

# LGen: Program Generator for Linear Algebra



Daniele Spampinato  
Markus Püschel



Computer Science  
**ETH** zürich



**Kalman Filter**

Predict

$$x_k = Ax_{k-1} + Bu_k$$
$$P_k = AP_{k-1}A^T + Q$$

Update

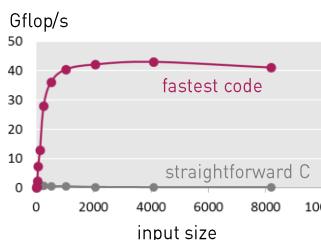
$$x_k = x_k + P_k H^T (HP_k H^T + R)^{-1} (z_k - Hx_k)$$
$$P_k = P_k - P_k H^T (HP_k H^T + R)^{-1} H P_k$$

Fast code needed

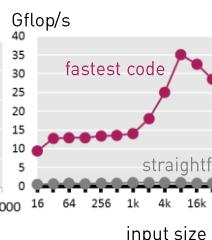
For example, commonly used in robotics  
Let's assume 13 states

## Three Performance Plots

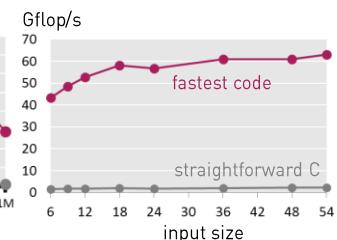
Matrix multiplication



Fast Fourier transform



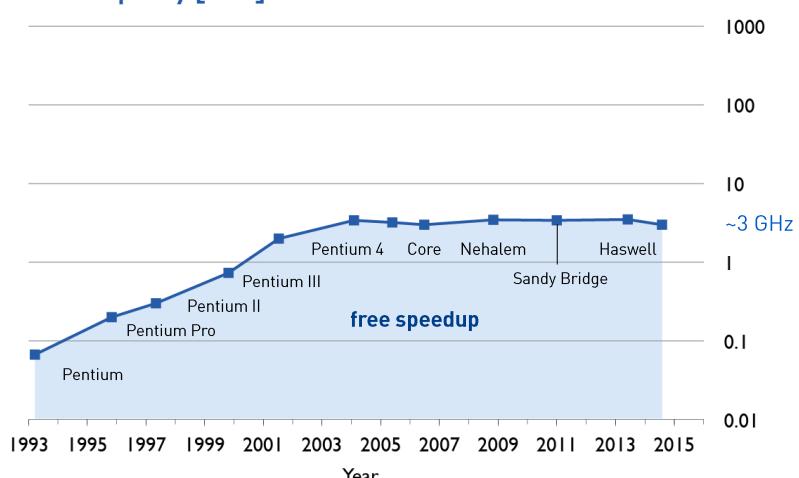
WiFi Receiver



Same op count  
Best compiler + optimization flags

## Evolutions of Processors (Intel)

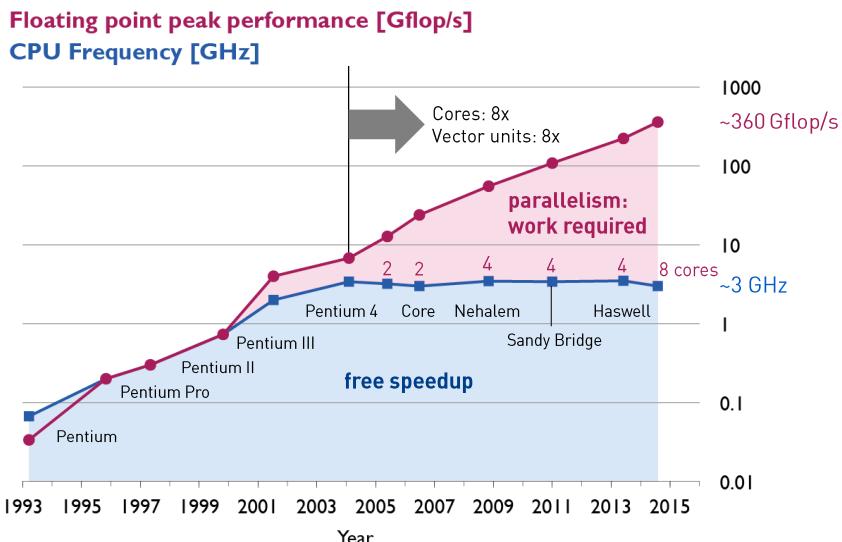
CPU Frequency [GHz]



Source: Wikipedia/Intel/PCGuide

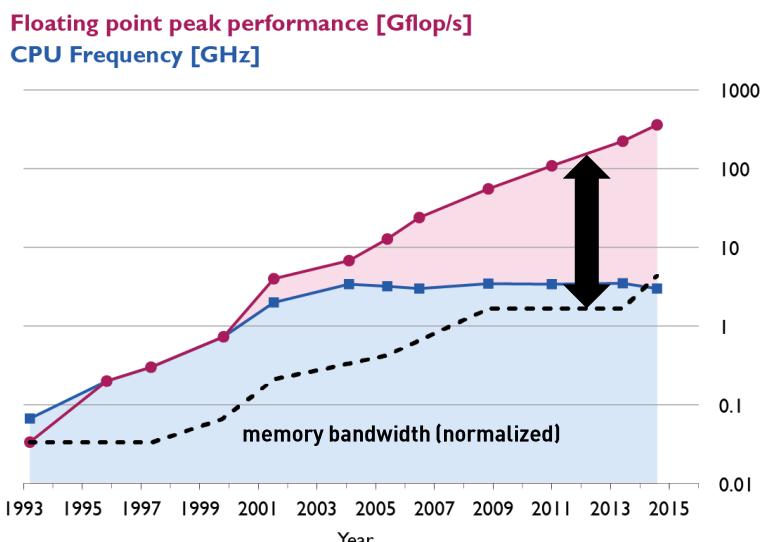
4

## Evolutions of Processors (Intel)



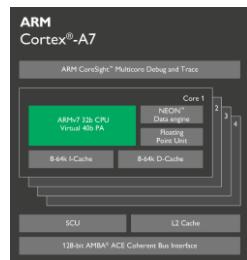
5

## Evolutions of Processors (Intel)

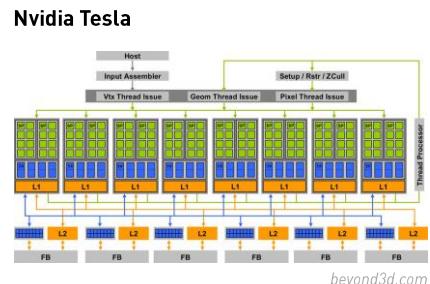


6

## And there is Processor Variety ...

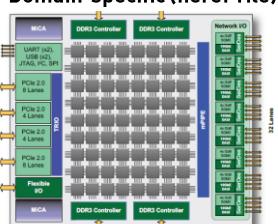


arm.com



beyond3d.com

### Domain-specific (here: Tile)



[mellanox.com](http://mellanox.com)

FPGA accelerators



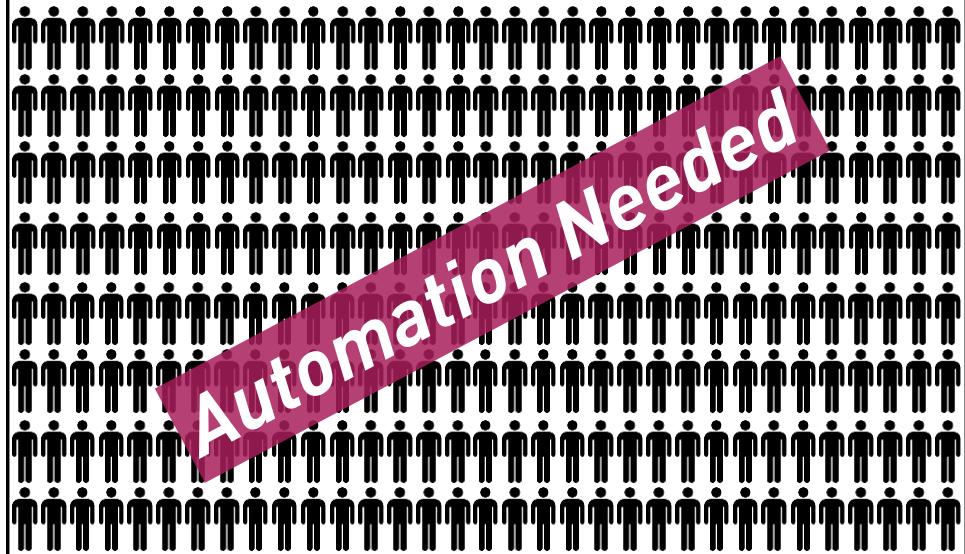
7

**Fast code = good algorithm + code style + locality + vectorization + parallelization**

## LTE Viterbi Decoder (scalar code)

## LTE Viterbi Decoder (vector code)

**Current practice:** Thousands of programmers tune performance-critical code to processors. This is redone for every new processor and for every new processor generation.



# Goal: Program Generation for Linear Algebra

Generate *highest performance code* for linear algebra computations directly from a mathematical description

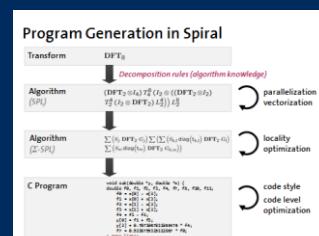
## Approach

Mathematical DSLs

*Rewriting systems for difficult optimizations*

Compiler

### *Learning and search for fine-tuning*

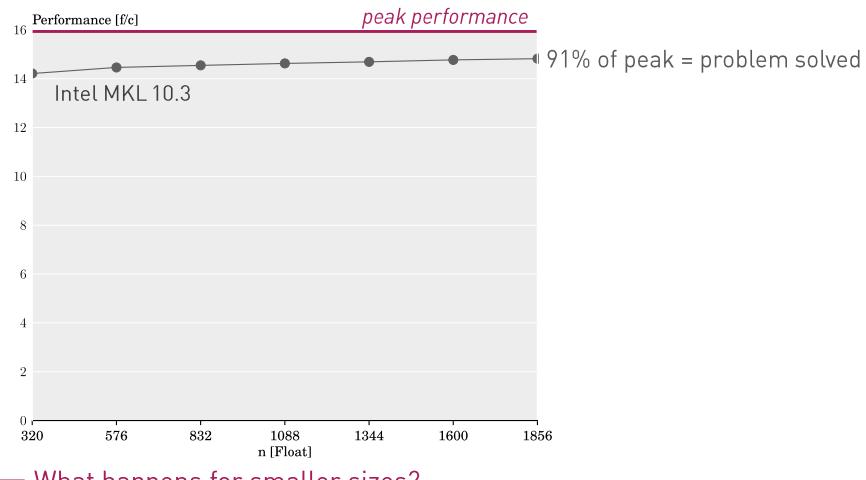


Example: Linear Transforms  
www.spiral.net

## Start: Basic linear algebra

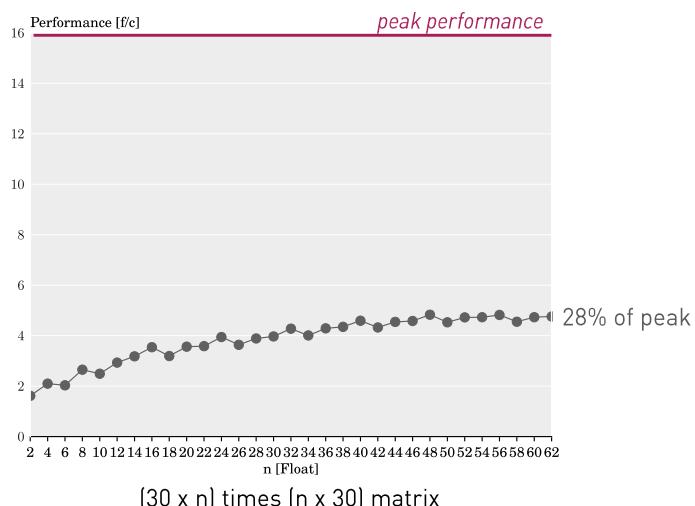
## Library performance sgemm ( $C += AB$ )

Intel Core i7-2600 CPU @ 3.40GHz



## A closer look at small problem sizes

Intel Core i7-2600 CPU @ 3.40GHz



## Are small problems so important?

Required by many performance-critical applications:

*Optimization algorithms*

*Kalman filters*

*Geometric transformations*

*Real-time localization and mapping*

Often for specific input sizes

Do not necessarily comply with standard interface (e.g., BLAS)

Of special interest for a variety of embedded systems

*Reduced HW and SW resources*

## Basic Linear Algebra Computations (BLACs)

*Examples:*

$$y = Ax$$

$$C = \alpha AB^T + \beta C$$

$$\gamma = x^T(A + B)y + \delta$$

*Composed of:*

Scalars, vectors, and matrices

Operators:

*Addition*

*Scalar multiplication*

*Matrix multiplication*

*Transposition*

*Assumption:* All input and output vectors and matrices have a fixed size

## Our Goal: From (any) BLAC to fast code

$$\gamma = \mathbf{x}^T (\mathbf{A} + \mathbf{B}) \mathbf{y} + \delta \quad \leftarrow \mathbf{A} \text{ is } 2 \times 3, \mathbf{x} \text{ is } 3 \times 1, \dots$$

**LGen**

```
void f(double const * A, double const * x, double * y) {
    double t0, ...;

    t0 = x[0];
    t1 = x[1];
    ...
    t9 = t3 * t0;
    t10 = t6 * t0;
    t11 = t4 * t1;
    t12 = t9 + t11;
    ...
    y[0] = t16;
    y[1] = t18;
}
```

## Our Goal: From (any) BLAC to fast code

$$\gamma = \mathbf{x}^T (\mathbf{A} + \mathbf{B}) \mathbf{y} + \delta \quad \leftarrow \mathbf{A} \text{ is } 2 \times 3, \mathbf{x} \text{ is } 3 \times 1, \dots$$

**LGen**

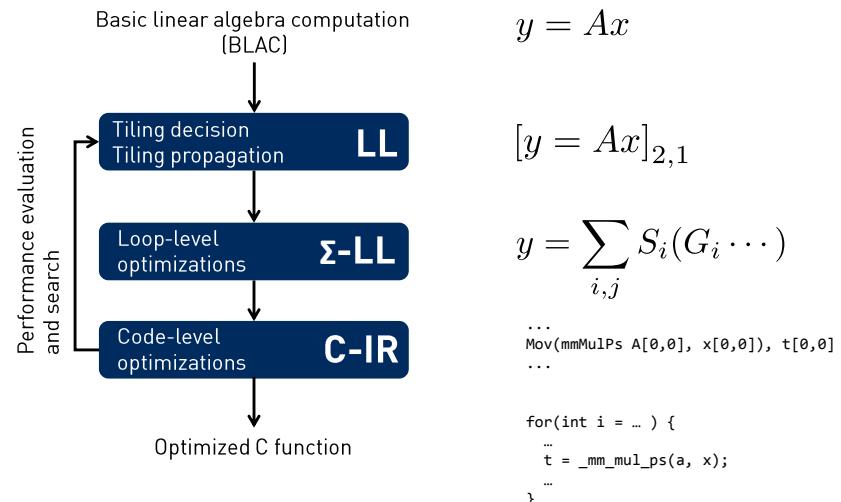
```
void f(double const * A, double const * x, double * y) {
    __m128d t0, ...;

    t0 = _mm_loadu_pd(A);
    t1 = _mm_load_sd(A + 2);
    ...
    t6 = _mm_hadd_pd(_mm_mul_pd(t0, t4), _mm_mul_pd(t2, t4));
    t7 = _mm_shuffle_pd(t1, t3, 0);
    t8 = _mm_mul_pd(t7, _mm_shuffle_pd(t5, t5, 0));
    t9 = _mm_add_pd(t6, t8);

    _mm_storeu_pd(y, t9);
}
```

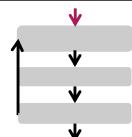
# Architecture of LGen

Design similar to Spiral



## Scalar code generation

The diagram shows the equation  $y = Ax + y$  being generated from a loop iteration. It consists of four terms: a large square block labeled with a brace under it (4), a small vertical rectangle, a plus sign, and another small vertical rectangle. The entire expression is enclosed in a brace labeled 4, indicating four iterations of the loop.



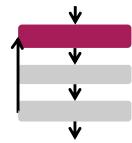
## Tiling in LL

$$\boxed{\quad} = \boxed{\quad} + \boxed{\quad}$$

$$[y = Ax + y]_{r,c}$$

*Tiling decision  
for equation*

$$r = 2, c = 1$$



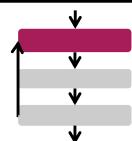
*Task:* Tiling decision for equation  $\rightarrow$  tiling decision for operands

## Tiling in LL

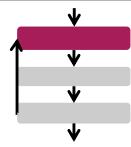
$$\boxed{\quad} = \boxed{\quad} + \boxed{\quad}$$

$$[y = Ax + y]_{2,1}$$

$$[y]_{2,1} = [Ax + y]_{2,1}$$



## Tiling in LL



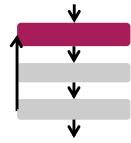
$$\begin{array}{c|c} \hline & \\ \hline \end{array} = \begin{array}{c|c} \hline & \\ \hline \end{array} + \begin{array}{c|c} \hline & \\ \hline \end{array}$$

$$[y = Ax + y]_{2,1}$$

$$[y]_{2,1} = [Ax + y]_{2,1}$$

$$[y]_{2,1} = [Ax]_{2,1} + [y]_{2,1}$$

## Tiling in LL



$$\begin{array}{c|c} \hline & \\ \hline \end{array} = \begin{array}{c|c|c|c} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} + \begin{array}{c|c} \hline & \\ \hline \end{array}$$

$$[y = Ax + y]_{2,1}$$

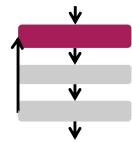
$$[y]_{2,1} = [Ax + y]_{2,1}$$

$$[y]_{2,1} = [Ax]_{2,1} + [y]_{2,1}$$

$$[y]_{2,1} = [A]_{2,k}[x]_{k,1} + [y]_{2,1}$$

Choice that can be used for search

## Tiling in LL



$$\begin{array}{c} \text{---} \\ | \quad | \\ \text{---} \end{array} = \begin{array}{c} \text{---} \\ | \quad | \\ \text{---} \end{array} + \begin{array}{c} \text{---} \\ | \quad | \\ \text{---} \end{array}$$

$$[y = Ax + y]_{2,1}$$

$$[y]_{2,1} = [Ax + y]_{2,1}$$

$$[y]_{2,1} = [Ax]_{2,1} + [y]_{2,1}$$

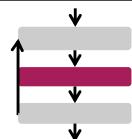
$$[y]_{2,1} = [A]_{2,2}[x]_{2,1} + [y]_{2,1}$$

## $\Sigma$ -LL: Basics

Extension of Sigma-SPL (Franchetti et al., PLDI 2005)

Gathers:  $G_L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$        $G_R = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$

Scatters:  $S_L = G_R$        $S_R = G_L$



Extracting a block

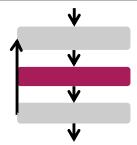
$$A = \begin{array}{|c|c|} \hline B & \text{---} \\ \hline \text{---} & \text{---} \\ \hline \end{array} \quad B = A(0 : 1, 0 : 1) = G_L A G_R$$

Expanding a block

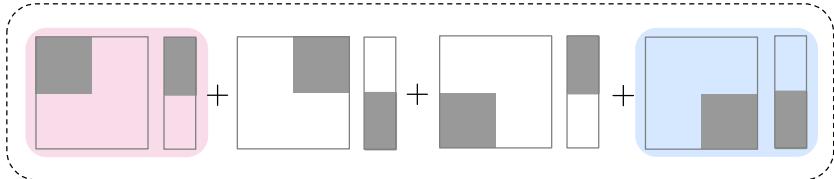
$$C = \begin{array}{|c|c|} \hline B & 0 \\ \hline 0 & 0 \\ \hline \end{array} \quad C = S_L B S_R$$

**Gathers** and **scatters** make data accesses explicit

## LL to $\Sigma$ -LL



$$[y]_{2,1} = \boxed{[A]_{2,2}[x]_{2,1}} + [y]_{2,1}$$

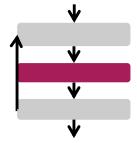


$$S_0(G_0AG_0)S_0 \cdot S_0(G_0x) + \dots + S_2(G_2AG_2)S_2 \cdot S_2(G_2x)$$

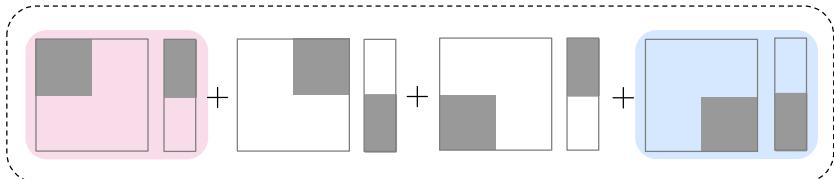
$$= \sum_{\iota=0,2}^3 \sum_{j=0,2}^3 S_\iota(G_\iota AG_j)(G_jx)$$

$$= \sum_{\iota=0,2}^3 \sum_{j=0,2}^3 S_\iota \sum_{\iota'=0}^1 \sum_{j'=0}^1 S_{\iota'}(G_{\iota'}G_\iota AG_jG_{j'})(G_{j'}G_jx)$$

## LL to $\Sigma$ -LL



$$[y]_{2,1} = \boxed{[A]_{2,2}[x]_{2,1}} + [y]_{2,1}$$

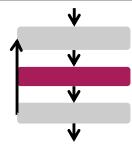


$$S_0(G_0AG_0)S_0 \cdot S_0(G_0x) + \dots + S_2(G_2AG_2)S_2 \cdot S_2(G_2x)$$

$$= \sum_{\iota=0,2}^3 \sum_{j=0,2}^3 S_\iota(G_\iota AG_j)(G_jx)$$

$$= \sum_{\iota,j,\iota',j'} S_{\iota+\iota'}(G_{\iota+\iota'}AG_{j+j'})(G_{j+j'}x)$$

## LL to $\Sigma$ -LL: Loop fusion



$$[y]_{2,1} = [A]_{2,2}[x]_{2,1} + [y]_{2,1}$$



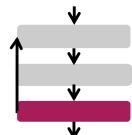
$$\begin{cases} t = \sum_{\iota, j, \iota', j'} S_{\iota+\iota'} (G_{\iota+\iota'} A G_{j+j'}) (G_{j+j'} x) \\ y = \sum_{i, i'} S_{i+i'} (G_{i+i'} t + G_{i+i'} y) \end{cases}$$

$\downarrow$

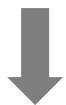
$$y = \sum_{i, i', j, j'} S_{i+i'} [(G_{i+i'} A G_{j+j'}) (G_{j+j'} x) + G_{i+i'} y]$$



## $\Sigma$ -LL to C-IR



$$y = \sum_{i, i', j, j'} S_{i+i'} [(G_{i+i'} A G_{j+j'}) (G_{j+j'} x) + G_{i+i'} y]$$



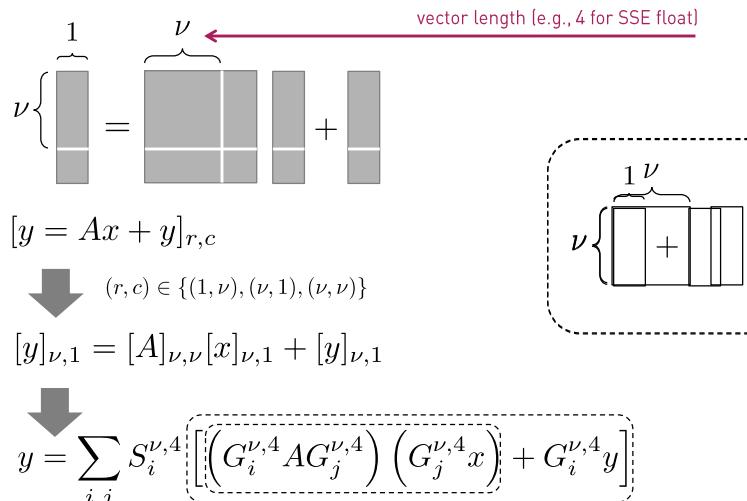
Loop unrolling

```
...
Mov (Mul A[0,0], x[0,0]), t[0,0]
Mov (Mul A[0,1], x[1,0]), t[1,0] ← Scalar replacement
Mov (Mul A[0,2], x[2,0]), t[2,0] ← SSA normalization
...

```

Peephole optimizations

## Vector code generation: Basic Idea



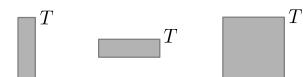
*Computation expressed in terms of  $\nu$ -BLACS*

## v-BLACs

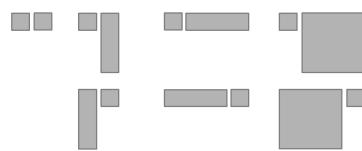
Addition (3 v-BLACs)



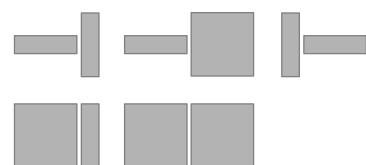
Transposition (3 v-BLACs)



Scalar Multiplication (7 v-BLACs)



Matrix Multiplication (5 v-BLACs)



*18 cases implemented once for every vector ISA*

## Performance evaluation & search

Search on tiling strategies

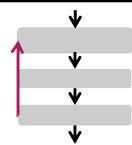
$$r \{ \overbrace{\text{---}}^c \} = \overbrace{\begin{array}{|c|c|} \hline k & \\ \hline \end{array}} + \overbrace{\text{---}}$$

Other degrees of freedom: currently model

Current search methods:

*exhaustive search*

*random search (in the following: 10 samples)*



## Experiments

Hardware details

*Intel Xeon X5680 (Westmere EP) @ 3.3 GHz*

*32 kB L1 D-cache*

*SSE 4.2 (theoretical peak 8 flops/cycle)*

*Intel's SpeedStep and Turbo Boost disabled*

Software details

*RHEL Server 6 – kernel v. 2.6.32*

*icc v. 13.1 with flags: -O3 -xHost -fargument-noalias -fno-alias -ipo*

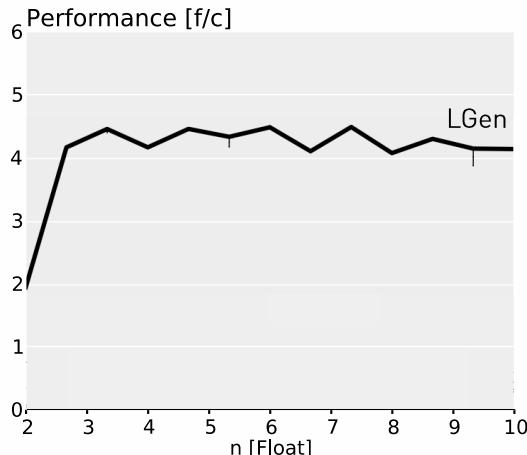
Comparisons

*Handwritten naïve code: Fixed and general size*

*Libraries: Intel MKL v. 11, Intel IPP v. 7.1*

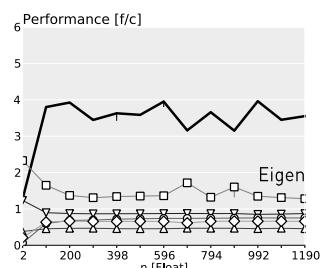
*Generators: Eigen v.3.1.3, BTO v.1.3*

## Plotting

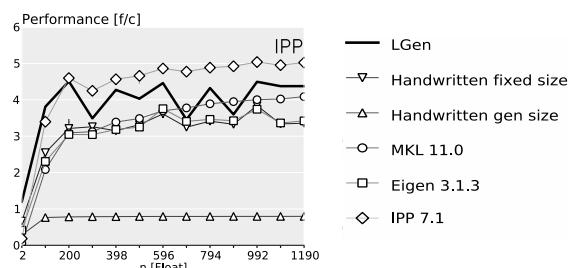


## Case 1: Simple BLACs

$$y = Ax$$



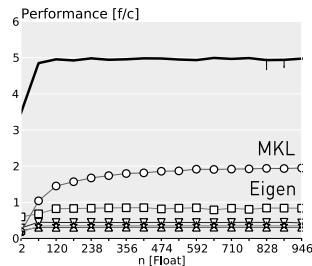
$$A \in \mathbb{R}^{n \times 4}$$



$$A \in \mathbb{R}^{4 \times n}$$

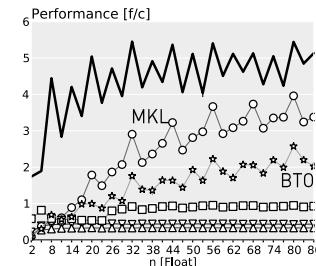
## Case 2: BLACs closely matching BLAS

$$C = \alpha AB + \beta C$$



$$A \in \mathbb{R}^{n \times 4}$$

$$B \in \mathbb{R}^{4 \times 4}$$



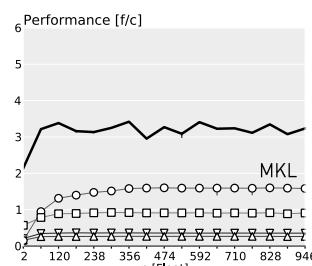
$$A \in \mathbb{R}^{n \times 4}$$

$$B \in \mathbb{R}^{4 \times n}$$

- LGen
- ▽— Handwritten fixed size
- △— Handwritten gen size
- MKL 11.0
- Eigen 3.1.3
- ◇— IPP 7.1

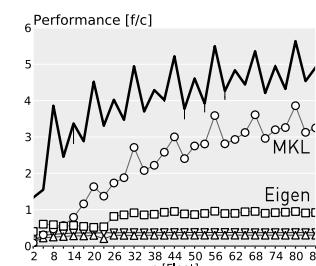
## Case 3: More than one BLAS call

$$C = \alpha(A_0 + A_1)^T B + \beta C$$



$$A_0 \in \mathbb{R}^{4 \times n}$$

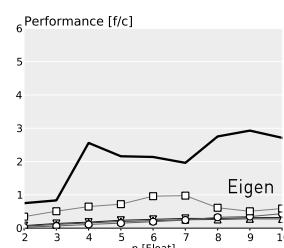
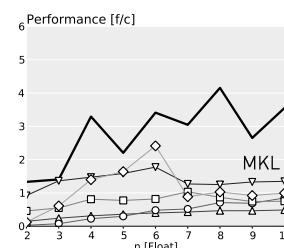
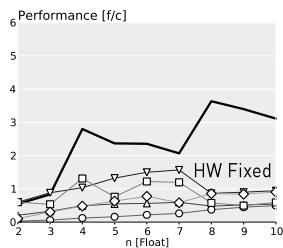
$$B \in \mathbb{R}^{4 \times 4}$$



- LGen
- ▽— Handwritten fixed size
- △— Handwritten gen size
- MKL 11.0
- Eigen 3.1.3
- ◇— IPP 7.1

## Case 4: Micro BLACs

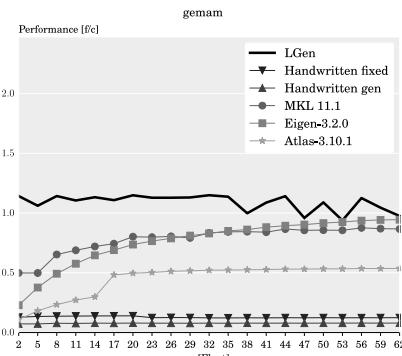
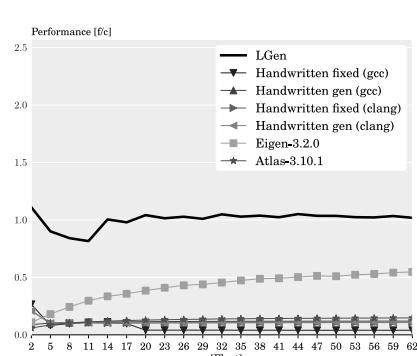
- LGen
- ▽— Handwritten fixed size
- △— Handwritten gen size
- MKL 11.0
- Eigen 3.1.3
- ◇— IPP 7.1



## On Embedded Processors

Work by Nikos Kyrtatas

$$C = \alpha(A_0 + A_1)^T B + \beta C$$



# Challenge: Alignment Analysis

unaligned load/stores only

```

for( size_t i2 = 0; i2 < 400; i2+=16 ) {
    for( size_t j3 = 0; j3 < 112; j3+=4 ) {
        for( size_t ii4 = 0; ii4 < 16; ii4+=4 ) {
            t0_7_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3);
            t0_8_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3 + 115);
            t0_9_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3 + 230);
            t0_4_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3 + 345);
            t0_3_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3);
            t0_2_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3 + 115);
            t0_1_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3 + 230);
            t0_0_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3 + 345);

            // 4x4 -> 4x4 - Incompact
            t0_8_0 = __mm_add_ps(t0_7_0, t0_3_0);
            t0_9_0 = __mm_add_ps(t0_6_0, t0_2_0);
            t0_10_0 = __mm_add_ps(t0_5_0, t0_1_0);
            t0_11_0 = __mm_add_ps(t0_4_0, t0_0_0);

            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3, t0_8_1);
            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3 + 115, t0_9_1);
            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3 + 230, t0_10_1);
            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3 + 345, t0_11_1);
        }
    }
}

```

with aligned load/stores

```

for( size_t i2 = 0; i2 < 400; i2+=16 ) {
    for( size_t j3 = 0; j3 < 112; j3+=4 ) {
        for( size_t ii4 = 0; ii4 < 16; ii4+=4 ) {
            t0_7_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3);
            t0_6_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3 + 115);
            t0_5_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3 + 230);
            t0_4_0 = __mm_loadu_ps(A + 115*i2 + 115*ii4 + j3 + 345);
            t0_3_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3);
            t0_2_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3 + 115);
            t0_1_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3 + 230);
            t0_0_0 = __mm_loadu_ps(B + 115*i2 + 115*ii4 + j3 + 345);

            // 4x4 -> 4x4 - Incompact
            t0_8_0 = __mm_add_ps(t0_7_0, t0_3_0);
            t0_9_0 = __mm_add_ps(t0_6_0, t0_2_0);
            t0_10_0 = __mm_add_ps(t0_5_0, t0_1_0);
            t0_11_0 = __mm_add_ps(t0_4_0, t0_0_0);

            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3, t0_8_1);
            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3 + 115, t0_9_1);
            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3 + 230, t0_10_1);
            _mm_storeu_ps(C + 115*i2 + 115*ii4 + j3 + 345, t0_11_1);
        }
    }
}

```

goal

# Solution: Abstract interpretation

assume B[] is aligned:

```

for( size_t j5 = 0; j5 < 80; j5+=16 ) {
    ...
    for( size_t k4 = 8; k4 < 48; k4+=8 ) {           // k4->([8,40], 0+8Z)
        for( size_t kk7 = 0; kk7 < 8; kk7+=4 ) {     // kk7->([0,4], 0+4Z)
            for( size_t jj8 = 0; jj8 < 16; jj8+=4 ) { // jj8->([0,12], 0+4Z)
                ...
                t2219900=__mm_madd_ps(B + j5 + jj8 + 81*k4 + 81*kk7);
                // Eval(B + j5 + jj8 + 81*k4 + 81*kk7) = ([ -oo,+oo], 0+4Z) + ([0,64], 0+16Z) +
                // ([0,12], 0+4Z) + ([81,81], 81+0Z) * ([8,40], 0+8Z) + ([81,81], 81+0Z) * ([0,4], 0+4Z) =
                // = ([ -oo,+oo], 0+gcd(4,16,4,648,324)Z) = ([ -oo,+oo], 0+4Z)
            }
        }
    }
}

interval ↓           congruence ↓
// B->([-oo,+oo], 0+4Z)
// jj8->([0,12], 0+4Z)
// kk7->([0,4], 0+4Z)
// k4->([8,40], 0+8Z)
// Eval(B + j5 + jj8 + 81*k4 + 81*kk7) = ([ -oo,+oo], 0+4Z) + ([0,64], 0+16Z) +
// ([0,12], 0+4Z) + ([81,81], 81+0Z) * ([8,40], 0+8Z) + ([81,81], 81+0Z) * ([0,4], 0+4Z) =
// = ([ -oo,+oo], 0+gcd(4,16,4,648,324)Z) = ([ -oo,+oo], 0+4Z)

aligned ↑

```

Analysis is sound and precise

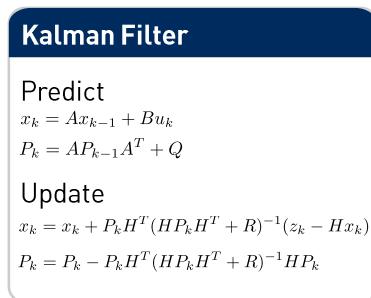
## What's Next?

Step 1: BLACs – [CGO 14, DATE 15]

Step 2: structured BLACs – [CGO 16]

Step 3: higher level linear algebra – collaboration

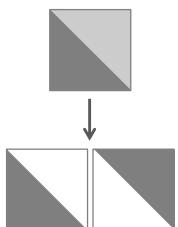
Step 4: applications (e.g., Kalman filter) – not yet



fast code

## Cl1ck: Synthesis of Linear Algebra Algorithms

Cholesky factorization



Cl1ck

```

Algorithm:  $A := \text{CHOL\_BLK\_VAR}_k(A)$ 
Algorithm:  $A := \text{CHOL\_L\_BLK\_VAR}_k(A)$ 

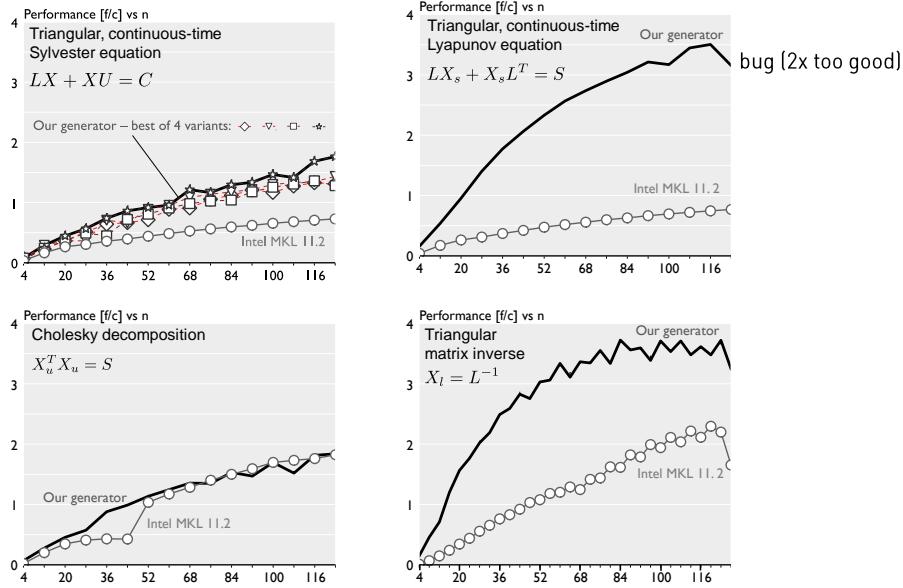
Partition  $A \rightarrow \begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix}$ 
where  $A_{TL}$  is  $0 \times 0$ 
while  $m(A_{TL}) < m(A)$  do
  Determine block size  $b$ 
  Repartition
     $\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \rightarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$ 
    where  $A_{11}$  is  $b \times b$ 
     $A_{11} := A_{11} - A_{10}A_{10}^T \quad (\text{SYRK})$ 
     $A_{21} := A_{21} - A_{20}A_{10}^T \quad (\text{GEMM})$ 
     $A_{11} := \text{CHOL}(A_{11}) \quad (\text{CHOL})$ 
     $A_{21} := A_{21}A_{11}^T \quad (\text{TRSM})$ 
  Continue with
     $\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \leftarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$ 
endwhile

```

### More on Cl1ck

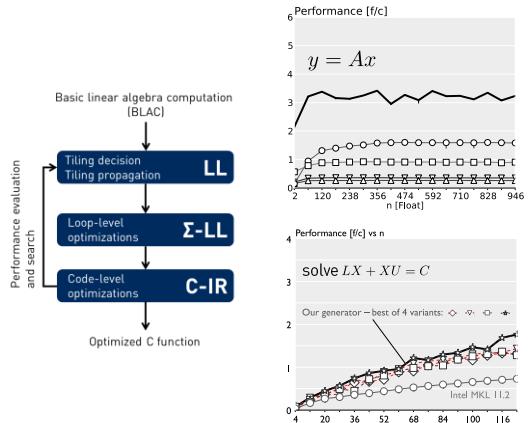
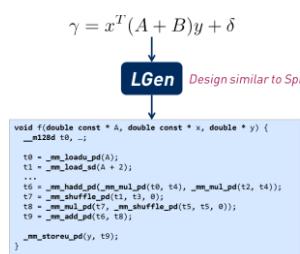
- D. Fabregat-Traver and P. Bientinesi. Knowledge-Based Automatic Generation of Partitioned Matrix Expressions. CASC, 2011.
- D. Fabregat-Traver and P. Bientinesi. Automatic Generation of Loop-invariants for Matrix Operations. ICCSA, 2011.

## Preliminary results



## Summary

Goal: Automatically from linear algebra to fast code



## Soon ...

### Kalman Filter

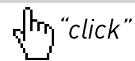
#### Predict

$$x_k = Ax_{k-1} + Bu_k$$
$$P_k = AP_{k-1}A^T + Q$$

#### Update

$$x_k = x_k + P_k H^T (HP_k H^T + R)^{-1} (z_k - Hx_k)$$
$$P_k = P_k - P_k H^T (HP_k H^T + R)^{-1} H P_k$$

Generate Code



For example, commonly used in robotics  
Let's assume 13 states

More info: <http://spiral.net/software/lgen.html>