**UNIVERSITÀ DEGLI STUDI DI TRIESTE**

# Distribuited Stencil Computation

A Case Study with OpenMP and MPI

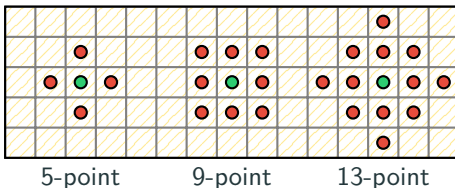Andrea Spinelli
25/07/2025

University of Trieste

# Introduction

## Problem Statement: Stencil Computation

The objective of this project is to implement a parallel **5-point stencil** computation utilizing OpenMP and MPI.

Stencil computations have diverse applications including the numerical solution of partial differential equations, such as the **heat equation**.

- Each grid point $x_{ij}(t)$ is updated iteratively based on its own value at the previous time step, $x_{ij}(t-1)$, and the values of its **neighbors**.
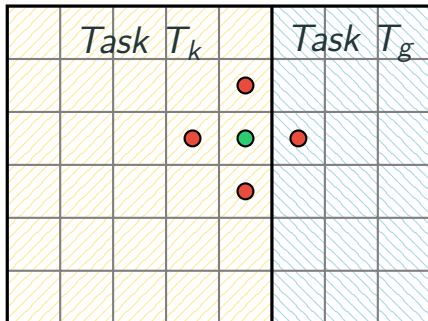


5-point          9-point          13-point

In this work, we focus on the 5-point stencil.

For large grids, the computational cost becomes significant, motivating the need for **parallelization**.

## Challenges: Data Dependencies & Communication

To enable parallel execution, the computational grid can be partitioned into subdomains, with each subdomain assigned to a separate process.



However, a key challenge arises: each grid point depends on the values of its neighbors, including those located in adjacent subdomains.

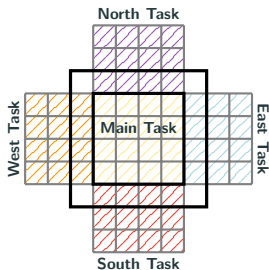This introduces data dependencies across process boundaries and necessitates inter-process communication.

## Challenges: Efficient Communication

Since processes must exchange border data at each iteration, we need the processes to communicate with each other.

The idea is to enlarge the local patches to include the **halo regions** (the border of the local patch).

Communication can be expenive if the problem is not well implemented:

- efficient communication pattern
  $\rightarrow$ Non-Blocking Communication
- efficient memory management
  $\rightarrow$ Avoid unnecessary allocations and copies
  $\rightarrow$ Cache locality among processes
- Computation / Communication Overlap
  (we will see further)



North Task

West Task

Main Task

East Task

South Task

# Implementation

## Parallelization Approach: Hybrid MPI + OpenMP

The implementation tackles both inter-node and intra-node parallelism:

**Distributed-memory parallelism (MPI):**

- The grid is decomposed into patches, each assigned to a process.
- At each iteration, processes exchange halo data with their neighbors using non-blocking communication (`MPI_Isend`, `MPI_Irecv`).
- Collective operations (`MPI_Reduce`, `MPI_Bcast`, `MPI_Gather`) are used for global statistics and data distribution.

**Shared-memory parallelism (OpenMP):**

- Within each MPI process, OpenMP parallelizes the main update loops (`#pragma omp parallel for`).
- Reductions (e.g., for total energy) are performed in parallel.

## Optimization: Overlapping Communication and Computation

In stencil codes, processes must wait for halo data from neighbors before updating boundary points. We can potentially improve performance by **overlapping communication with computation**.

We split the update loop into two parts, *interior* and *border*:

## Optimization: Overlapping Communication and Computation

In stencil codes, processes must wait for halo data from neighbors before updating boundary points. We can potentially improve performance by **overlapping communication with computation**.

We split the update loop into two parts, *interior* and *border*:

1. **Initiate Non-Blocking Communication**: Start sending and receiving halo data using MPI_Isend and MPI_Irecv.

## Optimization: Overlapping Communication and Computation

In stencil codes, processes must wait for halo data from neighbors before updating boundary points. We can potentially improve performance by **overlapping communication with computation**.

We split the update loop into two parts, *interior* and *border*:

1. **Initiate Non-Blocking Communication**: Start sending and receiving halo data using `MPI_Isend` and `MPI_Irecv`.
2. **Compute the Interior**: While the halo data is in transit over the network, compute the new values for the *interior points* of the local grid. These points do not depend on the halo data.

## Optimization: Overlapping Communication and Computation

In stencil codes, processes must wait for halo data from neighbors before updating boundary points. We can potentially improve performance by **overlapping communication with computation**.

We split the update loop into two parts, *interior* and *border*:

1. **Initiate Non-Blocking Communication**: Start sending and receiving halo data using `MPI_Isend` and `MPI_Irecv`.
2. **Compute the Interior**: While the halo data is in transit over the network, compute the new values for the *interior points* of the local grid. These points do not depend on the halo data.
3. **Synchronize and Complete**: Use `MPI_Waitall` to ensure all halo communications have completed.

## Optimization: Overlapping Communication and Computation

In stencil codes, processes must wait for halo data from neighbors before updating boundary points. We can potentially improve performance by **overlapping communication with computation**.

We split the update loop into two parts, *interior* and *border*:

1. **Initiate Non-Blocking Communication**: Start sending and receiving halo data using `MPI_Isend` and `MPI_Irecv`.
2. **Compute the Interior**: While the halo data is in transit over the network, compute the new values for the *interior points* of the local grid. These points do not depend on the halo data.
3. **Synchronize and Complete**: Use `MPI_Waitall` to ensure all halo communications have completed.
4. **Compute the Border**: Once the halo data is available, compute the new values for the *border points* of the local grid.

5

# Overlapping: Consequences

Contrary to the hypothesis, this optimization yielded no significant improvements because the workload is heavily **compute-bound**.
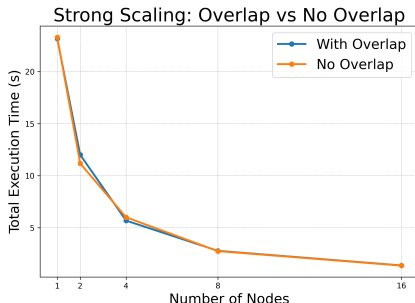
The time spent in calculations vastly outweighs communication time.

For example, in an 8-node run, the breakdown was:

**Computation Time:**  ∼**96%**
**Communication Time:**  ∼**4%**

Hiding the negligible communication latency provides no visible benefit.



Strong Scaling: Overlap vs No Overlap

Overlapping is effective only when communication latency is a significant performance bottleneck.

## Implementation Details: Memory and Data Management

Other techniques were used to ensure performance and correctness.

**Efficient Buffer Management**

To minimize memory allocation overhead and avoid unnecessary copies:

- Communication buffers are *pre-allocated* once at initialization and reused throughout the simulation.
- For contiguous boundaries (North/South), direct pointer arithmetic is used to send/receive data *in-place* without intermediate buffering.
- Only non-contiguous boundaries require temporary buffers.

**NUMA-Aware Memory Allocation**

The *first-touch policy* on NUMA architectures is exploited. By parallelizing the initialization of the grid with OpenMP, we ensure that each thread allocates the memory pages for its assigned rows on the NUMA node it is running on. This minimizes remote memory access during the computation phase.

## Build Process and Compilation

Compilation is handled by a `Makefile` that leverages the `mpicc` wrapper.

**Compiler Flags**

Key flags used for the final performance runs:

- `-Ofast`: Aggressive optimization, enables `-O3` optimizations and more
- `-flto`: Link-time optimization for improved performance.
- `-fopenmp`: Enables OpenMP support for parallel regions.
- `-march=native`: Generates code optimized for the build machine's architecture.
- `-Wall -Wextra -Wpedantic -Wshadow -Wuninitialized`: Enables a wide range of compiler warnings for code quality.

**Note on Compilation Environment**

Compile on a compute node to ensure `-march=native` targets the correct architecture (login node may differ).

## Verbose Output and Instrumentation

The code uses preprocessor directives to control the level of runtime output via compile-time flags:

- VERBOSE_LEVEL=1 (verbose1 target): Prints neighbor information, timing statistics, and energy statistics at each step.
- VERBOSE_LEVEL=2 (verbose2 target): In addition to level 1, prints the full grid state at each iteration for detailed inspection.



This approach ensures that verbose output and measurement code are excluded from the final performance binaries unless explicitly requested, avoiding unnecessary overhead.

## Reproducibility

A key goal is to ensure the parallel code is **reproducible and bitwise compatible with the serial version**, allowing reliable verification.

This is challenging, especially for source generation. In a naive parallel approach, each process generates sources independently, which cannot match the serial version.

To ensure compatibility:

- All sources are generated by the rank 0 process and then distributed to the other processes.
- A -s (seed) flag sets the random seed, so the same input always produces the same results.

# Results

## Overview of Scaling Tests

To comprehensively evaluate performance, three principal types of scalability tests were conducted:

- **OpenMP thread scaling**:
  single MPI task,
  varying threads (1, 2, 4, 8, 16, 32, 56, 84, 112)
- **MPI strong scaling**:
  fixed problem size,
  increasing nodes (1, 2, 4, 8, 16)
- **MPI weak scaling**:
  fixed workload per node,
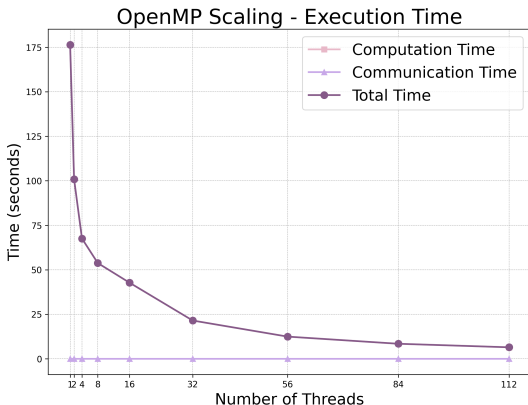  increasing nodes (1, 2, 4, 8, 16) and grid size accordingly

Batch scripts were used to automate runs and collect timing data.

For each configuration, *total*, *computation*, and *communication* times were measured.
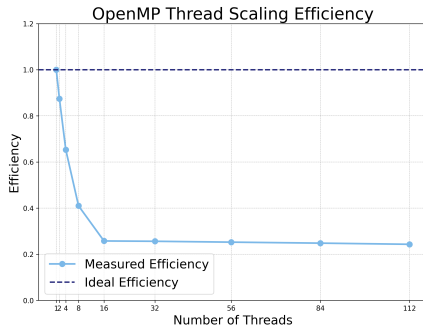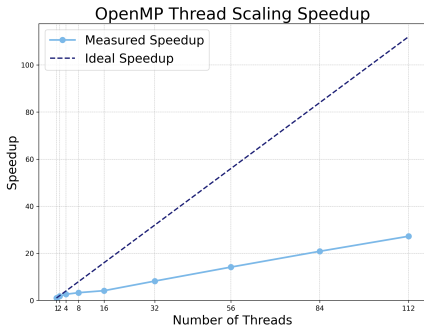
# OpenMP Scaling: Execution Time

**OpenMP thread scaling**: single MPI task, varying number of threads.

Measures how execution time changes as more threads are used



OpenMP Scaling - Execution Time

Execution time decreases with more threads, but the benefit diminishes at high thread counts due to memory bandwidth and thread overhead.
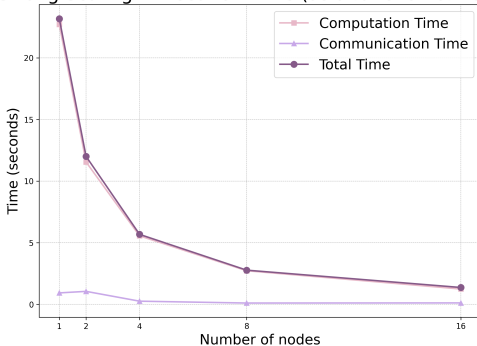
# OpenMP Scaling: Speedup and Efficiency



- Speedup is sublinear for large thread counts; best performance is achieved with moderate thread numbers.
- Efficiency drops as the number of threads increases, mainly due to memory bandwidth limitations and thread management overhead.

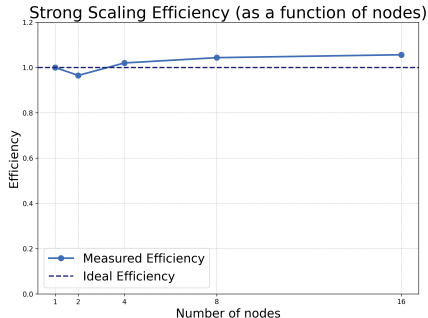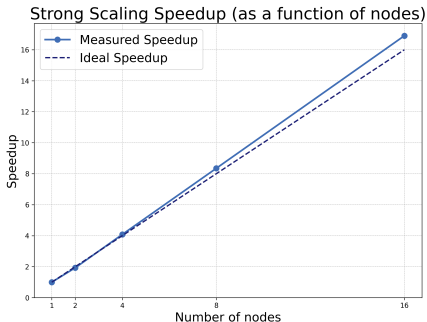**Strong scaling**: fixed problem size, increasing number of nodes.

Measures how execution time decreases as more resources are used for the same total work.



Strong Scaling - Execution Time (as a function of nodes)

Execution time drops as more nodes are used.

# MPI Strong Scaling: Speedup and Efficiency



Strong Scaling Speedup (as a function of nodes)

Strong Scaling Efficiency (as a function of nodes)
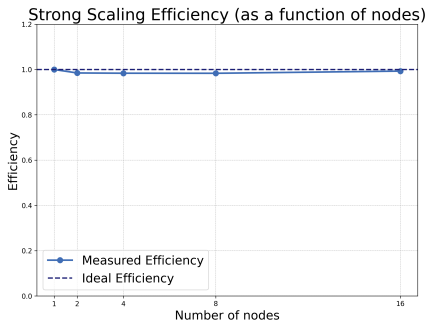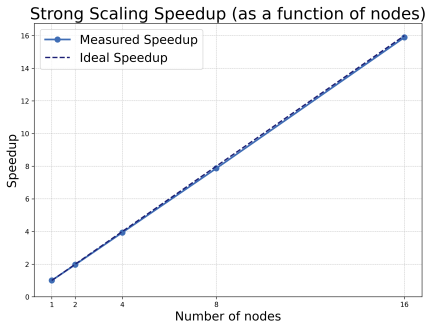
Ideally, speedup should be linear.

Results show a **superlinear speedup**, with efficiency greater than 1.

This means the parallel code is more than $N$ times faster on $N$ nodes, which seems to violate Amdahl's Law. This is not an error, but (probably) a real phenomenon caused by **cache effects**.

# MPI Strong Scaling: Superlinear Speedup Explained

Superlinear speedup can occur due to improved cache utilization as the local problem size per node decreases.

To verify this, strong scaling was repeated with only one MPI task per node, removing the cache benefit from multiple tasks sharing a node.
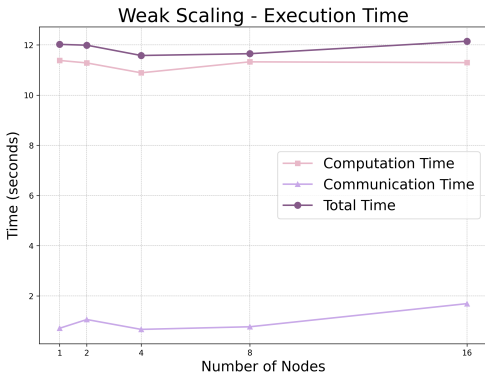


Strong Scaling Speedup (as a function of nodes)

Strong Scaling Efficiency (as a function of nodes)

In this configuration, the superlinear effect disappears and efficiency returns to below 1, confirming the cache explanation.
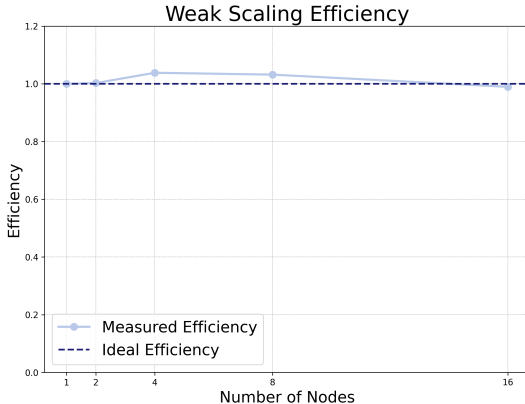
**Weak scaling**: fixed workload per node, increasing number of nodes.

Measures if execution time remains constant as the problem size grows proportionally with resources.



Weak Scaling - Execution Time

Ideally, execution time should remain constant. Observed times are nearly flat, with minor variations due to system effects.

Weak Scaling Efficiency

- Efficiency stays near 1, as expected for good weak scaling.
- Observed efficiency slightly above 1 in some cases, likely due to cache effects or measurement noise.

# Conclusions

## Key Findings and Conclusions

The parallel stencil solver, using a hybrid MPI+OpenMP approach, demonstrated strong scalability and high efficiency on a modern HPC.

**Main Achievements**

- **Strong and Weak Scaling:**
  The code achieved (over) ideal strong and weak scaling, confirming its effectiveness for large-scale workloads.
- **Hybrid Parallelism:**
  Combining OpenMP (for shared-memory parallelism within nodes) and MPI (for distributed-memory parallelism across nodes) provided robust performance and flexibility.

**Key Takeaway**

For compute-bound applications on modern HPC systems, maximizing data locality and cache usage can deliver performance gains that exceed predictions from simple parallel scaling laws.

## Future Work and Potential Enhancements

Building on this solid foundation, several avenues for further optimization and exploration can be pursued:

**Algorithmic & Communication**

- **Implement MPI Derived Datatypes**: Replace manual packing/unpacking of non-contiguous halo data (East/West boundaries) with `MPI_Type_vector`.
- **Advanced Domain Decomposition**: Implement dynamic load balancing strategies, such as recursive bisection, to handle non-uniform workloads (e.g., simulations with localized sources).

**Architecture & Performance Tuning**

- **GPU Acceleration**: Port the `update_plane` kernel to a GPU using OpenMP or CUDA to exploit massive data parallelism.
- **Deeper NUMA Analysis**: Experiment with different thread/task placement strategies using SLURM's binding options.

# Thank You!