  $< \epsilon$ ) A code size bound  $\mathscr{C}_2$ 

A metric d

#### Output:

Code to compute  $A \cdot B$  s.t.  $\mathscr{C}_1$  and  $\mathscr{C}_2$  are satisfied, or no code otherwise

#### Algorithm:

1:  $\mathscr{S}_{\Delta} \leftarrow \{A_0, \ldots, A_{n-1}\}$ 2:  $\mathscr{S}_{B} \leftarrow \{B_{0}, \ldots, B_{n-1}\}$ 3: while C1 is satisfied do  $(u_{\Delta}, v_{\Delta}), d_{\Delta} \leftarrow findClosestPair(\mathcal{S}_{\Delta}, d)$ 4:  $(u_B, v_B), d_B \leftarrow findClosestPair(\mathcal{S}_B, d)$ 5: 6: if  $d_A \leq d_B$  then 7: remove( $u_{\Delta}, v_{\Delta}, \mathscr{S}_{\Delta}$ ) 8:  $insert(u_{A} \cup v_{A}, \mathscr{S}_{A})$ 9: else 10. remove(up, vp, Sp) 11: insert( $u_B \cup v_B, \mathscr{S}_B$ ) 12: end if 13: for  $(A_i, B_i) \in \mathscr{S}_A \times \mathscr{S}_B$  do 14.  $DPSynthesis(A_i, B_i)$ 15: end for 16: end while





 $\mathscr{C}_1$  is satisfied

17: /\* Revert the last merging step, and check the bound  $\mathscr{C}_2.$  \*/

#### Input:

Two matrices  $A \in \mathbb{F}ix^{n \times n}$  and  $B \in \mathbb{F}ix^{n \times n}$ An accuracy bound  $\mathscr{C}_1$  (ex. the average error bound is  $< \epsilon$ ) A code size bound  $\mathscr{C}_2$ 

A metric d

#### Output:

Code to compute  $A \cdot B$  s.t.  $\mathscr{C}_1$  and  $\mathscr{C}_2$  are satisfied, or no code otherwise

#### Algorithm:

1:  $\mathscr{S}_{\Delta} \leftarrow \{A_0, \ldots, A_{n-1}\}$ 2:  $\mathscr{S}_{B} \leftarrow \{B_{0}, \ldots, B_{n-1}\}$ 3: while C1 is satisfied do  $(u_{\Delta}, v_{\Delta}), d_{\Delta} \leftarrow findClosestPair(\mathcal{S}_{\Delta}, d)$ 4:  $(u_B, v_B), d_B \leftarrow findClosestPair(\mathcal{S}_B, d)$ 5: 6: if  $d_A \leq d_B$  then 7: remove( $u_{\Delta}, v_{\Delta}, \mathscr{S}_{\Delta}$ ) 8:  $insert(u_{A} \cup v_{A}, \mathscr{S}_{A})$ 9: else 10. remove(up, vp, Sp) 11. insert( $u_B \cup v_B, \mathscr{S}_B$ ) 12: end if 13: for  $(A_i, B_i) \in \mathscr{S}_A \times \mathscr{S}_B$  do 14.  $DPSynthesis(A_i, B_i)$ 15: end for 16: end while

17: /\* Revert the last merging step, and check the bound  $\mathscr{C}_2.$  \*/



#### $\mathscr{C}_1$ is no longer satisfied

#### Input:

Two matrices  $A \in \mathbb{F} ix^{n \times n}$  and  $B \in \mathbb{F} ix^{n \times n}$ An accuracy bound  $\mathscr{C}_1$  (ex. the average error bound is  $< \epsilon$ ) A code size bound  $\mathscr{C}_2$ 

A metric d

#### Output:

Code to compute  $A \cdot B$  s.t.  $\mathscr{C}_1$  and  $\mathscr{C}_2$  are satisfied, or no code otherwise

#### Algorithm:

1:  $\mathscr{S}_{\Delta} \leftarrow \{A_0, \ldots, A_{n-1}\}$ 2:  $\mathscr{S}_{B} \leftarrow \{B_{0}, \ldots, B_{n-1}\}$ 3: while C1 is satisfied do  $(u_{\Delta}, v_{\Delta}), d_{\Delta} \leftarrow findClosestPair(\mathcal{S}_{\Delta}, d)$ 4:  $(u_B, v_B), d_B \leftarrow findClosestPair(\mathcal{S}_B, d)$ 5: 6: if  $d_A \leq d_B$  then 7: remove( $u_{\Delta}, v_{\Delta}, \mathscr{S}_{\Delta}$ ) 8:  $insert(u_{A} \cup v_{A}, \mathscr{S}_{A})$ 9: else 10. remove(up, vp, Sp) 11. insert( $u_B \cup v_B, \mathscr{S}_B$ ) 12: end if 13: for  $(A_i, B_i) \in \mathscr{S}_A \times \mathscr{S}_B$  do 14.  $DPSynthesis(A_i, B_i)$ 15: end for 16: end while

17: /\* Revert the last merging step, and check the bound  $\mathscr{C}_2.$  \*/



#### $\mathscr{C}_1$ is satisfied

 $\stackrel{\longrightarrow}{\rightarrow} \text{Revert the last merging step and} \\ \text{check if } \mathscr{C}_2 \text{ is satisfied}$ 

## **Benchmark matrices**



## Efficiency of the distance-based heuristic: 6 × 6 matrix product

- Generating the 41 209 different algorithms: 2h15 per benchmark
- Running our algorithm: 10 seconds per benchmark



### Impact of the metric on the trade-off strategy



## Linear algebra basic blocks



## Cholesky decomposition and triangular matrix inversion

### Cholesky decomposition

v

$$b_{i,j} = \begin{cases} \sqrt{c_{i,i}} & \text{if } i = j \\\\ \frac{c_{i,j}}{b_{j,j}} & \text{if } i \neq j \end{cases}$$
with  $c_{i,j} = m_{i,j} - \sum_{k=0}^{j-1} b_{i,k} \cdot b_{j,k}$ 

### Triangular matrix inversion

$$n_{i,j} = \begin{cases} \frac{1}{b_{i,i}} & \text{if } i = j \\ \frac{-c_{i,j}}{b_{i,i}} & \text{if } i \neq j \end{cases}$$
where  $c_{i,j} = \sum_{k=j}^{i-1} b_{i,k} \cdot n_{k,j}$ 

## Cholesky decomposition and triangular matrix inversion

### Cholesky decomposition

$$b_{i,j} = \begin{cases} \sqrt{c_{i,i}} & \text{if } i = j \\\\ \frac{c_{i,j}}{b_{j,j}} & \text{if } i \neq j \end{cases}$$
  
with  $c_{i,j} = m_{i,j} - \sum_{k=0}^{j-1} b_{i,k} \cdot b_{j,k}$ 

### Triangular matrix inversion

$$n_{i,j} = \begin{cases} \frac{1}{b_{i,i}} & \text{if } i = j \\ \frac{-c_{i,j}}{b_{i,i}} & \text{if } i \neq j \end{cases}$$
where  $c_{i,j} = \sum_{k=j}^{i-1} b_{i,k} \cdot n_{k,j}$ 





Dependencies of the coefficient  $b_{4,2}$  in the decomposition and inversion of a 6 × 6 matrix.

W

## FPLA (Fixed-Point Linear Algebra)



## Impact of the output format of division

Different functions to set the output format of division

1.  $f_1(i_1, i_2) = t$ ,

2. 
$$f_2(i_1, i_2) = \min(i_1, i_2) + t$$

3. 
$$f_3(i_1, i_2) = \max(i_1, i_2) + t$$
,

4. 
$$f_4(i_1, i_2) = \lfloor (i_1 + i_2)/2 \rfloor + t$$
,

 $i_1$  and  $i_2$ : integer parts of the numerator and denominator and  $t \in [-2, 8]$ 



Maximum errors with various functions used to determine the output formats of division.

## How fast is generating triangular matrix inversion codes?

### • We use $f_4(i_1, i_2) = \lfloor (i_1 + i_2)/2 \rfloor + 1$ to set the output format of division



Generation time for the inversion of triangular matrices of size 4 to 40.

## How fast is generating triangular matrix inversion codes?

### • We use $f_4(i_1, i_2) = \lfloor (i_1 + i_2)/2 \rfloor + 1$ to set the output format of division



Error bounds and experimental errors for the inversion of triangular matrices of size 4 to 40.

## Decomposing some well known matrices

- 2 ill-conditioned matrices: Hilbert and Cauchy
- 2 well-conditioned matrices: KMS and Lehmer



## Decomposing some well known matrices

- 2 ill-conditioned matrices: Hilbert and Cauchy
- 2 well-conditioned matrices: KMS and Lehmer



- Ill-conditioned matrices tend to overflow more often
  - similar behaviour in floating-point arithmetic
- The decompositions of KMS and Lehmer are highly accurate

## Outline of the talk

### A framework for certified fixed-point code synthesis

- Formalization and implementation of an arithmetic model
  - allows certification

▶ handles √ and /

### A framework for certified fixed-point code synthesis

- Formalization and implementation of an arithmetic model
  - allows certification
- Adaptation of the CGPE tool to the model:
  - generates code for fine grained expressions

▶ handles √ and /

instruction selection

### A framework for certified fixed-point code synthesis

- Formalization and implementation of an arithmetic model
  - allows certification
- Adaptation of the CGPE tool to the model:
  - generates code for fine grained expressions
- Development of FPLA:
  - automated and certified code synthesis for linear algebra basic block
    - → matrix multiplication: accuracy vs. code size trade-offs,
    - → Cholesky decomposition and triangular matrix inversion: study of divisions' impact

- handles \sqrt{ and }
- instruction selection

### A framework for certified fixed-point code synthesis

handles ,/ and /

instruction selection

- Formalization and implementation of an arithmetic model
  - allows certification
- Adaptation of the CGPE tool to the model:
  - generates code for fine grained expressions
- Development of FPLA:
  - automated and certified code synthesis for linear algebra basic block
    - → matrix multiplication: accuracy vs. code size trade-offs,
    - → Cholesky decomposition and triangular matrix inversion: study of divisions' impact

#### Publications

-	DASIP14: Toward the synthesis of fixed-point code for matrix inversion based on Cholesky decomposition	[?]
	SYNASC14: Automated Synthesis of Target-Dependent Programs for Polynomial Evaluation in Fixed-Point Arithmetic	[?]
	PECCS14: Code Size and Accuracy-Aware Synthesis of Fixed-Point Programs for Matrix Multiplication	[?]
	DASIP12: Design of Fixed-point Embedded Systems (DEFIS) French ANR Project.	[?]
	SCAN12: Approach based on instruction selection for fast and certified code generation	[?]

Integrate the matrix inversion flow



Integrate the matrix inversion flow



Extend the code size vs. accuracy trade-off strategy to Cholesky decomposition and triangular matrix inversion

Integrate the matrix inversion flow



- Extend the code size vs. accuracy trade-off strategy to Cholesky decomposition and triangular matrix inversion
- Extend the arithmetic model to support complex arithmetic
  - Objective: inverting co-variance matrices for Space Time Adaptive Processing

Integrate the matrix inversion flow



- Extend the code size vs. accuracy trade-off strategy to Cholesky decomposition and triangular matrix inversion
- Extend the arithmetic model to support complex arithmetic
  - Objective: inverting co-variance matrices for Space Time Adaptive Processing
- Target hardware implementations
  - Suggest an arithmetic model for fully custom word-length variables
  - Code size vs. accuracy → chip area vs. accuracy



Ε

[Wil98] H. Keding, M. Willems, M. Coors, and H. Meyr. Fridge: a fixed-point design and simulation environment.

[IEEE754] IEEE 754. IEEE Standard for Floating-Point Arithmetic.

- [MCCS02] Daniel Menard, Daniel Chillet, François Charot, and Olivier Sentieys. Automatic floating-point to fixed-point conversion for DSP code generation.
- [IBMK10] Ali Irturk, Bridget Benson, Shahnam Mirzaei, and Ryan Kastner. GUSTO: An Automatic Generation and Optimization Tool for Matrix Inversion Architectures.

[LHD12] Benoit Lopez, Thibault Hilaire, and Laurent-Stéphane Didier. Sum-of-products evaluation schemes with fixed-point arithmetic, and their application to IIR filter implementation.

- [FRC03] Claire F. Fang, Rob A. Rutenbar, and Tsuhan Chen. Fast, accurate static analysis for fixed-point finite-precision effects in dsp designs.
- [MRS12] Daniel Ménard, Romuald Rocher, Olivier Sentieys, Nicolas Simon, Laurent-Stéphane Didier, Thibault Hilaire, Benoît Lopez, Eric Goubault, Sylvie Putot, Franck Vedrine, Amine Najahi, Guillaume Revy, Laurent Fangain, Christian Samoyaeu, Fabrice Lemonnier, and Christophe Clienti.

Design of Fixed-Point Embedded Systems (defis) French ANR Project.

[LHD14] Benoit Lopez, Thibault Hilaire, and Laurent-Stéphane Didier. Formatting bits to better implement signal processing algorithms. [Rev09] Guillaume Revy.

Implementation of binary floating-point arithmetic on embedded integer processors - Polynomial evaluation-based algorithms and certified code generation.

[MNR12] Christophe Mouilleron, Amine Najahi, and Guillaume Revy.

Approach based on instruction selection for fast and certified code generation.

- [MR11] Christophe Mouilleron and Guillaume Revy. <u>Automatic Generation of Fast and Certified Code for</u> Polynomial Evaluation.
- [KG08] David R. Koes and Seth C. Goldstein. Near-optimal instruction selection on DAGs.
- [MNR14b] Matthieu Martel, Amine Najahi, and Guillaume Revy. <u>Toward the synthesis of fixed-point code for matrix</u> inversion based on cholesky decomposition.
- [MNR14c] Christophe Mouilleron, Amine Najahi, and Guillaume Revy.

Automated Synthesis of Target-Dependent Programs for Polynomial Evaluation in Fixed-Point Arithmetic.

- [MNR14a] Matthieu Martel, Amine Najahi, and Guillaume Revy. Code Size and Accuracy-Aware Synthesis of Fixed-Point Programs for Matrix Multiplication.
  - CG09] Jason Cong, Karthik Gururaj, Bin Liu 0006, Chunyue Liu, Zhiru Zhang, Sheng Zhou, and Yi Zou. Evaluation of static analysis techniques for fixed-point precision optimization.
- [LV09] Dong-U Lee and John D. Villasenor. Optimized custom precision function evaluation for embedded processors.