



Concepts and Models of Parallel and Data-centric Programming

Parallel Algorithms III

Lecture, Summer 2020

Dr. Christian Terboven <terboven@itc.rwth-aachen.de>

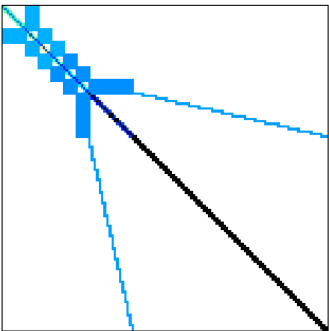
Outline

- 0. Organization
- 1. Foundations
- 2. Shared Memory
- 3. GPU Programming
- 4. Bulk-Synchronous Parallelism
- 5. Message Passing
- 6. Distributed Shared Memory
- 7. Parallel Algorithms**
 - a. Berkeley DWARFS
 - b. Dense Linear Algebra
 - c. Sparse Linear Algebra
 - d. Monte Carlo Methods
 - e. Graph Traversal
- 8. Parallel I/O
- 9. MapReduce
- 10. Apache Spark

Sparse Linear Algebra

Motivation

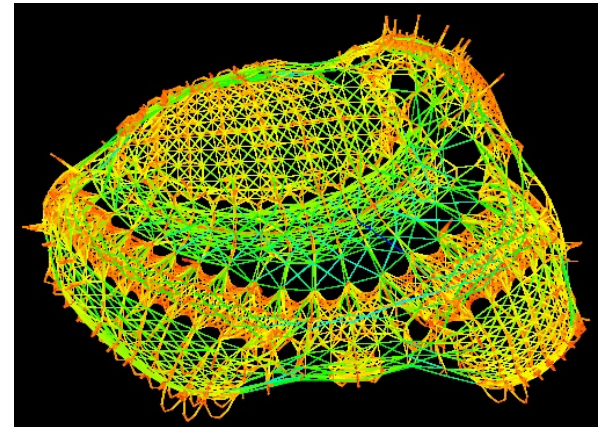
- Sparse Linear Algebra
 - Sparse Linear Equation Systems occur in many scientific disciplines.
 - Sparse matrix-vector multiplications (SpMxV) are the dominant part in many iterative solvers (like CG) for such systems.
 - number of non-zeros $\ll n \cdot n$



Beijing Botanical Garden

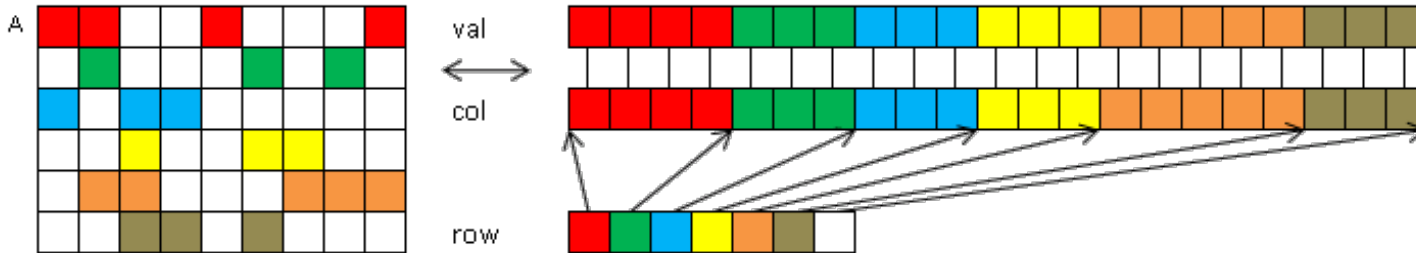
Upper right: Original building
Lower right: Model
Lower left: Matrix

(Source: Beijing Botanical Garden and University of Florida, Sparse Matrix Collection)



Representation of a sparse matrix

- CRS: Compressed Row Storage



- A concrete example

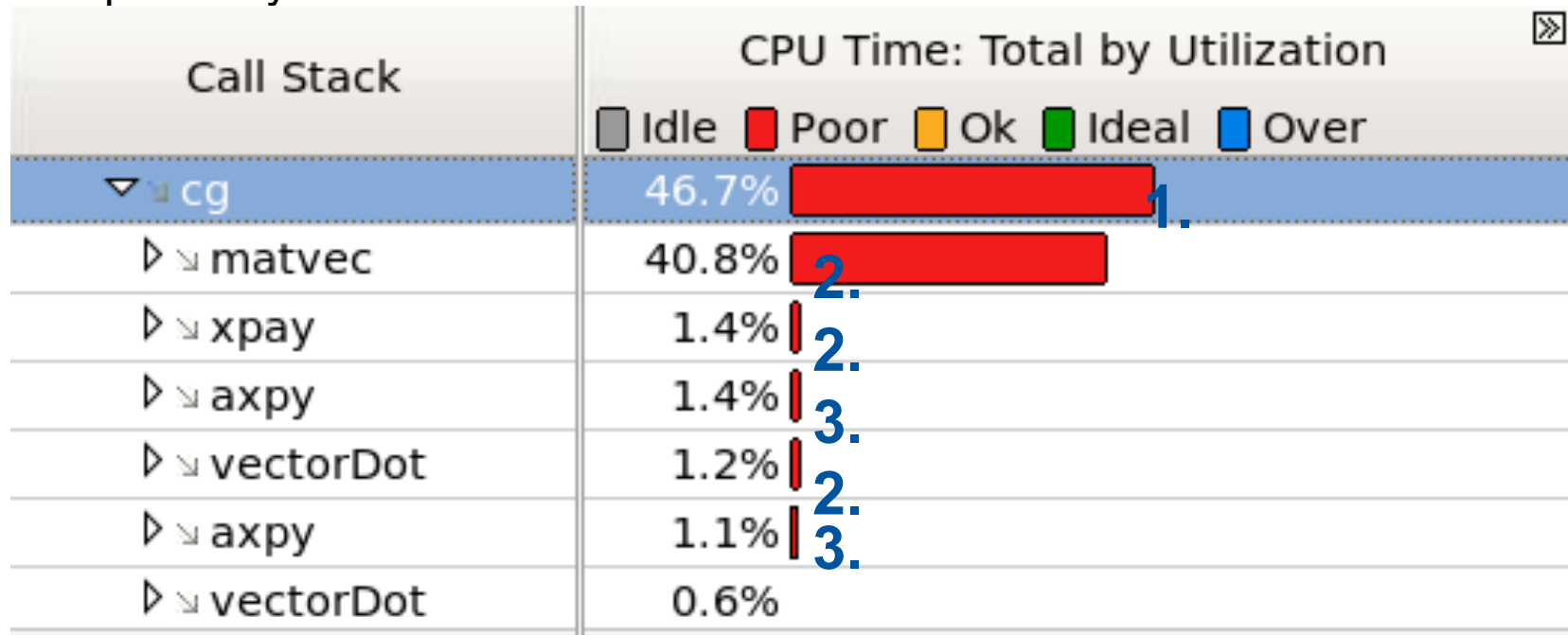
$$A = \begin{pmatrix} 10 & 0 & 0 & 12 & 0 \\ 0 & 0 & 11 & 0 & 13 \\ 0 & 16 & 0 & 0 & 0 \\ 0 & 0 & 11 & 0 & 13 \end{pmatrix}$$

Non-zero elements are labeled with their row and column indices: (0,0), (0,3), (1,2), (1,4), (2,1), (3,2), (3,4).

$$\begin{aligned} val &= \begin{pmatrix} 10 & 12 & 11 & 13 & 16 & 11 & 13 \\ (0,0) & (0,3) & (1,2) & (1,4) & (2,1) & (3,2) & (3,4) \end{pmatrix} \\ colInd &= \begin{pmatrix} 0 & 3 & 2 & 4 & 1 & 2 & 4 \\ (0) & & (1) & & (2) & (3) & \end{pmatrix} \\ rowPtr &= \begin{pmatrix} 0 & 2 & 4 & 5 & 7 \\ (0) & (1) & (2) & (3) & (4) \end{pmatrix} \end{aligned}$$

Case Study CG / 1

Hotspot analysis of the serial code:



Hotspots are:

1. matrix-vector multiplication
2. scaled vector additions
3. dot product

- SpMXV: sparse matrix vector multiplication
- Algorithm ($y = A * x$) adopted for the CRS format:

```
for i = 1, n
```

```
    y(i) = 0
```

```
    for j = row(i), row(i+1) - 1
```

```
        y(i) = y(i) + val(j) * x(col(j))
```

```
    end;
```

```
end;
```

- What is the intuitive approach to parallelize such a loop?

Case Study CG / 2

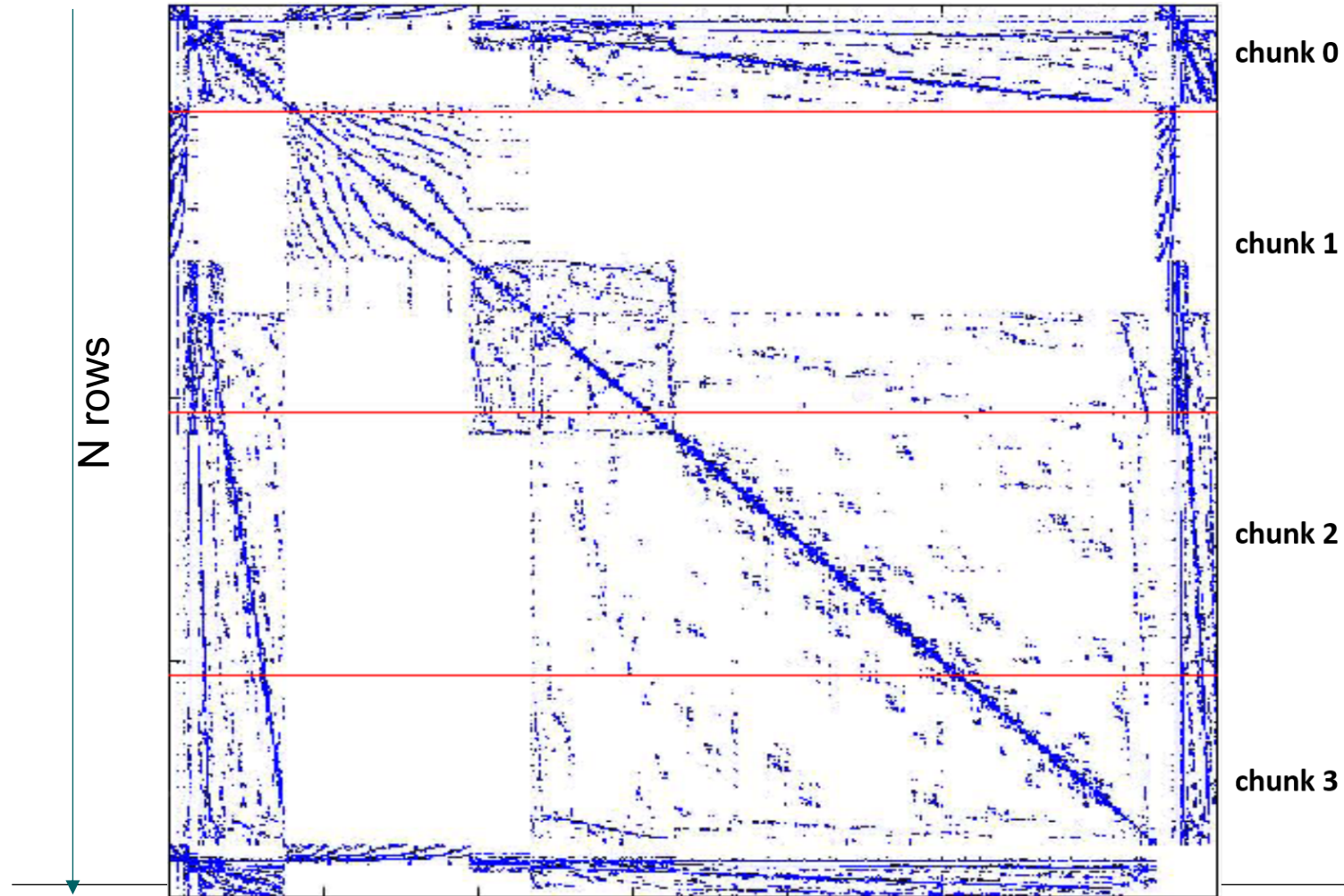
- Analyzing load imbalance in the concurrency view:

So.. Line	Source	CPU Time: Total by... Idle Poor Ok	Ove... and...
49	void matvec(const int n, const int nnz,		
50	int i,j;		
51	#pragma omp parallel for private(j)	22.462s	10.612s
52	for(i=0; i<n; i++){	0.050s	0s
53	y[i]=0;	0.060s	0s
54	for(j=ptr[i]; j<ptr[i+1]; j++){	1.741s	0s
55	y[i]+=value[j]*x[index[j]];	9.998s	0s

- 10 seconds out of ~35 seconds are overhead time
- other parallel regions which are called the same amount of time only produce 1 second of overhead

Matrix structure (spy plot)

- Chunking illustrated for four threads



Solution in OpenMP

- Manual computation of the work distribution for load balancing
- Exploitation of the load balancing in the computation:

```
#pragma omp parallel private(i) num_threads(tOptions->iNumThreads)
{
    for (i = start; i < end; i++) {
        double sum = 0.0;

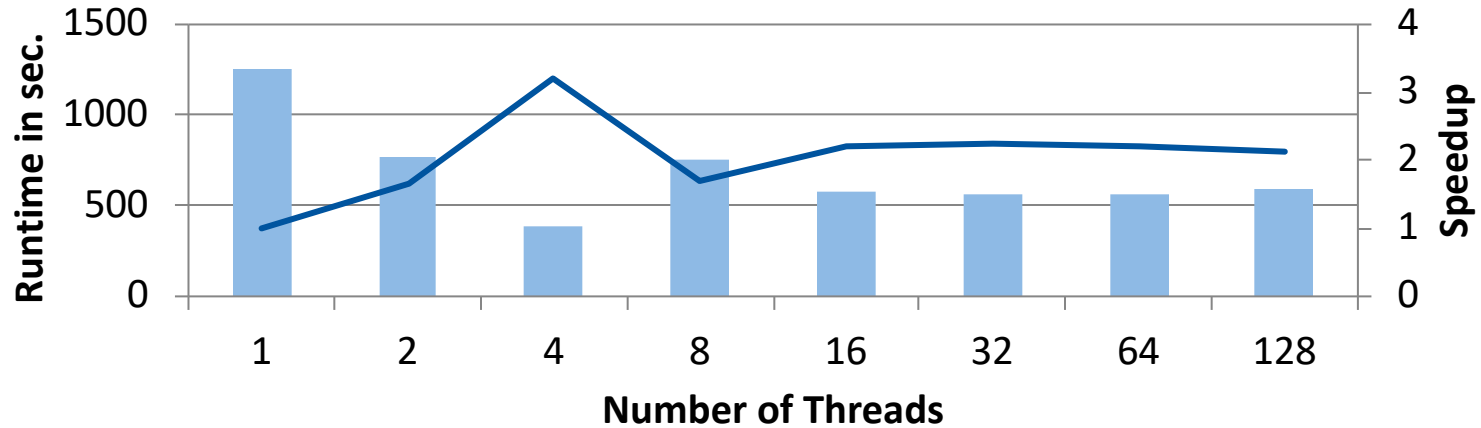
        int rowbeg = Arow[i]; int rowend = Arow[i+1];
        int nz;

        #pragma omp simd reduction(+:sum)
            for (nz = rowbeg; nz < rowend; nz++) {
                sum += Aval[nz] * x[ Acol[nz] ];
            }

        y[i] = sum;
    } }
```

SpMXV within CG

- Naive approach:



- Optimal approach:

