



Management of Deep Memory Hierarchies - Recursive Blocking and Hybrid Data Structures for Dense Matrix Computations

Bo Kågström
Dept of Computing Science & HPC2N
Umeå University, Sweden

Föreläsning #F8 vt 09: Design och analys av algoritmer för paralleldatorsystem



High Performance Computing Center North (HPC2N)



HPC2N - "HPC to North"

- National center for Scientific and Parallel Computing
- Competence network with 5 partners:
 - Luleå University of Technology
 - Mid Sweden University
 - Swedish Institute of Space Physics
 - Swedish University of Agricultural Sciences - SLU
 - Umeå University
- Funded by the Swedish Research Council and its meta-center SNIC

HPC2N Computer Systems

Computational-intensive jobs, DMM

- **sarek**: HPC2N Opteron Cluster, 384 proc., 04-06
- **ingrid**: HPC2N Swegrid Cluster, 101 proc. 03-10
- **seth**: HPC2N Super Cluster, 240 proc., 02-05

A **new Swegrid cluster** (432 → 512 cores, 864 GB - 1TB RAM, 20 TB disk storage)

A **new scalable cluster** (~ April 10, 2008):

- > 50 Tflops/s peak (> 5000 core) 1 Tera = 10^{12}
- > 10 TB RAM
- > 100 TB disk

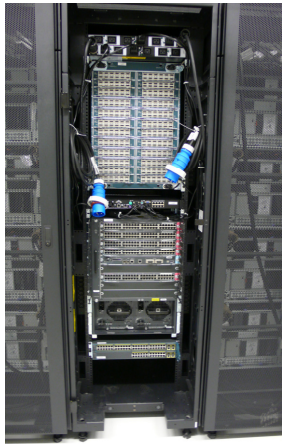
Akka - the "green" HPC2N cluster



- Gridcore
- IBM
- Intel
- Microsoft

- 672 dual Intel L5420 quad core nodes (5376 cores)
- L5420: low power 50W cpu (80W/chip is normal)
- dual boot (Linux, Windows HPC 2008)

Akka - switches and cables



- 288 ports Cisco Infiniband DDR
- GigE

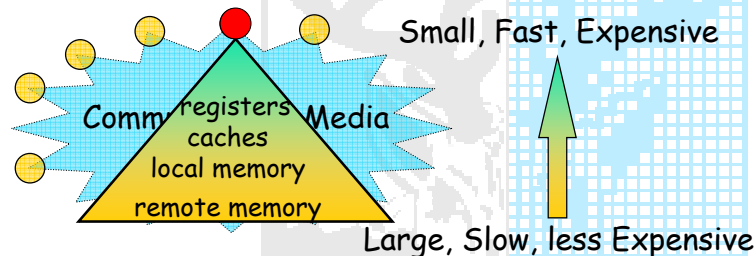


Matrix Computations

- Fundamental and ubiquitous in computational science and its vast application areas
- Library software - optimized building blocks for fundamental operations
 - BLAS, (Sca)LAPACK, SLICOT (see also NETLIB)
 - ESSL and other vendors
 - Portability and efficiency
- Continuing demand for new and improved algorithms and software along with the computer evolution

"Data transport" in memory hierarchies

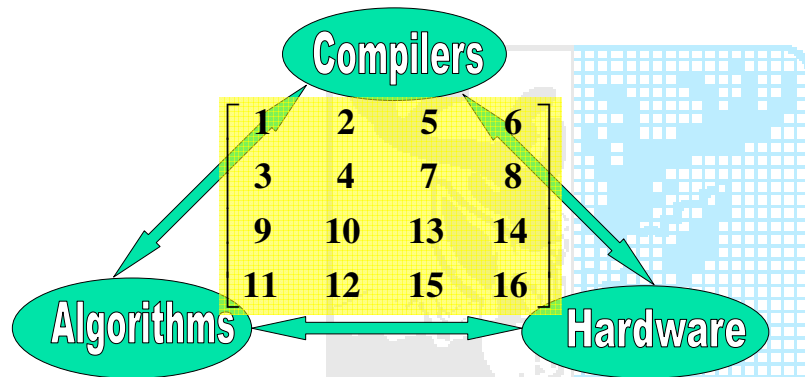
- of today's computer systems
 - PC - cluster - supercomputer



Management of deep memory hierarchies

- Architecture evolution: HPC systems with multiple SMP nodes, several levels of caches, more functional units per CPU
- Key to performance: understand the algorithm and architecture interaction
- Hierarchical blocking

The fundamental AHC triangle



Outline

- Hierarchical blocking: motivation and implications
- Recursive blocked templates and algorithms
- Recursive blocked data structures
- Case studies:
 - General matrix multiply and add (GEMM)
 - Packed Cholesky factorization
 - QR factorization and linear systems
 - Triangular matrix equations and condition estimation
- Some related and complementary work
- Work in progress: periodic matrix equations
- Concluding remarks

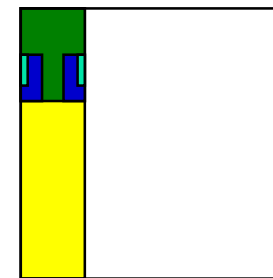
Recursive Blocked Algorithms and Hybrid Data Structures for Dense Matrix Library Software*

Joint work with:

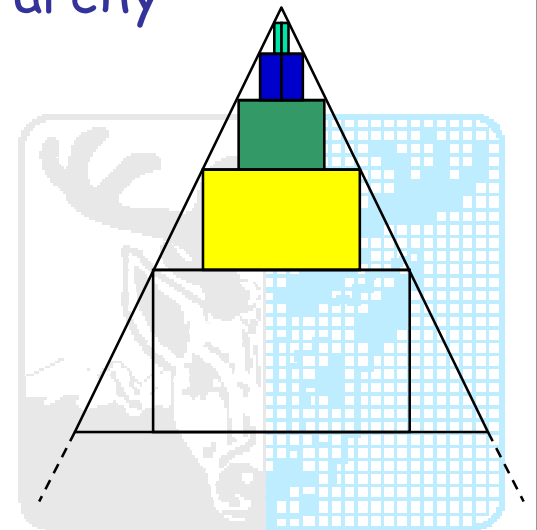
Erik Elmroth[†]
Fred Gustavson[‡]
Isak Jonsson[†]
Bo Kågström[†]

Abstract. Matrix computations are both fundamental and ubiquitous in computational science and its vast application areas. Along with the development of more advanced computer systems with complex memory hierarchies, there is a continuing demand for new algorithms and library software that efficiently utilize and adapt to new architecture features. This article reviews and details some of the recent advances made by applying the paradigm of recursion to dense matrix computations on today's memory-tiered computer systems. Recursion

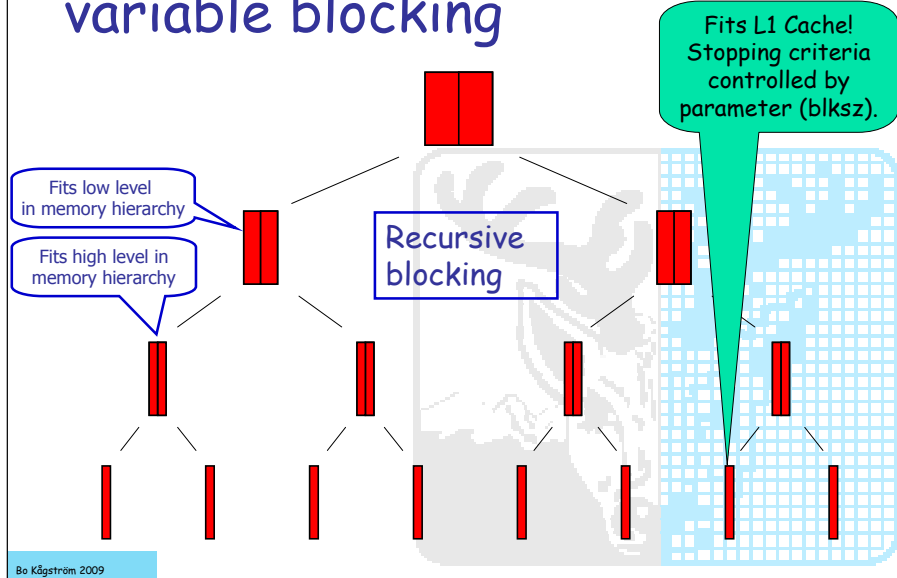
Blocking for a memory hierarchy



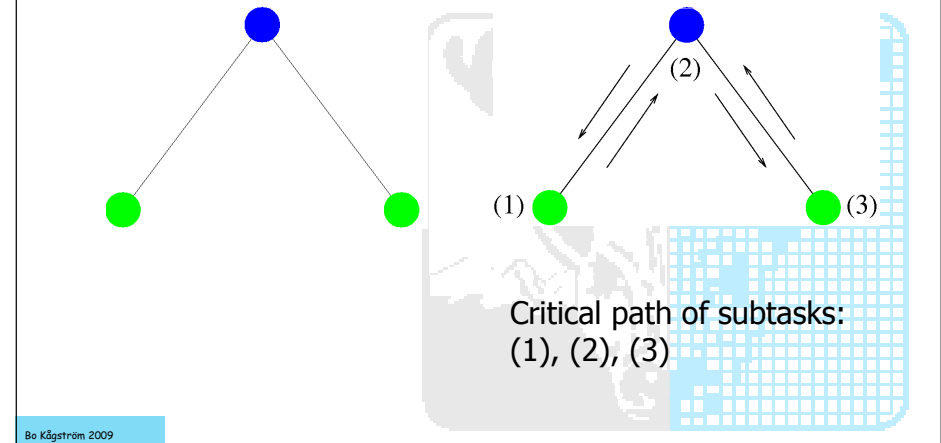
Explicit multi-level blocking



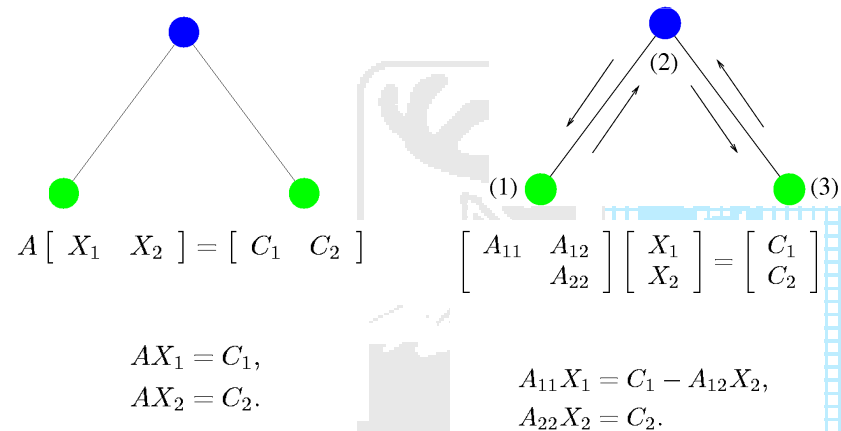
Recursion leads to automatic variable blocking



Splittings defining independent and dependent tasks

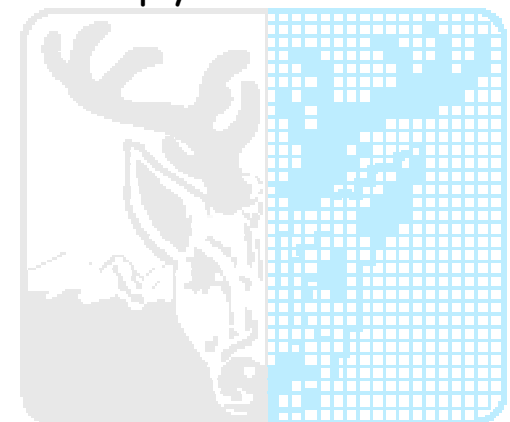


TRSM Operation: $AX = C$, A mxm upper triangular, C/X mxn



Case Study 1

General matrix multiply and add
(GEMM)



Recursive splittings for GEMM:

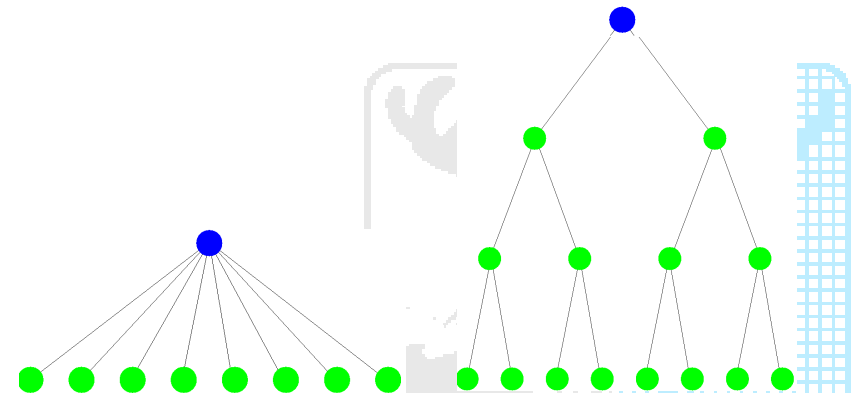
$$C \leftarrow \beta \text{op}(C) + \alpha \text{op}(A)\text{op}(B)$$

Split $m \times n$ $m \times k$ $k \times n$

$$\begin{aligned}
 m, n, k & \quad \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} + \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \\
 m & \quad = \begin{bmatrix} [C_{11} \ C_{12}] + [A_{11} \ A_{12}] \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \\ [C_{21} \ C_{22}] + [A_{21} \ A_{22}] \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \end{bmatrix} = \\
 n & \quad = \begin{bmatrix} [C_{11}] + [A_{11} \ A_{12}] \begin{bmatrix} B_{11} \\ B_{21} \end{bmatrix}, [C_{12}] + [A_{11} \ A_{12}] \begin{bmatrix} B_{12} \\ B_{22} \end{bmatrix} \\ [C_{21}] + [A_{21} \ A_{22}] \begin{bmatrix} B_{11} \\ B_{21} \end{bmatrix}, [C_{22}] + [A_{21} \ A_{22}] \begin{bmatrix} B_{12} \\ B_{22} \end{bmatrix} \end{bmatrix} = \\
 k & \quad = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} + \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} [B_{11} \ B_{12}] + \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} [B_{21} \ B_{22}]
 \end{aligned}$$

Bo Kågström 2009

Recursive splitting - by breadth or by depth



Bo Kågström 2009

GEMM recursive blocked template - splitting by depth

```

◦ C = rgemm( A, B, C, blksz)
  If m, n, and k <= blksz
    C = opt_gemm( A, B, C) % optimized GEMM kernel!
  elseif m = max(m, n, k) % split m: m2 = m/2
    C(1:m2, :) = rgemm( A(1:m2, :), B, C(1:m2, :), blksz)
    C(m2+1:m, :) = rgemm( A(m2+1:m, :), B, C(m2+1:m, :), blksz)
  elseif n = max(n, k) % split n: n2 = n/2, k
    C(:, 1:n2) = rgemm( A, B(:, 1:n2), C(:, 1:n2), blksz)
    C(:, n2+1:n) = rgemm( A, B(:, n2+1:n), C(:, n2+1:n), blksz)
  else % split k: k2 = k/2,
    C = rgemm(A(:, 1:n2), B(1:m2, :), C, blksz)
    C = rgemm(A(:, n2+1:n), B(m2+1:m, :), C, blksz)
  
```

When to end the recursive splitting?

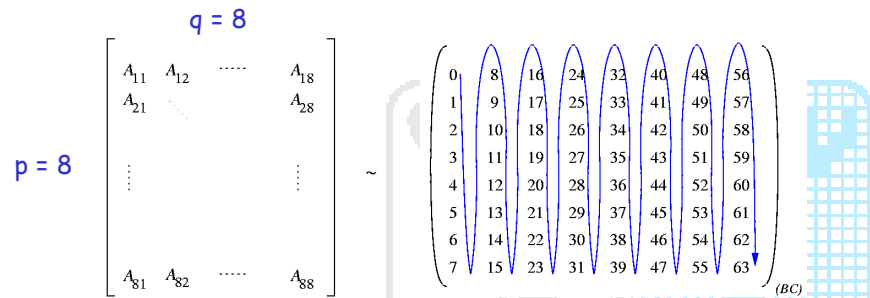
Bo Kågström 2009

Locality of reference

- Recursive blocked algorithms mainly improve on the **temporal locality**
- Further performance improvements by matching the data structure with the algorithm (and vice versa)
- Recursive blocked data structures improve on the **spatial locality**

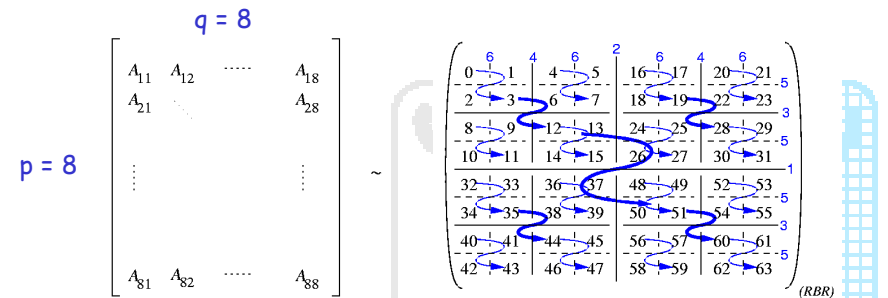
Bo Kågström 2009

Blocked data formats



Blocks A_{ij} of size $mb \times nb$ can be ordered in $(pq)!$ different ways

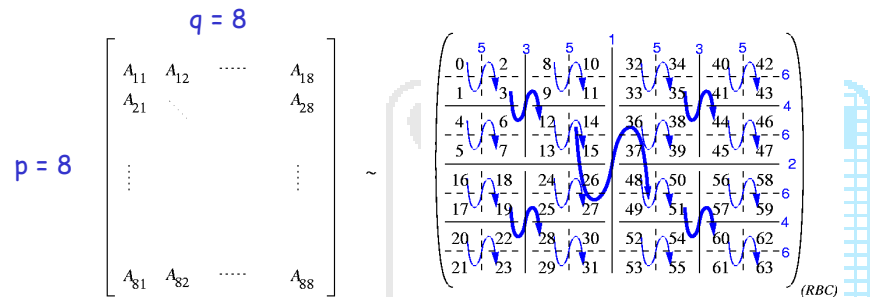
Recursive blocked row format



Recursive ordering: a 1-dim tour through a 2-dim object (Hilbert space filling heuristics)

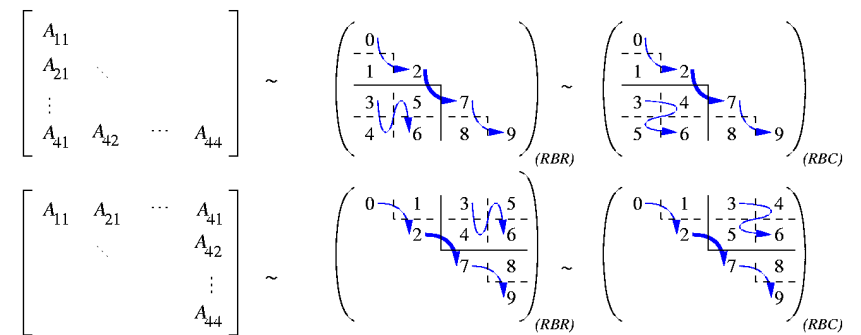
RBR \leftrightarrow Z-Morton ordering

Recursive blocked column format



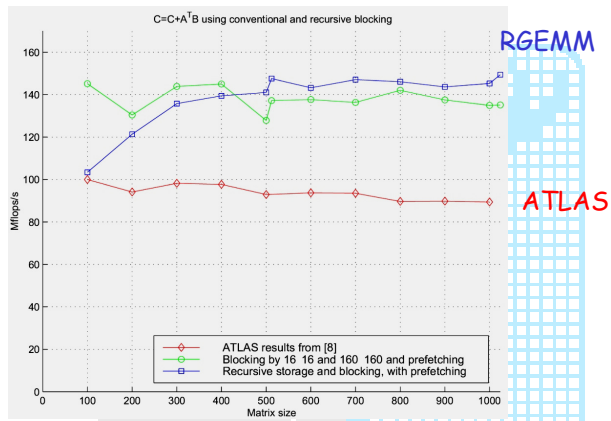
RBC \leftrightarrow reflected-N-Morton space filling ordering

Triangular recursive data format



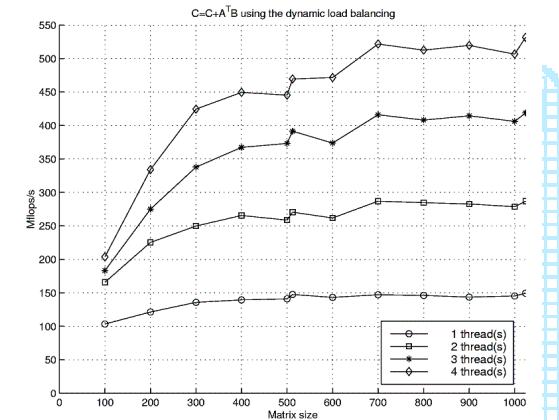
Recursive GEMM: multi-level vs. recursive blocking

IBM PPC604,
112 MHz

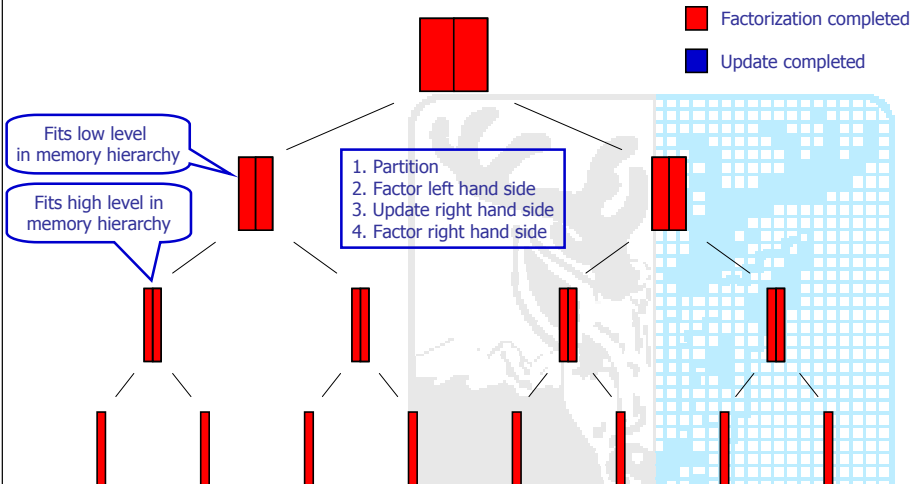


Recursive blocked GEMM and SMP parallelism via threads

IBM PPC604, 4 proc

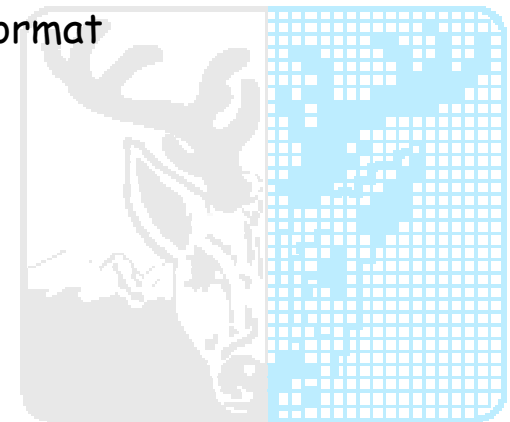


Recursion template for one-sided matrix factorization



Case Study 2

Cholesky factorization for matrices in packed format



Packed Cholesky factorization

$$A \equiv \begin{bmatrix} A_{11} & A_{21}^T \\ A_{21} & A_{22} \end{bmatrix} = LL^T \equiv \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} L_{11}^T & L_{21}^T \\ 0 & L_{22}^T \end{bmatrix}$$

Standard approach (typified by LAPACK):

- Packed storage -> cannot use standard level 3 BLAS (e.g., DGEMM)
- Possible to produce packed level 3 BLAS routines at a great programming cost
- Run at level 2 performance, i.e., much below full storage routines.
- Minimum storage: $1/2n(n+1)$ elements

Packed recursive blocked data

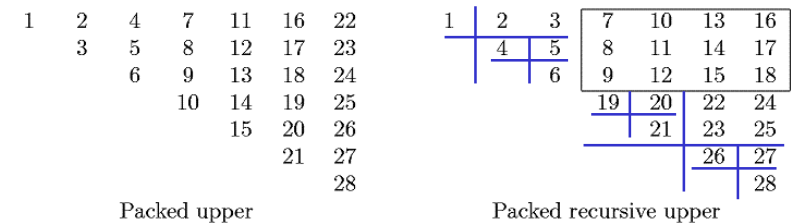


FIG. 3.1. Memory indices for 7×7 upper triangular matrix stored in traditional packed format and recursive packed format.

- Divide into two isosceles triangles T1, T2 and rectangle R
- Divide triangles recursively down to element level
- Store in order: T1, R, T2
- Rectangles stored in full format → Possible to use full storage level 3 BLAS

Cholesky recursive blocked template

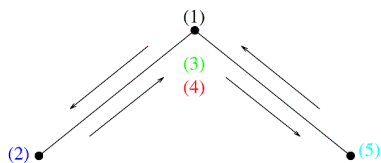
$$A = \begin{pmatrix} A_{11} & A_{21}^T \\ A_{21} & A_{22} \end{pmatrix} = LL^T = \begin{pmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{pmatrix} \begin{pmatrix} L_{11}^T & L_{21}^T \\ 0 & L_{22}^T \end{pmatrix} \quad (1)$$

$$\text{Factor : } A_{11} = L_{11}L_{11}^T. \quad (2)$$

$$\text{TRSM : } L_{21}L_{11}^T = A_{21}. \quad (3)$$

$$\text{SYRK : } \tilde{A}_{22} = A_{22} - L_{21}L_{21}^T \quad (4)$$

$$\text{Factor : } \tilde{A}_{22} = L_{22}L_{22}^T \quad (5)$$

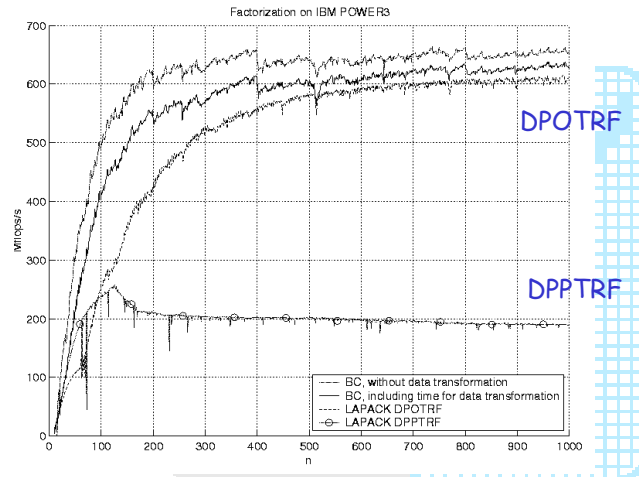


Packed recursive blocked Cholesky highlights

- Recursive blocked algorithm + recursive packed data layout => can make use of high performance level 3 BLAS routines (e.g., DGEMM)
- Use minimal storage for matrix A
- Temporary workspace = $1/8n^2$ elements (~25%)
- Leaf problems ($< \text{blksz}$) are solved using superscalar kernels (Cholesky, TRSM, SYRK)

Recursive blocked Cholesky vs. LAPACK - (rec.) packed format

Runs at level 3 performance - at least!



Case Study 3

QR factorization and linear systems



Recursive blocked QR factorization

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} = Q \begin{pmatrix} R_{11} & R_{12} \\ 0 & R_{22} \end{pmatrix}$$

1. Divide A $m \times n$ in two parts (left & right)

Stopping criteria:
if $n < 4$ use
standard algorithm

2. Factorize left hand side by a

#flops grows cubically with
Householder transformations
being aggregated (compact WY)!

$$Q_1 \begin{pmatrix} R_{11} \\ 0 \end{pmatrix} = \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$$

3. Update right hand side

$$\begin{pmatrix} R_{12} \\ \tilde{A}_{22} \end{pmatrix} \leftarrow Q_1^T \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix}$$

4. Factorize by a recursive call

$$Q_2 R_{22} = \tilde{A}_{22}$$

Recursive blocked QR highlights

- Recursive splitting controlled by nb (splitting point = $\min(nb, n/2)$, $nb = 32-64$)
- Level 3 algorithm for generating $Q = I - YTY^T$ (compact WY) within the recursive blocked algorithm (T triangular of size $\leq nb$)
- Replaces LAPACK level 2 and 3 algorithms

Recursive QR vs. LAPACK

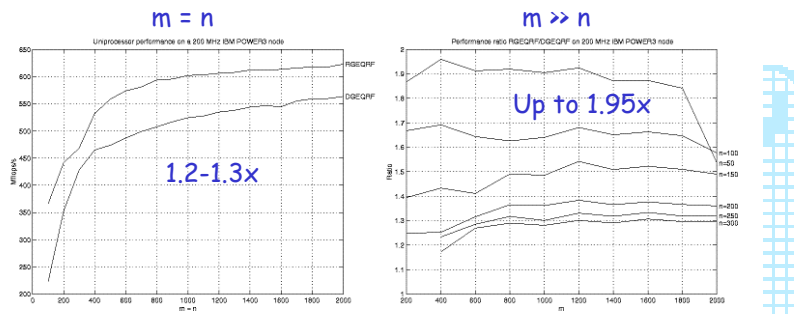
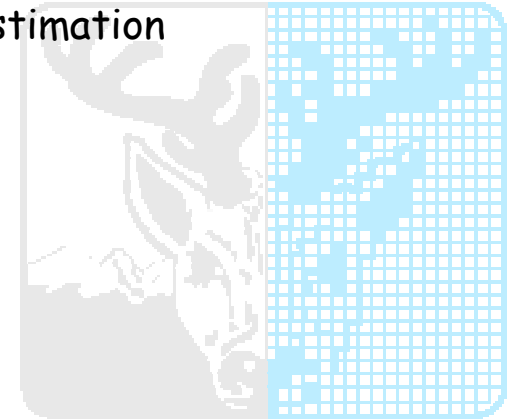


Fig. 4.1 Performance results in Mflops/s for square matrices (left) and performance ratio for tall, thin matrices (right) for the recursive algorithm RGEQRF and DGEQRF of LAPACK on the 200 MHz IBM Power3.

Case Study 4

Triangular matrix equations and condition estimation



Matrix equations

Name	Matrix equation	Acronym
Standard Sylvester (CT)	$AX - XB = C$	SYCT
Standard Lyapunov (CT)	$AX + XA^T = C$	LYCT
Generalized coupled Sylvester	$(AX - YB, DX - YE) = (C, F)$	GCSY
Standard Sylvester (DT)	$AXB^T - X = C$	SYDT
Standard Lyapunov (DT)	$AXA^T - X = C$	LYDT
Generalized Sylvester	$AXB^T - CXD^T = E$	GSYL
Generalized Lyapunov (CT)	$AXE^T + EXA^T = C$	GLYCT
Generalized Lyapunov (DT)	$AXA^T - EXE^T = C$	GLYDT

One-sided (top) and two-sided (bottom)

Separation of two matrices

$$\text{Sep}[A, B] = \inf_{\|X\|_F=1} \|AX - XB\|_F = \sigma_{\min}(Z),$$

$$\text{where } Z = I_n \otimes A - B^T \otimes I_m.$$

Computing $\text{Sep}[A, B]$ costs $O(m^3n^3)$ - impractical!

Reliable Sep -estimates of cost $O(m^2n + mn^2)$:

$$\frac{\|x\|_2}{\|y\|_2} = \frac{\|X\|_F}{\|C\|_F} \leq \|Z^{-1}\|_2 = \frac{1}{\sigma_{\min}(Z)} = \text{Sep}^{-1},$$

$$(mn)^{-1/2} \|Z^{-1}\|_1 \leq \|Z^{-1}\|_2 \leq \sqrt{mn} \|Z^{-1}\|_1.$$

Matrix equation Sep-functions

Z-matrix	Sep-function = $\sigma_{\min}(Z)$
$Z_{\text{SYCT}} = I_n \otimes A - B^T \otimes I_m$	$\inf_{\ X\ _F=1} \ AX - XB\ _F$
$Z_{\text{LYCT}} = I_n \otimes A + A \otimes I_n$	$\inf_{\ X\ _F=1} \ AX - X(-A^T)\ _F$
$Z_{\text{GCSY}} = \begin{bmatrix} I_n \otimes A & -B^T \otimes I_m \\ I_n \otimes D & -E^T \otimes I_m \end{bmatrix}$	$\inf_{\ (X,Y)\ _F=1} \ (AX - YB, DX - YE)\ _F$
$Z_{\text{SYDT}} = B \otimes A - I_n \otimes I_m$	$\inf_{\ X\ _F=1} \ AXB^T - X\ _F$
$Z_{\text{LYDT}} = A \otimes A - I_n \otimes I_n$	$\inf_{\ X\ _F=1} \ AXA^T - X\ _F$
$Z_{\text{GSYL}} = B \otimes A - D \otimes C$	$\inf_{\ X\ _F=1} \ AXB^T - CXD^T\ _F$
$Z_{\text{GLYCT}} = E \otimes A + A \otimes E$	$\inf_{\ X\ _F=1} \ AXE^T - EX(-A^T)\ _F$
$Z_{\text{GLYDT}} = A \otimes A - E \otimes E$	$\inf_{\ X\ _F=1} \ AXA^T - EXE^T\ _F$

$Zx = b$, Z is a Kronecker product representation

Sep-function = smallest singular value of Z

Recursive blocked SYCT template

Case 1: $1 \leq n \leq m/2$

$$\left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline & A_{22} \end{array} \right] \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] - \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] \left[\begin{array}{c|c} B_{11} & B_{12} \\ \hline & B_{22} \end{array} \right] = \left[\begin{array}{c|c} C_{11} & C_{12} \\ \hline C_{21} & C_{22} \end{array} \right]$$

Case 2: $1 \leq m \leq n/2$

$$\left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline & A_{22} \end{array} \right] \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] - \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] \left[\begin{array}{c|c} B_{11} & B_{12} \\ \hline & B_{22} \end{array} \right] = \left[\begin{array}{c|c} C_{11} & C_{12} \\ \hline C_{21} & C_{22} \end{array} \right]$$

Case 3: $n/2 < m < 2n$

$$\left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline & A_{22} \end{array} \right] \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] - \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] \left[\begin{array}{c|c} B_{11} & B_{12} \\ \hline & B_{22} \end{array} \right] = \left[\begin{array}{c|c} C_{11} & C_{12} \\ \hline C_{21} & C_{22} \end{array} \right]$$

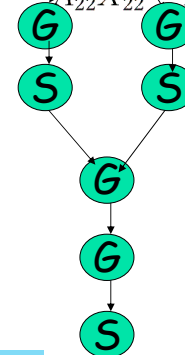
Recursive SYCT - Case 3

$$\left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline & A_{22} \end{array} \right] \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] - \left[\begin{array}{c|c} X_{11} & X_{12} \\ \hline X_{21} & X_{22} \end{array} \right] \left[\begin{array}{c|c} B_{11} & B_{12} \\ \hline & B_{22} \end{array} \right] = \left[\begin{array}{c|c} C_{11} & C_{12} \\ \hline C_{21} & C_{22} \end{array} \right]$$

$$\begin{aligned} A_{11}X_{11} - X_{11}B_{11} &= C_{11} - A_{12}X_{21} \\ A_{11}X_{12} - X_{12}B_{22} &= C_{12} - A_{12}X_{22} + X_{11}B_{12} \\ A_{22}X_{21} - X_{21}B_{11} &= C_{21} \\ A_{22}X_{22} - X_{22}B_{22} &= C_{22} + X_{21}B_{12} \end{aligned}$$

Recursive SYCT - Case 3

$$\begin{aligned} A_{11}X_{11} - X_{11}B_{11} &= C_{11} - A_{12}X_{21} \\ A_{11}X_{12} - X_{12}B_{22} &= C_{12} - A_{12}X_{22} + X_{11}B_{12} \\ A_{22}X_{21} - X_{21}B_{11} &= C_{21} \\ A_{22}X_{22} - X_{22}B_{22} &= C_{22} + X_{21}B_{12} \end{aligned}$$



1. SYLV('N', 'N', A_{22} , B_{11} , C_{21})
- 2a. GEMM('N', 'N', $\alpha = +1$, C_{21} , B_{12} , C_{22})
- 2b. GEMM('N', 'N', $\alpha = -1$, A_{12} , C_{21} , C_{11})
- 3a. SYLV('N', 'N', A_{22} , B_{22} , C_{22})
- 3b. SYLV('N', 'N', A_{11} , B_{11} , C_{11})
4. GEMM('N', 'N', $\alpha = -1$, A_{12} , C_{22} , C_{12})
5. GEMM('N', 'N', $\alpha = +1$, C_{11} , B_{12} , C_{12})
6. SYLV('N', 'N', A_{11} , B_{22} , C_{12})

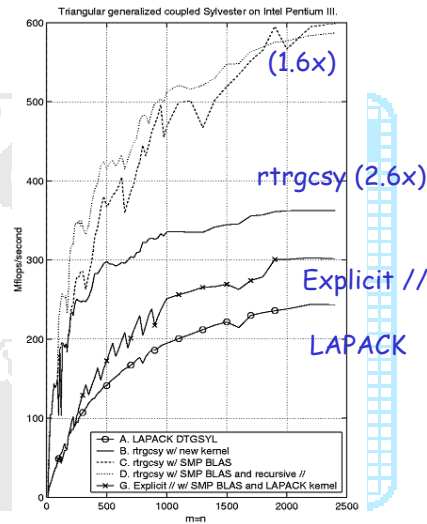
Triangular generalized coupled Sylvester equation - GCSY

$$AX - YB = C$$

$$DX - YE = F$$

(A, D) and (B, E) in generalized Schur form

Solution (X, Y) overwrites r.h.s. (C, F)



Two-sided matrix equation: GLYDT

- $AXA^T - EXE^T = C$
- $C = C^T$ $n \times n$; (A, E) $n \times n$ in gen. Schur form
- Unique sol'n $X = X^T \Leftrightarrow$
 λ_i of $A - \lambda E$ satisfy $\lambda_i \lambda_j \neq 1$
- Recursive splitting:

$$\begin{bmatrix} A_{11} & A_{12} \\ & A_{22} \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ & X_{22} \end{bmatrix} \begin{bmatrix} A_{11}^T & \\ & A_{22}^T \end{bmatrix} - \begin{bmatrix} E_{11} & E_{12} \\ & E_{22} \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ & X_{22} \end{bmatrix} \begin{bmatrix} E_{11}^T & \\ & E_{22}^T \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ & C_{22} \end{bmatrix}$$

Two-sided matrix product

$$C = \beta C + \alpha \text{op}(A) X \text{op}(B)^T$$

- A and/or B can be dense or triangular
- One or several of A, B and C can be symmetric
- Extra workspace - size of r.h.s.

Make use of symmetry, e.g., in GLYDT:

$$C_{11} = C_{11} - A_{12} X_{22} A_{12}^T \quad \text{and} \quad C_{11} = C_{11} + E_{12} X_{22} E_{12}^T$$

GLYDT performance with optional condition estimation

Table 5.3 Timings for solving unreduced two-sided matrix equations (GLYDT) with optional condition estimation. (Job = X, compute solution only; Job = X + Sep, compute solution and Sep-estimation.) Results from 375 MHz IBM Power3.

n	SG03AD using SG03AX		SG03AD using rtrglydt		Speedup	Job
	Total time	Solver part	Total time	Solver part		
50	0.0277	49.9 %	0.0185	20.1 %	1.50	X
100	0.180	51.2 %	0.0967	9.0 %	1.86	X
250	2.89	46.8 %	1.62	4.7 %	1.79	X
500	59.0	42.3 %	34.5	1.5 %	1.71	X
750	303.4	42.0 %	177.5	0.9 %	1.71	X
1000	646.6	44.6 %	361.8	1.0 %	1.79	X
50	0.117	87.6 %	0.0263	45.6 %	4.44	X+Sep
100	0.709	87.3 %	0.152	40.6 %	4.68	X+Sep
250	9.98	84.5 %	2.08	25.4 %	4.81	X+Sep
500	178.6	80.9 %	37.8	9.4 %	4.73	X+Sep
750	924.1	80.9 %	184.4	4.5 %	5.01	X+Sep
1000	2076.6	82.7 %	391.8	8.4 %	5.30	X+Sep

Up to 1.9x

Up to 5.3x

RECSY library

- Recursive blocked algorithms for solving reduced matrix equations
- Recursion implemented in F90
- SMP versions using OpenMP
- F77 wrappers for LAPACK and SLICOT routines
- www.cs.umu.se/research/parallel/recsy/
- Part of Isak Jonsson's PhD Thesis, Dec 2003

SCASY library

- ScalLAPACK-style software package of matrix equation solvers for distributed memory machines.
- Triangular solvers are used in implementing parallel condition estimators for each matrix equation.
- With Robert Granat, PhD Nov 2007

$op(A)X \pm Xop(B) = C$	SYCT	✓
$op(A)X + Xop(A^T) = C$	LYCT	✓
$op(A)Xop(B) \pm X = C$	SYDT	✓
$op(A)Xop(A^T) - X = C$	LYDT	✓
$op(A)X \pm Yop(B) = C,$ $op(D)X \pm Yop(E) = F$	GCSY	
$op(A)Xop(B) \pm op(D)Xop(E) = C$	GSYL	
$op(A)Xop(A^T) - op(E)Xop(E^T) = C$	GLYCT	
$op(A)X(E^T) + op(E)Xop(A^T) = C$	GLYDT	

Discrete-time periodic systems

$$\begin{aligned} x_{k+1} &= A_k x_k + B_k u_k \\ y_k &= C_k x_k + D_k u_k \end{aligned}$$

$$\left. \begin{aligned} A_k &\in \mathbb{R}^{n_{k+1} \times n_k}, & B_k &\in \mathbb{R}^{n_{k+1} \times m} \\ C_k &\in \mathbb{R}^{p \times n_k}, & D_k &\in \mathbb{R}^{p \times m} \end{aligned} \right\} - K\text{-periodic}$$

$$(A_{k+K} = A_k, B_{k+K} = B_k, \dots)$$

From discretization of **continuous-time periodic models**: revolving satellite; helicopter in forward flight; multi-rate sampled control systems \rightarrow time-varying dimensions

Periodic Sylvester equation

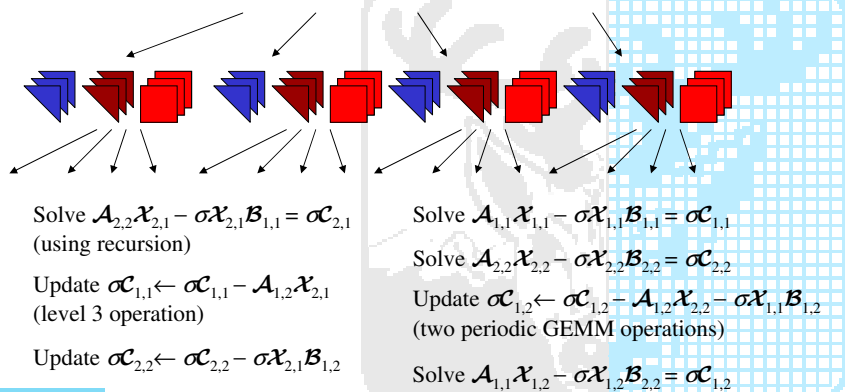
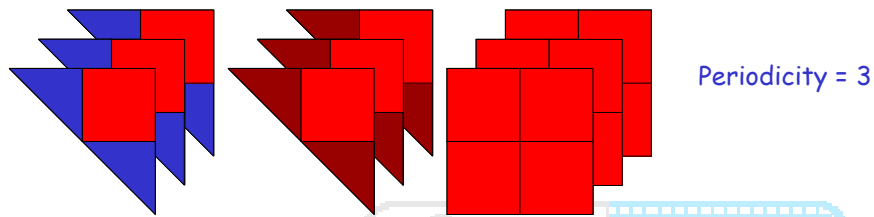
- $A(:, :, k) X(:, :, k) + X(:, :, k+1) B(:, :, k) = C(:, :, k+1)$, for $k = 1 : K-1$
- $A(:, :, K) X(:, :, K) + X(:, :, 1) B(:, :, K) = C(:, :, 1)$



Script notations: X_k K -periodic

$$X_k := \text{diag}(X_k, X_{k+1}, \dots, X_{k+K-1})$$

$$\sigma X_k := \text{diag}(X_{k+1}, \dots, X_{k+K-1}, X_k)$$



Solve $\mathcal{A}_{2,2}\mathcal{X}_{2,1} - \sigma\mathcal{X}_{2,1}\mathcal{B}_{1,1} = \mathcal{C}_{2,1}$
(using recursion)

Update $\mathcal{C}_{1,1} \leftarrow \mathcal{C}_{1,1} - \mathcal{A}_{1,2}\mathcal{X}_{2,1}$
(level 3 operation)

Update $\mathcal{C}_{2,2} \leftarrow \mathcal{C}_{2,2} - \sigma\mathcal{X}_{2,1}\mathcal{B}_{1,2}$

Solve $\mathcal{A}_{1,1}\mathcal{X}_{1,1} - \sigma\mathcal{X}_{1,1}\mathcal{B}_{1,1} = \mathcal{C}_{1,1}$

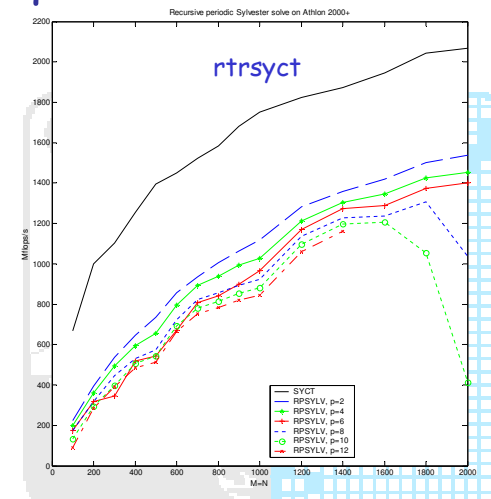
Solve $\mathcal{A}_{2,2}\mathcal{X}_{2,2} - \sigma\mathcal{X}_{2,2}\mathcal{B}_{2,2} = \mathcal{C}_{2,2}$

Update $\mathcal{C}_{1,2} \leftarrow \mathcal{C}_{1,2} - \mathcal{A}_{1,2}\mathcal{X}_{2,2} - \sigma\mathcal{X}_{1,1}\mathcal{B}_{1,2}$
(two periodic GEMM operations)

Solve $\mathcal{A}_{1,1}\mathcal{X}_{1,2} - \sigma\mathcal{X}_{1,2}\mathcal{B}_{2,2} = \mathcal{C}_{1,2}$

Recursive blocked periodic Sylvester equation

Some preliminary performance results



Recursive blocking ...

- creates new algorithms for linear algebra software
- expresses dense linear algebra algorithms entirely in terms of level~3 BLAS like matrix-matrix operations
- introduces an automatic variable blocking that targets every level of a deep memory hierarchy
- can also be used to define hybrid data formats for storing block-partitioned matrices (**general, triangular, symmetric, packed**) - L1, L2 and TLB misses are minimized for certain block sizes (Park-Hong-Prosana'03)

High-performance software

- implementations are based on data locality and superscalar optimization techniques
- recursive blocked algorithms improve on the temporal data locality
- hybrid data formats improve on the spatial data locality
- portable and generic superscalar kernels ensure that all functional units on the processor(s) are used efficiently

Acknowledgements

- Erik Elmroth, Isak Jonsson, Fred Gustavson (co-authors and co-workers)
- André Henriksson, Olov Gustavsson and Andreas Lindkvist (earlier MSc students)
- Bjarne Andersén, Jerzy Wasniewski (e.g., packed Cholesky)
- Robert Granat (PhD student)
- HPC and LA team at Umeå University
- **Community that do related and complementary work!** (see SIAM Rev. 2004)

- Thanks for your attention!



Some related and complementary work

- Recursive algorithms and hybrid data structures
 - **Winograd-Strassen'69**: Douglas etal'94, ESSL, Demmel-Higham'92 (stability)
 - **Quad- and octtrees**: Samet'84, Salman-Warner'94 (N-body, Barnes-Hut'84)
 - **Cache oblivious algorithms**: Leiserson etal'99 (sorting, FFT, AT)
 - **GEMM**: Chatterjee etal'02, Valsalam and Skjellum'02, ATLAS-project
 - **LU**: Toledo'97(dense), Dongarra, Eijkhout Luszczek'01 (sparse)
 - **QR**: Rabani and Toledo'01 (out-of-core), Frens and Wise'03 (Givens-based)

Some related and complementary work

- Automated generation of library software and compiler technology
 - **Empirical optimization**:
 - **PHiPAC** - Bilmes, Demmel etal'97,
 - **ATLAS** - Whaley, Petitet and Dongarra'00,
 - **Sparse kernels** - Vuduc, Demmel et al'03
 - **FLAME**: Gunnels, Goto, Van de Geijn etal'01, '02
 - **Compiler blockability**: Wolf and Lam'91 (loop transformations), Carr and Lehoucq'97
 - **Automatic generation of recursive codes**: Ahmed and Pingali'00 (iterative algorithms -> recursive), Yi, Adve and Kennedy'00 (convert loop nests into recursive form)

◦ Thanks for your attention!

