Automated synthesis of fixed-point programs:

the case of matrix multiplication and, some elements on matrix inversion

Amine Najahi

Advisers: M. Martel and G. Revy

Équipe-projet DALI, Univ. Perpignan Via Domitia LIRMM, CNRS: UMR 5506 - Univ. Montpellier 2











- Embedded systems are ubiquitous
 - microprocessors and/or DSPs 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 and/or DSPs 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 and/or DSPs 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 and/or DSPs 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

With the floating-point arithmetic, it is very easy to program !!

```
int main()
{
    tarting is a start of the interval of th
```

With the floating-point arithmetic, it is very easy to program !!

```
int main()
{
    f
    int i,j,k;
    float A[N][N]={...}, B[N][N]={...}, C[N][N]={0,...,0};
    for (i = 0; i < N; i++)
    for (j = 0; j < N; j++)
    for (k = 0; k < N; k++) /* This inner loop computes the dot product of row i and column j */
        C[i][j]+=A[i][k]*B[k][j];
}</pre>
```

What makes the problem harder in fixed-point?

Intermediate computations depend on the input variables range and computation scheme

With the floating-point arithmetic, it is very easy to program !!



What makes the problem harder in fixed-point?

Intermediate computations depend on the input variables range and computation scheme

Some works on linear algebra primitives in fixed-point

- Lee et al. (2006): 8 × 8 matrix-vector products for the computation of DCT's
 - The first matrix is constant: DCT coefficients
- Relies on some DCT properties
- Frantz et al. (2007): linear algebra routines (mostly matrix inversion) based on simulation
 - No strict guarantee on the error bounds

Based on lengthy simulations

1. Synthesizing fixed-point formulas: combinatorial and numerical issues

2. Efficient matrix multiplication in fixed-point arithmetic

3. State of the art on matrix inversion in fixed-point arithmetic

Outline of the talk

1. Synthesizing fixed-point formulas: combinatorial and numerical issues

2. Efficient matrix multiplication in fixed-point arithmetic

3. State of the art on matrix inversion in fixed-point arithmetic

Background on fixed-point arithmetic

- Main idea of fixed-point arithmetic:
 - Interpret bit packets as integers coupled with a scale factor: $z \cdot 2^{-n}$
 - Example with z = (10000010)₂ and n = 4



Background on fixed-point arithmetic

- Main idea of fixed-point arithmetic:
 - Interpret bit packets as integers coupled with a scale factor: $z \cdot 2^{-n}$
 - Example with z = (10000010)₂ and n = 4



Δ The scale factor (or fixed-point format) is implicit, only the programmer is aware of it

Background on fixed-point arithmetic

- Main idea of fixed-point arithmetic:
 - Interpret bit packets as integers coupled with a scale factor: $z \cdot 2^{-n}$
 - Example with z = (10000010)₂ and n = 4



 Δ The scale factor (or fixed-point format) is implicit, only the programmer is aware of it

We will denote by Q_{a,b} a fixed-point format with a integer bits and b fractional bits

Fixed-point arithmetic model (1/2)

Arithmetic model to track errors in fixed-point computations

- For each intermediate variable r_i , we store 2 intervals **val** (r_i) and **err** (r_i)
- For each basic operator, we have rules to compute $val(r_i)$ and $err(r_i)$

Fixed-point arithmetic model (1/2)

Arithmetic model to track errors in fixed-point computations

- For each intermediate variable r_i , we store 2 intervals $val(r_i)$ and $err(r_i)$
- For each basic operator, we have rules to compute $val(r_i)$ and $err(r_i)$

Addition:

The two variables have to be in the same fixed-point format



Fixed-point arithmetic model (2/2)

- Multiplication:
 - ► The product of a $Q_{a,b}$ variable by a $Q_{c,d}$ variable yields a $Q_{a+c,b+d}$ variable



Fixed-point arithmetic model (2/2)

- Multiplication:
 - The product of a Q_{a,b} variable by a Q_{c,d} variable yields a Q_{a+c,b+d} variable



Physical and virtual shifts:



Numerical issues in dot product generation

- The building block of matrix multiplication is the dot product operation
 - Let us consider a size 3 dot product: $(a_0 \times b_0) + (a_1 \times b_1) + (a_2 \times b_2)$ and the following input fixed-point formats:

	a ₀	b ₀	a ₁	<i>b</i> 1	a ₂	b ₂
Value	[0.1, 1.57]	[0, 1.98]	[0.01,0.87]	[1.1,1.86]	[0,15.4]	[2, 3.3]
Fixed-point format	Q _{1,7}	Q _{1,7}	Q _{0,8}	Q _{1,7}	Q _{4,4}	Q _{2,6}

Numerical issues in dot product generation

- The building block of matrix multiplication is the dot product operation
 - Let us consider a size 3 dot product: $(a_0 \times b_0) + (a_1 \times b_1) + (a_2 \times b_2)$ and the following input fixed-point formats:

	a ₀	b ₀	a ₁	b ₁	a ₂	b ₂
Value	[0.1, 1.57]	[0,1.98]	[0.01,0.87]	[1.1, 1.86]	[0,15.4]	[2, 3.3]
Fixed-point format	Q _{1,7}	Q _{1,7}	Q _{0,8}	Q _{1,7}	Q _{4,4}	Q _{2,6}

Let us focus on 2 different schemes to compute the sum of products:



Numerical issues in dot product generation

- The building block of matrix multiplication is the dot product operation
 - Let us consider a size 3 dot product: $(a_0 \times b_0) + (a_1 \times b_1) + (a_2 \times b_2)$ and the following input fixed-point formats:

	a ₀	b ₀	a ₁	b ₁	a ₂	b ₂
Value	[0.1, 1.57]	[0,1.98]	[0.01,0.87]	[1.1, 1.86]	[0,15.4]	[2, 3.3]
Fixed-point format	Q _{1,7}	Q _{1,7}	Q _{0,8}	Q _{1,7}	Q _{4,4}	Q _{2,6}

Let us focus on 2 different schemes to compute the sum of products:



Combinatorial issues in dot product generation

Number of dot product evaluation schemes

Given by the sequence A001147(n) in the OEIS and the formula: (2n-1)!!

Dot product size	3	5	10	16	20	
Number of schemes	3	105	$34459425 \approx 2^{25}$	$6190283353629375 \approx 2^{52}$	$8200794532637891559375\approx 2^{73}$	

Remarks

- Picking a scheme that minimizes the evaluation error is one of the difficulties of writing fixed-point code
 - Makes it hard to write fixed-point code by hand
 - Appeals for tools with strong heuristics to automate the process

The CGPE¹ library

- Initially developed by Revy and Mouilleron
 - With the aim of generating fast and certified C code for polynomial evaluation
- fast → selects schemes that reduce the evaluation latency on a given target, by using (as much as possible) the architectural features
- certified → produces a bound on the error entailed by the evaluation within the given target's arithmetic



¹Code Generation for Polynomial Evaluation

The CGPE¹ library

- Initially developed by Revy and Mouilleron
 - With the aim of generating fast and certified C code for polynomial evaluation
- fast → selects schemes that reduce the evaluation latency on a given target, by using (as much as possible) the architectural features
- certified → produces a bound on the error entailed by the evaluation within the given target's arithmetic



- Front-ends available so far: sum, dot product, univariate and bivariate polynomials
- Back-ends available so far: C code, VHDL code, GAPPA certificates

¹Code Generation for Polynomial Evaluation

Outline of the talk

1. Synthesizing fixed-point formulas: combinatorial and numerical issues

2. Efficient matrix multiplication in fixed-point arithmetic

3. State of the art on matrix inversion in fixed-point arithmetic

Defining the problem

Inputs:

▶ a black box (CGPE) that synthesises code for dot products in fixed-point arithmetic



2 fixed-point matrices A and B

Defining the problem

Inputs:

▶ a black box (CGPE) that synthesises code for dot products in fixed-point arithmetic



2 fixed-point matrices A and B

Output

C code that evaluates the product $M = A \cdot B$ in fixed-point arithmetic

Straightforward algorithms

Accurate product

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

AccurateProduct

Inputs:

Two fixed-point square matrices A and B

Outputs:

C code to compute the product AB

Steps:

- 1: for $1 < i \le n$ do
- 2: **for** $1 < j \le n$ **do**
- 3: $cgpeGenDotProduct(A_i, B_j);$
- 4: end for

5: end for

Compact product

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

CompactProduct

Inputs:

Two fixed-point square matrices A and B

Outputs:

C code to compute the product AB

Steps:

- 1: compute v such that $v = A_1 \cup A_2 \cup \cdots \cup A_n$
- 2: compute *w* such that $w = B_1 \cup B_2 \cup \cdots \cup B_n$
- 3: cgpeGenDotProduct(v,w);

Illustration through a toy example

We consider the multiplication of the following two fixed-point matrices:

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

Illustration through a toy example

We consider the multiplication of the following two fixed-point matrices:

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

Accurate product

- 1. Output format of *DotProduct*_{0,0}: *Q*_{26,6}
- 2. Output format of *DotProduct*_{0,1}: Q_{18,14}
- 3. Output format of *DotProduct*_{1,0}: *Q*_{15,17}
- 4. Output format of *DotProduct*_{1,1}: Q_{7,25}
- Certified errors bounds:

 $\begin{pmatrix} 0.03125 & 0.00012207 \\ (1.52588e - 05 & 5.96046e - 08 \end{pmatrix}$

Average error bound: $0.00784 \approx 2^{-7}$

Compact product

$$u = ([-1000, 1000] \quad [-3000, 3000])$$
$$v = ([-2000, 2000])$$
$$(-4000, 4000])$$

- 1. Output format of *DotProduct*_{U,V}: Q_{26,6}
- Certified errors bounds:

 $\begin{pmatrix} 0.03125 & 0.03125 \\ 0.03125 & 0.03125 \end{pmatrix}$

Average error bound: 0.03125 $\approx 2^{-5}$

Remarks

- Accurate product generates large code sizes (prohibitive in embedded systems)
- Compact product generates 1 dot product, to the expense of numerical accuracy

Remarks

- Accurate product generates large code sizes (prohibitive in embedded systems)
- Compact product generates 1 dot product, to the expense of numerical accuracy
- But are there interesting trade-offs to look for?

Remarks

- Accurate product generates large code sizes (prohibitive in embedded systems)
- Compact product generates 1 dot product, to the expense of numerical accuracy
- But are there interesting trade-offs to look for?



Accurate product



Remarks

- Accurate product generates large code sizes (prohibitive in embedded systems)
- Compact product generates 1 dot product, to the expense of numerical accuracy
- But are there interesting trade-offs to look for?

	a00	<i>a</i> 01	a ₀₂	a ₀₃	<i>a</i> 04)
	^a 10	a ₁₁	^a 12	^a 13	a ₁₄	
<i>A</i> =	^a 20	a ₂₁	a ₂₂	a ₂₃	a ₂₄	
	a ₃₀	<i>a</i> 31	a32	agg	a ₃₄	
	a40	a ₄₁	a ₄₂	a43	a ₄₄)

Compact product

	(b00	<i>b</i> 01	b ₀₂	b ₀₃	b ₀₄	
		b ₁₀	b ₁₁	b ₁₂	b ₁₃	b ₁₄	
B =		b ₂₀	b ₂₁	b ₂₂	b ₂₃	b ₂₄	
		b30	b31	b32	b33	b ₃₄	
	l	b40	b41	b42	b43	b44	

Remarks

- Accurate product generates large code sizes (prohibitive in embedded systems)
- Compact product generates 1 dot product, to the expense of numerical accuracy
- But are there interesting trade-offs to look for?





Idea: Merge certain rows/columns to reduce the number of the generated dot products

Remarks

- Accurate product generates large code sizes (prohibitive in embedded systems)
- Compact product generates 1 dot product, to the expense of numerical accuracy
- But are there interesting trade-offs to look for?





Idea: Merge certain rows/columns to reduce the number of the generated dot products

Number of ways to merge n vectors

Given by the *n*th Bell number

Number of vectors	3	5	10	16	20	
Number of schemes	5	52	$115975 \approx 2^{17}$	$10480142147\approx 2^{33}$	$51724158235372\approx 2^{46}$	

Distances

The Hausdorff distance d_H

$$d_{\underline{H}} : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$$
$$d_{H}(\underline{[a,\overline{a}]}, \underline{[b,\overline{b}]}) = max\{|\underline{a} - \underline{b}|, |\overline{a} - \overline{b}|\}$$

Example



Distances

The Hausdorff distance d_H

$$d_{\underline{H}} : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$$
$$d_{H}(\underline{[a,\overline{a}]}, \underline{[b,\overline{b}]}) = max\{|\underline{a} - \underline{b}|, |\overline{a} - \overline{b}|\}$$

Example

Let
$$A = [-3, 1]$$
 and $B = [2, 4]$ be two intervals in $I(\mathbb{R})$, we have:

$$\cup (A, B) = [-3, 4] \qquad \qquad \blacksquare \quad d_H(A, B) = 5$$

$$\xrightarrow{\begin{array}{c} a_1 \\ \hline a_1 \\ \hline \end{array}} \qquad \qquad \textcircled{\begin{array}{c} a_1 \\ \hline a_1 \\ \hline \end{array}} \qquad \qquad \textcircled{\begin{array}{c} b_1 \\ \hline a_H(A, B) \end{array}} \qquad \qquad \blacksquare \quad d_H(A, B) = 5$$

Another possible criterion

$$d_d : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$$
$$d_d \left([\underline{a}, \overline{a}], [\underline{b}, \overline{b}] \right) = diam \left([\underline{a}, \overline{a}] \cup [\underline{b}, \overline{b}] \right)$$

A. Najah

- Given, a distance and a collection of vectors:
 - we can write a routine that merges the closest pair of vectors.

- Given, a distance and a collection of vectors:
 - we can write a routine that merges the closest pair of vectors.



Closest pair: A0 and A3

- Given, a distance and a collection of vectors:
 - we can write a routine that merges the closest pair of vectors.







Closest pair: A1 and A4

- Given, a distance and a collection of vectors:
 - we can write a routine that merges the closest pair of vectors.







$$A = \begin{pmatrix} \frac{a_{00}}{a_{01}} & \frac{a_{02}}{a_{02}} & \frac{a_{03}}{a_{03}} & \frac{a_{04}}{a_{04}} \\ \frac{a_{10}}{a_{11}} & \frac{a_{12}}{a_{12}} & \frac{a_{13}}{a_{13}} & \frac{a_{14}}{a_{14}} \\ \frac{a_{20}}{a_{21}} & \frac{a_{22}}{a_{22}} & \frac{a_{23}}{a_{23}} & \frac{a_{24}}{a_{24}} \end{pmatrix} \begin{pmatrix} A_0' = A_0 \cup A_3 \\ A_1' = A_1 \cup A_4 \\ A_2 \end{pmatrix}$$

Algorithm 1 Dynamic Closest Pair algorithm

Inputs:

```
Two square matrices A \in \mathbb{F}ix^{n \times n} and B \in \mathbb{F}ix^{n \times n}
a criterion \mathscr{C}
```

Outputs:

C code to compute the product AB s.t. C is satisfied

Steps:

```
1: S_1 = \{A_1, \dots, A_n\}
 2: \mathcal{S}_2 = \{B_1, \dots, B_n\}

 while C is satisfied do

         (u_1, u_2, d_{u_1, u_2}) = findClosestPair(\mathcal{S}_1)
 4·
 5:
        (v_1, v_2, d_{v_1, v_2}) = findClosestPair(\mathcal{S}_2)
 6:
         if d_{U_1,U_2} \le d_{V_1,V_2} then
             remove(u_1, \bar{\mathcal{S}}_1); remove(u_2, \mathcal{S}_1); insert(u_1 \cup u_2, \mathcal{S}_1)
 7.
 8:
         else
 9:
             remove(v_1, \mathcal{S}_2); remove(v_2, \mathcal{S}_2); insert(v_1 \cup v_2, \mathcal{S}_2)
10.
          end if
11:
          for A_i \in \mathcal{S}_1 do
12:
               for B_i \in \mathscr{S}_2 do
                   cgpeGenDotProduct(A_i, B_i)
13
14:
               end for
15
          end for
16: end while
```

Algorithm 2 Dynamic Closest Pair algorithm

Inputs:

Two square matrices $A \in \mathbb{F}ix^{n \times n}$ and $B \in \mathbb{F}ix^{n \times n}$ a criterion \mathscr{C} : the average error bound is $\leq \epsilon$

Outputs:

C code to compute the product AB s.t. C is satisfied

Steps:

```
1: S_1 = \{A_1, \dots, A_n\}
 2: \mathcal{S}_2 = \{B_1, \dots, B_n\}
 3: while average error \leq \epsilon do
          (u_1, u_2, d_{u_1, u_2}) = findClosestPair(\mathcal{S}_1)
 4:
 5:
        (v_1, v_2, d_{v_1, v_2}) = findClosestPair(\mathcal{S}_2)
 6:
         if d_{U_1,U_2} \le d_{V_1,V_2} then
             remove(u_1, \bar{\mathcal{S}}_1); remove(u_2, \mathcal{S}_1); insert(u_1 \cup u_2, \mathcal{S}_1)
 7.
 8:
         else
 9:
             remove(v_1, \mathcal{S}_2); remove(v_2, \mathcal{S}_2); insert(v_1 \cup v_2, \mathcal{S}_2)
10.
          end if
11:
          for A_i \in \mathcal{S}_1 do
12:
              for B_i \in \mathscr{S}_2 do
                   cgpeGenDotProduct(A_i, B_i)
13
14:
              end for
15
          end for
16: end while
```

Benchmarks

- 1. Weight matrices with dynamic range $2^{\lfloor \frac{n}{2} \rfloor 1}$
- 2. Normally distributed random matrices (generated by matlab)
- 3. We took the Hadamard product of both matrices H
- 4. The matrices fed to the algorithm are midrad(H, 1)











Demo

Outline of the talk

1. Synthesizing fixed-point formulas: combinatorial and numerical issues

2. Efficient matrix multiplication in fixed-point arithmetic

3. State of the art on matrix inversion in fixed-point arithmetic

Matrix inversion in THALES STAP benchmark

```
void Mat_Invert(int ntt, int nsa, Cplfloat In[ntt*nsa][ntt*nsa], Cplfloat Out[ntt*nsa][ntt*nsa]) {
 double inv[nsa*ntt][2*nsa*ntt][2];
 double pivot[2], coef[2];
 float re. im:
 int i=0, j=0, k=0, l=0;
                                                                                                                    -
 for (i=0; i<ntt*nsa; i++)</pre>
                                                                                                                    Initialization
    for (j=0; j<ntt*nsa; j++)</pre>
     inv[i][j][0] = (double) In[i][j].re; inv[i][j][1] = (double) In[i][j].im;
      if(i == i) \{
        inv[i][j+nsa*ntt][0] =1.0; inv[i][j+nsa*ntt][1] =0.0;
        inv[i][j+nsa*ntt][0]=0.0; inv[i][j+nsa*ntt][1] =0.0;
 for (i=0; i<nsa*ntt; i++) {</pre>
   pivot[0]=inv[i][i][0]; pivot[1]=inv[i][i][1];
    if (pivot[0] == 0.) {
      printf("\n Pivot nul re = %f , im = %f\n", pivot[0], pivot[1]);
                                                                                                                    Iterate on number
      exit(0);
    for (j=i; j<2*nsa*ntt; j++)</pre>
      re = inv[i][i][0];
                            im = inv[i][j][1];
      inv[i][j][0] = (re * pivot[0] + im * pivot[1])/(pivot[0] * pivot[0] + pivot[1] * pivot[1]);
      inv[i][i][1] = (im * pivot[0] - re * pivot[1])/(pivot[0] * pivot[0] + pivot[1] * pivot[1]);
    for (k=0; k<nsa*ntt; k++) {</pre>
      if (i!=k) {
        coef[0] = inv[k][i][0]; coef[1] = inv[k][i][1];
                                                                                                                    of rows
        for (l=i; l<2*nsa*ntt; l++) {
          inv[k][1][0] -= (coef[0] * inv[i][1][0] - coef[1] * inv[i][1][1]);
          inv[k][1][1] -= (coef[0] * inv[i][1][1] + coef[1] * inv[i][1][0]);
 for (i=0; i<nsa*ntt; i++) {
    for (j=0; j<nsa*ntt; j++) {</pre>
                                                                                                                    in Out
      Out[i][i].re = (float) inv[i][i+nsa*ntt][0]: Out[i][i].im = (float) inv[i][i+nsa*ntt][1]:
```

the pivot

pivot's

row the

to the rest

Copying

Applying

Assign.

Norm.

Initialization:

r	In _{0,0} In _{1,0}	^{In} 0,1 ^{In} 1,1	···· ···	In _{0,n-1} In _{1,n-1}	1 0	 1	···· ···	0
	:	1	1. I.	:	0	1	- 1.	0
L	In _{n-1,0}	In _{n-1,1}		In _{n-1,n-1}	0			1)

Step 1: Picking the pivot

	In _{0,0} In _{1,0}	In _{0,1} In _{1,1}	 	In _{0,n-1} In _{1,n-1}	1 0	 1	···· ···	0 0	`
	:		÷.	:	0	1	1.	0	
l	In _{n-1,0}	In _{n-1,1}		In _{n-1,n-1}	0			1	,

Step 1: Normalizing the pivot's row

1 In _{1,0}	<u>In_{0,1} In_{0,0} In_{1,1}</u>	 	<u>In_{0,n-1}</u> In _{0,0} In _{1,n-1}	1 1n _{0,0} 0	 1	 	0
: In _{n-1,0}	In _{n-1,1}	·	In _{n-1,n-1}	0	: : 	·	0

Step 1: Putting zeros in the pivot's column





Some remarks on this implementation

- It is based on row reduction
- The implementation is optimized
- Computes the inverse in place
 - avoids backward as well as forward substitution

Which parts are easily convertible to fixed-point?

- Initialization
- 2. Picking the pivot 🗸
- Normalizing the pivot's row (Division by the pivot) X

Issues with the division operator

Consider the analytic method to invert a matrix: $A^{-1} = \frac{com(A)^T}{det(A)}$

For a 2 × 2 matrix
$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$
, $A^{-1} = \frac{1}{det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$

Example

Let $a, b, c, d \in [1, 2]$:

- Suppose the format of *a*, *b*, *c* and *d* is *Q*_{3,29}
- Using interval arithmetic, $det(A) = ad bc \in [-3,3] \rightsquigarrow Q_{3,29}$

• Now consider the matrix
$$M = \begin{pmatrix} 1 + 2^{-28} & 1 \\ 1 & 1 \end{pmatrix} \rightsquigarrow det(M) = 2^{-28}$$

- Then what should be the format of $\frac{a}{det(A)}$?
 - ► To avoid overflow, it should be at least Q_{31,1}

Available works on linear algebra routines

Works that exploits simulation-based methods

- Frantz et al. (2007):
 - Decompositions investigated: SVD, Cholesky, LU and QR
 - Matrices size: up to 30
 - A posteriori error estimation
- Irturk et al. (2006): GUSTO tool
 - Decompositions investigated: Cholesky, LU, QR and analytic
 - Matrices size: up to 8
 - A posteriori error estimation

Works with proven and certified error bounds

• ?

DEFIS WP3 meeting Perpignan, July 9 th, 2013

Automated synthesis of fixed-point programs:

the case of matrix multiplication and, some elements on matrix inversion

Amine Najahi

Advisers: M. Martel and G. Revy

Équipe-projet DALI, Univ. Perpignan Via Domitia LIRMM, CNRS: UMR 5506 - Univ. Montpellier 2









