

AMATH 483/583

High Performance Scientific Computing

Lecture 9:

Strassen's Algorithm

Sparse Matrix Computation

Andrew Lumsdaine

Northwest Institute for Advanced Computing

Pacific Northwest National Laboratory

University of Washington

Seattle, WA

Overview

- Review: Locality and optimization strategies
- Sparsity
- Coordinate format (COO)
- Compressed sparse row (CSR)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

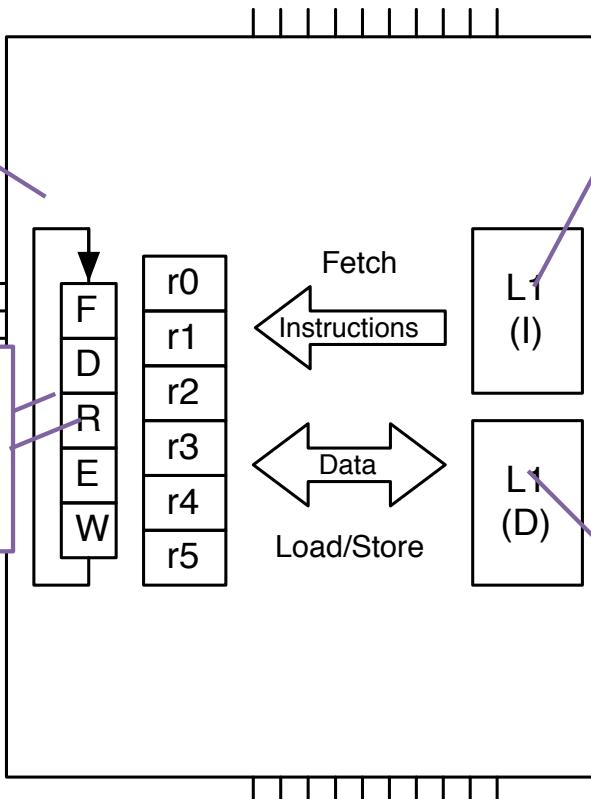
W
UNIVERSITY OF
WASHINGTON

Locality → Strategy

The next operand may be "near" the last

It could be "near" in time or space

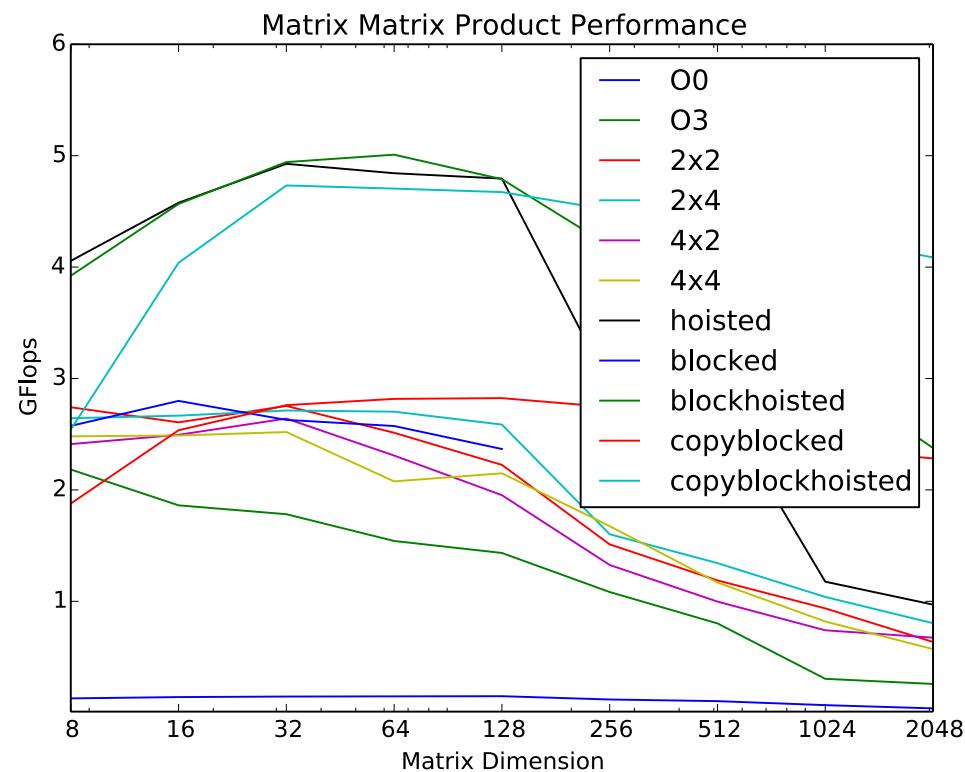
Clock
... → ← ...
cycle



Near in time
(temporal locality):
the next operand is a previous operand

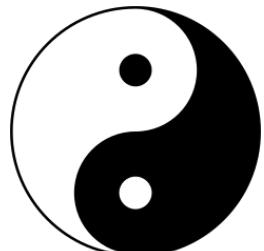
Near in space (**spatial locality**): the next operand is in a nearby memory location to a previous operand

Blocking and Tiling and Hoisting and Copying



What Else Can We Do for Performance

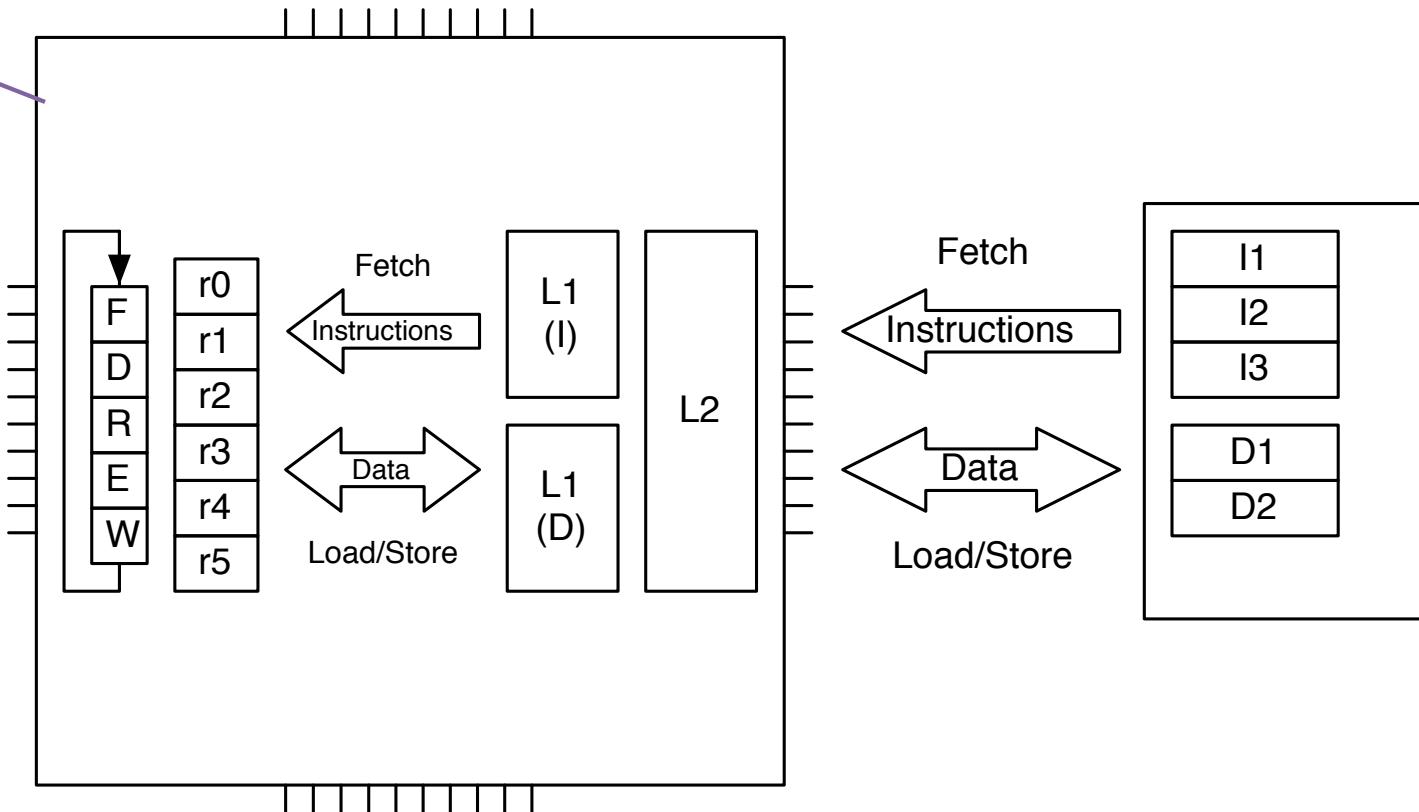
Exploit features
that make
hardware fast



Clock

... ...

cycle



General Performance Principles

- Work harder
 - Faster core
- Work smarter
 - Branch predictions, etc
 - Better compilation
 - Better algorithm
 - Better implementation
- Get help

Dennard scaling
(ended 2005)

What
about this?

We did this

Another Way to Work Smarter

(Work less)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

W
UNIVERSITY OF
WASHINGTON

Strassen's Algorithm

Volker Strassen.

Gaussian Elimination is not Optimal.

Numer Math, Vol 13, No.4, Aug 1969.

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

$$C_{00} = A_{00}B_{00} + A_{01}B_{10}$$

$$C_{01} = A_{00}B_{01} + A_{01}B_{11}$$

$$C_{10} = A_{10}B_{00} + A_{11}B_{10}$$

$$C_{11} = A_{10}B_{01} + A_{11}B_{11}$$

Eight multiplies

If these are matrix
blocks: Eight
matrix multiplies

Strassen's Algorithm

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

$$T_0 = (A_{00} + A_{11})(B_{00} + B_{11})$$

$$T_1 = (A_{10} + A_{11})(B_{00})$$

$$T_2 = (A_{00})(B_{01} - B_{11})$$

$$T_3 = (A_{11})(B_{10} - B_{00})$$

$$T_4 = (A_{00} + A_{01})(B_{11})$$

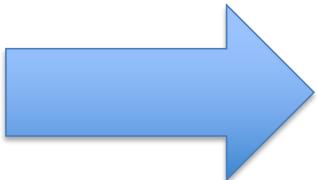
$$T_5 = (A_{10} - A_{00})(B_{00} + B_{01})$$

$$T_6 = (A_{01} - A_{11})(B_{10} + B_{11})$$

Seven
multiplies

Seven
matrix
multiplies

Recurse



$$C_{00} = T_0 + T_3 - T_4 + T_6$$

$$C_{01} = T_2 + T_4$$

$$C_{10} = T_1 + T_4$$

$$C_{11} = T_0 - T_1 + T_2 + T_5$$

Many adds
and subtracts

Strassen's Algorithm

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

Seven
matrix
multiplies

Recurse

$$T_0 = (A_{00} + A_{11})(B_{00} + B_{11})$$

$$T_1 = (A_{10} + A_{11})(B_{00})$$

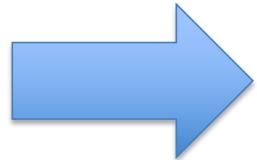
$$T_2 = (A_{00})(B_{01} - B_{11})$$

$$T_3 = (A_{11})(B_{10} - B_{00})$$

$$T_4 = (A_{00} + A_{01})(B_{11})$$

$$T_5 = (A_{10} - A_{00})(B_{00} + B_{01})$$

$$T_6 = (A_{01} - A_{11})(B_{10} + B_{11})$$



$$C_{00} = T_0 + T_3 - T_4 + T_6$$

$$C_{01} = T_2 + T_4$$

$$C_{10} = T_1 + T_4$$

$$C_{11} = T_0 - T_1 + T_2 + T_5$$

$O(N^3)$ work vs $O(N^2)$ data

Multiply

Add

Strassen's Algorithm

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

$$\begin{aligned} T_0 &= (A_{00} + A_{11})(B_{00} + B_{11}) \\ T_1 &= (A_{10} + A_{11})(B_{00}) \\ T_2 &= (A_{00})(B_{01} - B_{11}) \\ T_3 &= (A_{11})(B_{10} - B_{00}) \\ T_4 &= (A_{00} + A_{01})(B_{11}) \\ T_5 &= (A_{10} - A_{00})(B_{00} + B_{01}) \\ T_6 &= (A_{01} - A_{11})(B_{10} + B_{11}) \end{aligned}$$

Divide and Conquer

Recurse

$$\begin{aligned} C_{00} &= T_0 + T_3 - T_4 + T_6 \\ C_{01} &= T_2 + T_4 \\ C_{10} &= T_1 + T_4 \\ C_{11} &= T_0 - T_1 + T_2 + T_5 \end{aligned}$$

$O(N^3)$ work vs $O(N^2)$ data

Seven matrix multiplies

Each block is size $\frac{N}{2}$

$$\left(\frac{N}{2}\right)^3 = \frac{N^3}{8}$$

$$\frac{7}{8}N^3$$

Strassen's Algorithm

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

$$T_0 = (A_{00} + A_{11})(B_{00} + B_{11})$$

$$T_1 = (A_{10} + A_{11})(B_{00})$$

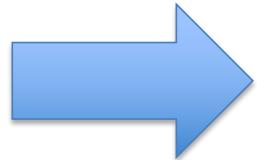
$$T_2 = (A_{00})(B_{01} - B_{11})$$

$$T_3 = (A_{11})(B_{10} - B_{00})$$

$$T_4 = (A_{00} + A_{01})(B_{11})$$

$$T_5 = (A_{10} - A_{00})(B_{00} + B_{01})$$

$$T_6 = (A_{01} - A_{11})(B_{10} + B_{11})$$



$$\frac{7}{8} \frac{7}{8} \cdots \frac{7}{8}$$

How many
of these

Divide and
conquer

$$C_{00} = T_0 + T_3 - T_4 + T_6$$

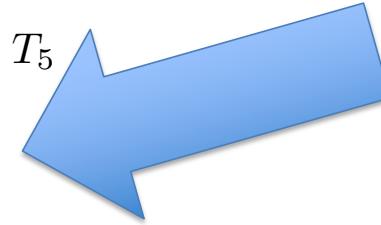
$$C_{01} = T_2 + T_4$$

$$C_{10} = T_1 + T_4$$

$$C_{11} = T_0 - T_1 + T_2 + T_5$$

$\log_2(N)$

$O(N^{\log_2 7})$

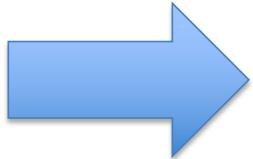


$$O(N^{\log_2 7}) \ll O(N^{\log_2 8}) = O(N^3)$$

Strassen's Algorithm

$$\begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

$$\begin{aligned} T_0 &= (A_{00} + A_{11})(B_{00} + B_{11}) \\ T_1 &= (A_{10} + A_{11})(B_{00}) \\ T_2 &= (A_{00})(B_{01} - B_{11}) \\ T_3 &= (A_{11})(B_{10} - B_{00}) \\ T_4 &= (A_{00} + A_{01})(B_{11}) \\ T_5 &= (A_{10} - A_{00})(B_{00} + B_{01}) \\ T_6 &= (A_{01} - A_{11})(B_{10} + B_{11}) \end{aligned}$$



Limit?

$O(N^{2.38})$

Better algorithms

$$\begin{aligned} C_{00} &= T_0 + T_3 - T_4 + T_6 \\ C_{01} &= T_2 + T_4 \\ C_{10} &= T_1 + T_4 \\ C_{11} &= T_0 - T_1 + T_2 + T_5 \end{aligned}$$

Require large N

Limit Unknown, Biggest open question in numerical linear algebra

Another Way to Work Smarter

(Work less)

NORTHWEST INSTITUTE for ADVANCED COMPUTING

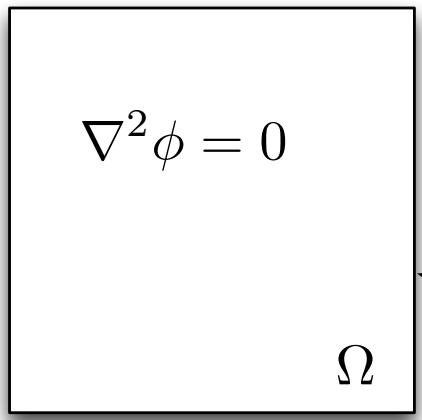
AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

W
UNIVERSITY OF
WASHINGTON

In Practice

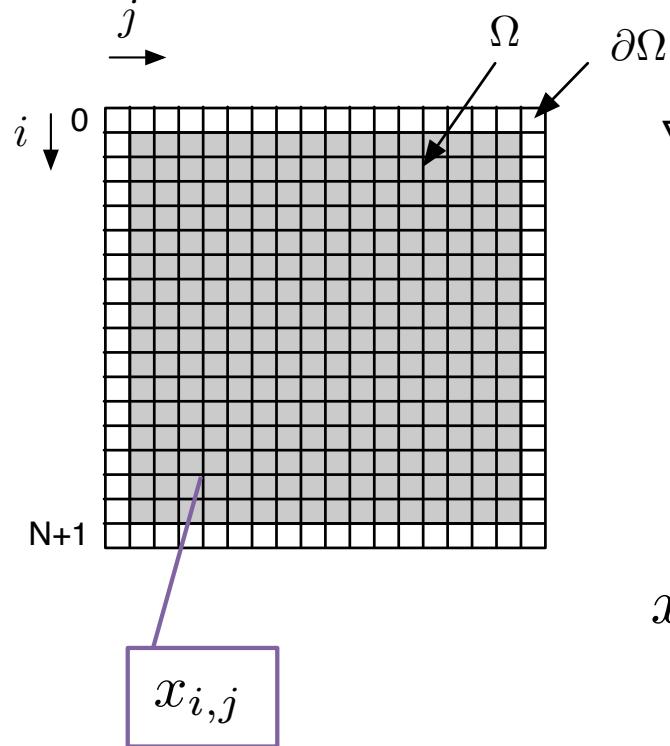
- Many scientific applications are based on solving systems of partial differential equations that model physical phenomena
- Laplace's equation on unit square is prototypical PDE

$$\nabla^2 \phi = 0 \quad \text{on } \Omega$$


A diagram of a unit square with vertices at (-1, -1), (1, -1), (1, 1), and (-1, 1). The interior of the square is labeled $\nabla^2 \phi = 0$. The bottom edge of the square is labeled Ω . An arrow points from the text $\partial\Omega$ to the bottom edge of the square.

$$\nabla \phi = f \quad \text{on } \partial\Omega$$

Laplace's Equation on a Regular Grid



$$\begin{aligned}\nabla^2 \phi &= 0 \quad \text{on } \Omega \\ \phi &= f \quad \text{on } \partial\Omega\end{aligned}$$

$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \cdots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \ddots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \cdots & -1 & 4 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \end{bmatrix}$$

Discretization

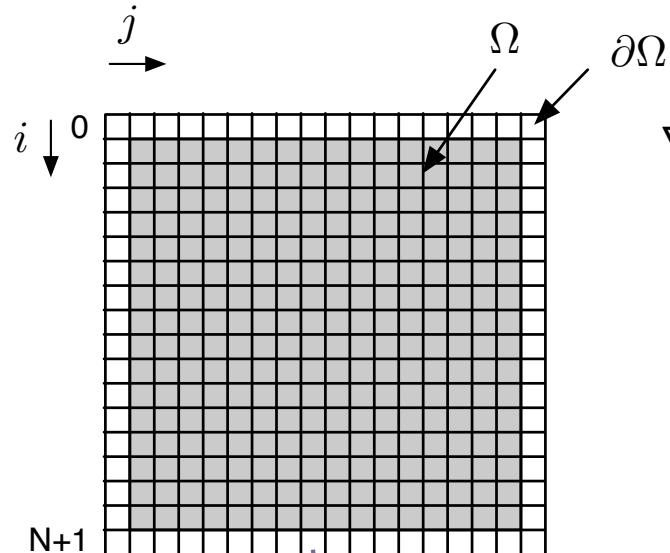
$$x_{i-1,j} + x_{i+1,j} + x_{i,j-1} + x_{i,j+1} - 4x_{i,j} = 0$$

$$x_{i,j} = (x_{i-1,j} + x_{i+1,j} + x_{i,j-1} + x_{i,j+1})/4$$

The value of each point on the grid

The average of its neighbors

Laplace's Equation on a Regular Grid



The boundary
is non-zero

$$x_{i,j} = (x_{i-1,j} + x_{i+1,j} + x_{i,j-1} + x_{i,j+1})/4$$

The value of each
point on the grid

The average of
its neighbors

Why isn't 0
the solution?

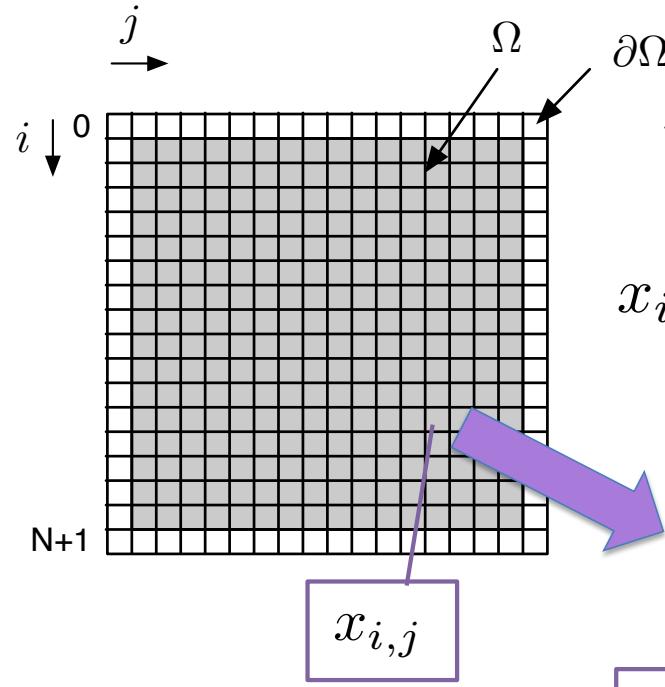
$$\begin{aligned}\nabla^2 \phi &= 0 \quad \text{on } \Omega \\ \phi &= f \quad \text{on } \partial\Omega\end{aligned}$$

Discret
The boundary
is non-zero

$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \cdots & -1 & & \\ -1 & \ddots & \ddots & \ddots & \ddots & \\ \vdots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \ddots & \ddots & \ddots & \ddots & -1 & \\ -1 & \cdots & -1 & 4 & & \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ \vdots \end{bmatrix}$$

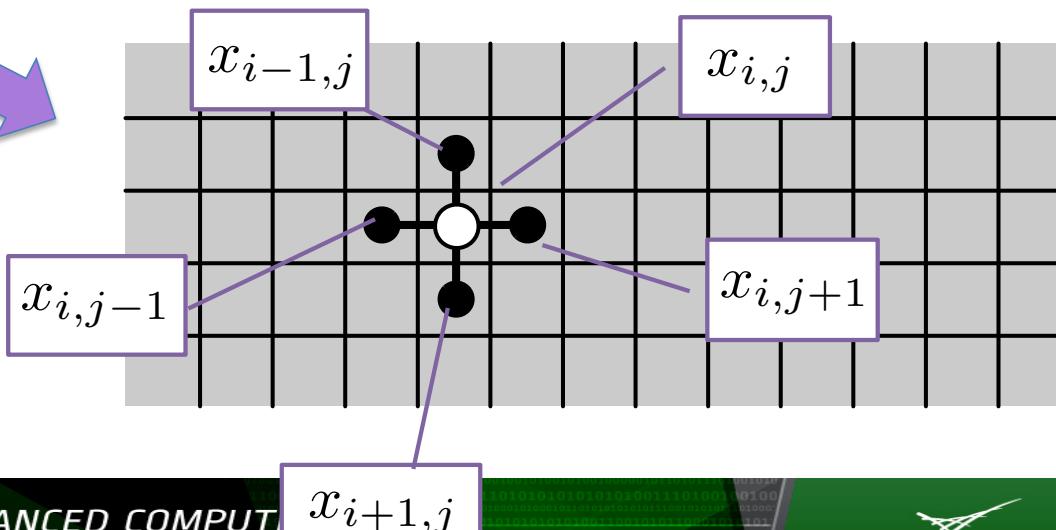
Non-zeros in
here due to
boundary

Laplace's Equation on a Regular Grid

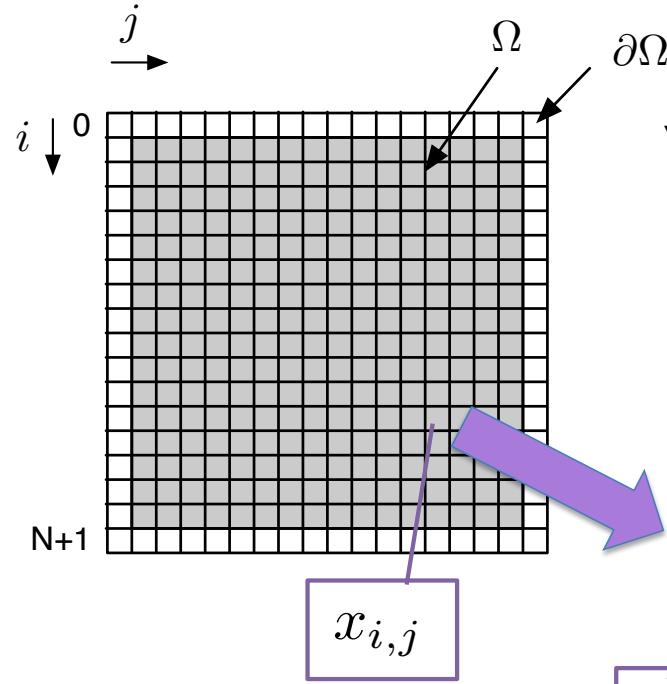


$$\begin{aligned}\nabla^2 \phi &= 0 \quad \text{on } \Omega \\ \phi &= f \quad \text{on } \partial\Omega\end{aligned}$$

$$x_{i,j} = (x_{i-1,j} + x_{i+1,j} + x_{i,j-1} + x_{i,j+1})/4$$

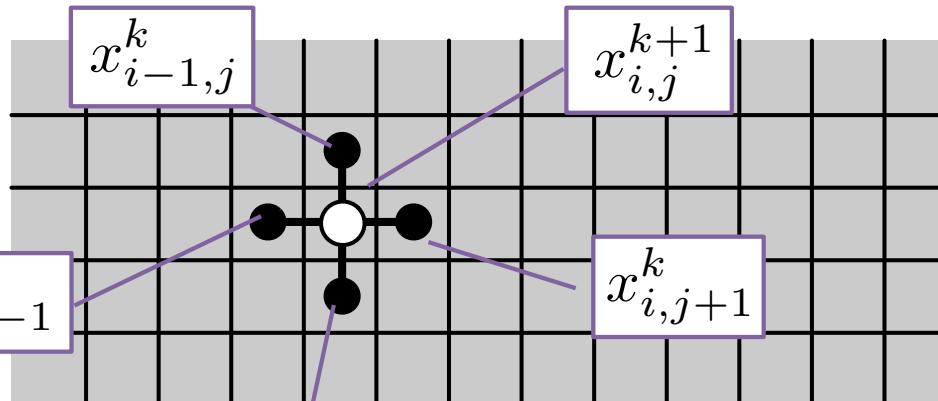


Iterating for a solution



$$\begin{aligned}\nabla^2 \phi &= 0 \quad \text{on } \Omega \\ \phi &= f \quad \text{on } \partial\Omega\end{aligned}$$

$$x_{i,j}^{k+1} = (x_{i-1,j}^k + x_{i+1,j}^k + x_{i,j-1}^k + x_{i,j+1}^k)/4$$



Approximation at iteration $k+1$

Average of approximation at iteration k

Iterating for a solution

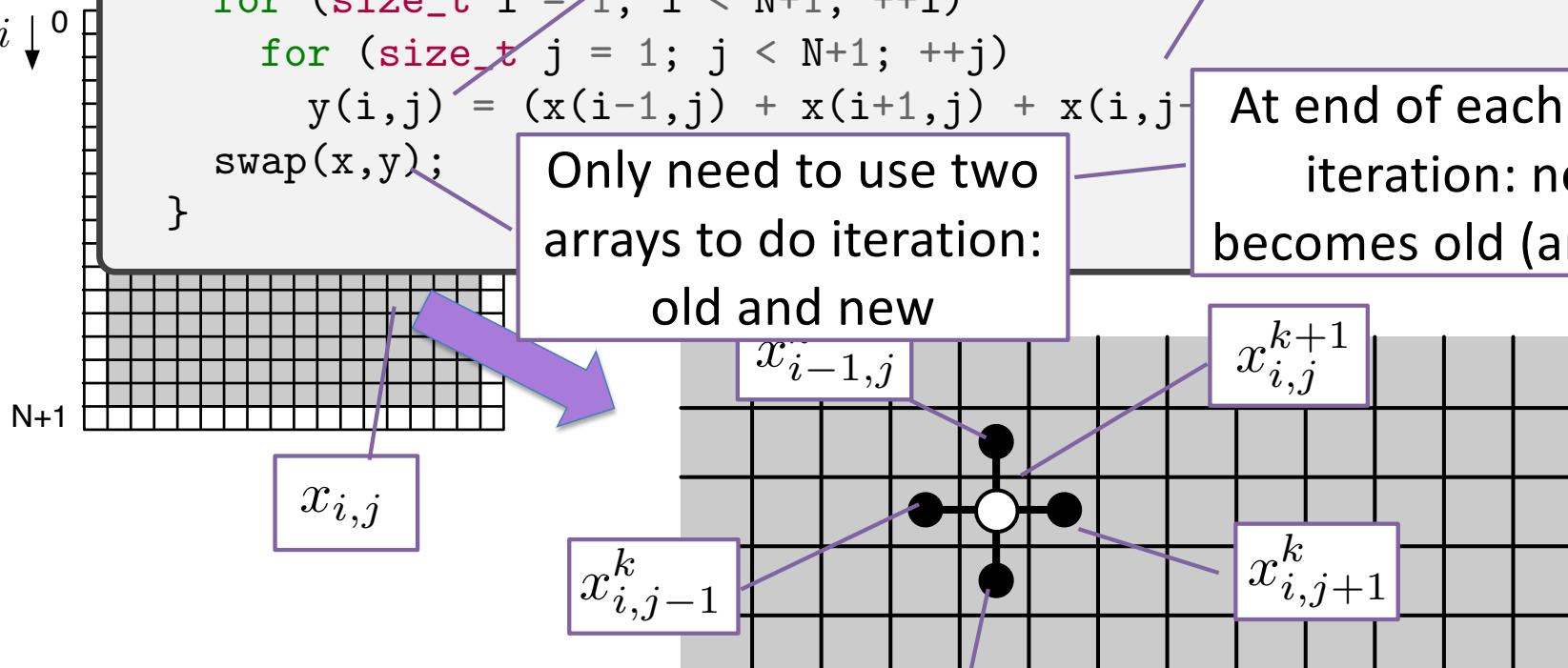
```
while (! converged())
    for (size_t i = 1; i < N+1; ++i)
        for (size_t j = 1; j < N+1; ++j)
            y(i,j) = (x(i-1,j) + x(i+1,j) + x(i,j))
            swap(x,y);
    }
```

Approximation at iteration k+1

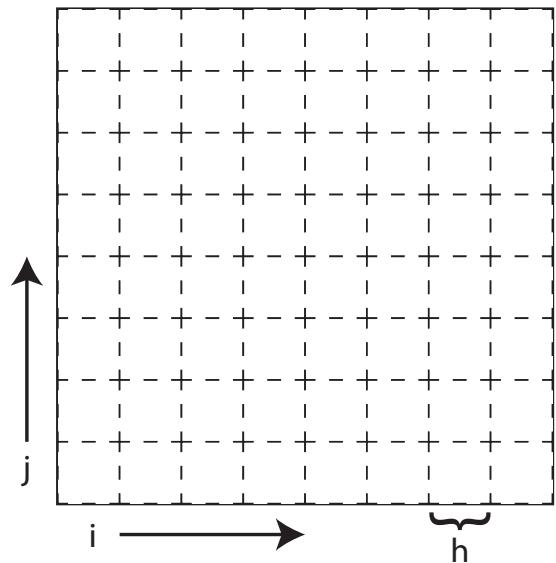
Average of approximation at iteration k

Only need to use two arrays to do iteration:
old and new

At end of each outer iteration: new becomes old (and v.v.)



Discretized



- Del operator $\nabla\phi = \frac{\partial\phi}{\partial x} + \frac{\partial\phi}{\partial y}$
 $\nabla^2\phi = \frac{\partial^2\phi}{\partial x^2} + \frac{\partial^2\phi}{\partial y^2}$
- Finite difference approximation to derivative

$$\begin{aligned}\frac{dx}{dt}(t_0) &\approx \frac{x(t_0+h)-x(t_0)}{h} \\ \frac{d^2x}{dt^2}(t_0) &\approx \frac{\frac{dx}{dt}(t_0+h)-\frac{dx}{dt}(t_0)}{h} \\ &= \frac{x(t_0+h+h)-x(t_0+h)-x(t_0+h)+x(t_0)}{h^2} \\ &= \frac{x(t_0+2h)-2x(t_0+h)+x(t_0)}{h^2} \\ &= \frac{x(t_0+h)-2x(t_0)+x(t_0-h)}{h^2}\end{aligned}$$

- Finite difference approximation to del

$$\frac{\phi_{i+1,j} + \phi_{i-1,j} + \phi_{i,j-1} + \phi_{i,j+1} - 4\phi_{i,k}}{h^2} = 0$$

Matrix Formulation

- Lexicographically order unknowns (note some will be boundary values)

$$\frac{x_{i+1} + x_{i-1} + x_{i+N} + x_{i-N} - 4x_i}{h^2} = 0$$

- Formulate as a matrix problem:

$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \dots & -1 & & \\ -1 & \ddots & \ddots & \ddots & \ddots & \\ \vdots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & \vdots \\ & \ddots & \ddots & \ddots & \ddots & -1 \\ & & -1 & \dots & -1 & 4 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ \vdots \end{bmatrix}$$

Linear System Solution

```
void multiply(const Matrix& A, const Matrix& B, Matrix& C) {
    for (size_t i = 0; i < A.num_rows(); ++i) {
        for (size_t j = 0; j < B.num_cols(); ++j) {
            for (size_t k = 0; k < A.num_cols(); ++k) {
                C(i, j) += A(i, k) * B(k, j);
            }
        }
    }
}
```

Work Smarter!
Don't multiply and
add zero to zero

Multiplying and
adding zero to zero

Matrix-matrix
product is kernel
operation

What happens with
the Laplacian
matrix?

Solution?

```
void multiply(const Matrix& A, const Matrix& B, Matrix& C) {  
    for (size_t i = 0; i < A.num_rows(); ++i) {  
        for (size_t j = 0; j < B.num_cols(); ++j) {  
            for (size_t k = 0; k < A.num_cols(); ++k) {  
                if(A(i,k) != 0.0 && B(k,j) != 0) {  
                    C(i, j) += A(i, k) * B(k, j);  
                }  
            }  
        }  
    }  
}
```

Avoid zeros

But we still touch every element

And that's what expensive

Solution?

```
void multiply(const Matrix& A, const Matrix& B, Matrix& C) {  
    for (size_t i = 0; i < A.num_rows(); ++i) {  
        for (size_t j = 0; j < B.num_cols(); ++j) {  
            for (size_t k = 0; k < A.num_cols(); ++k) {  
                if(A(i,k) != 0.0 && B(k,j) != 0) {  
                    C(i, j) += A(i, k) * B(k, j);  
                }  
            }  
        }  
    }  
}
```

We need to avoid zeros

Without looking to see if there is a zero

Solution: Sparse Matrices

In order to
avoid zeros

Don't store
zeros

A zero is a
null op

Use data structures
and algorithms
accordingly

$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \cdots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \ddots & \ddots & \ddots & \ddots & -1 & -1 \\ -1 & \cdots & -1 & -1 & 4 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ \vdots \end{bmatrix}$$

Sparse matrix
techniques

Solving Sparse Systems

- Work only with non-zeros
- Direct methods
 - Perform LU factorization on sparse matrix
 - Create non-zeros during elimination process
 - Pre-order (using heuristics) to minimize the amount of fill
 - Very sequential
 - Fill can be quite significant
- Iterative methods
 - Successively create better approximations to x
 - Relaxation methods (e.g., Jacobi) – very very very slow to converge
 - Krylov subspace methods (e.g., conjugate gradient)
 - Good preconditioning often required

$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \cdots & -1 \\ -1 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & -1 \\ & \ddots & \ddots & \ddots & \ddots & \ddots & -1 \\ & & -1 & \cdots & -1 & 4 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ \vdots \end{bmatrix}$$

A zero
here

Turns into
a non-zero

Need to create
new space (fill)

Conjugate Gradient Algorithm

Initial $r^{(0)} = b - Ax^{(0)}$

For $i=1, 2, \dots$

solve $Mz^{(i-1)} = r^{(i-1)}$

$\rho_{i-1} = r^{(i-1)T} z^{(i-1)}$

If $i=1$

$p^{(1)} = z^{(0)}$

Else

$\beta_{i-1} = \rho_{i-1} / \rho_{i-2}$

$p^{(i)} = z^{(i-1)} + \beta_{i-1} p^{(i-1)}$

Endif

$q^{(i)} = Ap^{(i)}$

$\alpha_i = \rho_{i-1} / p^{(i)T} q^{(i)}$

$x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}$

$r^{(i)} = r^{(i-1)} - \alpha_i q^{(i)}$

Check convergence

end

Key operation

```
mult(A, scaled(x, -1.0), b, r);
while (! iter.finished(r)) {
    solve(M, r, z);
    rho = dot_conj(r, z);

    if ( iter.first() )
        copy(z, p);
    else {
        beta = rho / rho_1;
        add(z, scaled(p, beta), p);
    }
    mult(A, p, q);
    alpha = rho / dot_conj(p, q);
    add(x, scaled(p, alpha), x);
    add(r, scaled(q, -alpha), r);
    rho_1 = rho;
    ++iter;
}
```

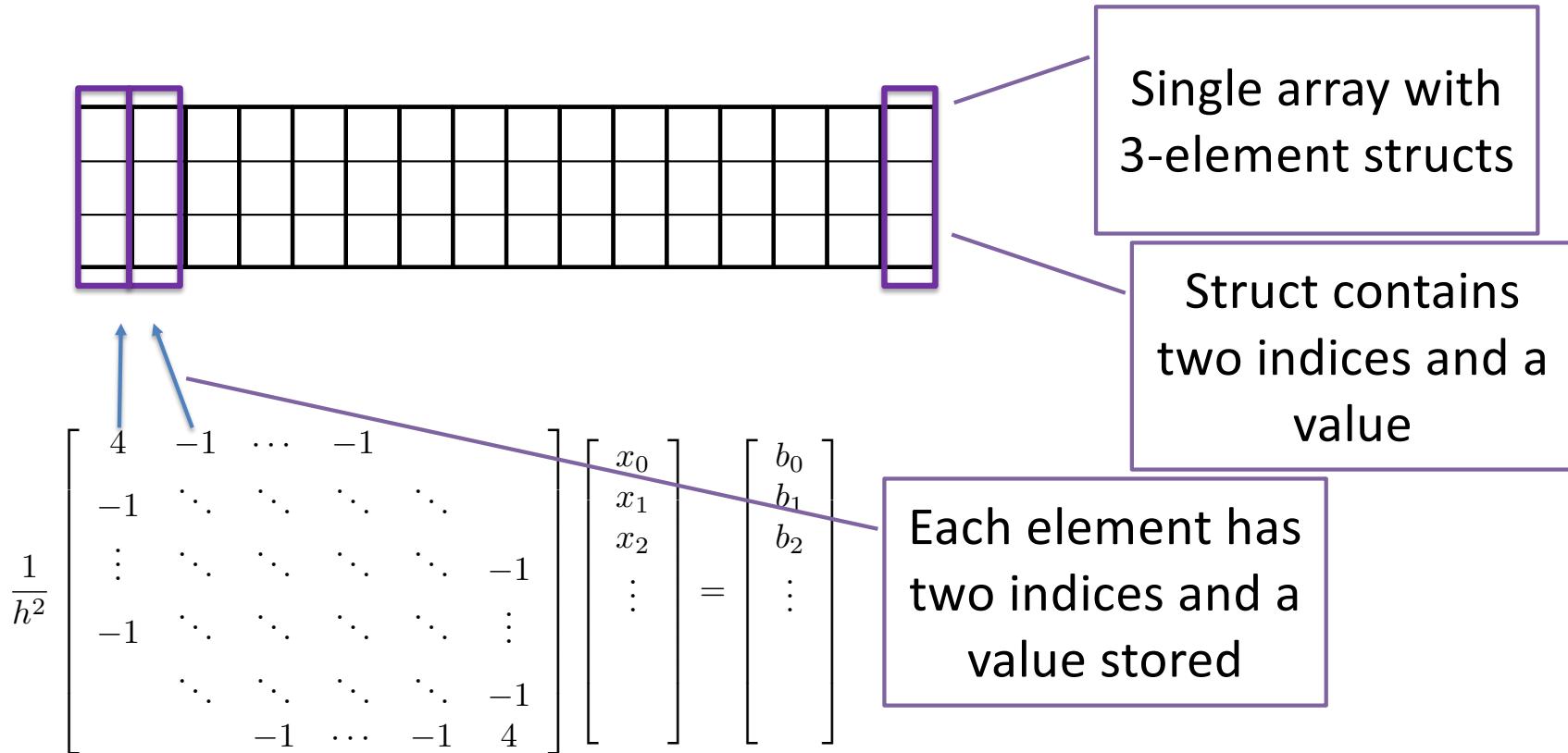
Sparse Storage

- A matrix is map from two indices to a value

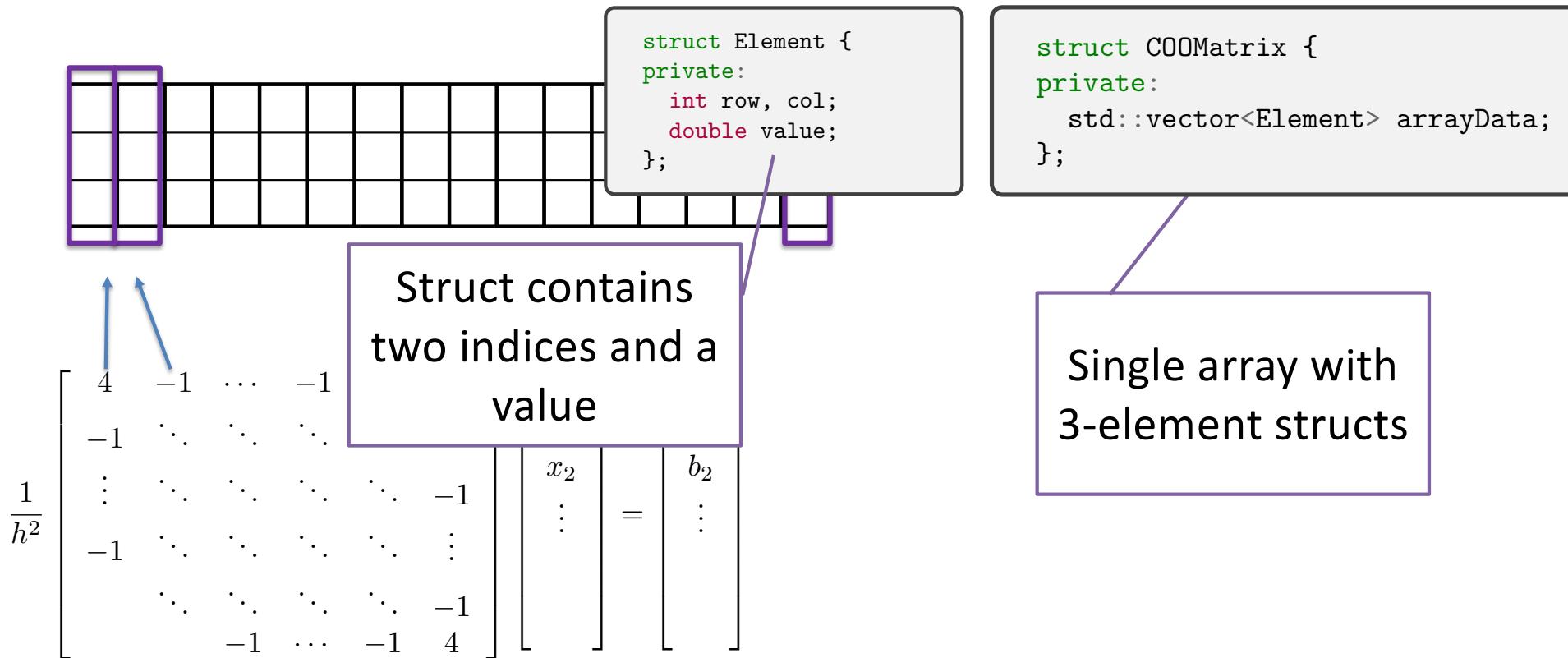
$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \cdots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & \vdots \\ & \ddots & \ddots & \ddots & \ddots & -1 \\ & & -1 & \cdots & -1 & 4 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \end{bmatrix}$$

- So if we want to store just elements that are not zero (the “non-zeros”)
- We need to store the two indices and the value

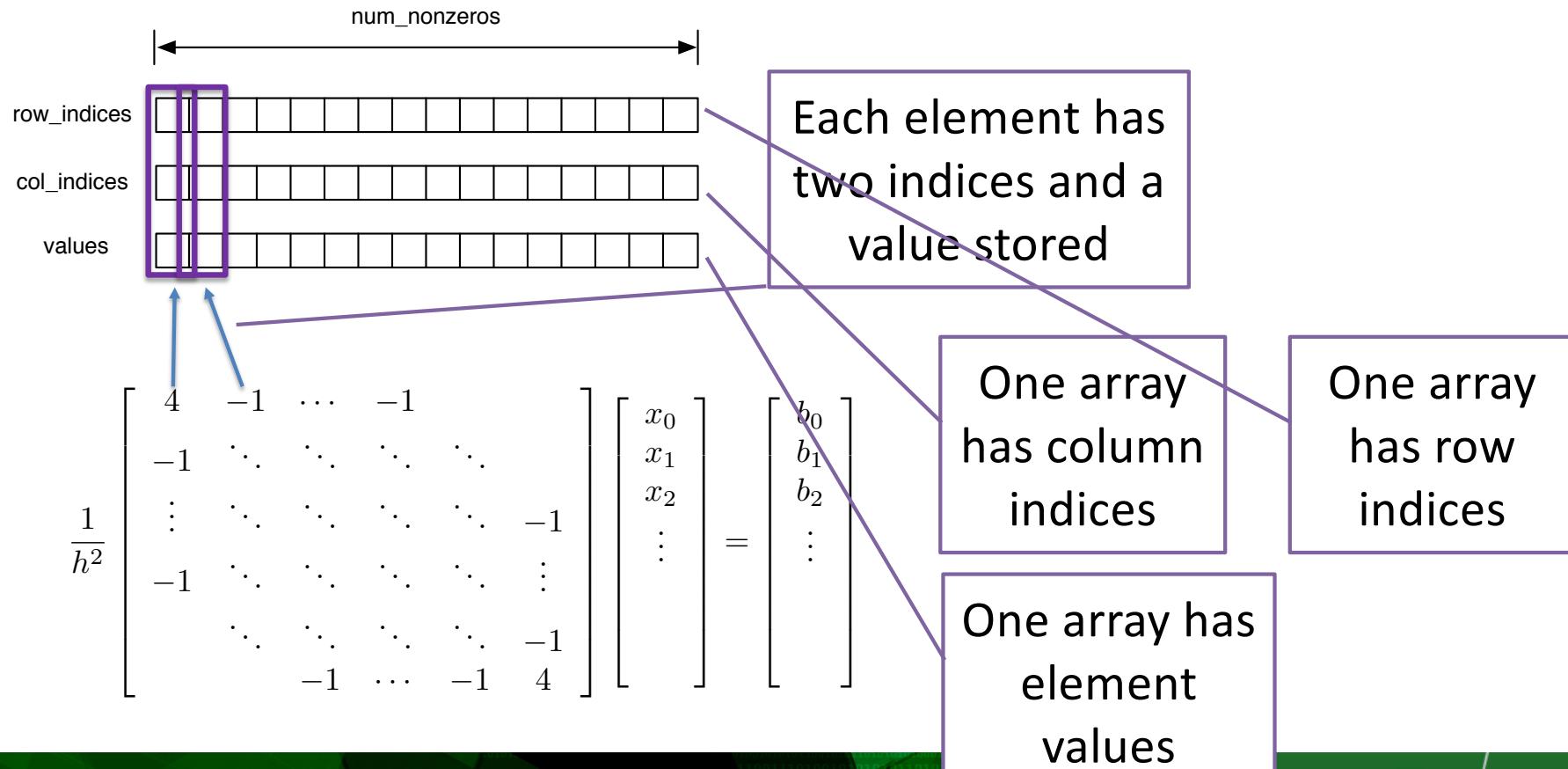
Coordinate Storage (Array of Structs)



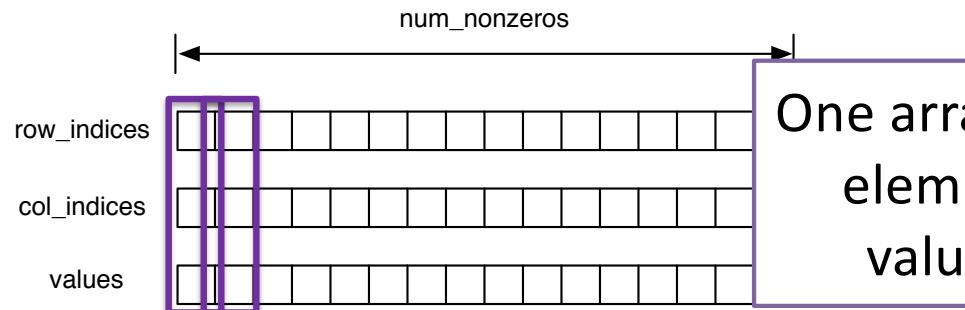
Coordinate Storage (Array of Structs)



Coordinate Storage (Struct of Arrays)



Coordinate Storage (Struct of Arrays)



```
struct COOMatrix {  
private:  
    std::vector<size_t> row_indices_;  
    std::vector<size_t> col_indices_;  
    std::vector<double> values_;  
};
```

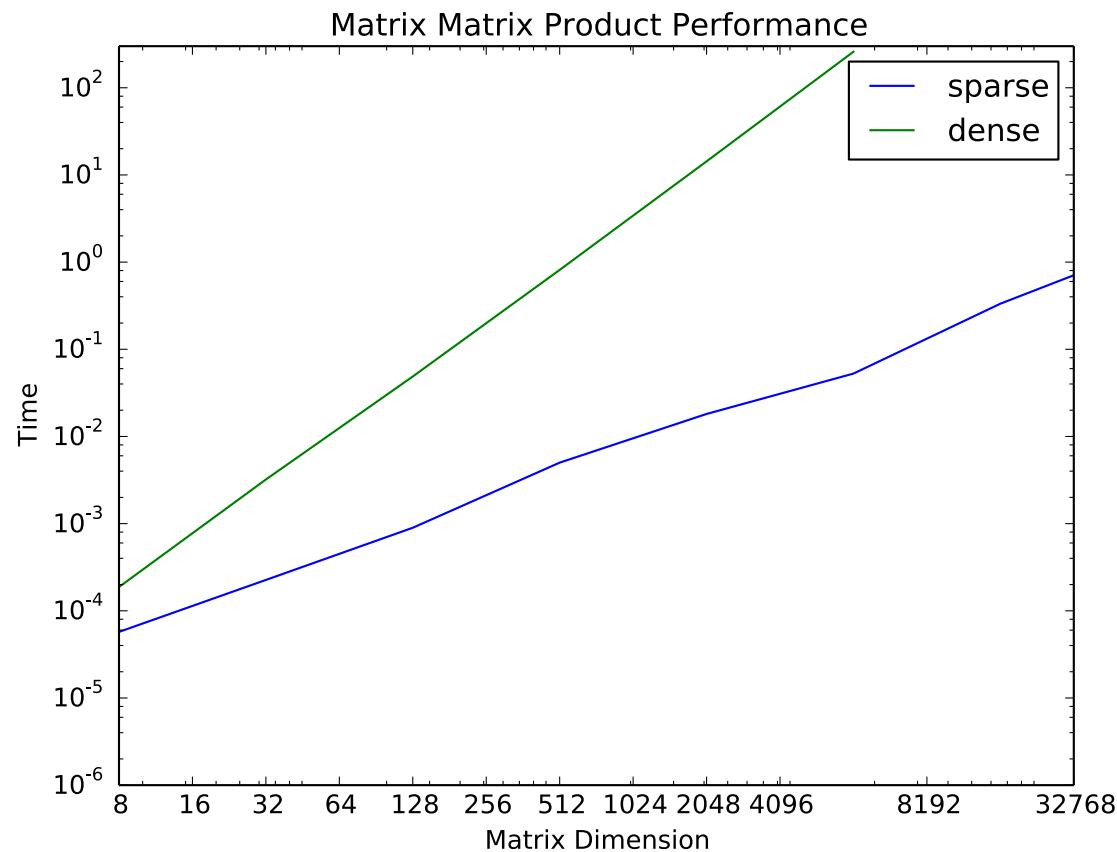
$$\frac{1}{h^2} \begin{bmatrix} 4 & -1 & \cdots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots \\ \vdots & \ddots & \ddots & \ddots & \ddots & -1 \\ -1 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \ddots & \ddots & \ddots & \ddots & -1 & \\ -1 & \cdots & -1 & 4 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \end{bmatrix}$$

One array has column indices

One array has row indices

Conventional Wisdom:
Struct of Arrays is faster

Performance Comparison



NORTHWEST INSTITUTE for ADVANCED COMPUTING

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Partially Operated by Battelle
for the U.S. Department of Energy

W
UNIVERSITY OF
WASHINGTON

What's the Catch?

```
class Matrix {  
public:  
    Matrix(size_t M, size_t N) : num_rows_(M), num_cols_(N), storage_(num_rows_ * num_cols_) {}  
  
    double& operator()(size_t i, size_t j) { return storage_[i * num_cols_ + j]; }  
    const double& operator()(size_t i, size_t j) const { return storage_[i * num_cols_ + j]; }  
  
    size_t num_rows() const { return num_rows_; }  
    size_t num_cols() const { return num_cols_; }  
  
private:  
    size_t num_rows_, num_cols_;  
    std::vector<double> storage_;  
};
```

In fact, it's a reference, so we can modify it

Provide indices, get back value

In constant time

Uh...

```
class COOMatrix {  
public:  
    COOMatrix(size_t M, size_t N) : num_rows_(M),  
        num_cols_(N)  
    {  
        row_indices_.resize(M);  
        col_indices_.resize(N);  
        storage_.reserve(M * N);  
    }  
  
    size_t num_rows() const { return num_rows_; }  
    size_t num_cols() const { return num_cols_; }  
  
private:  
    size_t num_rows_, num_cols_;  
    std::vector<size_t> row_indices_, col_indices_;  
    std::vector<double> storage_;  
};
```

How do we get
to a value (in
constant time)?

We can't

Next Problem

```
void matvec(const Matrix& A, const Vector& x, Vector& y) {  
    for (size_t i = 0; i < A.num_rows(); ++i) {  
        for (size_t j = 0; j < A.num_cols(); ++j) {  
            y(i) += A(i, j) * x(j);  
        }  
    }  
}
```

Nice external
function using
operator()()

```
void matvec(const COOMatrix& A, const Vector& x, Vector& y) {  
    // ??  
}
```

No operator()()
no external
function

Coordinate Matvec

```
void matvec(const Matrix& A, const Vector& x, Vector& y) {
    for (size_t i = 0; i < A.num_rows(); ++i) {
        for (size_t j = 0; j < A.num_cols(); ++j) {
            y(i) += A(i, j) * x(j);
        }
    }
}
```

This is the
row index

This is the
value

This is the
column
index

Coordinate Matvec

```
void matvec(const Matrix& A, const Vector& x, Vector& y) {  
    for (size_t i = 0; i < A.num_rows(); ++i) {  
        for (size_t j = 0; j < A.num_cols(); ++j) {  
            y(i) = A(i, j) * x(j);  
        }  
    }  
}
```

Index into y with row index

Multiply by the corresponding value

Index into x with column index

We have these three things in coordinate format

Coordinate Matrix Mat Vec

```
class COOMatrix {  
public:  
    COOMatrix(size_t M, size_t N) : num_rows_(M), num_cols_(N)  
  
    void matvec(const Vector& x, Vector& y) const {  
        for (size_t k = 0; k < storage_.size(); ++k) {  
            y(row_indices_[k]) += storage_[k] * x(col_indices[k]);  
        }  
    }  
  
private:  
    int num_rows_;  
    std::vector<int> row_indices_;  
    std::vector<double> storage_;  
};
```

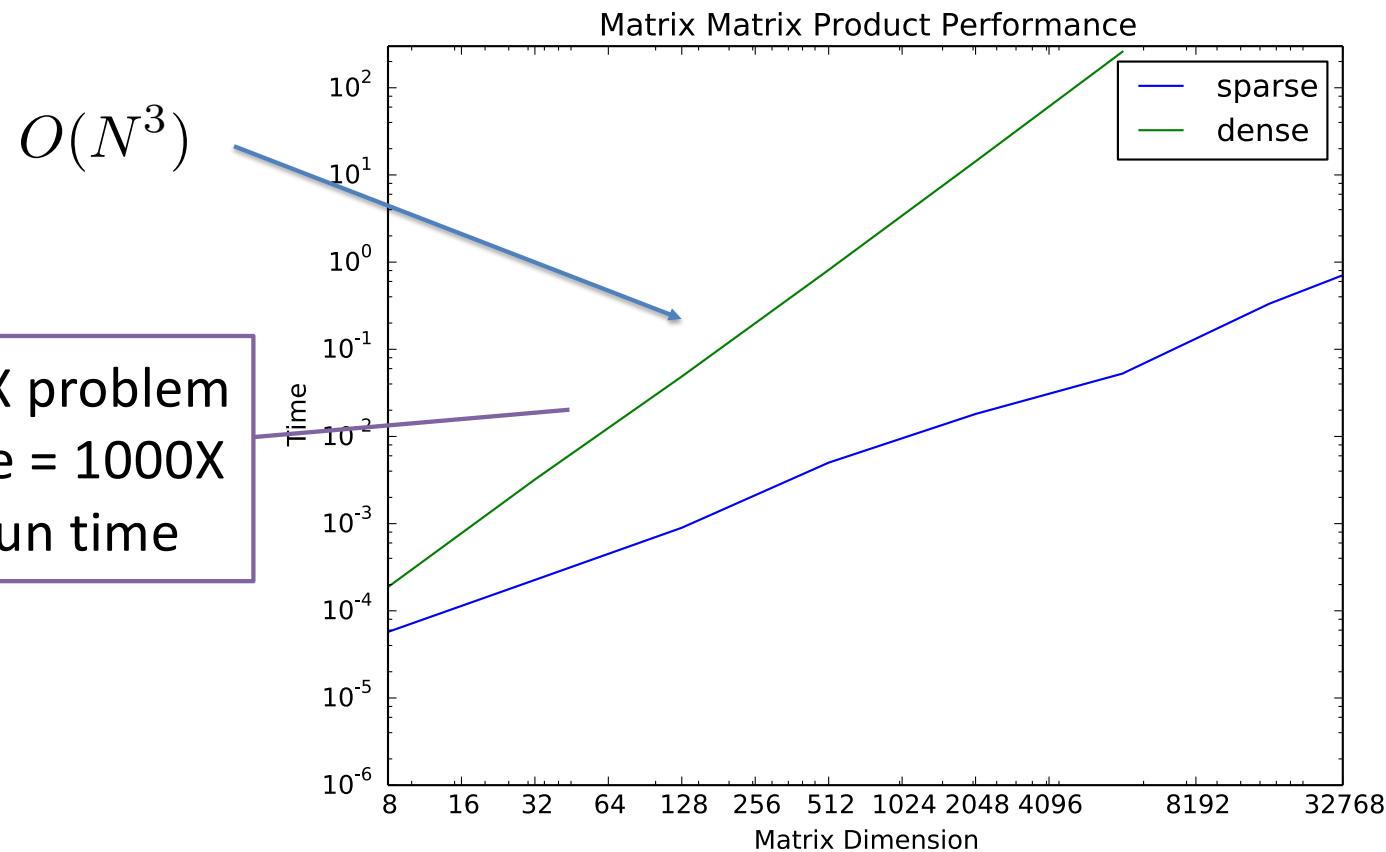
Meditate on
this

Index into y
with row
index

Multiply by
corresponding
value

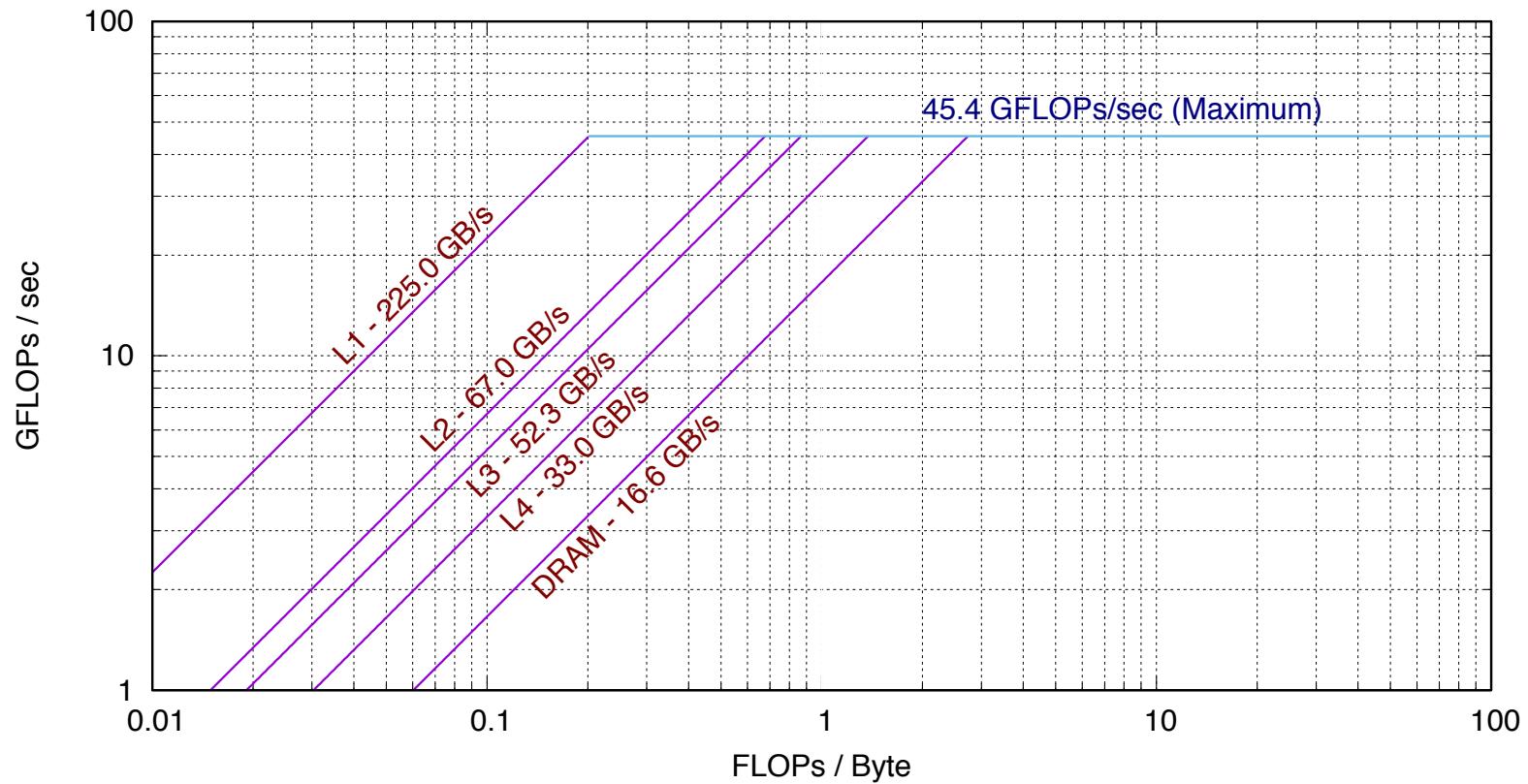
Index into x
with column
index

Performance Comparison



Roofline

Empirical Roofline Graph (Results.WE31821/Run.004)



Numerical Intensity

```
void matvec(const Vector& x, Vector& y) const {
    for (size_type k = 0; k < arrayData.size(); ++k) {
        y(rowIndices[k]) += arrayData[k] * x(rowIndices[k]);
    }
}
```

Two flops

Three doubles + 2 ints
= 32 bytes? (36 bytes?)

2 NNZ Flops

5N

NNZ doubles
+2 NNZ indexes
+2N doubles

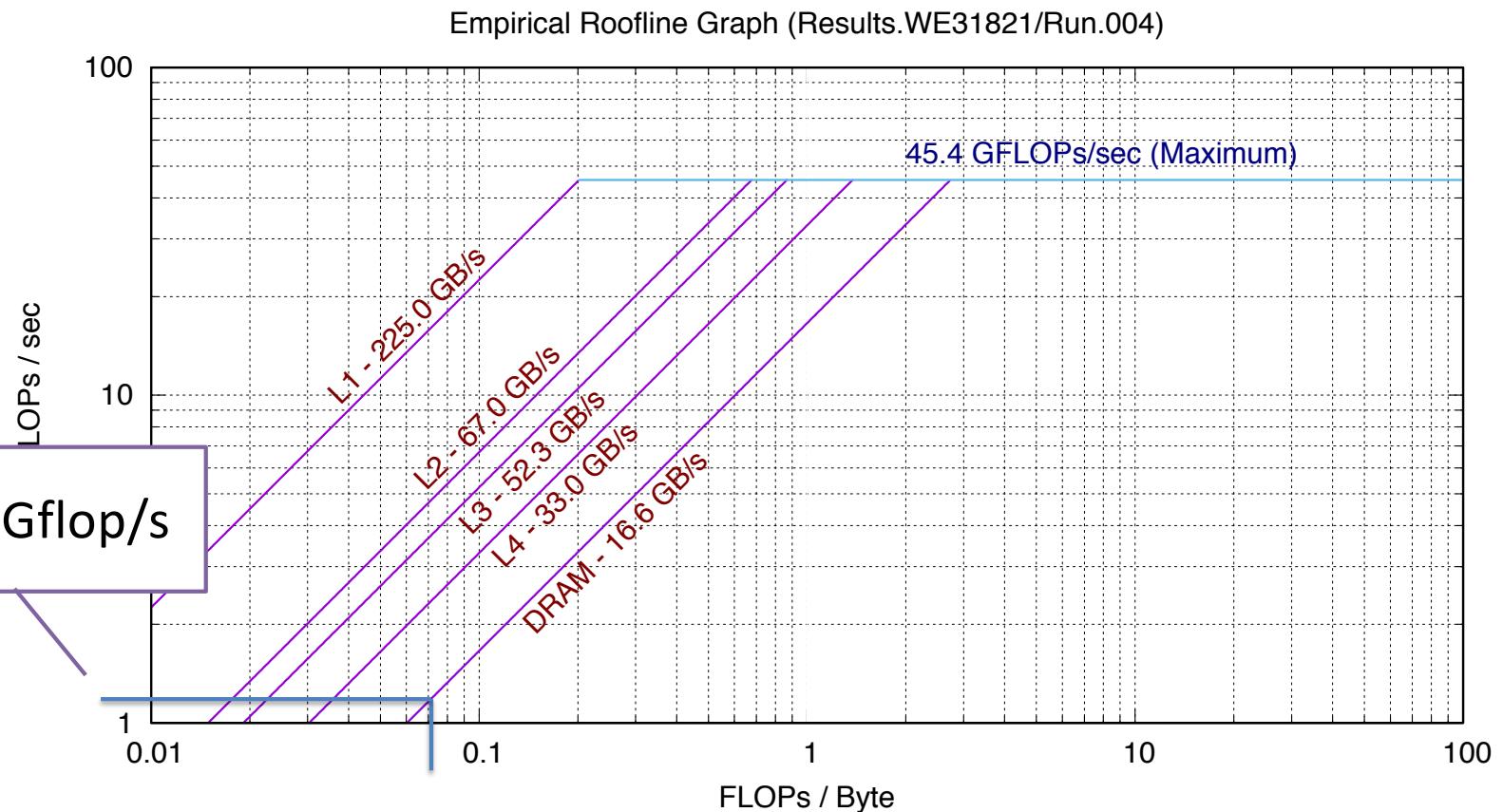
10N
Flops

7N doubles =
56 bytes

$\frac{1}{14}$ Flop
 $\frac{1}{14}$ byte

10N indexes =
40, 80 bytes

Measured



Coordinate Storage

```
class COOMatrix {  
public:  
    COOMatrix(size_t M, size_t N) : num_rows_(M), num_cols_(N) {}  
  
    void matvec(const Vector& x, Vector& y) const {  
        for (size_t k = 0; k < storage_.size(); ++k) {  
            y(row_indices_[k]) += storage_[k] * x(col_indices[k]);  
        }  
    }  
  
private:  
    int num_rows, num_cols;  
    std::vector<size_t> row_indices_, col_indices_;  
    std::vector<double> storage_;  
};
```

How do we initialize storage_?

In fact, how do we create a sparse matrix?

Filling a Sparse Matrix

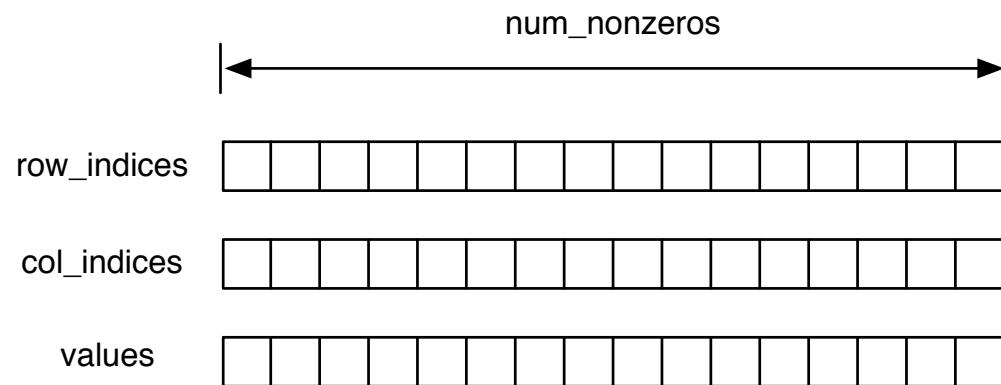
```
class COOMatrix {  
public:  
    COOMatrix(size_t M, size_t N) : num_rows_(M), num_col  
  
        void insert(sizet i, size_t j, double val) {  
            row_indices_.push_back(i);  
            col_indices_.push_back(j);  
            storage_.push_back(val);  
        }  
  
private:  
    size_t num_rows_, num_cols_;  
    std::vector<size_t> row_indices_, col_indices_;  
    std::vector<double> storage_;  
};
```

Often treated like variable initialization

Matrix is filled with something when created

Can also append elements (no ordering required)

There is Still Too Much Work



Coordinate Storage

- Storing elements in random order

$\begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 \\ 9 & 2 \\ 7 & 6 \end{bmatrix}$	row_indices	<table border="1"><tr><td>1</td><td>4</td><td>4</td><td>2</td><td>3</td><td>2</td><td>3</td><td>1</td><td>1</td></tr></table>	1	4	4	2	3	2	3	1	1
1	4	4	2	3	2	3	1	1			
	col_indices	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>2</td><td>1</td><td>3</td><td>3</td><td>4</td><td>2</td></tr></table>	1	2	4	2	1	3	3	4	2
1	2	4	2	1	3	3	4	2			
	values	<table border="1"><tr><td>3</td><td>7</td><td>6</td><td>1</td><td>9</td><td>5</td><td>2</td><td>4</td><td>1</td></tr></table>	3	7	6	1	9	5	2	4	1
3	7	6	1	9	5	2	4	1			

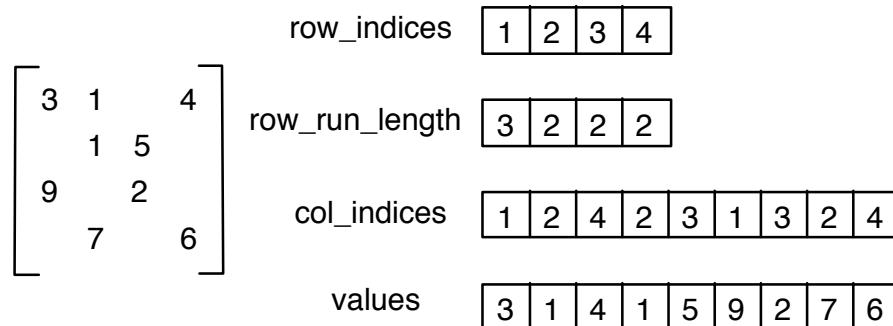
- Storing elements in sorted order (along rows)

$\begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 \\ 9 & 2 \\ 7 & 6 \end{bmatrix}$	row_indices	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>2</td><td>2</td><td>3</td><td>3</td><td>4</td><td>4</td></tr></table>	1	1	1	2	2	3	3	4	4
1	1	1	2	2	3	3	4	4			
	col_indices	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>2</td><td>3</td><td>1</td><td>3</td><td>2</td><td>4</td></tr></table>	1	2	4	2	3	1	3	2	4
1	2	4	2	3	1	3	2	4			
	values	<table border="1"><tr><td>3</td><td>1</td><td>4</td><td>1</td><td>5</td><td>9</td><td>2</td><td>7</td><td>6</td></tr></table>	3	1	4	1	5	9	2	7	6
3	1	4	1	5	9	2	7	6			

- Note there is significant redundancy

Run Length Encoding

- Row indices can be compressed with run-length encoding



```
size_t row_ptr = 0;
for (size_t i = 0; i < num_rows_; ++i) {
    for (size_t j = row_ptr; j < row_ptr + row_run_length[i]; ++j)
        y[row_indices[i]] += values[j] * x[col_indices[j]];
    row_ptr = row_ptr + row_ptr + row_run_length[i];
}
```

Compressed Sparse Row

- Store the running “row_ptr” value rather than computing it

	row_run_length	<table border="1"><tr><td>3</td><td>2</td><td>2</td><td>2</td></tr></table>	3	2	2	2													
3	2	2	2																
<table border="1"><tr><td>3</td><td>1</td><td>4</td></tr><tr><td>1</td><td>5</td><td></td></tr><tr><td>9</td><td>2</td><td></td></tr><tr><td>7</td><td>6</td><td></td></tr></table>	3	1	4	1	5		9	2		7	6		row_ptr	<table border="1"><tr><td>1</td><td>4</td><td>6</td><td>8</td><td>10</td></tr></table>	1	4	6	8	10
3	1	4																	
1	5																		
9	2																		
7	6																		
1	4	6	8	10															
	col_indices	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>2</td><td>3</td><td>1</td><td>3</td><td>2</td><td>4</td></tr></table>	1	2	4	2	3	1	3	2	4								
1	2	4	2	3	1	3	2	4											
	values	<table border="1"><tr><td>3</td><td>1</td><td>4</td><td>1</td><td>5</td><td>9</td><td>2</td><td>7</td><td>6</td></tr></table>	3	1	4	1	5	9	2	7	6								
3	1	4	1	5	9	2	7	6											

```
for (size_t i = 0; i < num_rows_; ++i) {
    for (size_t j = row_ptr[i]; j < row_ptr[i+1]; ++j)
        y[i] += values[j] * x[col_indices[j]];
}
```

Compressed Sparse Row

$\begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 \\ 9 & 2 \\ 7 & 6 \end{bmatrix}$	row_indices	<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td></tr></table>	1	2	3	4					
1	2	3	4								
	row_run_length	<table border="1"><tr><td>3</td><td>2</td><td>2</td><td>2</td></tr></table>	3	2	2	2					
3	2	2	2								
	col_indices	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>2</td><td>3</td><td>1</td><td>3</td><td>2</td><td>4</td></tr></table>	1	2	4	2	3	1	3	2	4
1	2	4	2	3	1	3	2	4			
	values	<table border="1"><tr><td>3</td><td>1</td><td>4</td><td>1</td><td>5</td><td>9</td><td>2</td><td>7</td><td>6</td></tr></table>	3	1	4	1	5	9	2	7	6
3	1	4	1	5	9	2	7	6			

Don't need
this

Starting
value

$\begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 \\ 9 & 2 \\ 7 & 6 \end{bmatrix}$	row_run_length	<table border="1"><tr><td>3</td><td>2</td><td>2</td><td>2</td></tr></table>	3	2	2	2					
3	2	2	2								
	row_ptr	<table border="1"><tr><td>1</td><td>4</td><td>6</td><td>8</td><td>10</td></tr></table>	1	4	6	8	10				
1	4	6	8	10							
	col_indices	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>2</td><td>3</td><td>1</td><td>3</td><td>2</td><td>4</td></tr></table>	1	2	4	2	3	1	3	2	4
1	2	4	2	3	1	3	2	4			
	values	<table border="1"><tr><td>3</td><td>1</td><td>4</td><td>1</td><td>5</td><td>9</td><td>2</td><td>7</td><td>6</td></tr></table>	3	1	4	1	5	9	2	7	6
3	1	4	1	5	9	2	7	6			

Instead of run
length values

Accumulate
them

Note extra
end value

CSR Implementation

```
class CSRMatrix {  
public:  
    CSRMatrix(size_t M, size_t N) : num_rows_(M), num_cols_(N), row_indices_(nu  
    size_t num_rows() const { return num_rows_; }  
    size_t num_cols() const { return num_cols_; }  
    size_t num_nonzeros() const { return storage_.size(); }  
  
private:  
    size_t num_rows_, num_cols_;  
    std::vector<size_t> row_indices_, col_indices_;  
    std::vector<double> storage_;  
};
```

Constructor

And
row_indices_

Note initial
value

Initialize
num_rows and
num_cols

Matrix size
accessors

Useful info for
sparse matrix

Private
implementation

CSR Implementation (Matrix Vector Multiply)

```
class CSRMatrix {  
  
public:  
    CSRMatrix(size_t M, size_t N) : num_rows_(M), num_cols_(N)  
    {}  
  
    void matvec(const Vector& x, Vector& y) const {  
        for (size_t i = 0; i < num_rows_; ++i) {  
            for (size_t j = row_indices_[i]; j < row_indices_[i+1]; ++j) {  
                y(i) += storage_[j] * x(col_indices_[j]);  
            }  
        }  
    }  
  
private:  
    size_t num_rows_, num_cols_;  
    std::vector<size_t> row_indices_, col_indices_;  
    std::vector<double> storage_;  
};
```

For each row

For each element in that row

Meditate on this

Row index

Matrix value

Column index

Building a CSR Matrix

```
class CSRMatrix {  
  
public:  
    void open_for_push_back() { is_open = true; }  
  
    void close_for_push_back() { is_open = false;  
        for (size_t i = 0; i < num_rows_; ++i) row_indices_[i+1] += row_indices_[i];  
        for (size_t i = num_rows_; i > 0; --i) row_indices_[i] = row_indices_[i-1];  
        row_indices_[0] = 0;  
    }  
  
    void push_back(size_t i, size_t j, double value) {  
        ++row_indices_[i];  
        col_indices_.push_back(j);  
        storage_.push_back(value);  
    }  
  
private:  
    bool is_open;  
    size_t num_rows_, num_cols_;  
    std::vector<size_t> row_indices_;  
    std::vector<double> storage_;  
};
```

When done pushing,
accumulate run lengths to
offsets

Should be
checked

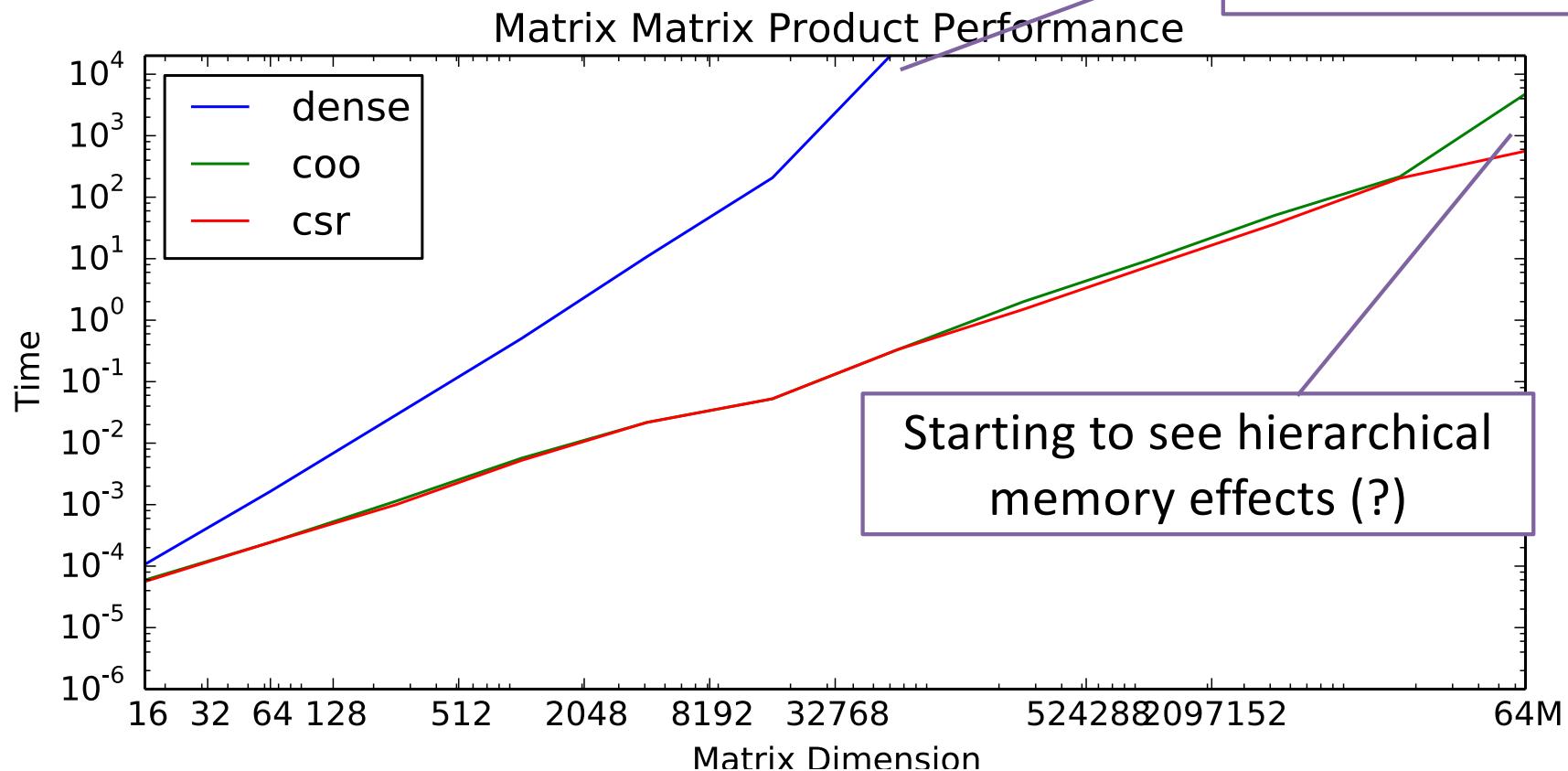
Push elements back
(similar to COO)

Accumulate run
row lengths

Push column
index and value

Rows **must** be
added in order and
contiguously

Performance



Review

- Explored variety of techniques for matching algorithm structure to hardware performance features (work smarter)
 - And we pushed this pretty far
- Strassen's algorithm (work way smarter)
- Sparse matrix representations and algorithms (don't do work you don't have to do)
- Get help



Last Chance for Questions Before we Leave the Sequential World



NORTHWEST INSTITUTE for ADVANCED COMPUTING

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

W
UNIVERSITY OF
WASHINGTON

Thank you!

NORTHWEST INSTITUTE for ADVANCED COMPUTING

AMATH 483/583 High-Performance Scientific Computing Spring 2019
University of Washington by Andrew Lumsdaine

Pacific Northwest
NATIONAL LABORATORY
Proudly Operated by Battelle
for the U.S. Department of Energy

W
UNIVERSITY OF
WASHINGTON

Creative Commons BY-NC-SA 4.0 License



© Andrew Lumsdaine, 2017-2019

Except where otherwise noted, this work is licensed under

<https://creativecommons.org/licenses/by-nc-sa/4.0/>

