# Sparse Matrix Multiplication

Alok Ranjan

241010011

Department of Aerospace Engineering
IIT Kanpur

April 21, 2025

# Serial: Core Mathematics

For matrices $A$ and $B$, compute:

$$C_{ij} = \sum_{k=1}^{80} A_{ik} \times B_{kj}$$

## Basic Optimization

Skip when $A_{ik} = 0$:

$$C_{ij} = \sum_{\substack{k=1 \\ A_{ik} \neq 0}}^{80} A_{ik} \times B_{kj}$$

▶ Operations Saved: 70% (from 512,000 to 153,600)
▶ Key Idea: Check before multiplying

# Serial Code Explained

```
1  for (int i = 0; i < 80; i++) {         // For each
       output row
2      for (int k = 0; k < 80; k++) {      // For each
           element
3          if (A[i][k] != 0) {             //  zero-check
4              for (int j = 0; j < 80; j++) {
5                  C[i][j] += A[i][k] * B[k][j];
6  }}}}
```

▶ Loop Order: i→k→j optimal for row-major memory
▶ Why This Works: Skips entire inner loops when $A_{ik} = 0$

# OpenMP: Parallel Mathematics

Extend serial approach by distributing rows:

$$C_{\text{thread}} = \sum_{\text{assigned rows}} A_{\text{thread}} \times B$$

## Implementation

- ▶ Divide $A$'s rows equally among threads
- ▶ Each thread computes its portion independently

# OpenMP Code

```
1  #pragma omp parallel for  //  parallel directive
2  for (int i = 0; i < 80; i++) {
3      // Same serial logic:
4      if (A[i][k] != 0) {          //  zero-check
5          for (int j = 0; j < 80; j++) {
6              C[i][j] += A[i][k] * B[k][j];
7  }}}
```

- ▶ How Distribution Works:
  - ▶ Thread 0: rows 0-19
  - ▶ Thread 1: rows 20-39
  - ▶ etc.
- ▶ Why Safe: Threads write to different $C[i]$ rows

# MPI: Data Distribution Mechanics

**Scattering Matrix A**:
- ▶ `MPI_Scatter` splits A by rows
- ▶ Example (4 processes):
  - ▶ Process 0: rows 0-19
  - ▶ Process 1: rows 20-39
  - ▶ Process 2: rows 40-59
  - ▶ Process 3: rows 60-79
- ▶ Each gets `localA[20][80]`

**Broadcasting Matrix B**:
- ▶ `MPI_Bcast` sends full B to all
- ▶ Each process gets `B[80][80]`
- ▶ Why? All need full B for multiplication

# MPI: Computation & Gathering

**Local Computation**:

- ▶ Each process calculates:

  $$localC = localA \times B$$

- ▶ Uses same zero-skipping as serial

- ▶ Works on its 20-row chunk

**Gathering Results**:

- ▶ `MPI_Gather` collects `localC` chunks
- ▶ Reassembles full C at Rank 0
- ▶ Like puzzle pieces coming together

# MPI Code

```
1   // Rank 0 divides A
2   MPI_Scatter(A, 20*80, MPI_INT, localA, 20*80, MPI_INT, 0,
        MPI_COMM_WORLD);
3
4   // All get full B
5   MPI_Bcast(B, 80*80, MPI_INT, 0, MPI_COMM_WORLD);
6
7   // local computation
8   for (i = 0; i < 20; i++) {
9       if (localA[i][k] != 0) {  // zero-check
10          for (j = 0; j < 80; j++) {
11              localC[i][j] += localA[i][k] * B[k][j];
12  }}}
13
14  // Rank 0 collects results
15  MPI_Gather(localC, 20*80, MPI_INT, C, 20*80, MPI_INT, 0,
        MPI_COMM_WORLD);
```

# Execution Time Comparison

| Implementation | 2 Threads | 4 Threads | 8 Threads |
|---|---|---|---|
| Serial | | 0.0066 s | |
| OpenMP | 0.002713 s | 0.001528 s | 0.004866 s |
| MPI | 0.002752 s | 0.002797 s | 0.000853 s |

Table: Execution times for different parallel implementations

- ▶ MPI with 8 threads showed the best performance (0.000853 s)
- ▶ OpenMP performance degraded with 8 threads (0.004866 s)
- ▶ Serial implementation was the slowest (0.0066 s)

# Parallel and Serial Ax=b linear Solvers
Implementing Jacobi, Gradient Descent and PETSc

# Jacobi Method - Algorithm

### Mathematical Formulation
Update each variable using previous iteration values:

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left( b_i - \sum_{j \neq i} a_{ij} x_j^{(k)} \right)$$

---

**Algorithm 1** Jacobi Iteration

---

1: Initialize $x^{(0)} = 0$
2: **for** $k = 1$ to max_iterations **do**
3:     **for** $i = 1$ to $n$ **do**
4:         $\sigma \leftarrow \sum_{j \neq i} a_{ij} x_j^{(k)}$
5:         $x_i^{(k+1)} \leftarrow (b_i - \sigma)/a_{ii}$
6:     **end for**
7:     **if** $\|x^{(k+1)} - x^{(k)}\| < \epsilon$ **then break**
8:

# Jacobi - Serial Code

```
1  // Main Jacobi loop
2  for (k = 0; k < MAX_ITER; k++) {
3      for (i = 0; i < N; i++) {
4          sum = 0;
5          for (j = 0; j < N; j++) {
6              if (j != i) sum += A[i][j] * x[j];
7          }
8          x_new[i] = (b[i] - sum) / A[i][i];
9      }
10     error = compute_error(x_new, x);
11     if (error < TOL) break;
12     // Update for next iteration
13     for (i = 0; i < N; i++) x[i] = x_new[i];
14 }
```

- ▶ Time: 0.000094 sec (16×16 matrix)
- ▶ Pros: Simple, easy to parallelize

# Gradient Descent - Algorithm

## Key Steps

1. Compute residual: $r = b - Ax$
2. Calculate step size: $\alpha = \frac{r^T r}{r^T A r}$
3. Update solution: $x = x + \alpha r$

---

**Algorithm 2** Gradient Descent

---

1: $x \leftarrow 0$
2: **repeat**
3:     $r \leftarrow b - Ax$
4:     $\alpha \leftarrow (r^T r)/(r^T A r)$
5:     $x \leftarrow x + \alpha r$
6: **until** $\|r\| < \epsilon = 0$

---

# Gradient Descent - Serial Code

```
1  for (iter = 0; iter < MAX_ITER; iter++) {
2      // Compute residual r = b - A*x
3      mat_vec_mult(A, x, Ax);
4      vector_sub(b, Ax, r);
5
6      // Compute step size
7      dot_rr = vector_dot(r, r);
8      mat_vec_mult(A, r, Ar);
9      alpha = dot_rr / vector_dot(r, Ar);
10
11     // Update solution
12     vector_add_scaled(x, alpha, r, x);
13
14     if (sqrt(dot_rr) < TOL) break;
15 }
```

- ▶ Time: 0.000604 sec (16×16 matrix)
- ▶ Pros: Faster convergence for large systems

# PETSc - Key Components

**Matrix Setup**

- ▶ `MatCreate()`: Creates matrix
- ▶ `MATAIJ`: Sparse format
- ▶ Preallocation: Critical for performance

**Solver Setup**

- ▶ `KSPCreate()`: Solver context
- ▶ `KSPCG`: Conjugate Gradient
- ▶ `KSPSetTolerances()`: Convergence

**Why These Choices?**

- ▶ Efficient for sparse systems
- ▶ Robust convergence
- ▶ Parallel-ready design

# PETSc - Code Implementation

```
 1  // 1. Matrix Setup
 2  MatCreate(PETSC_COMM_WORLD, &A);
 3  MatSetSizes(A, PETSC_DECIDE, PETSC_DECIDE, N, N);
 4  MatSetType(A, MATAIJ); // Sparse format
 5  MatMPIAIJSetPreallocation(A, 5, NULL, 5, NULL);
 6
 7  // 2. Vector Setup
 8  VecCreate(PETSC_COMM_WORLD, &x);
 9  VecSetSizes(x, PETSC_DECIDE, N);
10  VecDuplicate(x, &b);
11
12  // 3. Solver Configuration
13  KSPCreate(PETSC_COMM_WORLD, &ksp);
14  KSPSetType(ksp, KSPCG); // Conjugate Gradient
15  KSPSetTolerances(ksp, 1e-6, PETSC_DEFAULT, PETSC_DEFAULT,
        1000);
16
17  // 4. Solve System
18  KSPSetOperators(ksp, A, A);
19  KSPSolve(ksp, b, x);
```

# PETSc - Behind the Scenes

- **Matrix Assembly**:
  - PETSc distributes matrix across MPI processes
  - Each process owns portion of rows
- **Solver Workflow**:
  - Builds Krylov subspace (CG: conjugate directions)
  - Applies preconditioner (default: block Jacobi)
  - Checks convergence criteria
- **Parallel Communication**:
- MPI handles inter-process data exchange
- Overlaps computation/communication

# Execution Results

| Implementation | Time (s) |
|---|---|
| Serial Jacobi | 0.000094 |
| Serial Gradient | 0.000604 |
| PETSc Jacobi | – |
| PETSc Gradient | – |

- ▶ Jacobi: Fastest for small matrices
- ▶ Gradient: Better theoretical scaling
- ▶ PETSc: Implementation in progress

Thank You!