

# Introduction to OpenMP Programming Model

Dr. Ezhilmathi Krishnasamy

October 29, 2024



# Outline

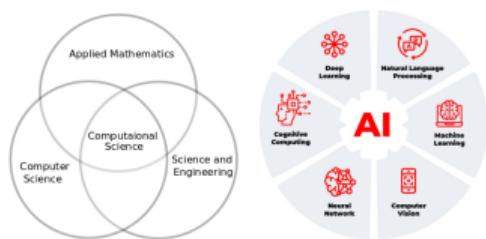
- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Outline for section 1

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Motivation: why we need parallel programming

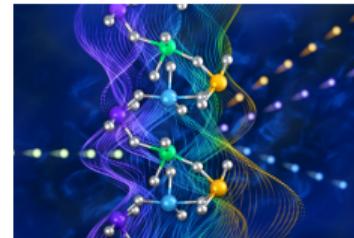
- ① In the fields of science, engineering, and artificial intelligence, we often engage in extensive arithmetic computations with numbers.
- ② Problems are typically defined by differential equations, particularly partial differential equations (PDEs).
- ③ To solve these PDEs, we convert them into a system of equations using various numerical methods, such as the finite difference method, finite element method, and finite volume method.
- ④ Ultimately, we must solve these systems of equations using either direct or indirect solvers to determine the values of unknown variables.
- ⑤ Similarly, in the realm of artificial intelligence (AI), we frequently work with matrices and vectors.



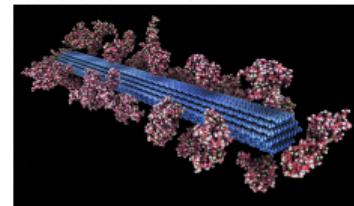
$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}.$$

## Applications and examples

- ① **Materials Science:** Materials science is an interdisciplinary field that combines chemistry, physics, and engineering to provide a nano-level understanding of materials. This knowledge enables the design of new materials with specific, predetermined properties.
  - ② **Proteins** are the fundamental building blocks of the human body. Proteins in cell membranes regulate the flow of substances into and out of cells. A deeper understanding of these proteins can lead to new treatments for diseases and potentially extend human lifespan.
  - ③ In **engineering**, disciplines such as computational fluid dynamics, computational heat transfer, and computational solid mechanics require significant computer power to simulate and design optimized systems. Traditionally, these processes relied on experiments and wind tunnel tests, which could be costly in terms of both finances and human resources.



Materials science: Credit:www.olcf.ornl.gov



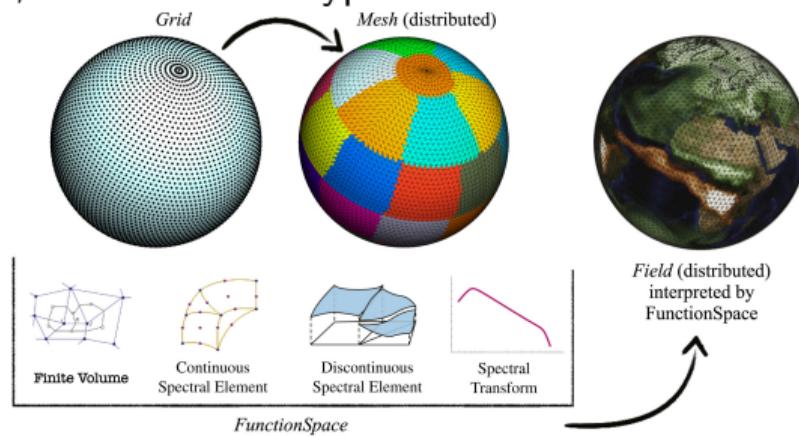
Biology: Credit:www.olcf.ornl.gov



Engineering: Credit: [www.olcf.ornl.gov](http://www.olcf.ornl.gov)

# Applications and examples

In weather forecasting, numerical methods are used to discretize the domain size. For instance, simulating a resolution of 3 kilometers (with 1.8 billion cells) on the IBM Blue Gene/Q Mira system at Argonne National Laboratory requires 786,000 cores. This simulation uses spectral element methods, which are one type of numerical method used in the field<sup>1</sup>.



Mesh and simulation settings <sup>2</sup>

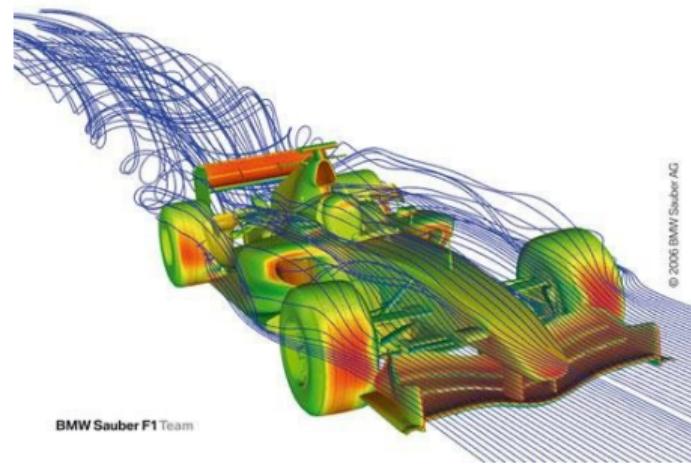
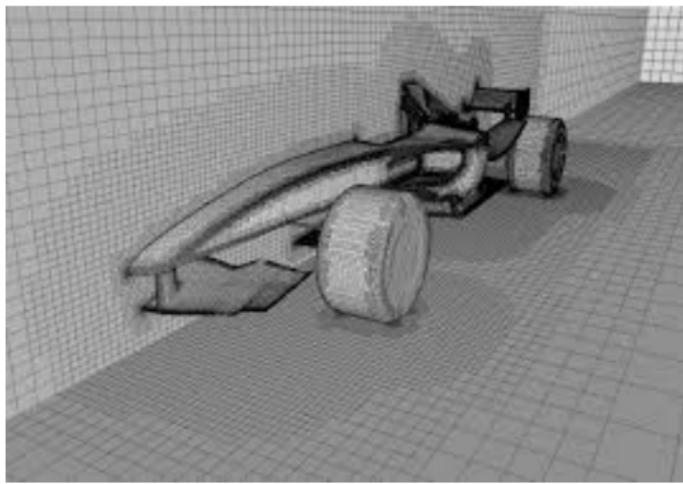


HPC for Weather Forecasting, John Michalakes

A library for numerical weather prediction and climate modelling, Willem Deconinck, et al.

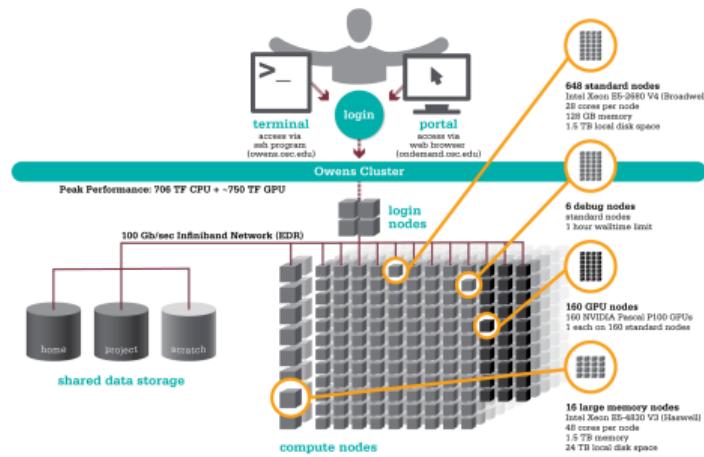
# Applications and examples

CERFACS in France utilized up to 250,000 cores with grid sizes ranging from 2 to 4 billion cells using Large Eddy Simulation, a numerical method strategy for computational fluid dynamics applications.



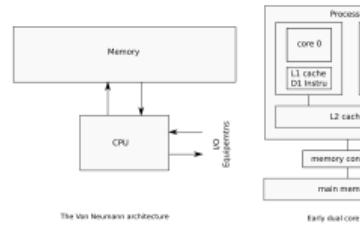
# What is supercomputer

- ① A large machine is composed of multiple processors and extensive memory.
- ② The processors (CPUs) and their memory are interconnected by network cables.
- ③ Multicore processors are included within the CPUs.
- ④ By employing parallel computing techniques, these machines can be utilized more efficiently.

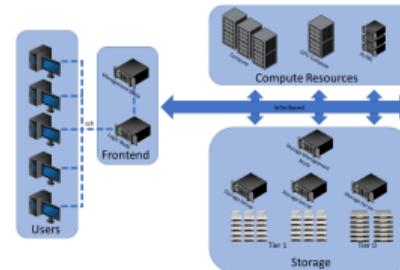
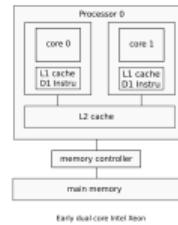


# Difference between laptop and supercomputer

Laptop	Supercomputer
few cores	lots of cores
limited memory	lots of memory
slower than supercomputer	faster than laptop
personalized Operating System (OS)	mostly Linux based
owned by you	owed by public sector and located somewhere
used by owner	many users



The Von Neumann architecture



# Computer speed

- ① FLOPS: Floating Point Operations Per Second - this metric measures the speed of a computer.
- ② The speed of a supercomputer is determined by running the [LINPACK](#) benchmark applications.
- ③ Typically, this involves solving N-dimensional linear systems using Gaussian elimination with partial pivoting.
- ④ The [TOP500](#) regularly lists the fastest supercomputers.
- ⑤ As of June 2024, the top 3 supercomputers are:
  - Frontier: HPE Cray EX235a, AMD EPYC 64C at 2GHz, AMD Instinct MI250X, Slingshot-11; DOE/SC/Oak Ridge, USA.
  - Aurora: HPE Cray EX, Intel Xeon CPU Max, Intel Data Center GPU Max, Slingshot-11; Argonne, USA.
  - Eagle: Microsoft NDv5, Intel Xeon Platinum 8480C, NVIDIA H100; Microsoft Azure, USA.

# Computer speed

Speed	Flops
1 Flop	$10^0 = 1$
1 Kflops	$10^3 = 1 \text{ Thousand}$
1 Mflops	$10^6 = 1 \text{ Million}$
1 Gflops	$10^9 = 1 \text{ Billion}$
1 Tflops	$10^{12} = 1 \text{ Trillion}$
1 Pflops	$10^{15} = 1 \text{ Quadrillion}$
1 Eflops	$10^{18} = 1 \text{ Quintillion}$
1 Zflops	$10^{21} = 1 \text{ Sextillion}$
1 Yflops	$10^{24} = \text{Septillion}$

Modern supercomputers primarily operate at the levels of petaflops (Pflops) and exaflops (Exaflops).

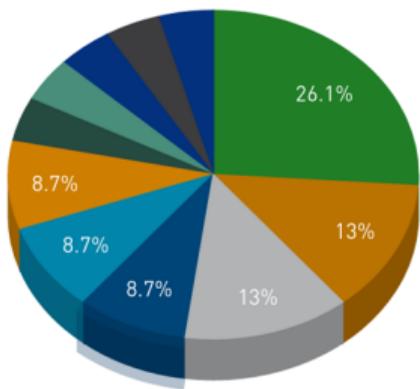
# Computing power trend



credit:top500; FLOPS(Floating Point Operations per Second)

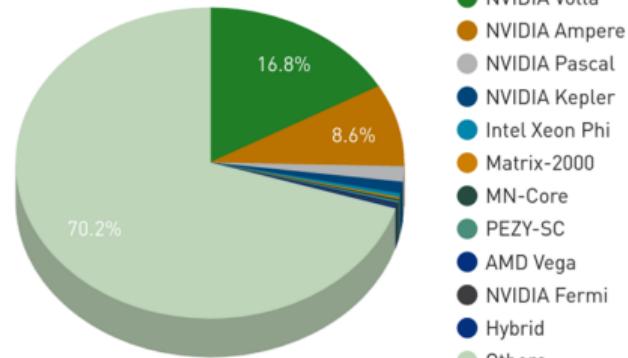
# How it is used in different research sectors

Application Area System Share



credit: top500

Accelerator/CP Family System Share



credit: top500

# Supercomputers in Luxembourg

- ① [MeluXina](#): Accelerator Module - BullSequana XH2000, featuring an AMD EPYC 7452 processor with 32 cores at 2.35 GHz, an NVIDIA A100 with 40 GB, and equipped with Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite, provided by Atos and LuxProvide in Luxembourg. [Ranked 36th in 2021 \(according to TOP500\)](#).
- ② [Aion](#) is a supercomputer developed by Atos/Bull/AMD, consisting of 318 compute nodes. It totals 40704 compute cores and 81408 GB of RAM, delivering a peak performance of approximately 1.70 PetaFLOP/s.
- ③ [Iris](#) is a Dell/Intel supercomputer comprised of 196 compute nodes, totaling 5824 compute cores and 52224 GB of RAM, with a peak performance of around 1.072 PetaFLOP/s.

# Supercomputers in Luxembourg

	Meluxina	Aion	Iris
<b>CPU Architecture</b>	AMD EPYC	AMD EPYC	Intel Xeon
<b>GPU Architecture</b>	Nvidia A100	-	Nvidia V100
<b>Performance RPEAK</b>	18 PFLOPS	1,70 PFLOPS	1,072 PFLOPS
<b>Host</b>	EuroHPC & LuxProvide	Luxembourg University	Luxembourg University
<b>Vendor</b>	AMD	AMD	Intel



# Outline for section 2

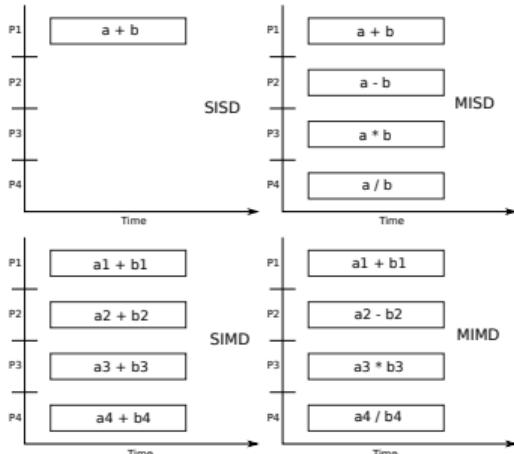
- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Classification of Parallel Computer Architecture

- ① Control Structure: Based on the instruction and data streams.
  - SISD - Single Instruction Stream, Single Data Stream
  - SIMD - Single Instruction Stream, Multiple Data Streams
  - MISD - Multiple Instruction Streams, Single Data Stream
  - MIMD - Multiple Instruction Streams, Multiple Data Streams
- ② Memory Organization: Shared Memory and Distributed Memory
- ③ Network Topology (Connectivity): For example, 3D Grid and Tree

# Control Structure - Flynn's Taxonomy

- ① SISD - Represents sequential execution, exemplified by conventional single processors in the von Neumann model.
- ② SIMD - Refers to array computers and traditional vector processors, where instructions can be executed either in a pipeline or in parallel.
- ③ MISD - This model is rarely used in practice.
- ④ MIMD - Operates with both shared and distributed memory models.



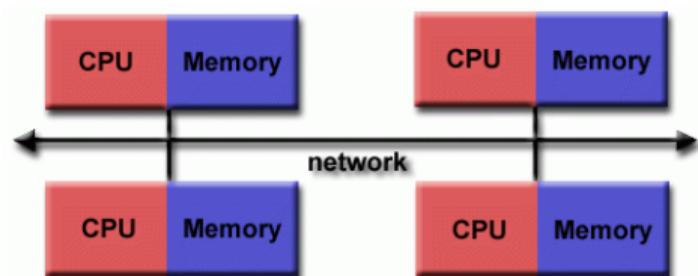
Instruction Streams	
one	many
one	<b>SISD</b> traditional von Neumann single CPU computer
many	<b>MISD</b> May be pipelined Computers
many	<b>SIMD</b> Vector processors fine grained data Parallel computers
many	<b>MIMD</b> Multi computers Multiprocessors

# Outline for section 3

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Distributed memory architecture

- ① A distributed memory architecture is designed using a network connection.
- ② There is no global address space; each compute node has its own local memory.
- ③ There is no cache coherence, meaning that if one processor updates memory, that change is not automatically visible to another processor at the hardware level.
- ④ It is up to the programmers to determine how and from where data should be accessed or transferred.



# Distributed memory architecture

## ① Advantages:

- It is scalable with respect to memory and processors, allowing for large-scale computations.
- Each processor within a node can quickly access its own memory without the need for global cache coherence.
- Networking enables access to multiple processors and memory resources.

## ② Disadvantages:

- Accessing local data is faster than accessing data across compute nodes via interconnect networks.
- It is the programmer's responsibility to manage read and write operations across multiple compute nodes and their remote memory.

# Outline for section 4

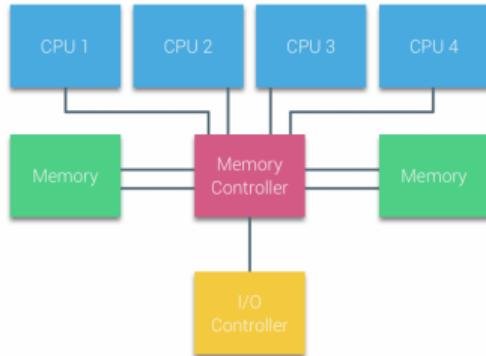
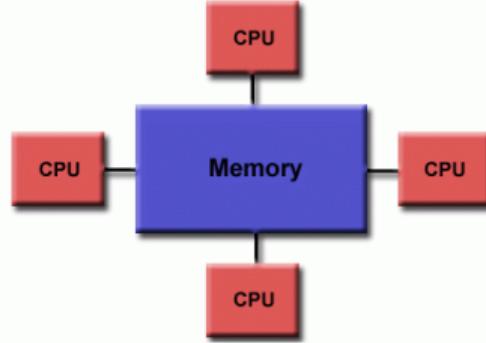
- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Shared memory architecture

- ① In shared memory computers, all processors have access to a global address space of memory.
- ② Multiple processors within a shared memory architecture can operate independently while sharing the same memory resources.
- ③ Changes made to a memory location by one processor are visible to all other processors.
- ④ Shared memory architecture is categorized into two types: Unified Memory Access (UMA) and Non-Unified Memory Access (NUMA), based on how memory is accessed.

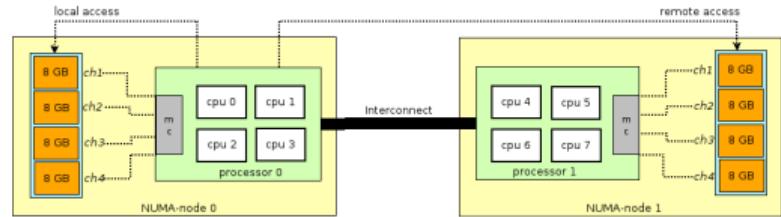
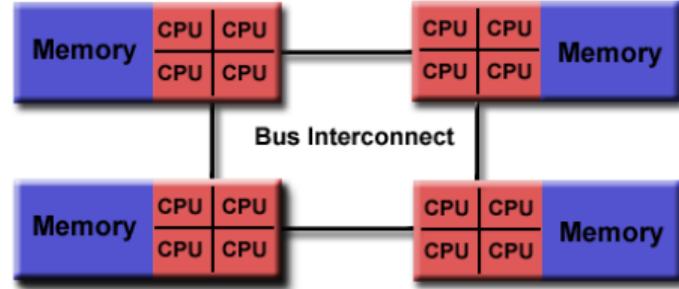
# Unified Memory Access (UMA)

- ① All processors access memory through an interconnection network, similar to how a single processor accesses its own memory.
- ② Memory access times for all processors are equal, meaning they take the same amount of time when accessing memory through the interconnection network (e.g., bus).
- ③ This architecture is often represented by Symmetric Multiprocessor (SMP) systems, which feature identical processors.
- ④ It is also referred to as Cache Coherent UMA (CC-UMA), where memory updates made by one processor are visible to other processors, achieved through hardware-level configurations.



# Non-Unified Memory Access (NUMA)

- ① Each set of processors (Symmetric Multiprocessing, or SMP) has its own dedicated memory.
- ② Processors from one SMP can access the memory of another SMP.
- ③ The time taken to access memory varies among different processors across the SMPs, leading to slower memory access when communicating between SMPs.
- ④ NUMA can also include Cache Coherent NUMA (CC-NUMA) through specific hardware-level configurations.



# Shared memory architecture

## ① Advantages

- Rapid and uniform data sharing between processes, due to the close proximity of memory to CPUs.
- The global address space provides a user-friendly perspective for memory access.
- No need for explicit specification of data communication between processes.
- Easy to learn and program (e.g., OpenMP).

## ② Disadvantages

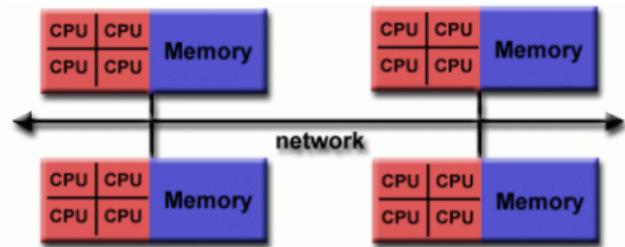
- Scalability is limited by memory and processor capabilities.
- Programmers must implement synchronization constructs (e.g., locks) to ensure access to the correct data in global memory.

# Outline for section 5

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Hybrid distributed-shared memory

- ① Modern supercomputers utilize both shared and distributed memory architectures.
- ② Shared memory components may include GPUs and other accelerators that can access this shared memory.
- ③ Distributed memory connects to shared memory through network connections, allowing memory access from one shared memory node to another via an interconnect network.
- ④ Scalability can be enhanced by increasing the number of compute nodes and memory—this is an advantage.
- ⑤ However, there is increased programmer complexity involving tasks, data management, and communication—this is a disadvantage.



# Outline for section 6

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# What is OpenMP

- ① **OpenMP** is an application programming interface (API) designed to facilitate parallel programming on shared memory architectures, including both UMA and NUMA.
- ② It uses a directive-based programming model to parallelize serial code written in **C/C++** and **Fortran**.
- ③ The **directives** instruct the compiler on how to execute code in parallel within the shared memory environment.
- ④ A single application can contain both serial and parallel versions of its code.
- ⑤ The latest specification of the **OpenMP API 5.2** was released in November 2021.

# OpenMP compilers

Vendor/Source	Compiler/Language
AMD	C/C++
ARM	C/C++/Fortran
GNU	GCC - C/C++/Fortran
HPE	CCE - C/C++/Fortran
IBM	XL - C/C++/Fortran
Intel	C/C++/Fortran
NVIDIA HPC compiler	C/C++/Fortran

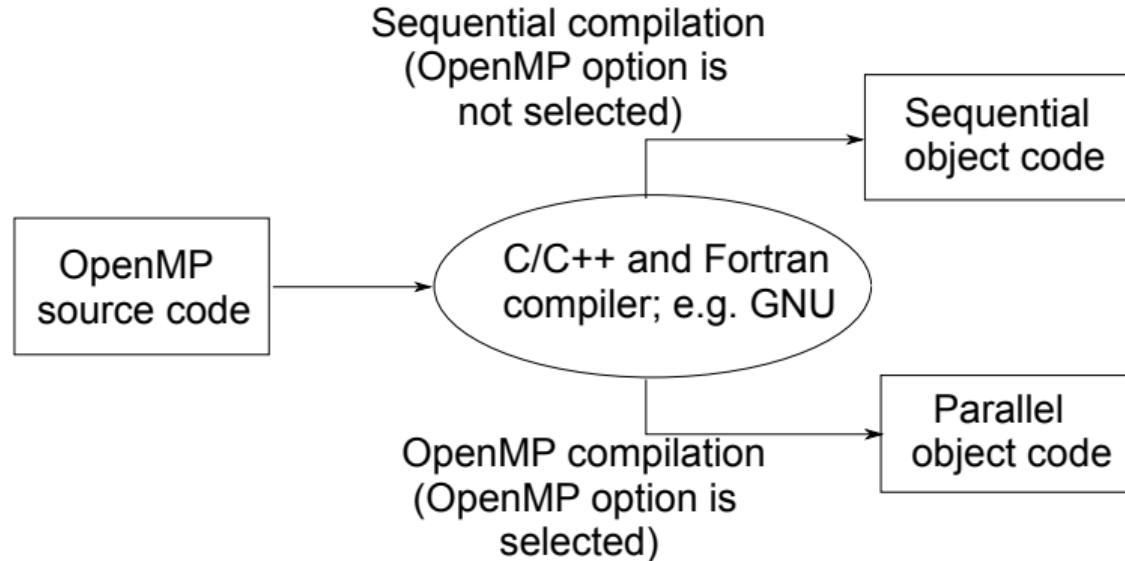
## GNU

```
g++ -fopenmp hello-world.cc -o hello-world  
gcc -fopenmp hello-world.c -o hello-world
```

## Intel

```
icpc -qopenmp hello-world.cc -o hello-world
```

# Compilation

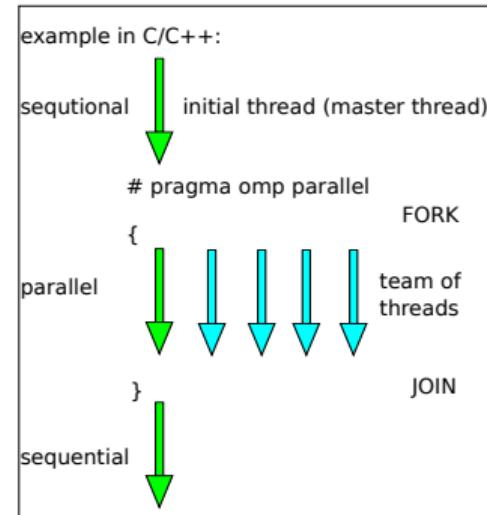


# Outline for section 7

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Parallel region

- ① **Fork-Join:** model of parallel execution.
- ② **FORK:** master thread creates a **team** of parallel threads (master thread is also part of a team of threads).
  - in C: **#pragma omp parallel** will create threads.
  - by default, thread number is defined by the given architecture.
- ③ **JOIN:** team threads complete statements in parallel regions.
  - then they synchronize & terminate, leaving only the master thread to continue.



# Parallel construct

## Example of parallel region

---

```
#include <iostream>
#include <omp.h>
using namespace std;
int main()
{
#pragma omp parallel /*-- create a parallel region "Fork" --*/
{
    cout << "Hello world from thread id "
        << omp_get_thread_num()
        << " from the team size of "
        << omp_get_num_threads()
        << endl;
} /*-- parallel region is closed "Join" --*/
return 0;
cout << "end of the program from the master thread" << endl;
}
```

---

```
g++ -fopenmp hello-world.cc -o hello-world
```

# Outline for section 8

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines**
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# OpenMP - environment routines

- ➊ `void omp_set_num_threads(int num_threads);`  
→ number of threads to be used in the parallel region
  
- ➋ `int omp_get_num_threads(void);`  
→ gets the number of threads in the current parallel region
  
- ➌ `int omp_get_max_threads(void);`  
→ get available maximum threads (system default)

# OpenMP - environment routines

- ① `int omp_get_thread_num(void);`  
→ gets the thread numbers (e.g., 1, 4, etc.)
  
- ② `int omp_get_num_procs(void);`  
→ gets the number of the available processor in the system
  
- ③ `int omp_in_parallel(void);`  
→ boolean variable to check if the parallel region is active (returns 1) or not (returns 0)

# OpenMP - environment routines

## OMP\_NUM\_THREADS

```
$ setenv OMP_NUM_THREADS 4
```

```
$ export OMP_NUM_THREADS=4
```

```
$ OMP_NUM_THREADS=4 ./omp_code.exe
```

# Outline for section 9

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Shared variable

- ① All threads can access the shared variable.
  - ② By default, in a parallel region, all variables are treated as shared.
  - ③ Shared variables must be managed with care to avoid race conditions in the program.
  - ④ In this context, all threads within the parallel region can access the elements of vector a.
- 

```
#pragma omp parallel for shared(a)
for (int i = 0; i < N; i++)
    a[i] += i;
/*-- End of parallel region --*/
```

---

# Private variable

- ① Each thread has its own copy of the private variable.
  - ② The private variable is accessible only within the parallel region, not outside of it.
  - ③ By default, loop iteration counters are considered private.
  - ④ Changes made by one thread are not visible to other threads.
- 

```
#pragma omp parallel for private(a)
for (int i = 0; i < N; i++)
    a = 1 + i;
/*-- End of parallel region --*/
```

---

```
default(shared | none)
#pragma omp parallel default(shared) private(var_list)
```

# Lastprivate

- ① **lastprivate:** is similar to a private clause.
  - ② Each thread will have an uninitialized copy of the variables specified as lastprivate.
  - ③ At the end of the parallel loop or sections, the value of the lastprivate variable will be the final value accessed by the last thread that modified it within the section or loop.
- 

```
int var = 5;  
#pragma omp parallel for lastprivate(var)  
for(int i = 0; i < n; i++)  
{  
    var += omp_get_thread_num();  
    cout << " lastprivate in parallel region " << var << endl;  
} /*-- End of parallel region --*/  
cout << "lastprivate after parallel region " << var << endl;
```

# Firstprivate

- ① **firstprivate**: behaves like a private clause.
  - ② Each thread receives an initialized copy of the variables passed as firstprivate.
  - ③ It is available for parallel constructs, loop, sections, and single constructs.
- 

```
int var = 5;  
#pragma omp parallel firstprivate(var)  
{  
    var += omp_get_thread_num();  
    cout << " firstprivate variable within the parallel region " << var << endl;  
} /*-- End of parallel region --*/  
cout << "firstprivate varibale after the parallel region " << var << endl;
```

---

# Outline for section 10

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Work-sharing constructs

- ① Loop construct
- ② Sections construct
- ③ Single construct
- ④ Master construct
- ⑤ Workshare construct (only available in Fortran)

# Loop construct

- ① **For:** The 'for' directive executes a loop in parallel.
- ② This directive should be placed within the parallel region using `#pragma omp parallel`.
- ③ By default, the loop counter index is private.
- ④ **Loop Scheduling:** This feature includes the following clauses:
  - static
  - dynamic
  - guided
  - auto
  - runtime

Functionality	Syntax in C/C++	Syntax in Fortran
Distribute iterations over the threads	<code>#pragma omp for</code>	<code>!\$omp do</code>

# Loop construct

## Example of work-sharing loop

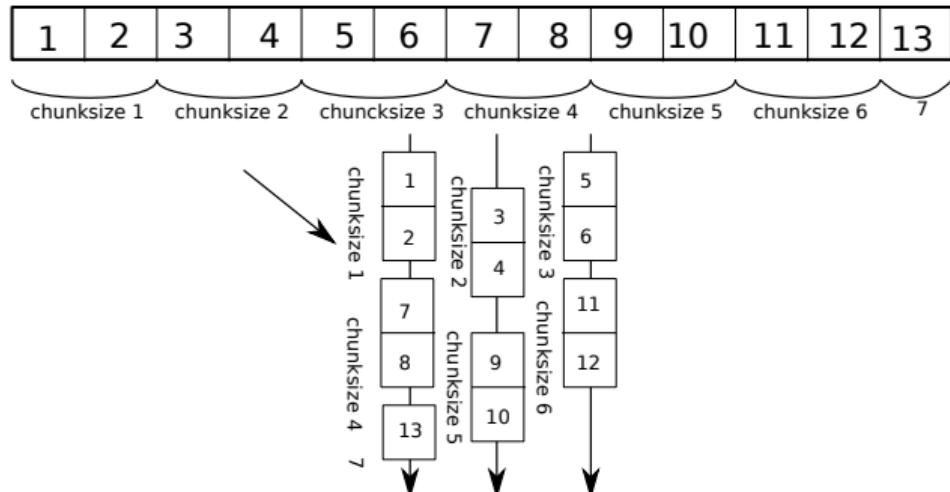
---

```
#include <iostream>
#include <omp.h>
using namespace std;
int main()
{
    int N = 10;
#pragma omp parallel num_threads(2)
    {
#pragma omp for
        for (int i = 0; i < N; i++)
        {
            cout << " Thread " << omp_get_thread_num()
                << " executes loop iteration " << i << endl;
        }
    } /*-- End of parallel region --*/
    return 0;
}
```



# Static clause

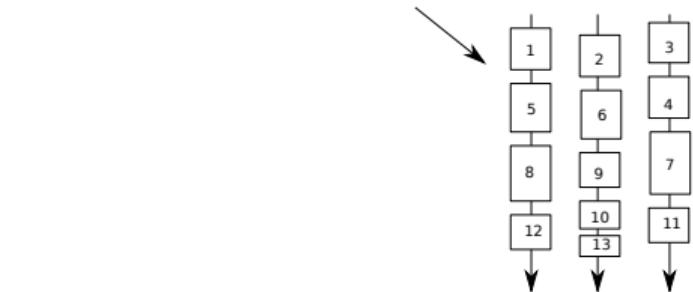
- ① Number of iterations is divided by chunksize (granularity of the workload)
- ② If the chunksize is not provided, a number of iterations will be divided by the size of the team of threads
  - e.g.,  $n=64$ ,  $\text{num\_threads}=4$ ; each thread will execute the 16 iterations in parallel
- ③ This approach is beneficial when the computational cost is similar for each iteration.



# Dynamic clause

- ① The number of iterations is divided by the chunk size.
- ② If the chunk size is not specified, the default value of 1 will be used.
- ③ The system will request a chunk of data until no more chunks are available.
- ④ There is no specific pattern for how different threads access the various chunks of data.
- ⑤ This approach is beneficial when the computational cost varies across iterations.
- ⑥ This method efficiently places chunks of data into the queue.

1	2	3	4	5	6	7	8	9	10	11	12	13
---	---	---	---	---	---	---	---	---	----	----	----	----



# Guided, Auto, and Runtime schedule clause

## ① Guided:

- Similar to dynamic scheduling, the number of iterations is divided into a "chunksize."
- However, the size of each chunk decreases proportionally to the number of unsigned iterations divided by the number of threads.
- If a "chunksize" is not specified, the system will use a default value of 1.
- This approach is useful when there is poor load balancing towards the end of the iterations.

## ② Auto:

- Compiler chooses the optimized chunksize for a number of iterations.

## ③ Runtime:

- Environmental variable OMP\_SCHEDULE is used to define parallel for loop scheduling.

# Schedule clause

## Example of schedule clause

---

```
#include <iostream>
#include <omp.h>
using namespace std;
int main()
{
    int i, N = 1000;
    int *A, *B, *C;
    A = new int[N]; B = new int[N]; C = new int[N];
#pragma omp parallel for schedule(runtime)
    for (i= 0; i < N; i++)
    {
        C[i] = A[i] + B[i];
    }
/*-- End of parallel region --*/
    return 0;
}
```

```
$ export OMP_SCHEDULE="DYNAMIC"
$ setenv OMP_SCHEDULE "GUIDED,4"
```

# Collapse clause

- ① Useful for managing nested loops.
  - ② The total number of iterations is divided among the available threads.
  - ③ For example, if the outer loop index equals the number of threads, the innermost loop (with a size larger than the number of threads) will have its work distributed among those threads.
  - ④ Only one collapse clause is permitted within the nested loop.
- 

```
int i, N = 1000;
int *A, *B, *C;
A = new int[N]; B = new int[N]; C = new int[N];
#pragma omp parallel for schedule(runtime) collapse(2)
for (int i = 0; i < N; i++)
    for (int j = 0; j < N; j++)
    {
        C[i]=A[i]*B[j];
    }
/* End of parallel and collapse region --*/
```

# Section(s) construct

- ① Different tasks can be executed in parallel.
  - ② Sections are useful when multiple tasks (functions or subroutines) need to be processed in parallel.
  - ③ Each thread will execute one task only once.
  - ④ Tasks within the sections should be independent of each other.
- 

```
#pragma omp parallel sections num_threads(9)
{
#pragma omp section
    cout << "section-1 from thread id " << omp_get_thread_num () << endl;
#pragma omp section
    cout << "section-2 from thread id " << omp_get_thread_num () << endl;
#pragma omp section
    cout << "section-3 from thread id " << omp_get_thread_num () << endl;
#pragma omp section
    cout << "section-4 from thread id " << omp_get_thread_num () << endl;
} /*-- End of parallel and sections region --*/
```

# Single construct

- ① The order in which the threads access the single construct is random.
  - ② Only one thread will execute the single construct, and only once.
  - ③ The single construct waits for all other constructs to finish their tasks.
- 

```
#pragma omp parallel num_threads(6)
{
#pragma omp single
{
    cout << "Thread id in single " << omp_get_thread_num()
        << " from team of threads " << omp_get_num_threads() << endl;
}
cout << "Parallel region thread id - "
    << omp_get_thread_num() << " from team of threads "
    << omp_get_num_threads() << endl;
} /*-- End of parallel and single region --*/
```

# Master construct

- ① Only the master thread executes the master construct.
  - ② It functions similarly to a single construct; however, it does not include an implied barrier.
  - ③ It must be used in conjunction with a barrier when there is a data dependency.
- 

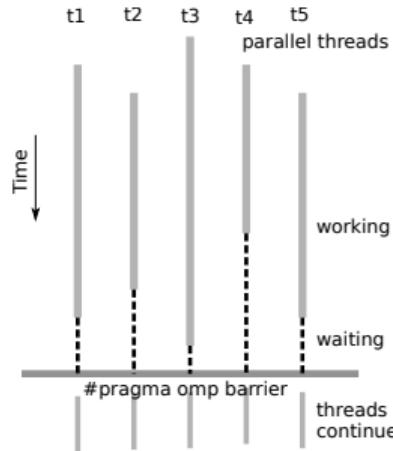
```
#pragma omp parallel shared(a,b)
{
    #pragma omp master
    {
        a = 1;
        cout << "master thread " << omp_get_thread_num() << endl;
    }
    #pragma omp barrier
    #pragma omp for
    for (int i = 0; i < n; i++)
        b[i] = a;
} /*-- End of the parallel region --*/
```

# Outline for section 11

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

# Synchronization constructs

- ① `#pragma omp barrier` will synchronize all the threads (waiting for other threads to finish their work) before it continues.
- ② Synchronization is useful to avoid data race in the application



```
int *x = new int[10];
#pragma omp parallel num_threads(10)
{
    int tid = omp_get_thread_num();
    x[tid] = my_task(tid);
    cout << " before barrier" << x[tid] << endl;
#pragma omp barrier
    cout << " after barrier" << x[tid] << endl;
}/*-- End of the parallel region --*/
```



# Critical construct

- ① It allows only one thread at a time within the critical section.
  - ② Multiple threads cannot access the critical section in parallel.
  - ③ This prevents race conditions.
- 

```
int local_sum, sum = 0;  
#pragma omp parallel num_threads(5) private(local_sum)  
{  
    for (int i = 0; i < n; i++)  
        local_sum += i;  
#pragma omp critical  
    {  
        sum += local_sum;  
    }  
} /*-- End of the parallel region --*/  
cout << "value of local sum is " << sum << endl;
```

# Atomic constructs

- ① This method is similar to critical sections but offers a performance advantage.
- ② It consists of a single statement followed by the atomic declaration.
- ③ This method safely updates a shared numeric variable and supports memory location updates.

```
int total = 0;  
#pragma omp parallel num_threads(5)  
{  
    for(int i = 0; i < 10; i++)  
    {  
        #pragma omp atomic  
        total++;  
    }  
} /*-- End of the parallel region --*/  
cout << "Total = " << total << endl;
```

# Reduction clause

- ① Useful for incrementing or summing the values of an array into a shared numerical variable.
- ② The syntax for reduction is: `reduction(operators: variable)`
  - Supported arithmetic reductions include: '+, \*, -, max, min'
  - Supported logical operator reductions include: '|, &, ^|'
- ③ This should be used instead of critical sections (which allow only one thread to update the variable at a time).

---

```
omp_set_num_threads(5);
int sum = 0;
#pragma omp parallel for schedule(runtime) reduction(+:sum)
for (int i = 0; i < N; i++)
    sum += a[i];
/*-- End of parallel region --*/
```

# Nowait clause

- ① The `nowait` clause disables the implicit barrier in worksharing constructs.
  - ② It allows for fine-tuning the program's performance.
  - ③ However, it will not affect the `#pragma omp barrier` directive.
- 

```
#pragma omp parallel num_threads(5)
{
#pragma omp for schedule(runtime) nowait
    for (int i = 0; i < 10; ++i)
    {
        cout << " no wait " << endl;
    }
    cout << " after no wait " << endl;
} /*-- End of the parallel region --*/
```

---

# Outline for section 12

- 1 Introduction to Parallel Computing
- 2 Parallel Computer Architecture
- 3 Distributed Memory Architecture
- 4 Shared Memory Architecture
- 5 Hybrid Distributed-Shared Memory
- 6 OpenMP and Compilers
- 7 OpenMP - Parallel Region
- 8 OpenMP - Environment Routines
- 9 Data Sharing
- 10 Work - Sharing Constructs
- 11 Synchronization Constructs
- 12 OpenMP - SIMD

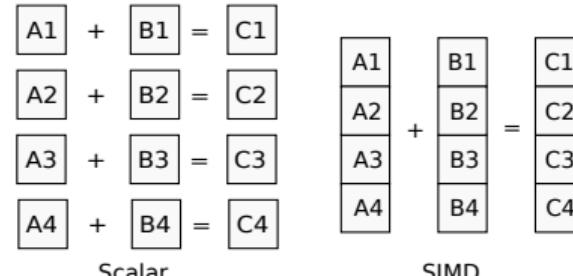


# Single Instruction Multiple Data

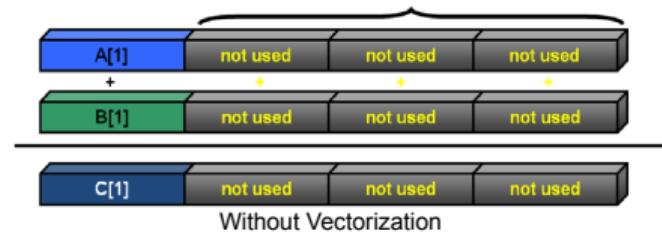
- 1 According to [Flynn's taxonomy](#), computer architecture can be classified into four categories:

- One of these categories is [Single Instruction Multiple Data \(SIMD\)](#),
- where the same instruction is executed on multiple data sets.

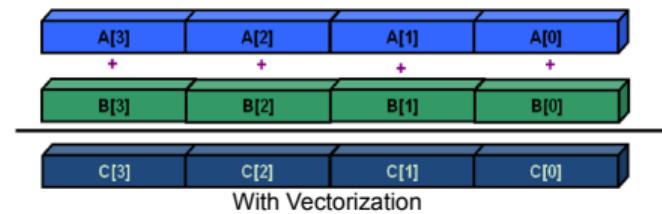
- 2 [OpenMP 4.x](#) directives support the generation of SIMD code generation



e.g. 3 x 32-bit unused integers

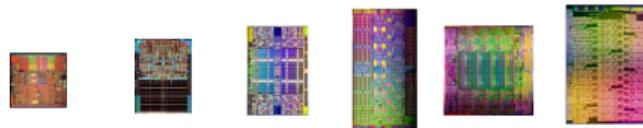


Without Vectorization



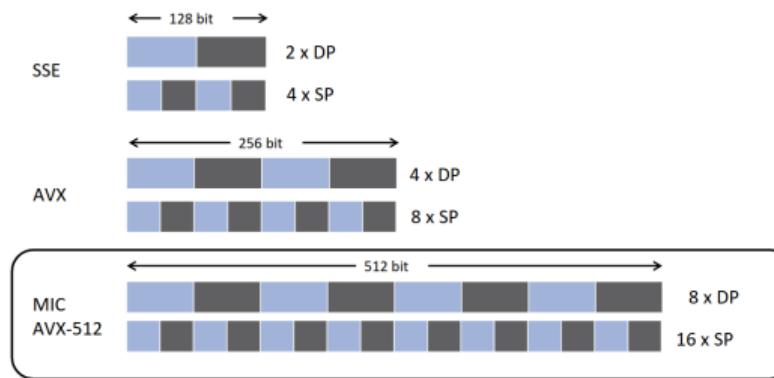
With Vectorization

# Intel's SIMD width



Images not intended to reflect actual die sizes

	64-bit Intel® Xeon® processor	Intel® Xeon® processor 5100 series	Intel® Xeon® processor 5500 series	Intel® Xeon® processor 5600 series	Intel® Xeon® processor E5-2600v2 series	Intel® Xeon Phi™ Co-processor 7120P
Frequency	3.6 GHz	3.0 GHz	3.2 GHz	3.3 GHz	2.7 GHz	1.238 MHz
Core(s)	1	2	4	6	12	61
Thread(s)	2	2	8	12	24	244
SIMD width	128 (2 clock)	128 (1 clock)	128 (1 clock)	128 (1 clock)	256 (1 clock)	512 (1 clock)



# SIMD and Collapse

## Example of SIMD and collapse

---

```
/*-- loop vectorization --*/
#pragma omp simd
for (int i = 0; i < N; i++) {
    a[i] = a[i] + b[i] * c[i];
} /*-- End of SIMD operation --*/
```

---

```
/*-- loop vectorization with scheduling --*/
#pragma omp parallel for simd schedule(static,4)
for (int i = 0; i < N; i++) {
    a[i] = a[i] + b[i] * c[i];
} /*-- End of the parallel region --*/
```

---

```
/*-- SIMD clause with Collapse --*/
#pragma omp simd collapse(2)
for (int i = 0; i < 10; ++i){
    for (int j = 0; j < 10; ++j){
        c[i][j] = a[i][j] / b[i][j];
    }
} /*-- End of the SIMD and collapse region --*/
```

# Reduction and Declare

## Example of reduction and declare

---

```
/*-- SIMD with Reduction --*/
int total = 0;
#pragma omp simd reduction(+:total)
for (int i = 0; i < size; ++i){
    total += a[i] + b[i];}
```

---

```
/*-- SIMD with Declare --/
#pragma omp declare simd
double square (double x){
    return x * x;
}
#pragma omp simd
for (int i = 0; i < n; ++i){
    c[i] = square(a[i]);
}
```

## Example for time measurement in OpenMP application

---

```
double start;
double end;
start = omp_get_wtime();
... work to be timed ...
end = omp_get_wtime();
printf("Work took %f seconds\n", end - start);
```

---

# Profiling and Performance Analysis

## ① ARM Forge

- Arm Forge assists with debugging, profiling, and analysis.
- It supports multiple programming models, including MPI, UPC, CUDA, and OpenMP, across various architectures and compilers.
- Both DDT and MAP offer a graphical user interface (GUI).

## ② Intel

- VTune: Identifies time-consuming sections of code, cache misses, and latency issues.
- Advisor: A collection of tools for gathering metrics and traces, which can be used for further code optimization. It can help analyze and explore ideas for efficient vectorization.
- Intel Inspector: Detects and identifies memory issues, deadlocks, and data races, such as memory access problems and leaks.
- ASP: Used to monitor CPU utilization, memory access efficiency, and vectorization.

## ③ AMD uProf employs a statistical sampling-based approach to collect profiling data, helping to identify performance bottlenecks in applications.

# Important - when using OpenMP

- ① Create a team of parallel threads
  - determine or have an idea about how many threads needed
- ② Properly share the work among the parallel threads
  - choosing appropriate scheduling clause for worksharing
- ③ Identify the shared variables and private variables
  - to avoid race condition
- ④ Finally, use synchronization constructs where it is necessary
  - avoid computational error