

Sparse Matrix Algorithms

combinatorics + numerical methods + applications
Math + X

Tim Davis
University of Florida

June 2013

contributions to the field
current work
vision for the future

- **Math+X**

- Math = [combinatorics + linear algebra + graph theory]
- X = [high-performance combinatorial scientific computing
+ many applications enabled by my contributions]

- **Roadmap: past, current, future work**

- **Sparse matrix algorithms**

- **Contributions to the field**

- from theory, to algorithms, to reliable software, to applications
- sparse Cholesky update/downdate (CHOLMOD)
- unsymmetric multifrontal LU (UMFPACK)
- multifrontal QR (SuiteSparseQR)

- **Current work**

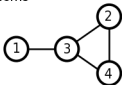
- highly concurrent methods (GPU or massive CPU core)
- NVIDIA Academic Partner

- **Future vision**

Math+X

X = high-performance combinatorial scientific computing + applications

combinatorics,
graph theory,
NP-hard problems



k , unmatched column

$i_1 \in \mathcal{A}_k$ $j_1 = \text{jmatch}[i_1]$

$i_2 \in \mathcal{A}_{j_1}$ $j_2 = \text{jmatch}[i_2]$

$i_3 \in \mathcal{A}_{j_2}$ $j_3 = \text{jmatch}[i_3]$

$i_4 \in \mathcal{A}_{j_3}$

(i_4 is unmatched)

$$\|b - Ax\|_2 = \|Q^T b - R x\|_2 = \left\| \begin{bmatrix} Q_1^T b - R_1 x \\ Q_2^T b \end{bmatrix} \right\|_2$$

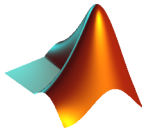
linear algebra,
scientific computing

$$Ax = b$$

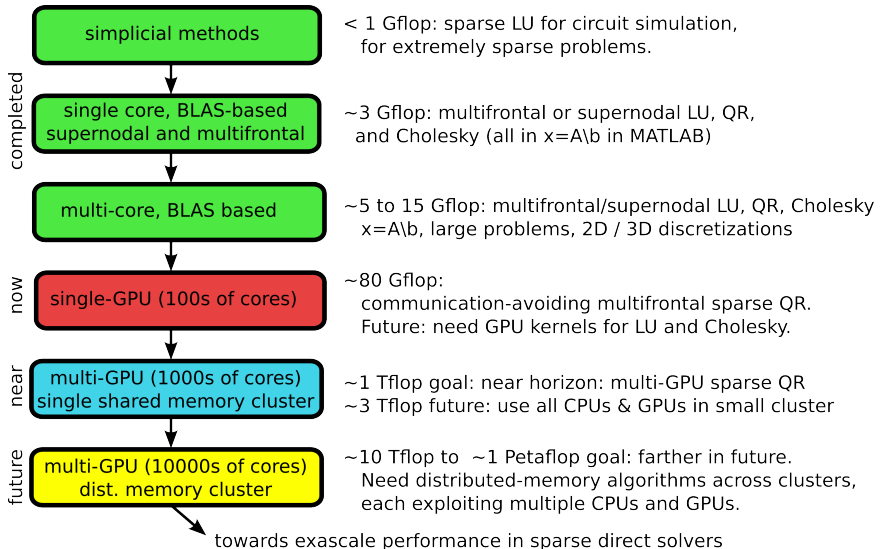
$$\begin{bmatrix} A_{11} & A_{13} \\ A_{22} & A_{23} \\ A_{13}^T & A_{23}^T & A_{33} \end{bmatrix}$$

**COMBINATORIAL
SCIENTIFIC
COMPUTING**

applications

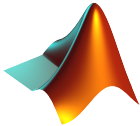


Roadmap: past, current, future work



Applications using SuiteSparse

General toolboxes for computational science



Applications using SuiteSparse

Linux distros



Ubuntu



Finite-element methods and differential equations



Applications using SuiteSparse

Circuit / power / semiconductor simulation



Computer graphics / computer vision / robotics



Applications using SuiteSparse

Mathematical optimization

C V X O P T
PYTHON SOFTWARE FOR CONVEX OPTIMIZATION

 **cvxpy**
A Python package for modeling convex optimization problems

 **CVX**
RESEARCH



mosek
<http://www.mosek.com>

Geophysical modeling



The logo for GEO MODELLING SOLUTIONS, featuring the text "GEO MODELLING SOLUTIONS" in a sans-serif font with a stylized arc above the text.



OpenSees

Financial simulation

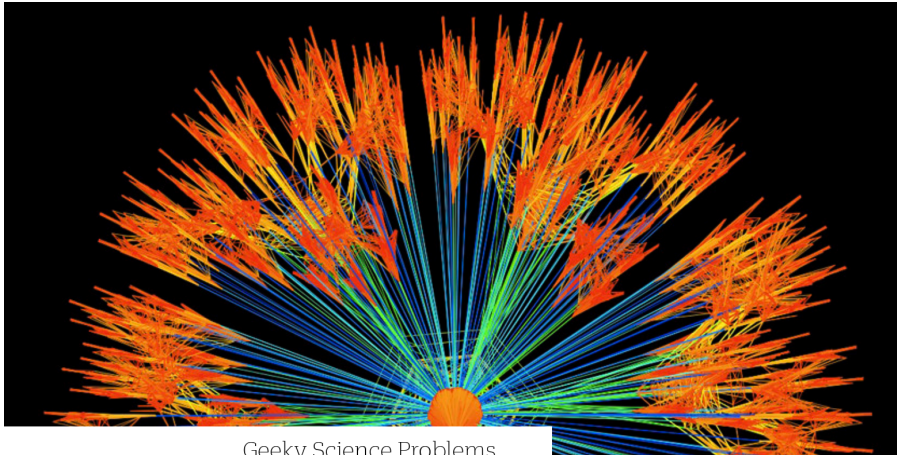
INTEX

Stellar evolution

MESA

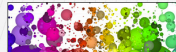
Applications using SuiteSparse

- **University of Florida Sparse Matrix Collection (with Hu)**
 - 2500+ sparse matrices from real applications



Geeky Science Problems
Double as Works of Art

SPARSE MATRICES: A LOT PRETTIER THAN THEY SOUND



Solve $Lx = b$ with L unit lower triangular; L, x, b are sparse

$x = b$

for $j = 0$ to $n - 1$ **do**

if $x_j \neq 0$

for each $i > j$ for which $l_{ij} \neq 0$ **do**

$x_i = x_i - l_{ij}x_j$

- non-optimal time $O(n + |b| + f)$, where $f = \text{flop count}$
- problem: outer loop and the test for $x_j \neq 0$
- solution: suppose we knew \mathcal{X} , the nonzero pattern of x
- optimal time $O(|b| + f)$, but how do we find \mathcal{X} ?
(Gilbert/Peierls)

Sparse matrix algorithms

Solve $Lx = b$ with L unit lower triangular; L, x, b are sparse

$x = b$

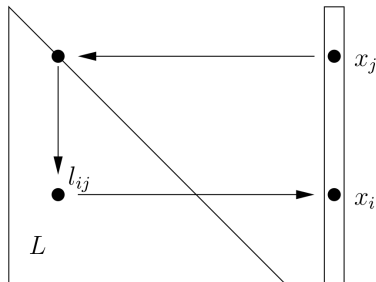
for $j = 0$ to $n - 1$ **do**

if $x_j \neq 0$

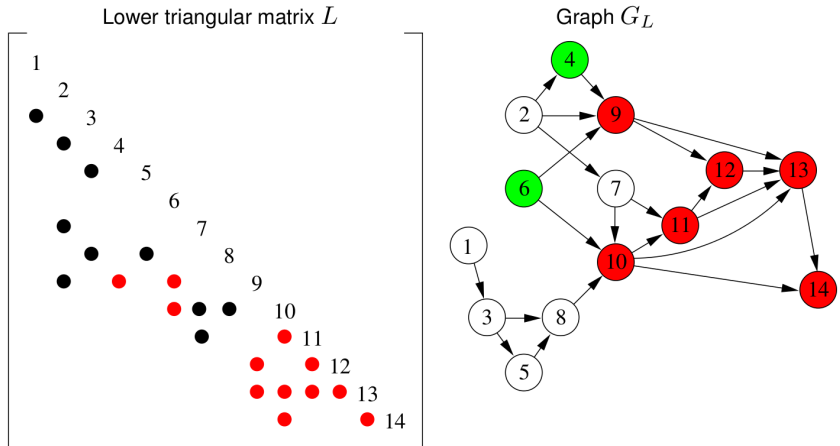
for each $i > j$ for which $l_{ij} \neq 0$ **do**

$x_i = x_i - l_{ij}x_j$

- if $b_i \neq 0$ then $x_i \neq 0$
- if $x_j \neq 0$ and $\exists i (l_{ij} \neq 0)$
 then $x_i \neq 0$
- start with pattern \mathcal{B} of b
- graph \mathcal{L} : edge (j, i) if $l_{ij} \neq 0$
- $\mathcal{X} = \text{Reach}_{\mathcal{L}}(\mathcal{B})$
(Gilbert/Peierls)



Sparse matrix algorithms



If $\mathcal{B} = \{4, 6\}$ then $\mathcal{X} = \{6, 10, 11, 4, 9, 12, 13, 14\}$

The update/downdate problem:

- Given $A = LL^T$
- A undergoes a low-rank change
- compute $\bar{L}\bar{L}^T = A \pm ww^T$
- arises in optimization, crack propagation, robotics, new data in least-squares, ...

Sparse Cholesky update/downdate

Rank-1 update, $LL^T + ww^T$ (Carlson)

$$\bar{\beta} = 1$$

for $j = 1$ to n

compute coefficients:

$$\alpha = w_j / l_{jj}$$

$$\beta = \bar{\beta}, \quad \bar{\beta} = \sqrt{\beta^2 + \alpha^2}$$

$$\gamma = \alpha / (\bar{\beta}\beta), \quad \delta = \beta / \bar{\beta}$$

update diagonal:

$$l_{jj} = \delta l_{jj} + \gamma w_j$$

$$w_j = \alpha$$

update below diagonal:

$$\mathbf{t} = \mathbf{w}_{j+1:n}$$

$$\mathbf{w}_{j+1:n} = \mathbf{w}_{j+1:n} - \alpha \mathbf{L}_{j+1:n,j}$$

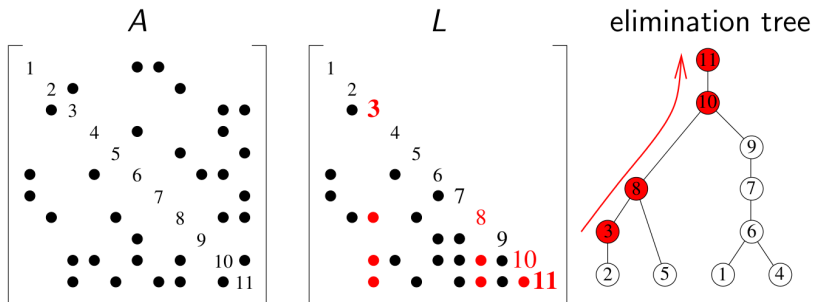
$$\bar{\mathbf{L}}_{j+1:n,j} = \delta \mathbf{L}_{j+1:n,j} + \gamma \mathbf{t}$$

end

Key observations:

- overwrites w with solution to $L\bar{w} = w$
- $w = L \backslash w$ in MATLAB notation
- j th column of L changes only if $\bar{w}_j = (L^{-1}w)_j \neq 0$,
- thus, need pattern of triangular solve
- ... but the Cholesky L can be pruned

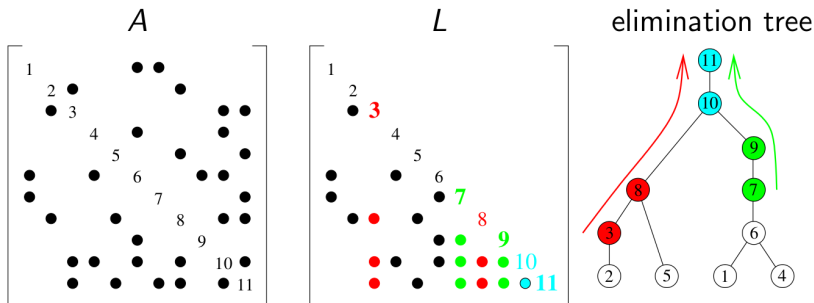
Sparse Cholesky update/downdate



Key results

- if \mathcal{L} doesn't change: columns in L that change = path from $\min \mathcal{W}$ to root of the etree
- if \mathcal{L} does change, follow the path in etree of \bar{L}
- Update/downdate in time proportional to the number of entries in L that change

Sparse Cholesky update/downdate



Multiple rank, with dynamic supernodes

- multiple rank: $A \pm WW^T$, follow multiple paths in the tree
- supernodes: adjacent columns of L with identical pattern; break apart and merge in update/downdate
- dynamic supernodes: find them as we go (cuts memory traffic)

Sparse Cholesky update/downdate

CHOLMOD update/downdate: key results / impact

- update/downdate faster than a $Lx = b$ solve for dense b
- example application: LPDASA (Hager and Davis)
 - maintains Cholesky factorization of $A_F A_F^T$ for basis set F
 - update/downdate as basis set changes
- example: g2o (Kümmerle et al)
 - robotics, simultaneous localization and mapping
 - builds a map of its surroundings
 - update/downdate as new images arrive
- example: crack propagation (Pais, Kim, Davis, Yeralan)
 - structural engineering problem: crack in aircraft fuselage
 - update/downdate as crack progresses through airframe

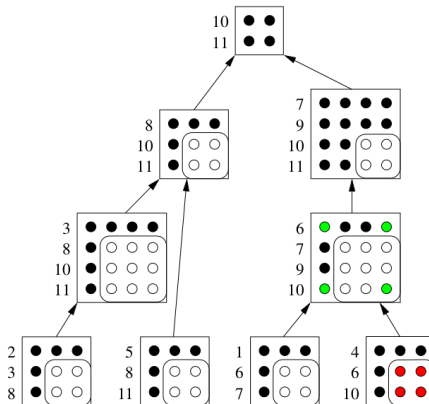
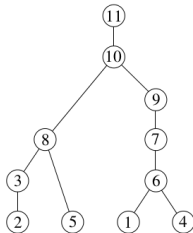
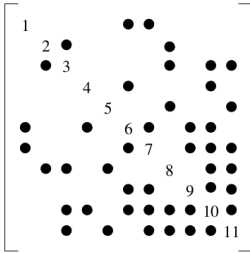
CHOLMOD: Supernodal Sparse Cholesky update/downdate

CHOLMOD: key results / impact

- sparse chol in MATLAB, and $x=A\backslash b$
- typical 10x speedup when incorporated into MATLAB
- peak performance over 50% of CPU theoretical peak
- example applications using CHOLMOD:
 - Google Street View, PhotoTours, and 3D Earth
 - Mentor Graphics: circuit simulation
 - Cadence: CAD design
 - VisionMap: satellite image processing
 - CVXOPT: convex optimization
 - ...

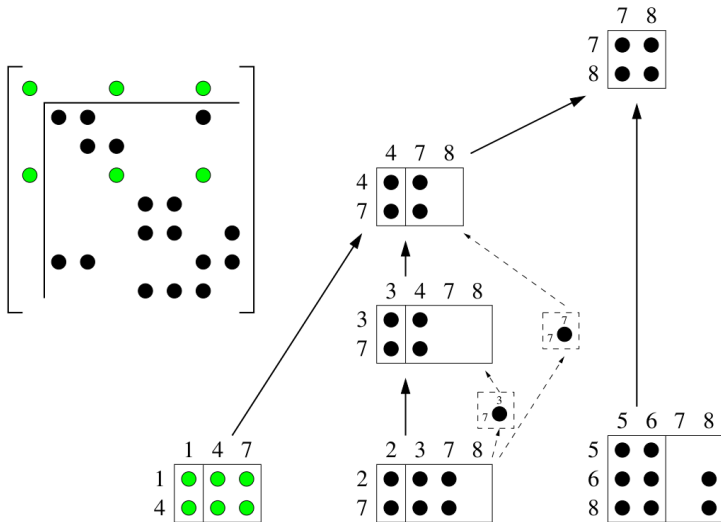
Multifrontal method

- Classic symmetric multifrontal method (Duff, Reid, others)
- Cliques + elimination tree = sequence of frontal matrices
- Dense factorization within a front; assemble data into parent



UMFPACK: unsymmetric multifrontal method

- Frontal matrices become rectangular
- Assemble data into ancestors, not just parents



UMFPACK: unsymmetric multifrontal method

Key results / impact

- high-performance via dense matrix kernels within each front
- symbolic preordering and analysis, followed by revised local pivot search with approximate unsymmetric degree update
- could extend to handle threshold rook pivoting
- widely used
 - sparse LU and $x=A\backslash b$ in MATLAB
 - Mathematica
 - IBM circuit simulation
 - finite-element solvers: NASTRAN, ANSYS, FEniCS, ...
 - CVXOPT
 - ...

SuiteSparseQR: multifrontal sparse QR factorization

Key results / impact

- rectangular fronts like UMFPACK, but simpler frontal matrix assembly
- multicore parallelism
- amenable to GPU implementation (in progress)
- on multicore CPU: up to 14 Gflops
- sparse qr in MATLAB, and $x=A\backslash b$

- Least squares problem: 2 million by 110 thousand

| Method | ordering | procs | time |
|--------------------------|----------|-------|-------------|
| prior $x=A \backslash b$ | COLMMD | 1 | ? |
| prior $x=A \backslash b$ | AMD | 1 | 11 days |
| MA49 | AMD | 1 | 3.5 hours |
| SuiteSparseQR | AMD | 1 | 1.5 hours |
| SuiteSparseQR | METIS | 1 | 45 minutes |
| SuiteSparseQR | METIS | 16 | 7.3 minutes |

- Algorithmic speedup vs prior $x=A \backslash b$: 375x
- Parallel speedup: 5.75x on 16 cores
- Total: 2,155x (14 Gflops on machine w/ 70 Gflops peak)
- Single core: 2.5 Gflop peak, same as LAPACK dense QR

Highly-concurrent heterogeneous computing

GPU computing: bulk synchronous model

- multiple Streaming Multiprocessors (SMs) on a GPU
- each SM with multiple SIMD threads
- one kernel launch: set of independent tasks

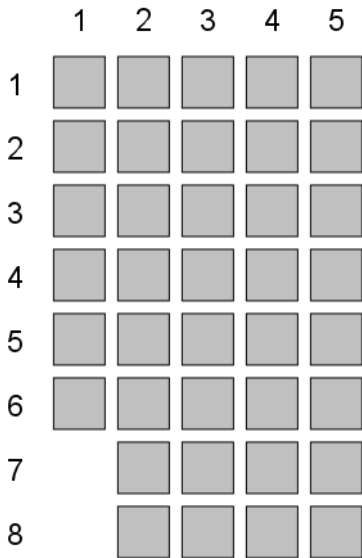
GPU-based sparse multifrontal QR

- symbolic ordering and analysis on the CPU (irregular)
- numerical factorization on the GPU (regular+irregular)

Multifrontal factorization and assembly

- Prior methods
 - one front at a time on the GPU
 - assembly on CPU
 - panel factorization on the CPU, applied on GPU
- Our multifrontal QR
 - many fronts on the GPU (entire subtree)
 - assembly on GPU
 - entire dense QR of each front on the GPU

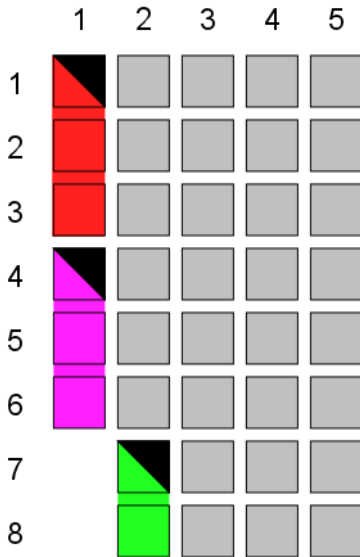
Highly-concurrent heterogeneous computing



Initial matrix:

- each tile a square submatrix
- panelsize: 3-by-1 tiles
- exploits initial staircase

Highly-concurrent heterogeneous computing



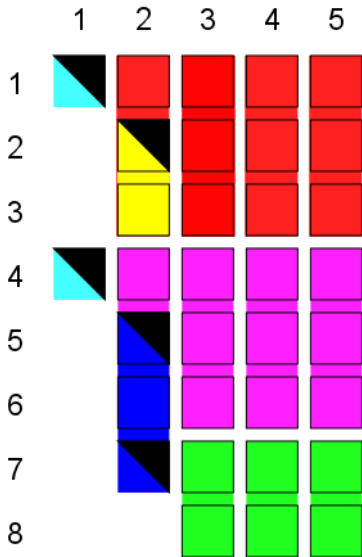
Householder bundles:

(1,2,3) new

(4,5,6) new

(7,8) new

Highly-concurrent heterogeneous computing



Householder bundles:

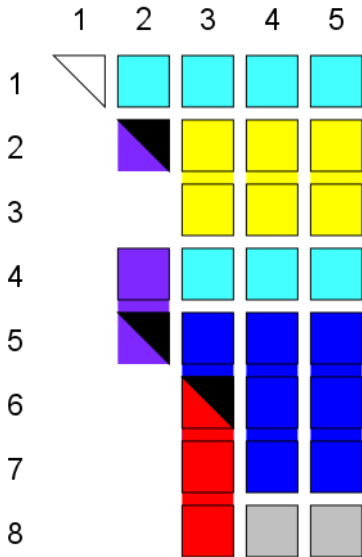
(1,2,3) applied, and
pipelined to (2,3)

(4,5,6) applied, and
pipelined to (5,6,7)

(7,8) applied

(1,4) new

Highly-concurrent heterogeneous computing



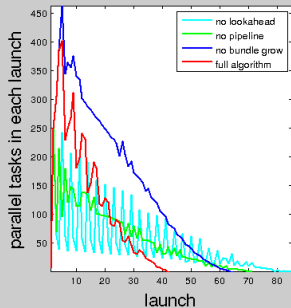
Householder bundles:

(2,3) applied

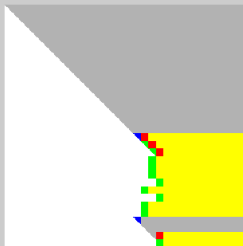
(5,6,7) applied, and
pipelined to (6,7,8)

(1,4) applied, and
pipelined to (2,4,5)

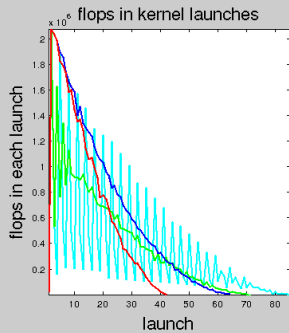
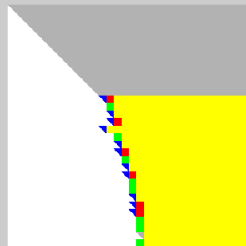
kernel launches



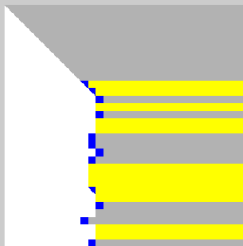
full pipelining



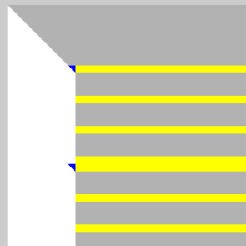
pipelining but no bundle growth



no pipelining



no pipelining, no lookahead

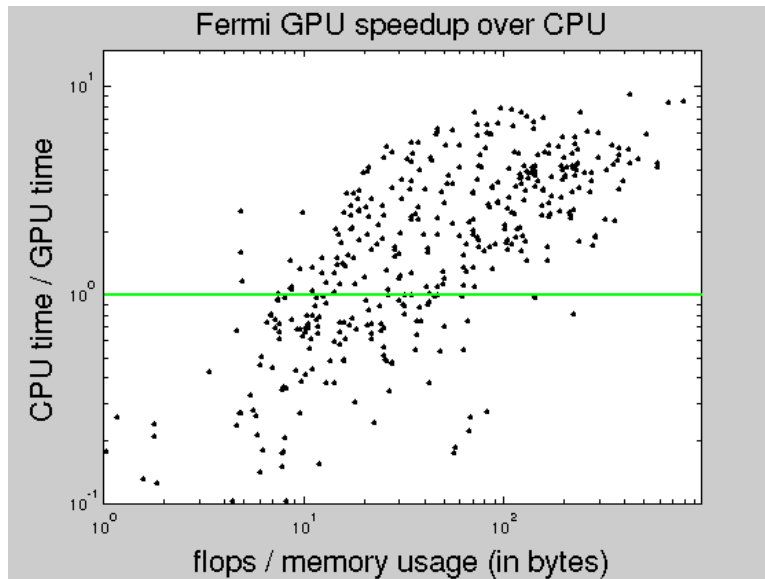


Highly-concurrent heterogeneous computing

- **Peak performance results on the GPU**

| | Fermi | Kepler |
|--------------------------|------------|---------------|
| GPU kernels: | | |
| apply block Householder | 183 Gflops | 260 Gflops |
| factorize 3 tiles | 27 Gflops | 20 Gflops |
| dense QR for large front | 107 Gflops | 120 Gflops |
| sparse QR on GPU | 80 Gflops | (in progress) |
| speedup over CPU | 9x | (in progress) |

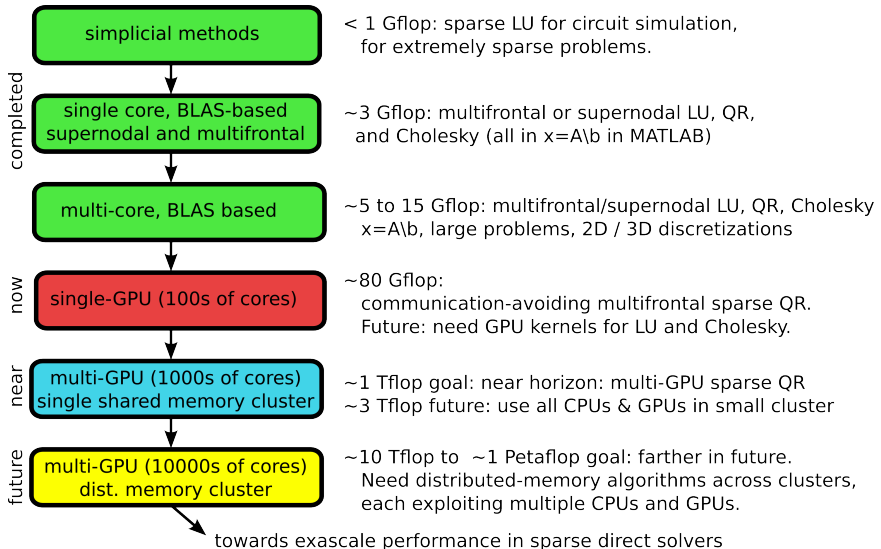
Highly-concurrent heterogeneous computing



Math + X

- Math = [combinatorics + linear algebra + graph theory]
- X = [high-performance combinatorial scientific computing
+ applications enabled by my contributions]
- computational mathematics: the future is heterogeneous
- driven by power constraints, need for parallelism
- high impact – getting it out the door:
 - novel algorithms: delivered in widely used robust software
 - Collected Algorithms of the ACM
 - enabling academic projects: Julia, R, Octave, FEnICS, ...
 - growing impact on industry, and industrial collaborations: Google, NVIDIA, Mentor Graphics, Cadence, MSC Software, Berkeley Design Automation, ...
 - applications: optimization, robotics, circuit simulation, computer graphics, computer vision, finite-element methods, geophysics, stellar evolution, financial simulation, ...

Roadmap: past, current, future work



- **Math+X**

- Math = [combinatorics + linear algebra + graph theory]
- X = [high-performance combinatorial scientific computing
+ many applications enabled by my contributions]

- **Roadmap: past, current, future work**

- **Sparse matrix algorithms**

- **Contributions to the field**

- from theory, to algorithms, to reliable software, to applications
- sparse Cholesky update/downdate (CHOLMOD)
- unsymmetric multifrontal LU (UMFPACK)
- multifrontal QR (SuiteSparseQR)

- **Current work**

- highly concurrent methods (GPU or massive CPU core)
- NVIDIA Academic Partner

- **Future vision** : individual and collaborative

- GPU / CPU heterogeneous parallel computing
- towards exascale
- continue creating algorithms with deep impact

Math + X, where X = [Poetry + Art]

Sea Fever, by Masefield (1902)

I must go down to the seas again, to the lonely sea and the sky,
And all I ask is a tall ship and a star to steer her by,
And the wheel's kick and the wind's song
 and the white sail's shaking,
And a grey mist on the sea's face and a grey dawn breaking.

C Fever, by T.D. (2010)

I must go code in both C and M, not only C and VI,
And all I ask is a Linux box and a mouse to steer her by,
And the while's break and the if then
 and the valgrind's shaking,
And a dash 0 so the C's fast and a switch case break-ing.

questions?

