

LGen

A Basic Linear Algebra Compiler

Daniele Spampinato
Markus Püschel

Computer Science
ETHzürich
SPIRAL 

Linear algebra: Central to many domains

Control systems

```

graph LR
    r((r)) --> Controller[Controller]
    Controller -- e --> Measurements[Measurements]
    ym((y_m)) --> Measurements
    Measurements -- u --> System[System]
    System -- y --> yout((y))
    disturbances[Disturbances] --> System
  
```

Optimization algorithms

Source: rain.aa.washington.edu

Computer Graphics

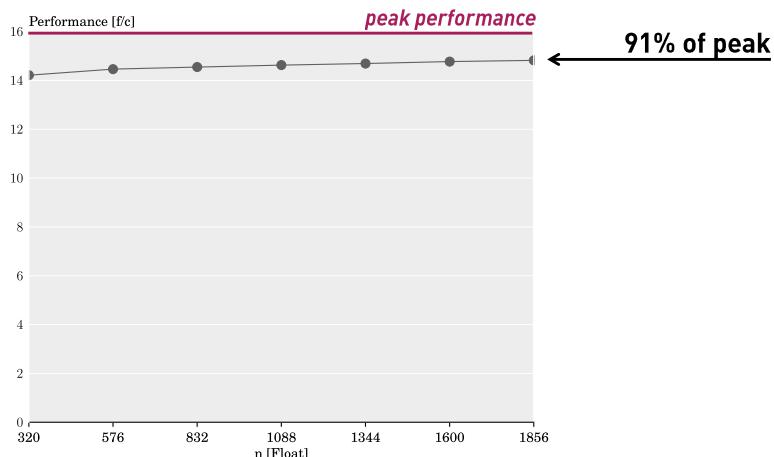
Source: disneyresearch.com

Quoting App performance section

"The frame rate of 12 FPS, however, is not enough to handle sudden motions. In future work, optimisation of implementation would improve the frame rate and provide better user experience."

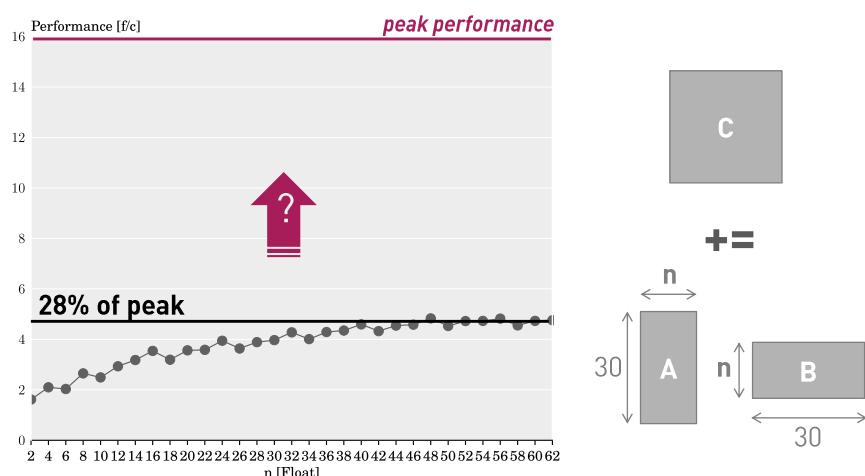
Library performance for sgemm ($C \doteq AB$)

Intel MKL on Intel Core i7 CPU (AVX)

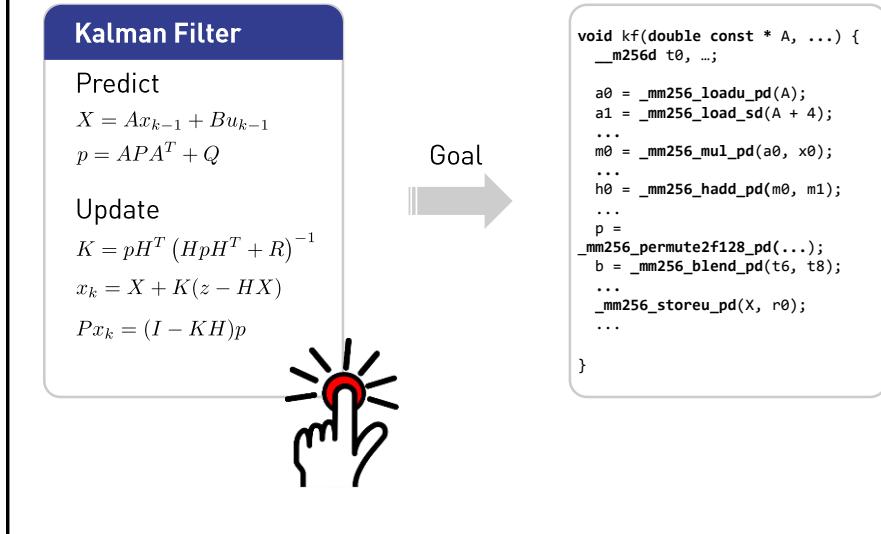


A closer look at small problem sizes

Intel MKL on Intel Core i7 CPU (AVX)



Bridging the gap with DSLs



LGen: Base System

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

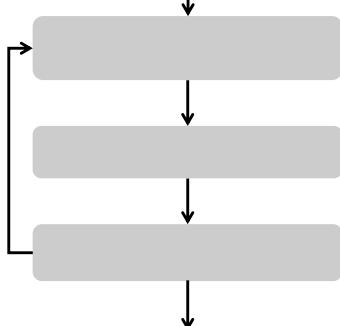
Transposition

All input and output vectors and matrices have a fixed size

LGen: A basic linear algebra compiler

Basic linear algebra computation
(BLAC)

$$y = Ax \quad \leftarrow \begin{array}{l} A \text{ is } 2 \times 3 \\ x \text{ is } 3 \times 1 \end{array}$$



LGen: A basic linear algebra compiler

Basic linear algebra computation (BLAC)

$$y = Ax \leftarrow \begin{array}{l} A \text{ is } 2 \times 3 \\ x \text{ is } 3 \times 1 \end{array}$$

Tiling decision
Tiling propagation LL

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

LL

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

LGen: A basic linear algebra compiler

Basic linear algebra computation (BLAC)

$$y = Ax \leftarrow \begin{array}{l} A \text{ is } 2 \times 3 \\ x \text{ is } 3 \times 1 \end{array}$$

Tiling decision
Tiling propagation LL

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

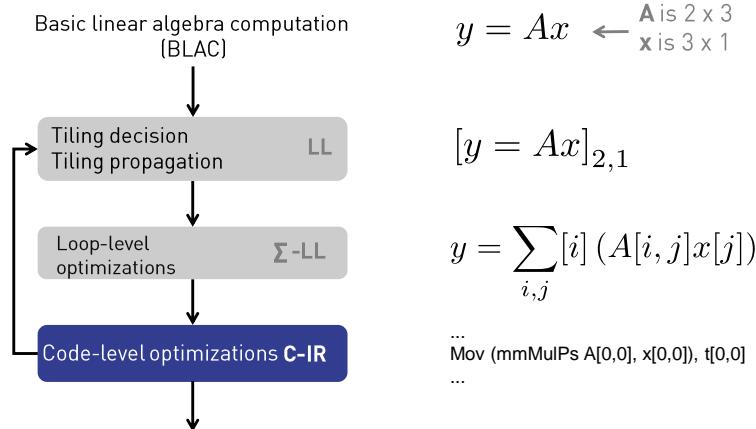
LL

Loop-level optimizations Σ -LL

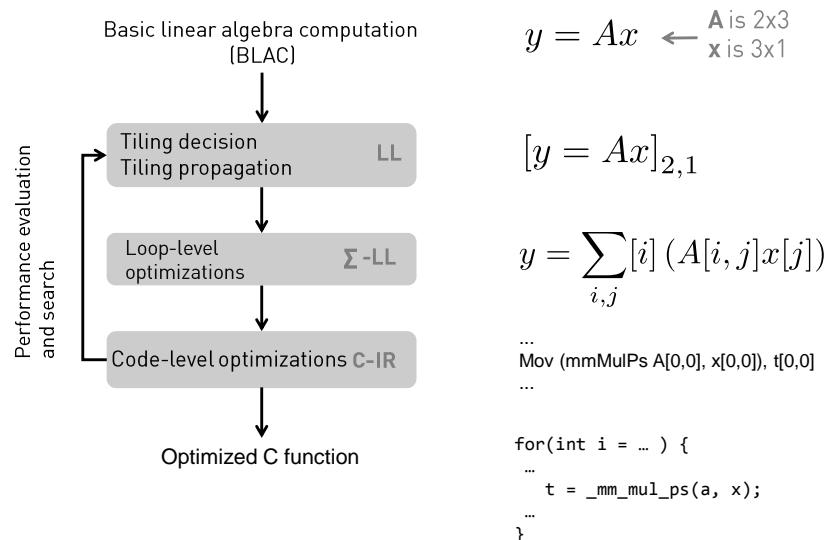
$$y = \sum_{i,j} [i] (A[i, j]x[j])$$

Σ -LL

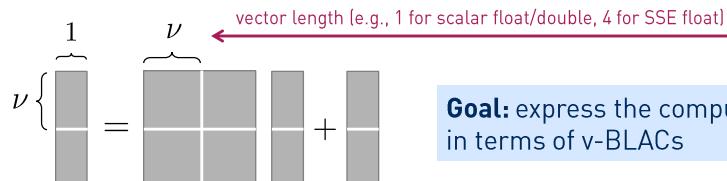
LGen: A basic linear algebra compiler



LGen: A basic linear algebra compiler



Vector code generation: Basic Idea

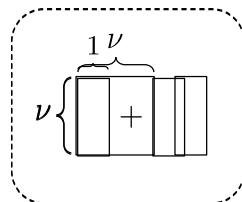


Goal: express the computation in terms of v-BLACs

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

\downarrow

$$y = \sum_{i,j} [i] \left([A[i,j]x[j]] + y[i] \right)$$

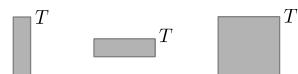


v-BLACs: Vectorization building blocks

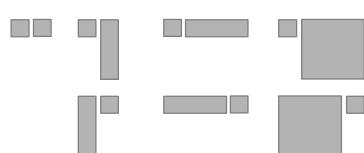
Addition (3 v-BLACs)



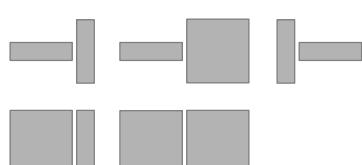
Transposition (3 v-BLACs)



Scalar Multiplication (7 v-BLACs)



Matrix Multiplication (5 v-BLACs)



18 cases implemented once for every ISA

The importance of structures

Kalman Filter

Predict

$$x_k = Ax_{k-1} + Bu_k$$

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

Update

$$K = P_k H^T (HP_k H^T + R)^{-1}$$

$$x_k = x_k + K(z_k - Hx_k)$$

$$P_k = (I - KH)P_k$$



The importance of structures

Kalman Filter

Predict

$$x_k = Ax_{k-1} + Bu_k$$

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

Update

$$K = P_k H^T (HP_k H^T + R)^{-1}$$

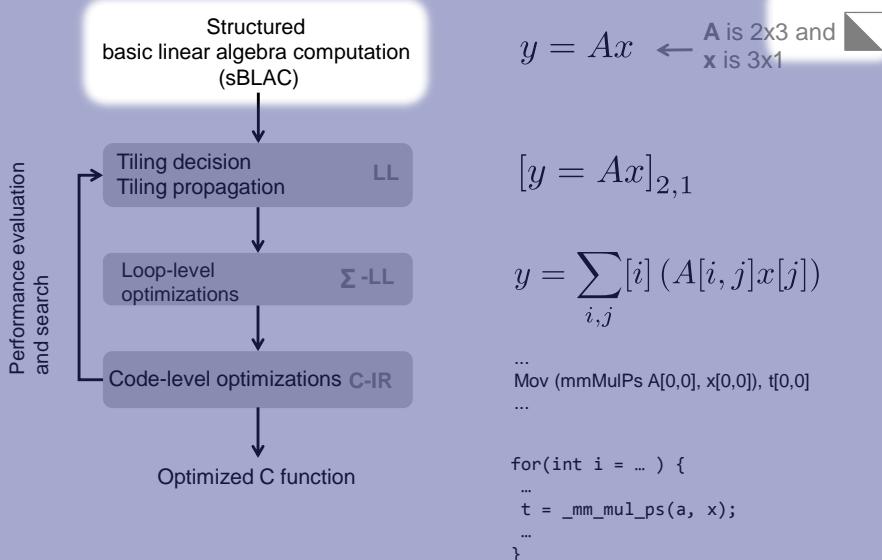
$$x_k = x_k + K(z_k - Hx_k)$$

$$P_k = (I - KH)P_k$$

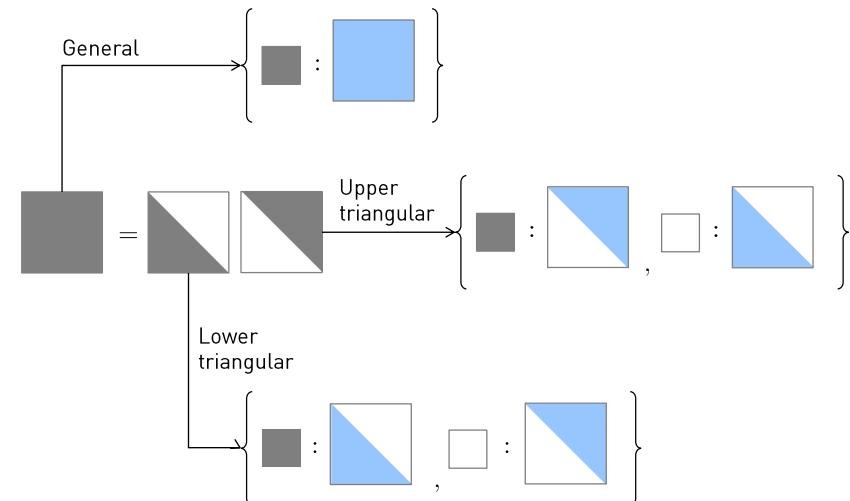


Extending LGen with Structures

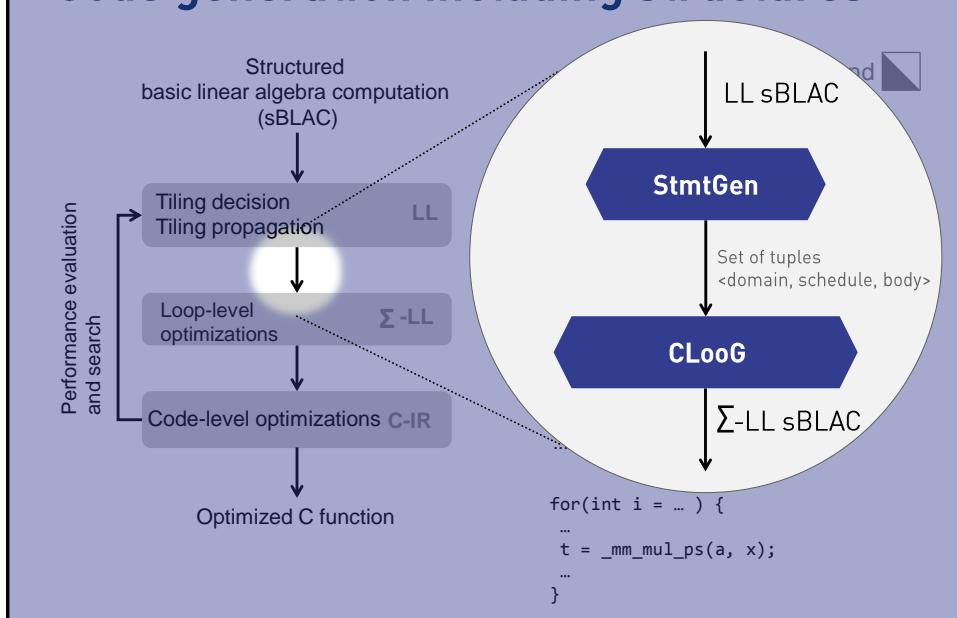
Code generation including structures



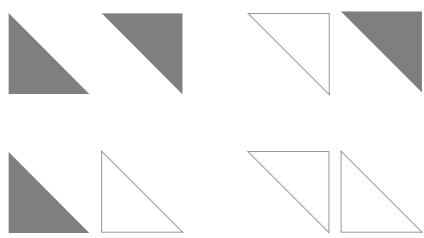
Structured matrices representation



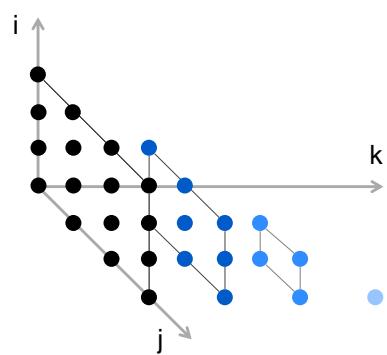
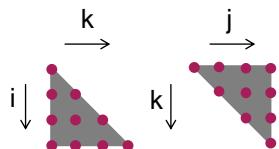
Code generation including structures



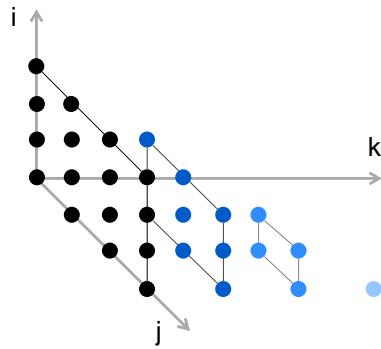
From LL to Σ -LL



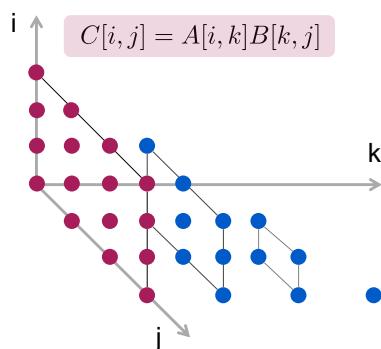
From LL to Σ -LL



From LL to Σ -LL



From LL to Σ -LL

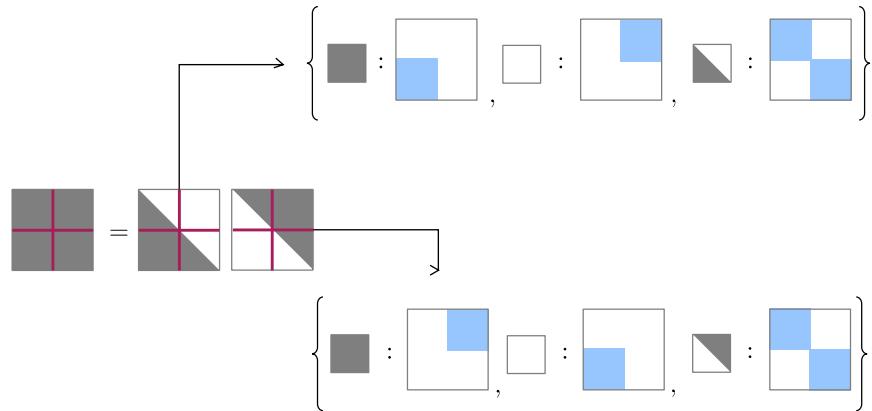


$$\text{CLooG} \Rightarrow C = \sum_{i=0}^3 \sum_{j=0}^3 [i, j] (A[i, 0]B[0, j]) + \sum_{k=1}^3 \sum_{i=k}^3 \sum_{j=k}^3 [i, j] (A[i, k]B[k, j])$$

$$C[i, j] = C[i, j] + A[i, k]B[k, j]$$

Scattering relation built based on optimization models (e.g., Goto model)

Representing structured tiled matrices



From LL to Σ -LL

$$\begin{matrix} \text{pink} & \text{yellow} \\ \text{green} & \text{blue} \end{matrix} = \begin{matrix} \text{dark gray} & \text{white} \\ \text{white} & \text{dark gray} \end{matrix} + \begin{matrix} \text{dark gray} & \text{white} \\ \text{white} & \text{dark gray} \end{matrix}$$

$$\text{pink} = \begin{matrix} \text{dark gray} & \text{white} \\ \text{white} & \text{dark gray} \end{matrix} \quad \text{yellow} = \begin{matrix} \text{dark gray} & \text{white} \\ \text{white} & \text{dark gray} \end{matrix}$$

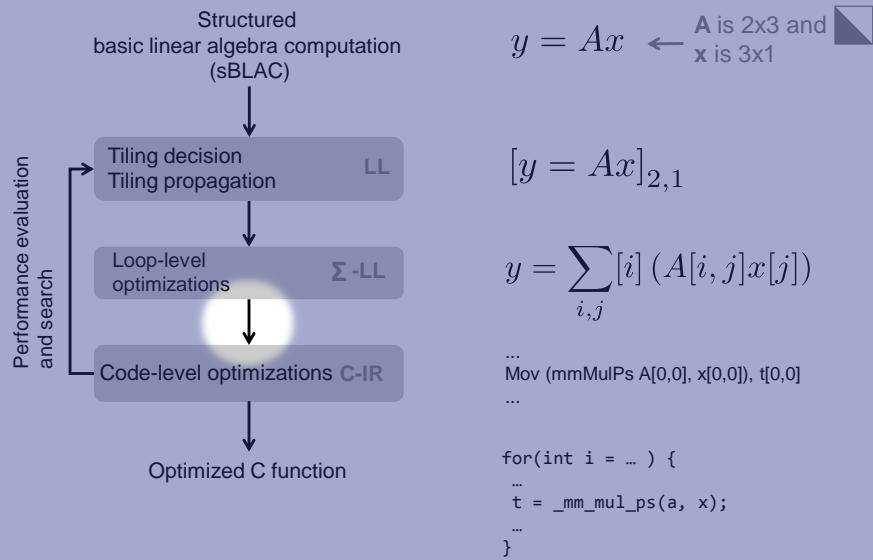
$$\text{green} = \begin{matrix} \text{dark gray} & \text{white} \\ \text{white} & \text{dark gray} \end{matrix}$$

$$\text{blue} = \begin{matrix} \text{dark gray} & \text{dark gray} \\ \text{dark gray} & \text{dark gray} \end{matrix} + \begin{matrix} \text{dark gray} & \text{white} \\ \text{white} & \text{dark gray} \end{matrix}$$

CLooG

$$\begin{aligned} C = & [0, 0](A[0, 0]B[0, 0]) \\ & + [0, 2](A[0, 0]B[0, 2]) \\ & + [2, 0](A[2, 0]B[0, 0]) \\ & + [2, 2](A[2, 0]B[0, 2] + A[2, 2]B[2, 2]) \end{aligned}$$

Code generation including structures



Experiments

Experimental settings

Intel Core i7 [Sandy Bridge]

Ubuntu 14.04 with Linux 3.13

Double computations with AVX

Warm cache scenario (32 kB L1 D-cache, 256 kB L2 cache)

Four competitors:

LGen w/ and w/o structures support

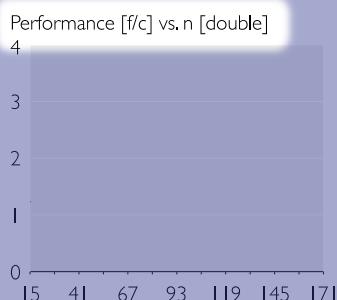
Intel MKL 11.2

Naïve code (non optimized double/triple loop code)

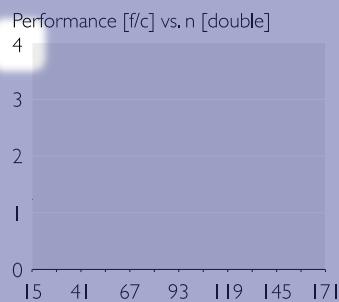
All kernels compiled with `icc 15` w/ flags:

`-O3 -xHost -fargument-noalias -fno-alias`

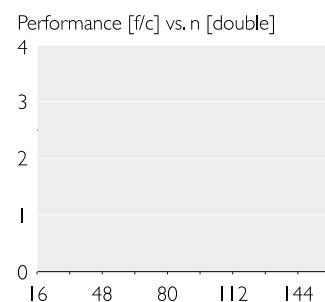
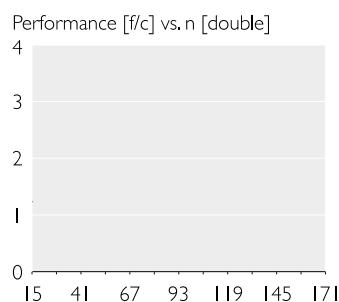
Plotting



Plotting

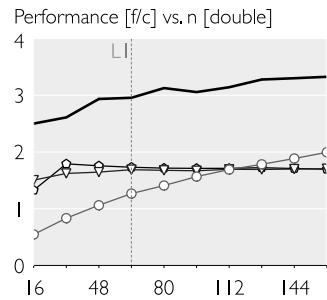
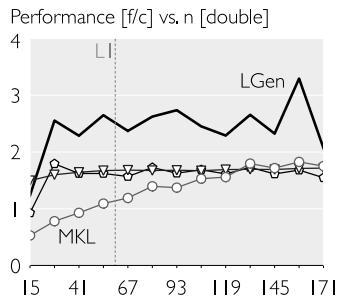


Plotting



BLAS category: dsyrk

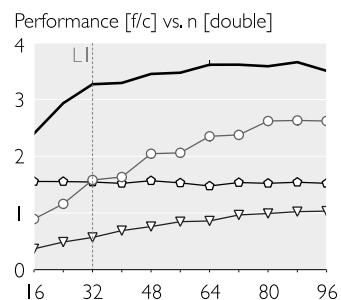
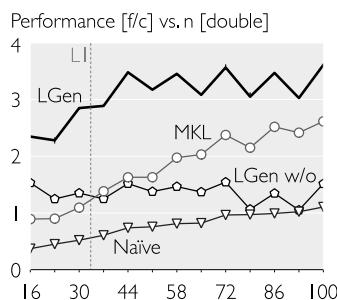
— LGen —○— Intel MKL 11.2 —▽— Naïve (icc 15) —◇— LGen w/o structures



$$S_u = AA^T + S_u, \quad A \in \mathbb{R}^{n \times 4}$$

BLAS-like category

— LGen —○— Intel MKL 11.2 —▽— Naïve (icc 15) —◇— LGen w/o structures

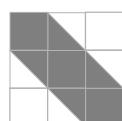


$$A = LU + S_l, \quad L, U \in \mathbb{R}^{n \times n}$$

More general structures & future work

Approach extensibility

Other important structures, e.g., banded matrices



Or other combined structures

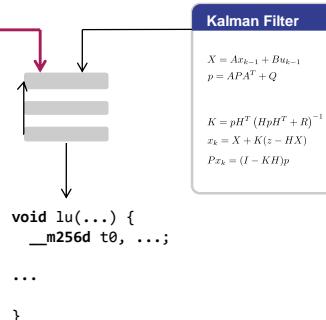


Support for higher-level functionalities

Algorithm: $[A] := \text{LU_BLK_VAR1}(A)$

$\text{Partition } A \rightarrow \left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$ where A_{TL} is 0×0 while $m(A_{TL}) < m(A)$ do Determine block size b Repartition $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c c c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right)$ where A_{11} is $b \times b$ <hr/> $A_{01} := U_{01} = L_{00}^{-1} A_{01}$ $A_{10} := L_{10} = A_{10} U_{00}^{-1}$ $A_{11} := \text{LU}(A_{11} - L_{10} U_{01})$ <hr/> Continue with $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c c c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right)$ endwhile

Source: FLAME Project - <http://www.cs.utexas.edu/~flame/web/>



Connecting with Cl1ck (Fabregat-Traver, Bientinesi)

Conclusion

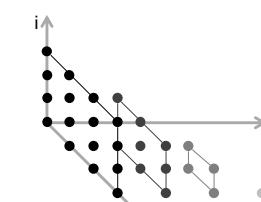
Our Approach

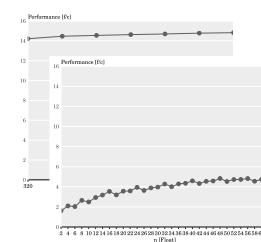
Kalman Filter

$$X = Ax_{k-1} + Bu_{k-1}$$
$$p = AP A^T + Q$$
$$K = p H^T (H p H^T + R)^{-1}$$
$$x_k = X + K(z - HX)$$
$$P x_k = (I - KH)p$$



```
void kernel(...) {  
    __m256d t0, ...;  
    ...  
}
```



Performance (flop)

n (Pivot)	Performance (flop)
2	~1.5
4	~2.5
6	~3.5
8	~4.5
10	~5.5
12	~6.5
14	~7.5
16	~8.5
18	~9.5
20	~10.5
22	~11.5
24	~12.5
26	~13.5
28	~14.5
30	~15.5
32	~16.5



void kernel(...)

```
__m256d t0, ...;  
...  
}
```

spiral.net/software/lgen.html