

Parallelism in Linnea

Henrik Barthels, Paolo Bientinesi

Introduction: Linnea

- How to compute the following expressions?

$$\mathbf{b} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{b} := (\mathbf{X}^T \mathbf{M}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{M}^{-1} \mathbf{y}$$

$$\mathbf{x} := \mathbf{W} (\mathbf{A}^T (\mathbf{A} \mathbf{W} \mathbf{A}^T)^{-1} \mathbf{b} - \mathbf{c})$$

$$\mathbf{x} := (\mathbf{A}^{-T} \mathbf{B}^T \mathbf{B} \mathbf{A}^{-1} + \mathbf{R}^T [\Lambda(\mathbf{R} \mathbf{z})] \mathbf{R})^{-1} \mathbf{A}^{-T} \mathbf{B}^T \mathbf{B} \mathbf{A}^{-1} \mathbf{y}$$

- High-level languages are easy to use, but performance is usually suboptimal.
- BLAS and LAPACK can be fast, but require a lot of expertise.

BLAS [DDC⁺90], **LAPACK** [AB⁺99]

- $\mathbf{y} \leftarrow \mathbf{A} \mathbf{x} + \mathbf{y}$
- $\mathbf{C} \leftarrow \mathbf{A} \mathbf{B} + \mathbf{C}$
- $\mathbf{B} \leftarrow \mathbf{A}^{-1} \mathbf{B}$
- ...

Introduction: Linnea

Input

$$z := (X^T S^{-1} X)^{-1} X^T S^{-1} y \quad S \text{ is symmetric positive definite}$$

Output

$L := \text{chol}(S)$	(potrf)
$U_1 := L^{-1} X$	(trsm)
$(Q, R) := \text{qr}(U_1)$	(geqrf)
$u_2 := L^{-1} y$	(trsv)
$u_3 := Q^T u_2$	(ormqr)
$z := R^{-1} u_3$	(trsv)

<https://github.com/HPAC/linnea>

Linear Algebra Knowledge

- Properties: trmm vs. gemm

- Inference of properties:  $\begin{matrix} \square & \square \\ \diagdown & \diagdown \end{matrix} \rightarrow \begin{matrix} \square \\ \diagdown \end{matrix}$

- Simplifications: $A^T \rightarrow A$ if $\text{Symmetric}(A)$

- Rewriting expressions:

$$\begin{aligned} X &:= A^T A + A^T B + B^T A && \rightarrow && Y := B + A/2 \\ &&& && X := A^T Y + Y^T A \end{aligned}$$

- Common subexpressions:

$$\begin{aligned} X &:= AB^{-T}C + B^{-1}A^T && \rightarrow && Z := AB^{-T} \\ &&& && X := ZC + Z^T \end{aligned}$$

- Matrix chains:

$$\begin{aligned} (AB)c & \quad \mathcal{O}(n^3) \\ A(Bc) & \quad \mathcal{O}(n^2) \end{aligned}$$

Problem 1: Generation of Parallel Code

Problem 2: Parallel Code Generation

Problem 1: Generation of Parallel Code

Problem 2: Parallel Code Generation

How to make use of parallelism?

- Threaded kernels.
- Kernels in parallel.
- Both.

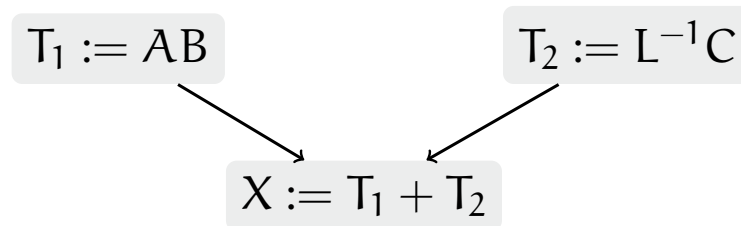
Generation of Parallel Code

Step 1: Parallelize a given sequence of kernels.

The Good

- Constructing dependency graph is easy.

$T_1 := AB$
 $T_2 := L^{-1}C$
 $X := T_1 + T_2$



- If operand sizes are known, amount of work is known (#FLOPs).

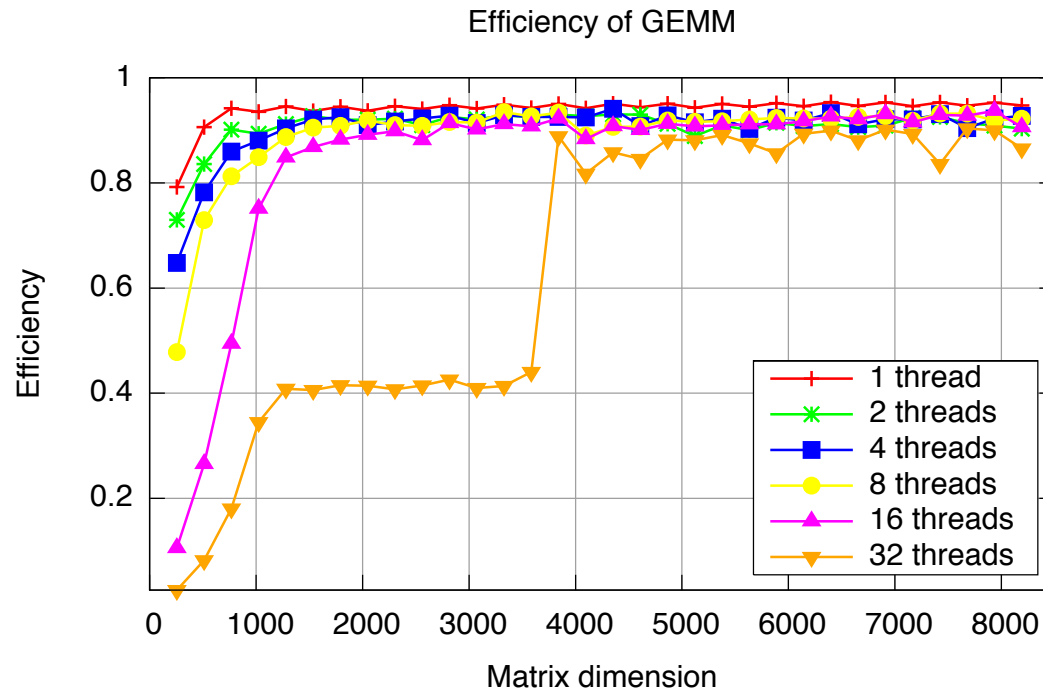
$A \in \mathbb{R}^{m \times k}, B \in \mathbb{R}^{k \times n} \rightarrow AB$ requires $2mnk$ FLOPs

- BLAS and LAPACK are parallelized.

Generation of Parallel Code

The Bad

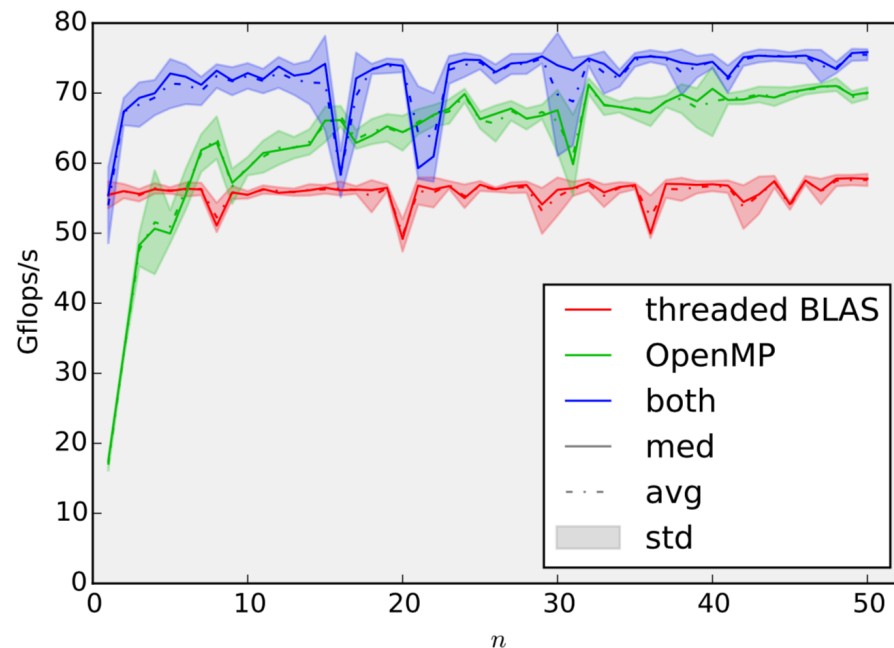
- FLOP count is not a good prediction for execution time.
- Performance modeling is a hard problem [PB12].
 - Performance is not composable.
 - Efficiency decreases with number of threads.



Generation of Parallel Code

The Bad

- FLOP count is not a good prediction for execution time.
- Performance modeling is a hard problem [PB12].
 - Performance is not composable.
 - Efficiency decreases with number of threads.
 - “Overbooking”: Parallelizing a sequence of n LU’s [PB15].



Generation of Parallel Code

Existing Tools

PaRSEC, OmpSs, StarPU, SuperGlue,...

- Built for large dependency graphs/large number of tasks.
- Only one thread per task.
- Do not consider cost (except StarPU [ATN09]).

We have:

- Small number of tasks.
- Multiple threads per task.
- Cost is (roughly) known.

Generation of Parallel Code

Step 2: Generate sequence that parallelizes well.

A good sequence of kernels for sequential execution may not be good for parallel execution.

Example: $X := ABCD$

$A, B, C \in \mathbb{R}^{n \times n}$

$D \in \mathbb{R}^{n \times m}$

$m < n$

$A(B(CD))$

- min #FLOPs
- dependencies
- threaded kernels still possible

$(AB)(CD)$

- more FLOPs
- fewer dependencies

Existing work: Matrix Chain Products on Parallel Systems [LKHL03]

Problem 1: Generation of Parallel Code

Problem 2: Parallel Code Generation

Parallel Code Generation

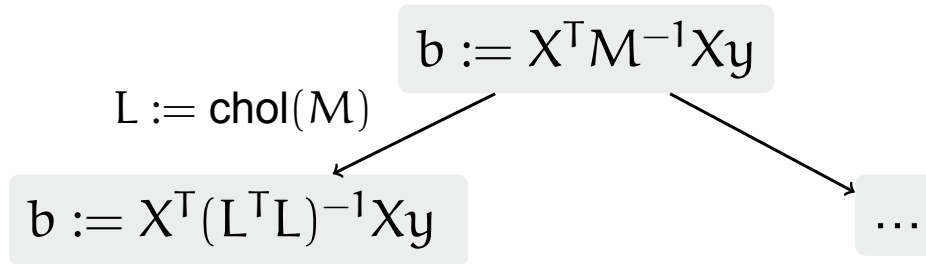
Motivation

- For sufficiently large matrices and/or enough runs of the program, generation time will be amortized.
- What about small computations?
- What about computations in interactive environment such as Matlab?

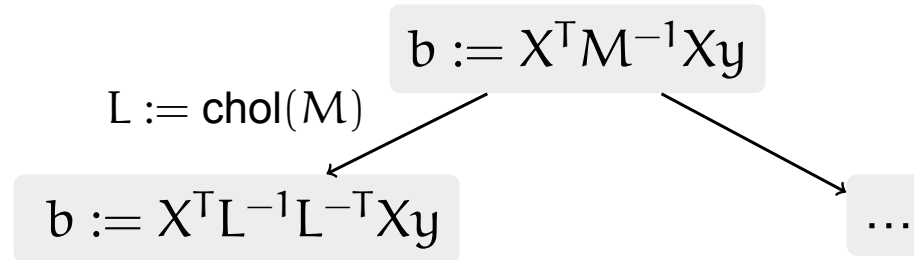
Parallel Code Generation

$$b := X^T M^{-1} X y$$

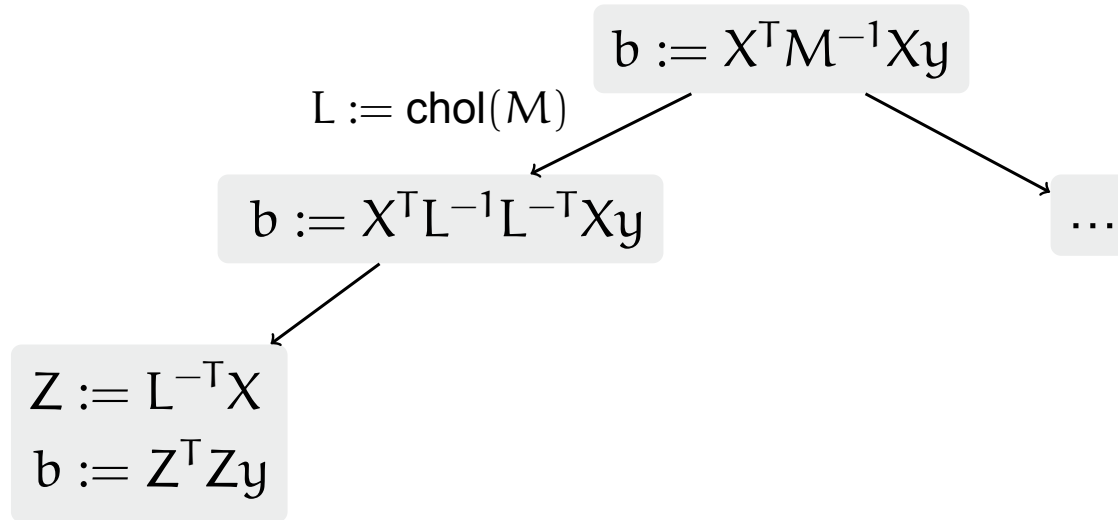
Parallel Code Generation



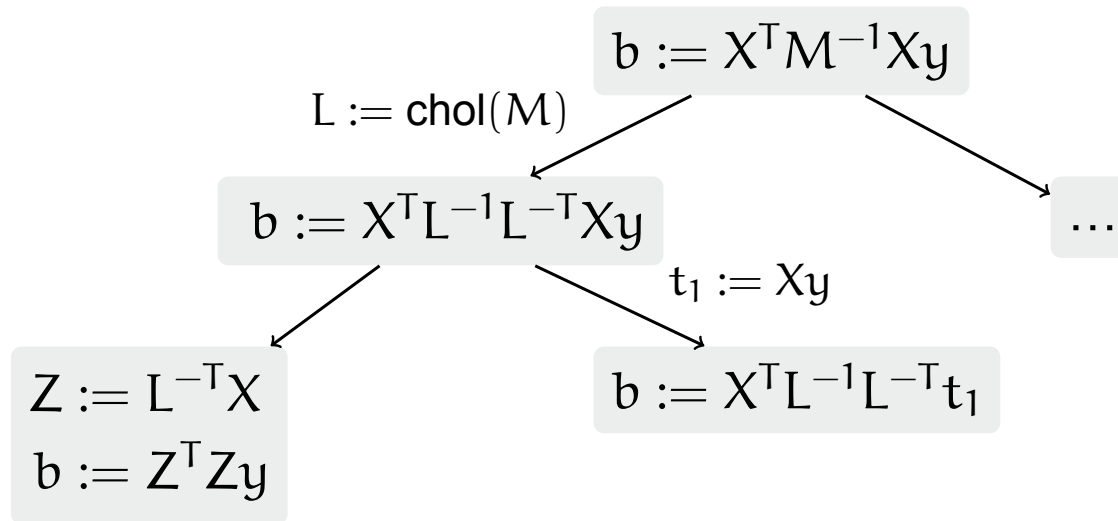
Parallel Code Generation



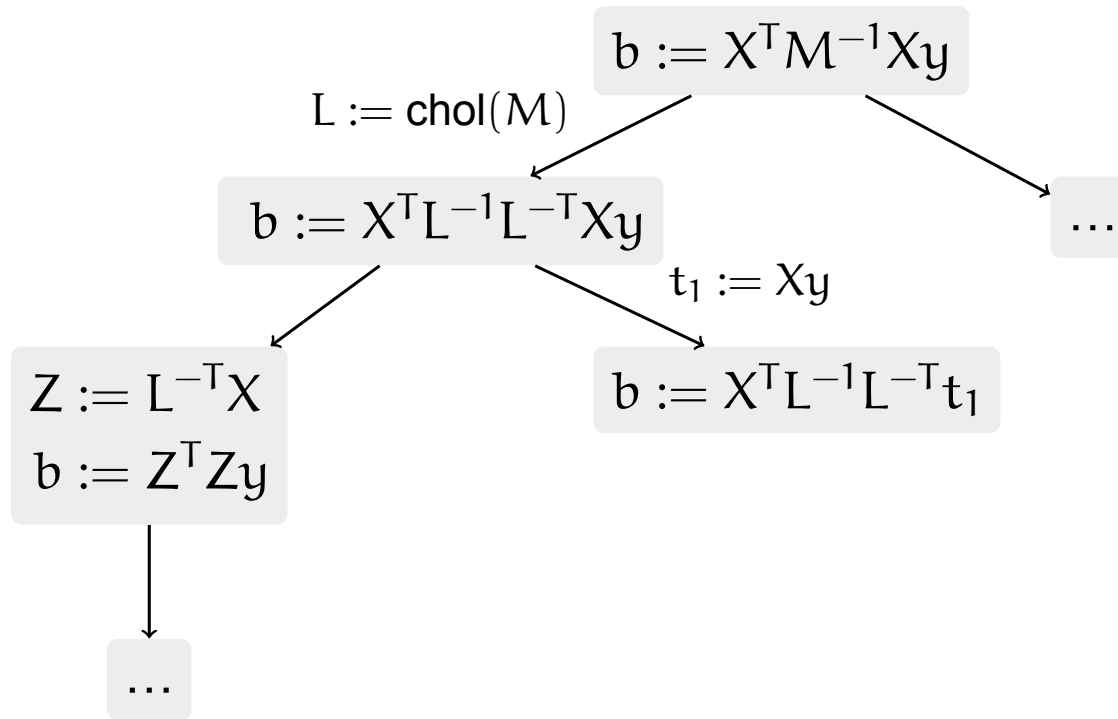
Parallel Code Generation



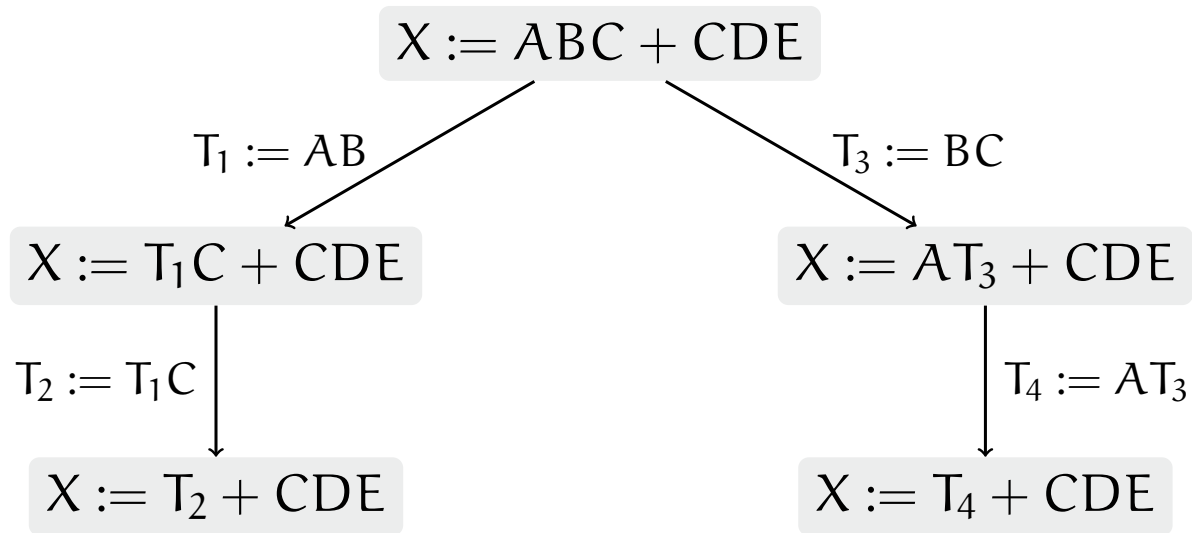
Parallel Code Generation



Parallel Code Generation



Reducing Redundancy



Parallel Code Generation

Reducing Redundancy

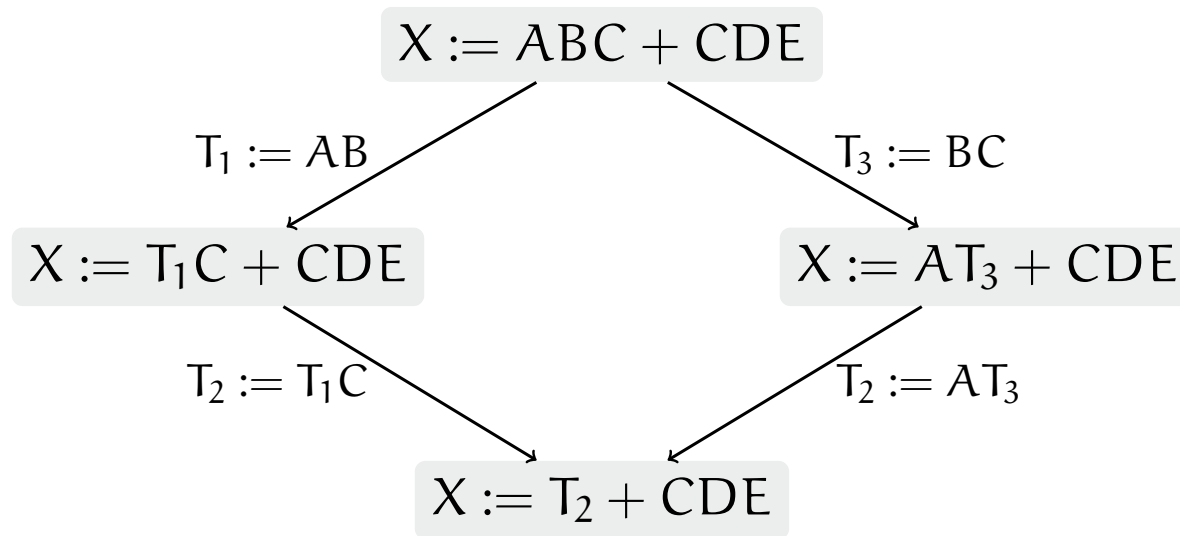
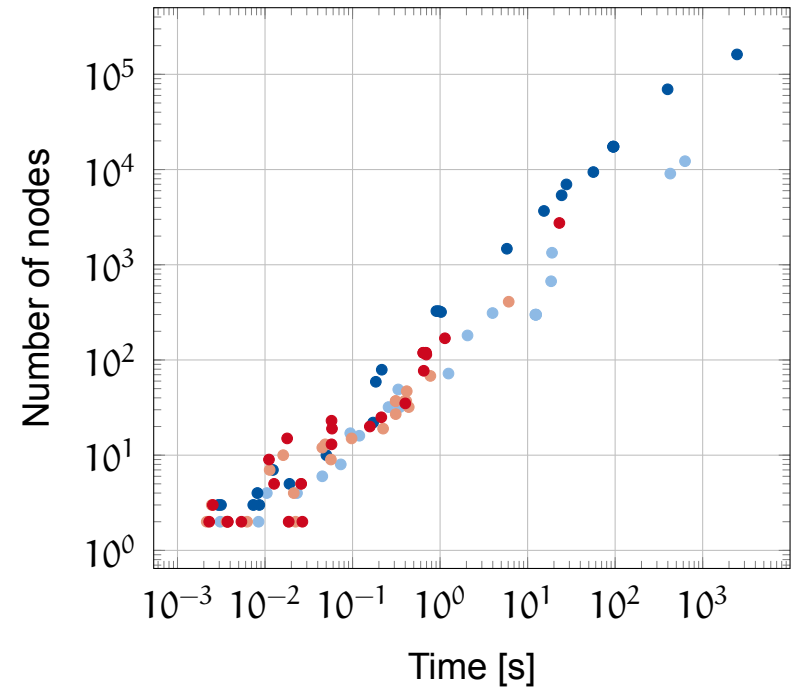
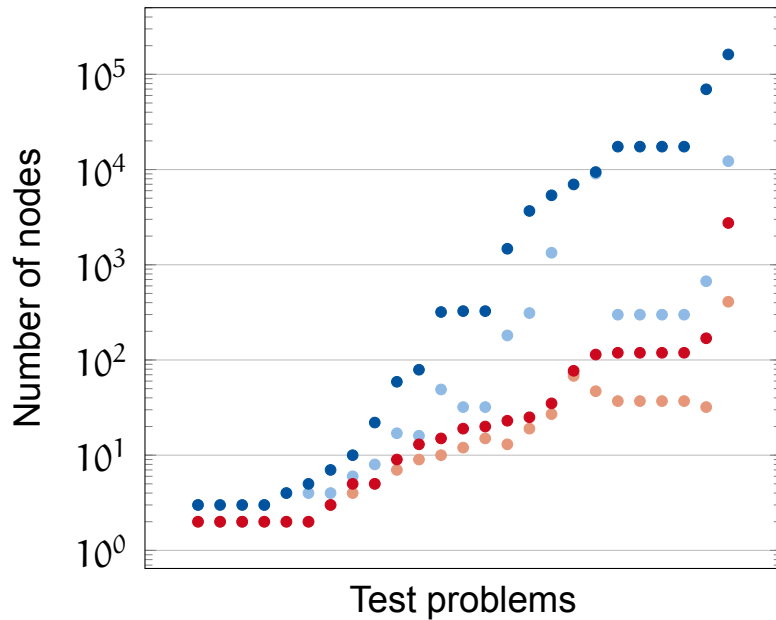


Table of expressions and intermediate operands:

tmp	expr
T_1	AB
T_2	ABC
T_3	BC

Parallel Code Generation

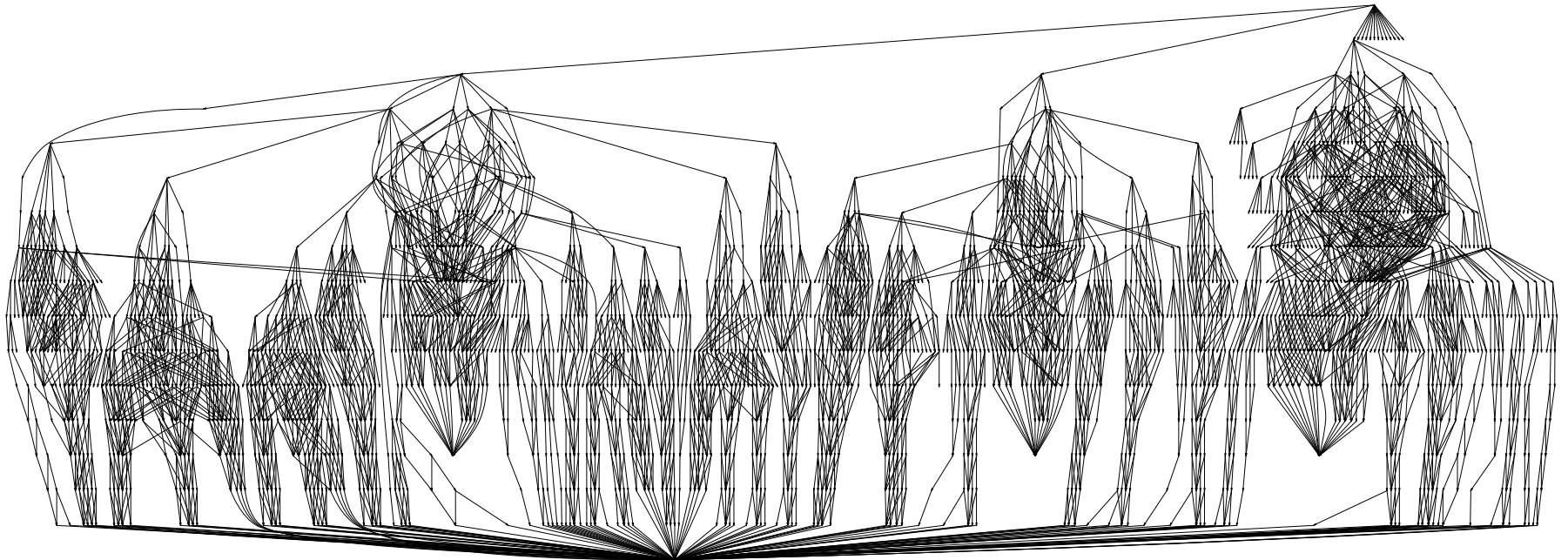
Reducing Redundancy



- exhaustive, merging
- exhaustive, no merging
- constructive, merging
- constructive, no merging

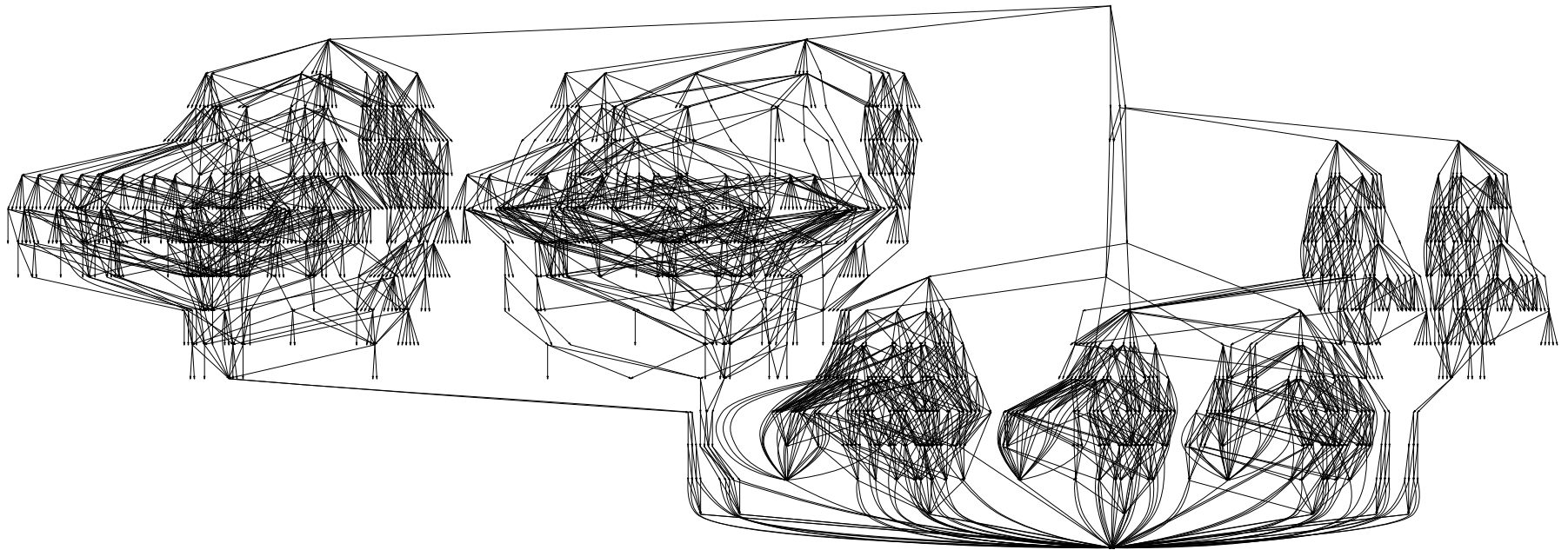
Parallel Code Generation

$$z := (X^T S^{-1} X)^{-1} X^T S^{-1} y \quad S \text{ is symmetric positive definite}$$



Parallel Code Generation

$$x := W(A^T(AWA^T)^{-1}b - c) \quad W \text{ is diagonal, diagonal elements are positive}$$



Parallel Code Generation

The Good

- Finding redundant nodes can be done efficiently using hash table.

The Bad

- Access to table of expressions & intermediates has to be protected.
- Merging nodes requires synchronization.
- Graph is not uniform at all.
- Graph is initially not known.

Possible solutions:

- Tradeoff between merging and redundancy?
- Parallelize computations on nodes?
- Nodes as tasks?

References

- [AB⁺99] Edward Anderson, Zhaojun Bai, et al. *LAPACK Users' guide*, volume 9. SIAM, 1999.
- [ATN09] Cédric Augonnet, Samuel Thibault, and Raymond Namyst. Automatic Calibration of Performance Models on Heterogeneous Multicore Architectures. *Euro-Par Workshops*, 6043(Chapter 9):56–65, 2009.
- [DDC⁺90] Jack J. Dongarra, Jeremy Du Croz, et al. A set of Level 3 Basic Linear Algebra Subprograms. *ACM TOMS*, 16(1):1–17, 1990.
- [LKHL03] Heejo Lee, Jong Kim, Sung Je Hong, and Sunggu Lee. Processor Allocation and Task Scheduling of Matrix Chain Products on Parallel Systems. *Parallel and Distributed Systems, IEEE Transactions on*, 14(4):394–407, 2003.
- [PB12] Elmar Peise and Paolo Bientinesi. Performance Modeling for Dense Linear Algebra. In *Proceedings of the 2012 SC Companion: High Performance Computing, Networking Storage and Analysis*

References

(*PMBS12*), SCC '12, pages 406–416, Washington, DC, USA, November 2012. IEEE Computer Society.

[PB15] Elmar Peise and Paolo Bientinesi. The ELAPS Framework - Experimental Linear Algebra Performance Studies. *CoRR*, cs.PF, 2015.