

Computing 101

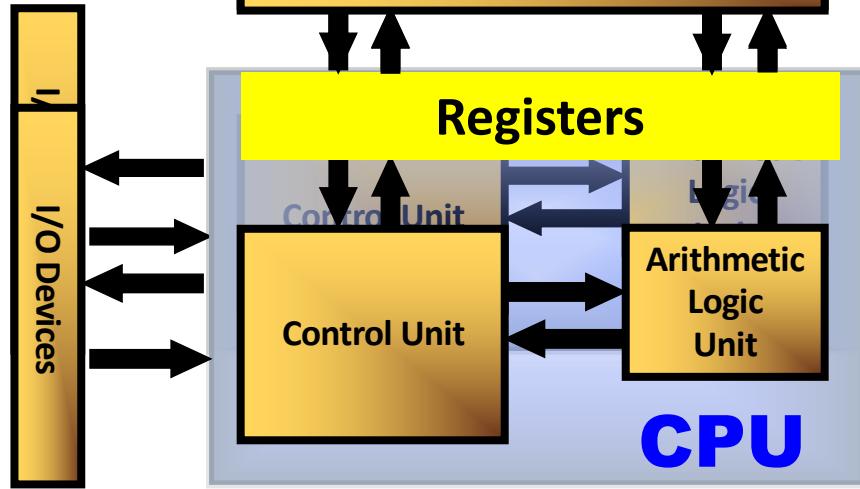
Ivan Girotto (ICTP)
igirotto@ictp.it





The Basic Model

John Von Neumann



```
int main( ){
    int c = 0;
    c = 2 + 3;
}
```

main.c

```
$gcc main.c -o main.x
```

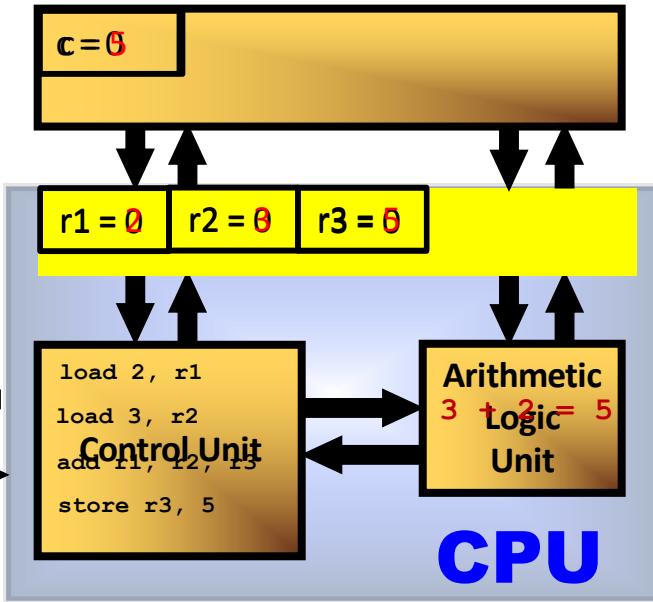
```
$. ./main.x
```

```
load 2, r1
load 3, r2
add r1, r2, r3
store r3, c
```



The Basic Model

John Von Neumann



```
int main( ){
    int c = 0;
    c = 2 + 3;
}
```

main.c

```
$gcc main.c -o main.x
```

```
$. ./main.x
```

```
load 2, r1
load 3, r2
add r1, r2, r3
store r3, c
```

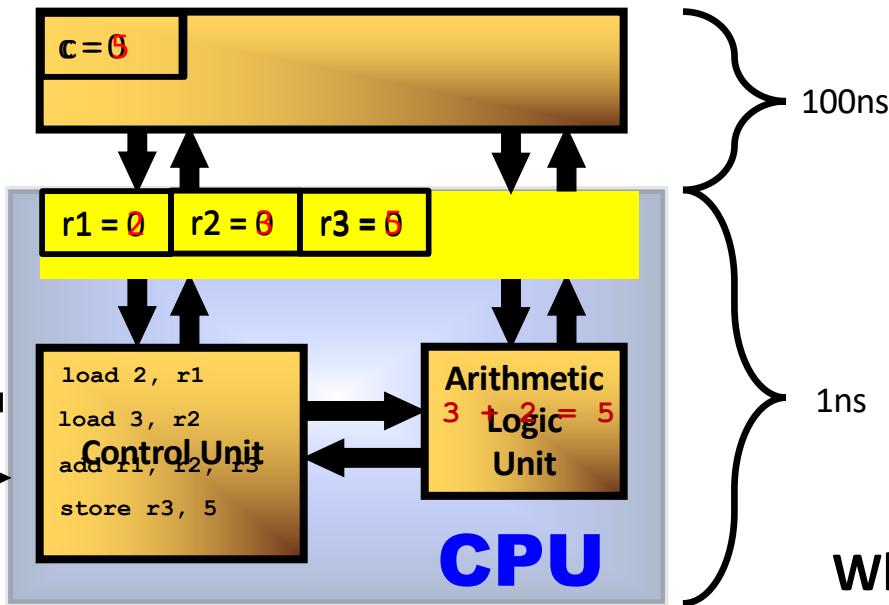
Performance assessment

- Let's suppose memory data transfer is for free ...
- ... and that our processor runs at 1×10^9 hertz (1GHz) => x1 cycle every ns
- ... and that for every clock cycle a single operation is performed
- ... to execute 4 operations, x4 ns (10^{-9}) are required!
- Meaning that at full speed, our processor runs our code executing 10^9 operations x second (100% efficiency)
- If operations are considered double precision FP (64bits), normally this number is referred as FLOP/s
- Unluckily this is not the case...



The Basic Model

John Von Neumann



```
int main( ){
    int c = 0;
    c = 2 + 3;
}
```

main.c

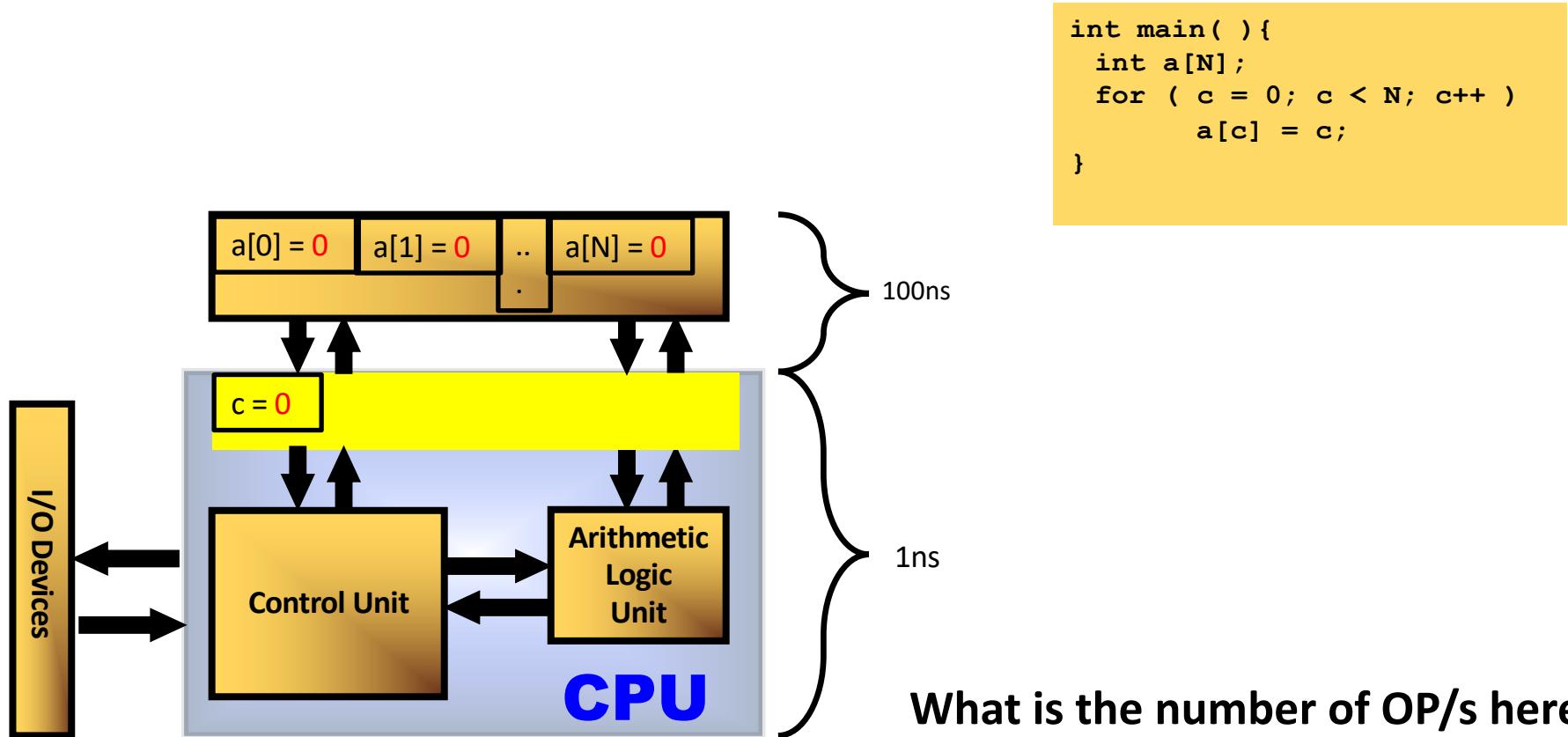
```
$gcc main.c -o main.x
```

```
./main.x
```

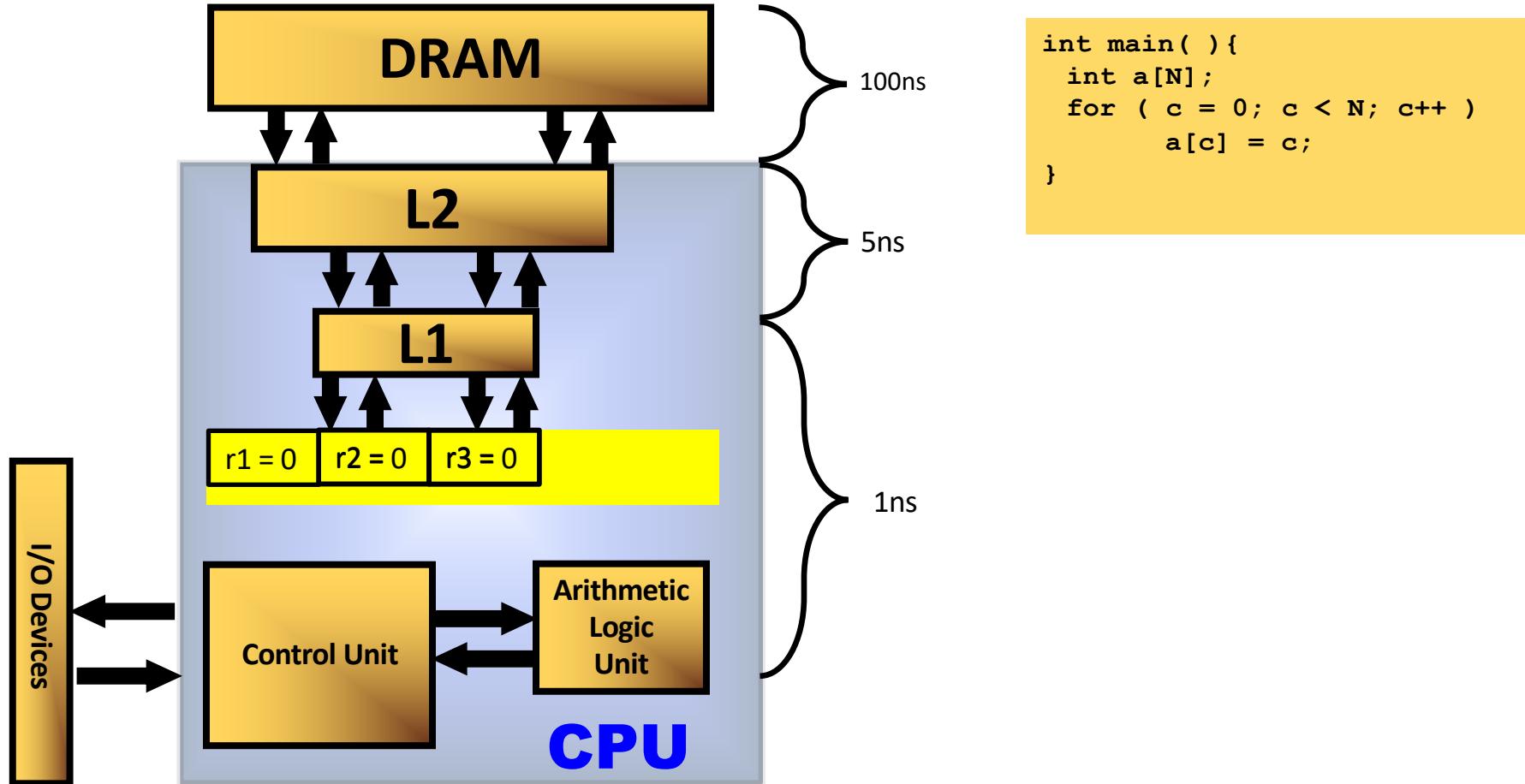
```
load 2, r1
load 3, r2
add r1, r2, r3
store r3, c
```

What is the number of OP/s here?

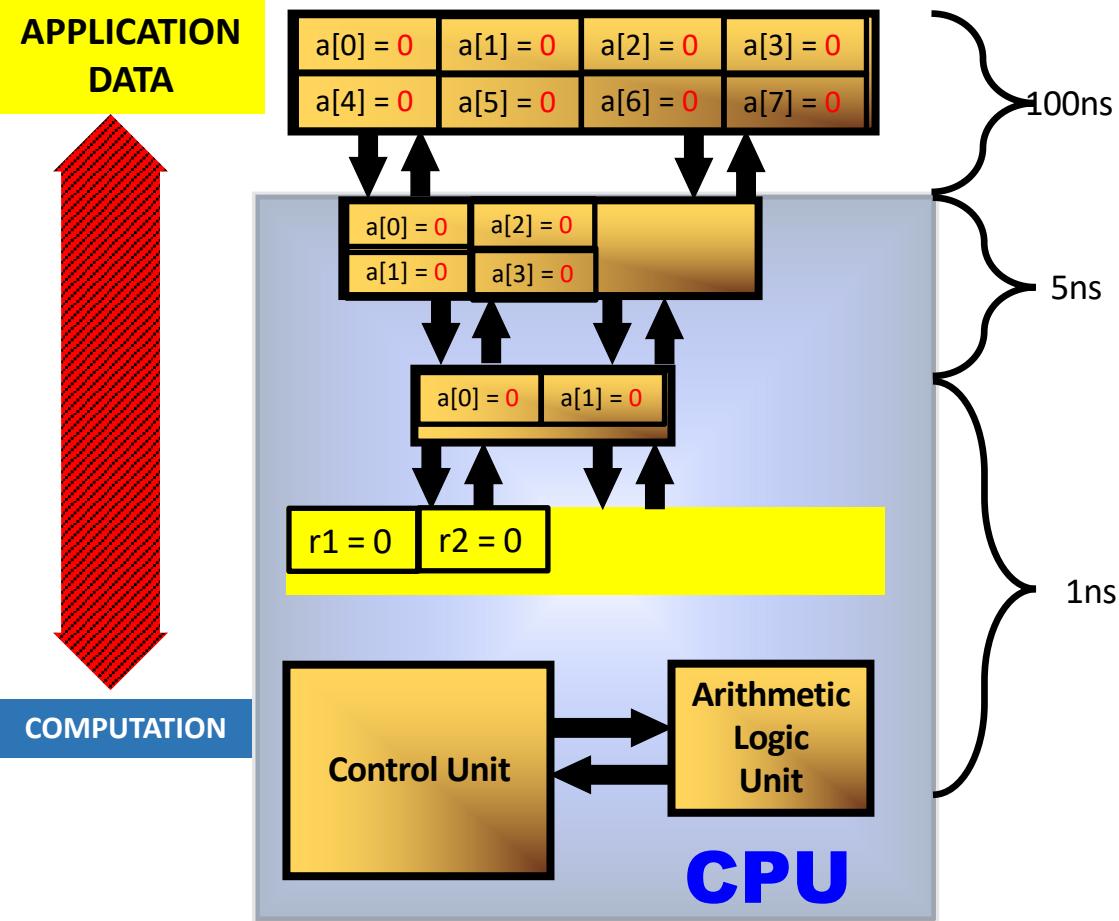
The Basic Model



The Memory Hierarchy



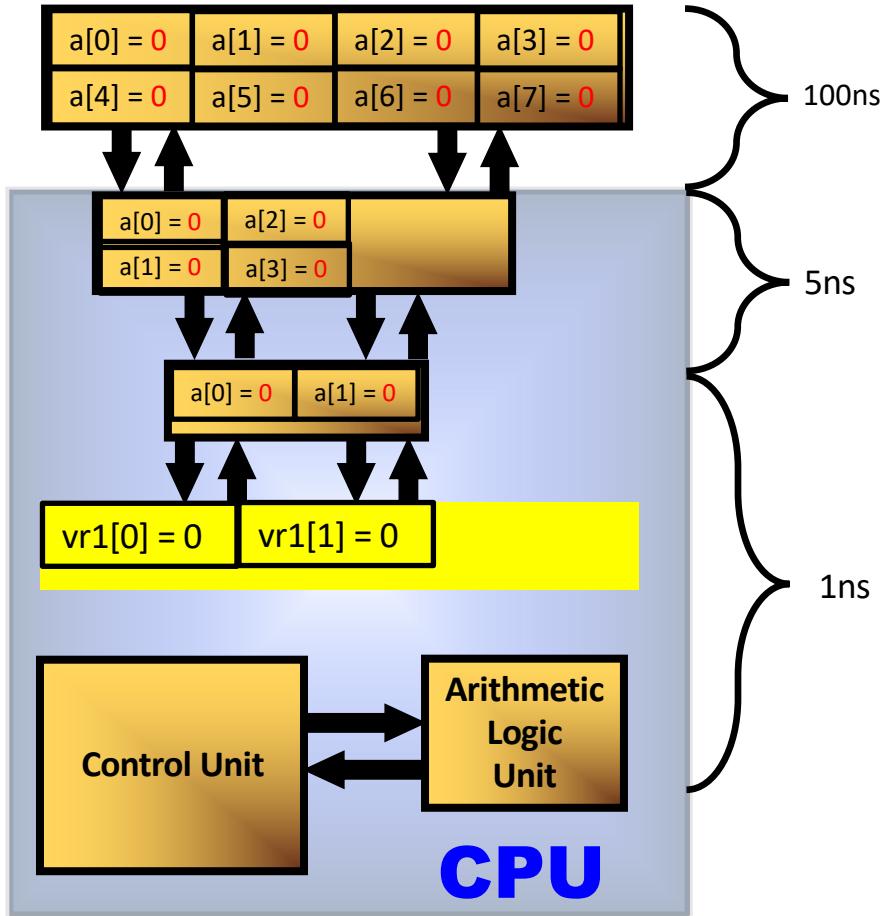
The Memory Hierarchy



```
int main( ) {
    int a[N];
    for ( c = 0; c < N; c++ )
        a[c] = c;
}
```

- Data is organized into blocks of fixed size, called *cache lines*.
- Operation of LOAD/STORE can lead to two different scenarios:
 - *cache hit*
 - *cache miss*

Vectorization



```
int main( ) {
    int a[N];
    for ( c = 0; c < N; c++ )
        a[c] = c;
}
```

Registers are capable to store more than one primitive type at a time

The vectorial ALU is capable to perform the same operation over multiple data of the same kind

Scalar Mode

$$\begin{array}{c} A \\ + \\ B \\ = \\ A+B \end{array}$$

Wrap up on modern computers

- Efficient data memory access patterns are crucial to best performances
 - design and development of computer algorithms that maximize data locality in time and space (data distance)
- Modern CPUs are equipped with vector registers and ALUs capable of multiple concurrent operations on multiple data
 - vectorization is aided by compilers, but requiring design and development of computer codes that maximize simple operations on contiguous data of the same type
- When all CPU component work at maximum speed that is called *peak of performance* (Floating point operations per seconds FLOP/s)
 - Tech-spec normally describe the theoretical peak
 - Benchmarks measure the real peak
 - Applications show the real performance value
- Real performance are also mostly related to the memory bandwidth (GBytes/s)

```
$gcc main.c -o main.x
```

```
./main.x & ./main.x
```



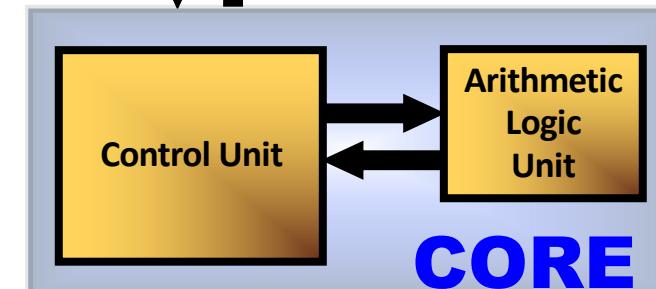
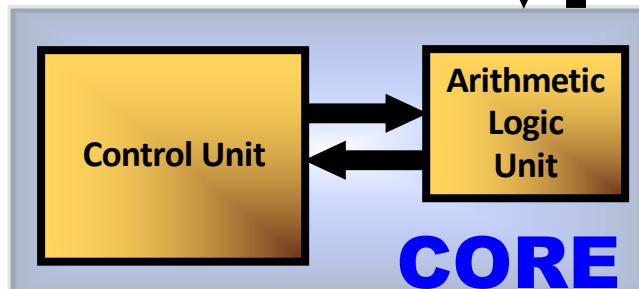
```
int main() {
    int a[N];
    for (c = 0; c < N; c++)
        a[c] = c;
}
```

```
$gcc main.c -o main.x -fopenmp
```

```
$export OMP_NUM_THREADS=2
$./main.x
```

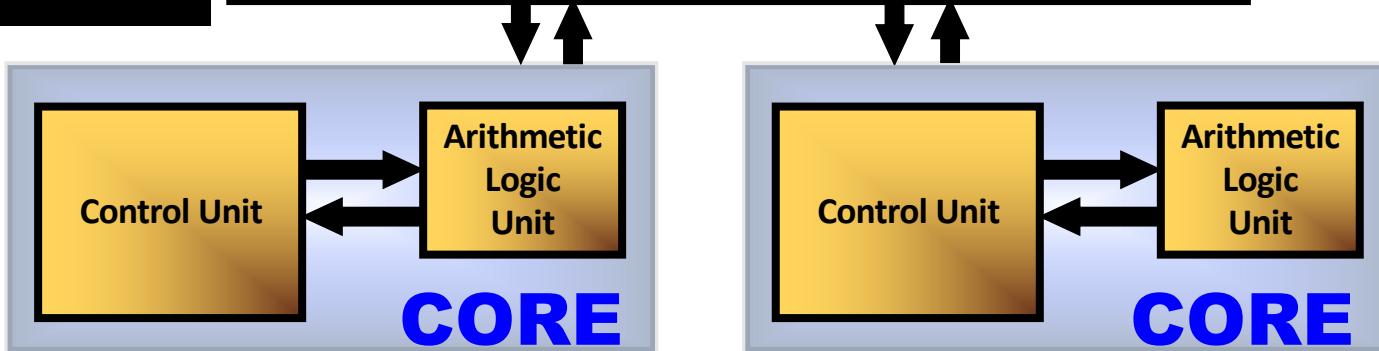


```
int main() {
    int a[N];
    #pragma omp parallel for
    for (c = 0; c < N; c++)
        a[c] = c;
}
```

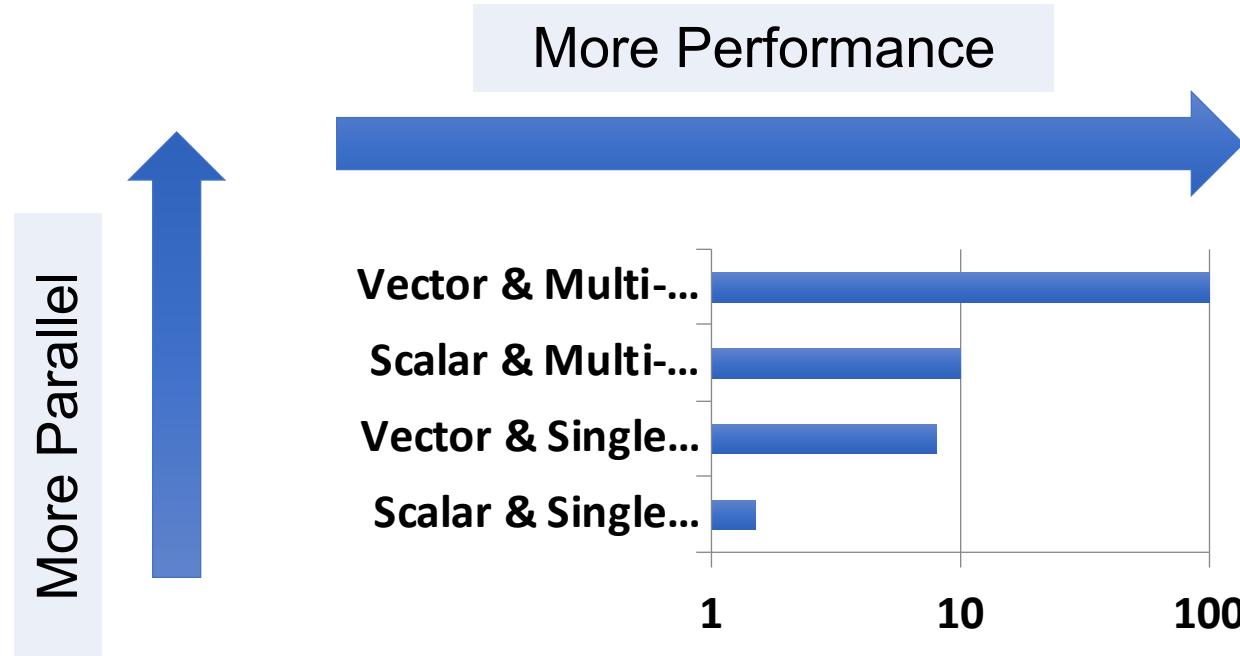


```
int main( ){
    int a[N], sum;
    #pragma omp parallel for
    for ( c = 0; c < N; c++ )
        sum += a[c];
}
```

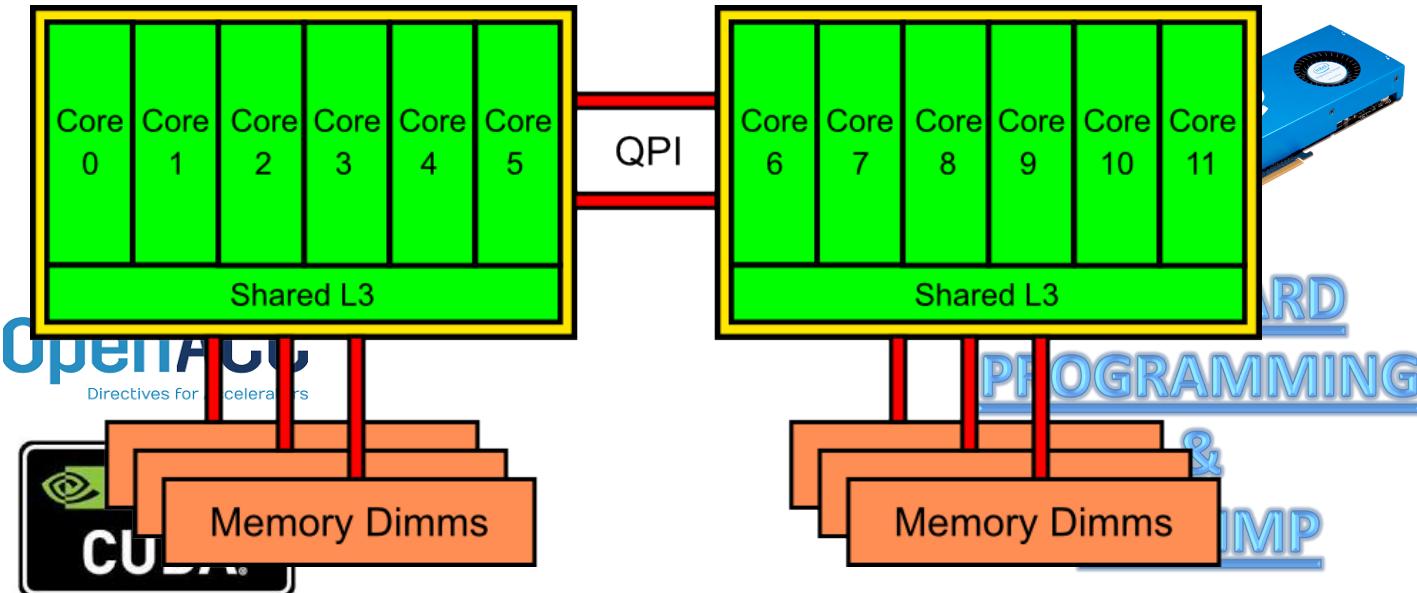
```
$gcc main.c -o main.x -fopenmp
$export OMP_NUM_THREADS=2
$./main.x
```



Threading and Vectorization

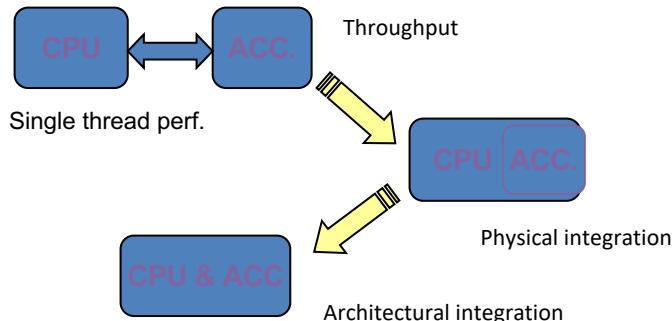


Multiple Socket CPUs + Accelerators

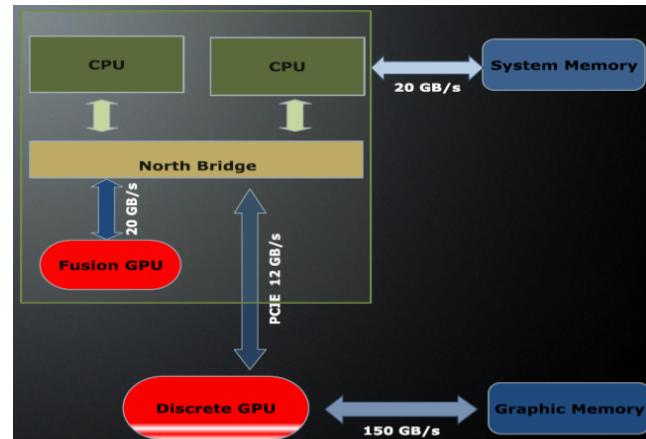


Accelerated co-Processors

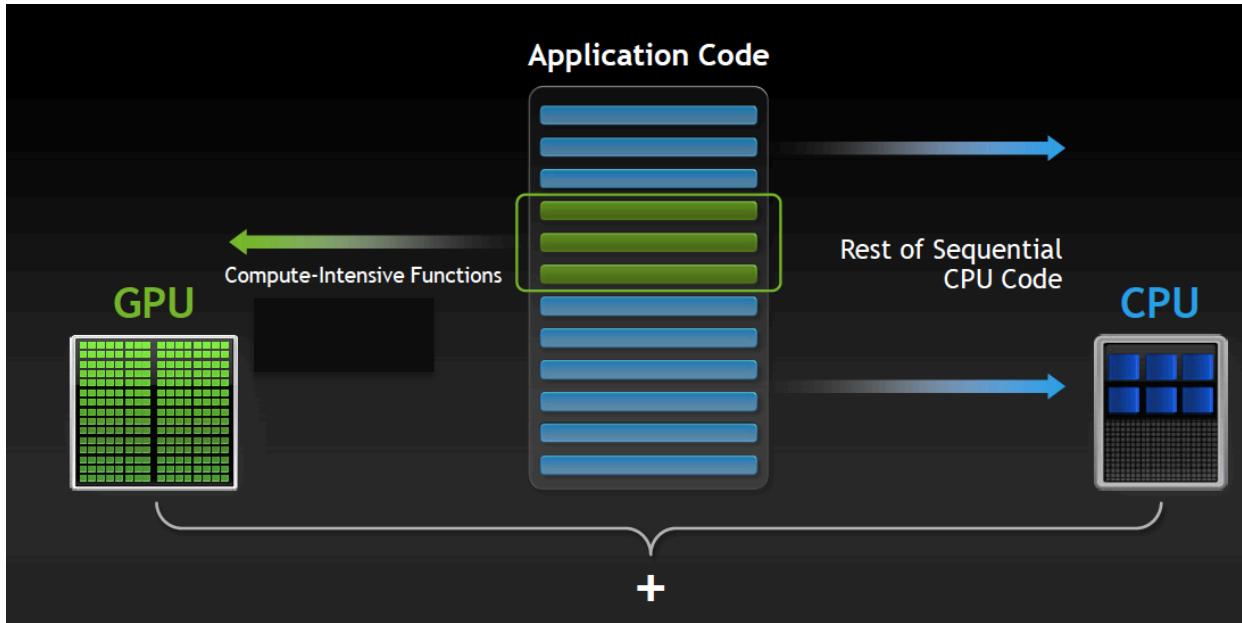
- A set of simplified execution units that can perform few operations (with respect to standard CPU) with very high efficiency. When combined with full featured CPU can accelerate the “nominal” speed of a system.

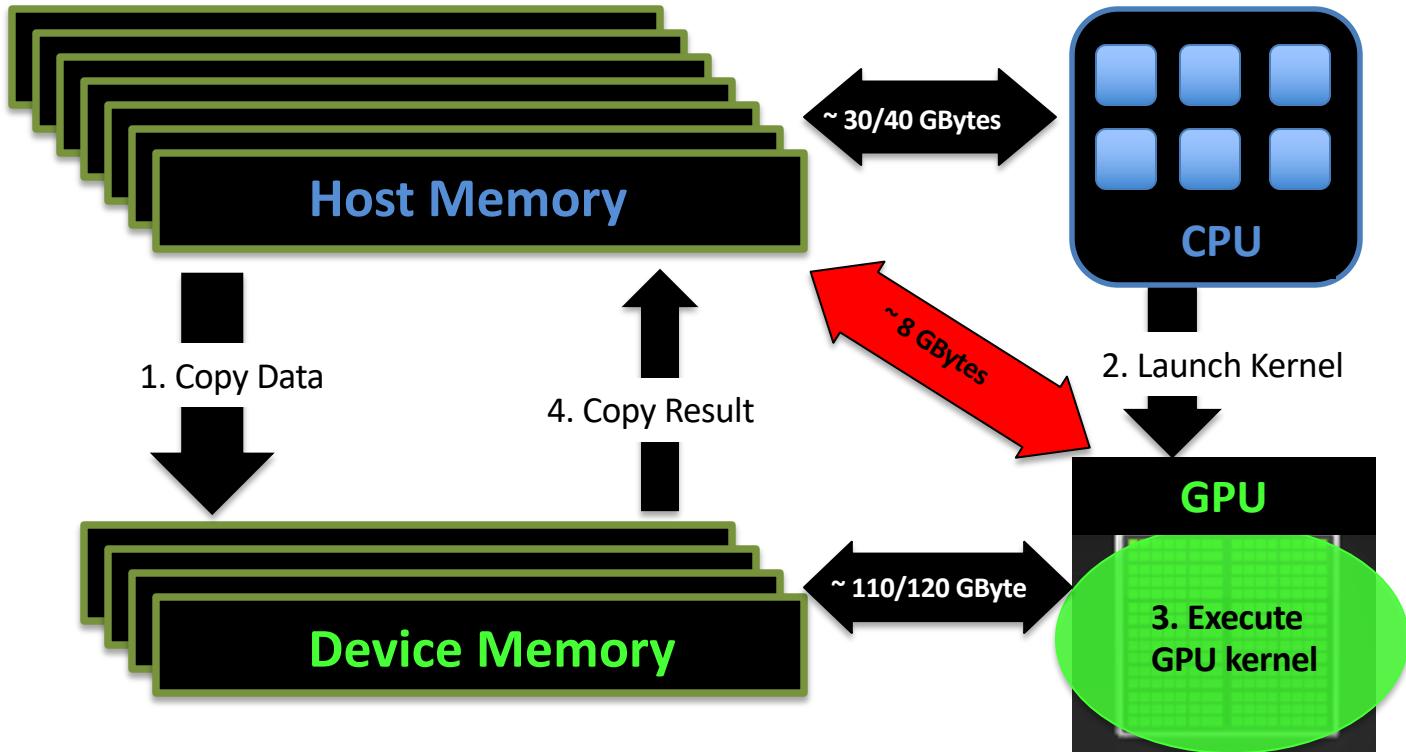


- Main approaches to accelerators:
 - Task Parallelism (MIMD) → MIC
 - Data Parallelism (SIMD) → GPU

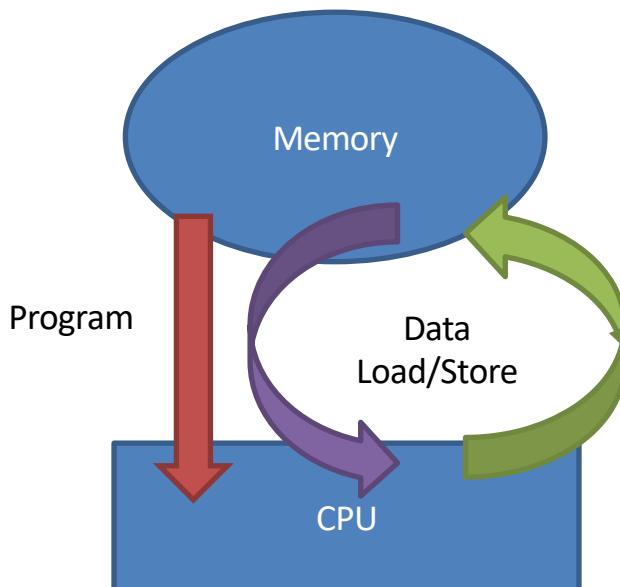


The General Concept of Accelerated Computing

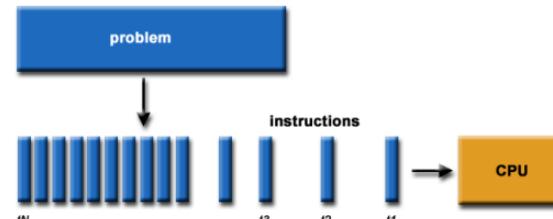




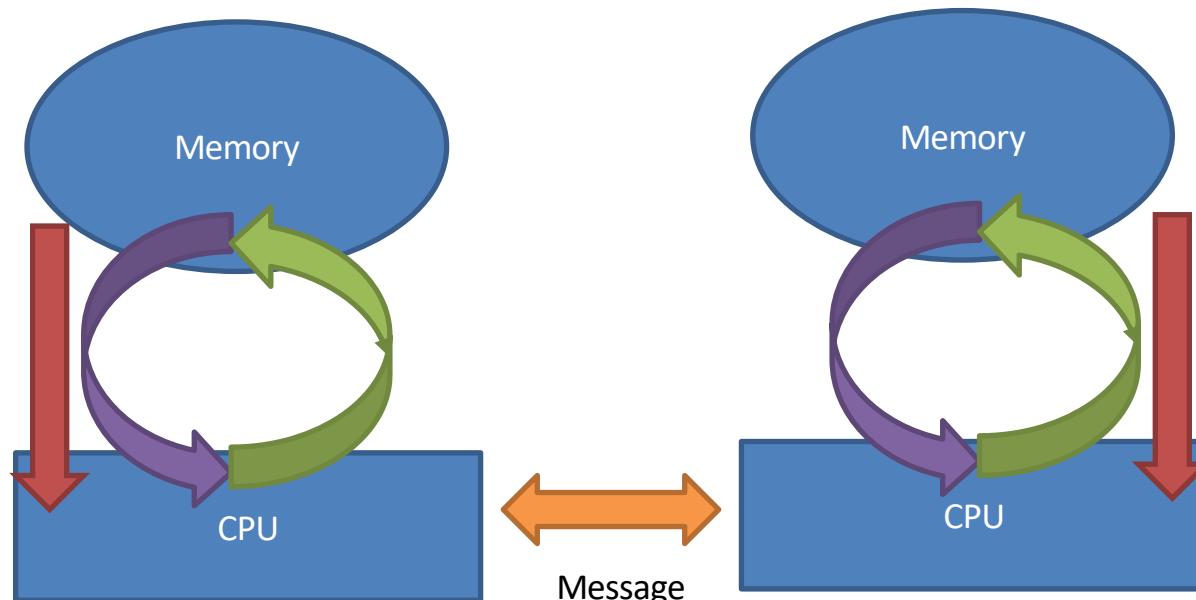
Serial Programming



A problem is broken into a discrete series of instructions.
Instructions are executed one after another.
Only one instruction may execute at any moment in time.

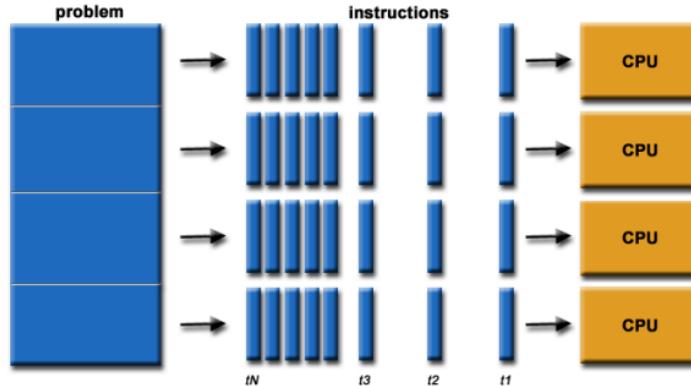


Parallel Programming



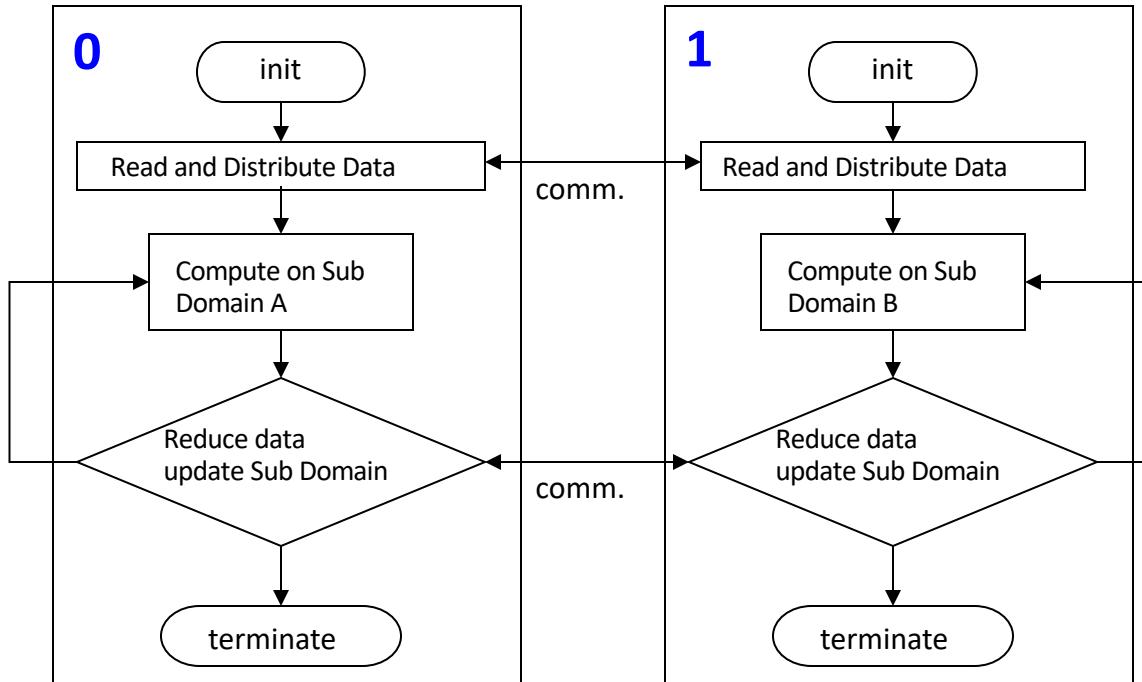
Concurrency

The first step in developing a parallel algorithm is to decompose the problem into tasks that can be executed concurrently



- A problem is broken into discrete parts that can be solved concurrently
- Each part is further broken down to a series of instructions
- Instructions from each part execute simultaneously on different processors
- An overall control / coordination mechanism is employed

What is a Parallel Program





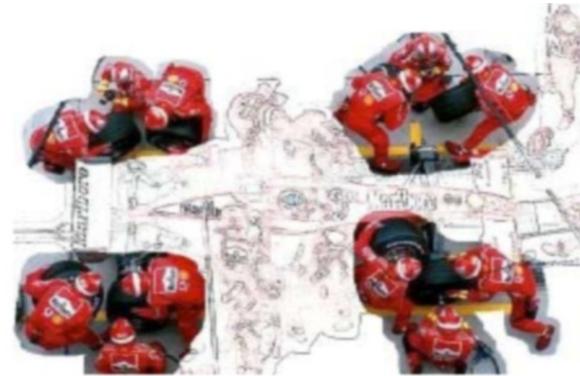
Fundamental Steps of Parallel Design

- Identify portions of the work that can be performed concurrently
- Mapping the concurrent pieces of work onto multiple processes running in parallel
- Distributing the input, output and intermediate data associated within the program
- Managing accesses to data shared by multiple processors
- Synchronizing the processors at various stages of the parallel program execution

Type of Parallelism

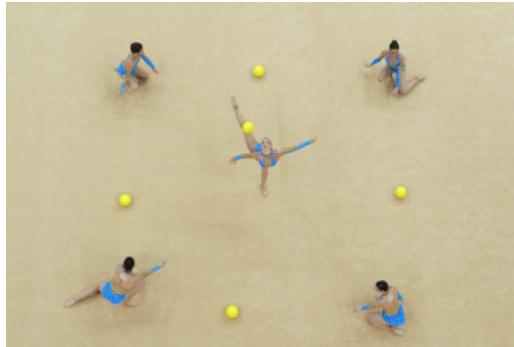
- **Functional (or task) parallelism:**
different people are performing different task at the same time

- **Data Parallelism:**
different people are performing the same task, but on different equivalent and independent objects



Process Interactions

- The effective speed-up obtained by the parallelization depend by the amount of overhead we introduce making the algorithm parallel
- There are mainly two key sources of overhead:
 1. Time spent in inter-process interactions (**communication**)
 2. Time some process may spent being idle (**synchronization**)

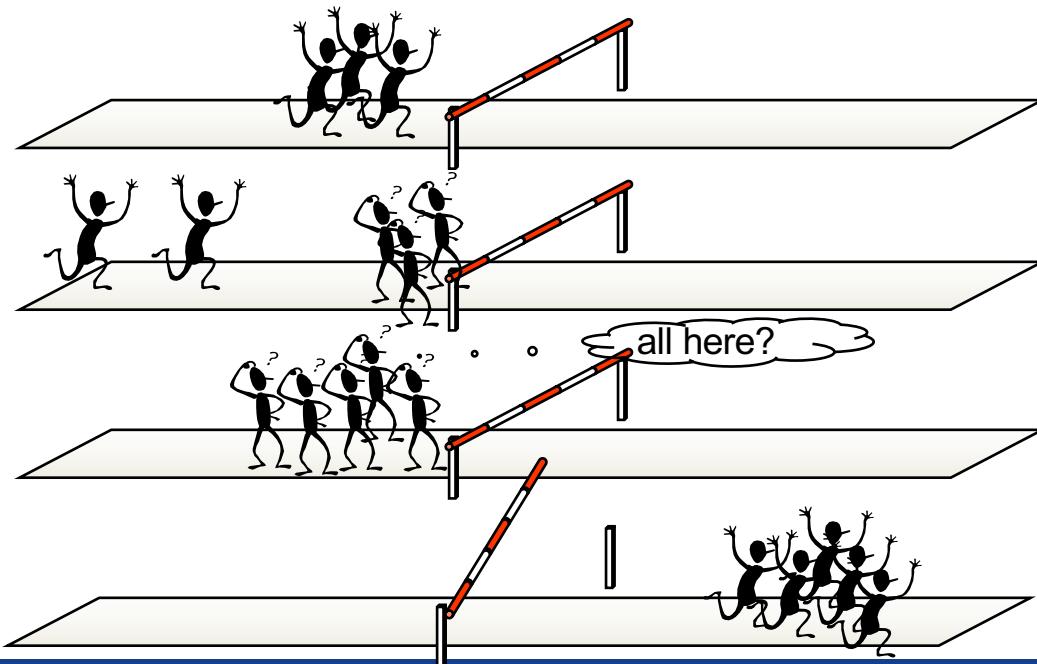




Load Balancing

- Equally divide the work among the available resource: processors, memory, network bandwidth, I/O, ...
- This is usually a simple task for the problem decomposition model
- It is a difficult task for the functional decomposition model

Effect of Load Unbalancing





Minimizing Communication

- When possible reduce the communication events:
 - group lots of small communications into large one
 - eliminate synchronizations as much as possible. Each synchronization levels off the performance to that of the slowest process



Overlap Communication and Computation

- When possible code your program in such a way that processes continue to do useful work while communicating
- This is usually a non trivial task and is afforded in the very last phase of parallelization
- If you succeed, you have done. Benefits are enormous

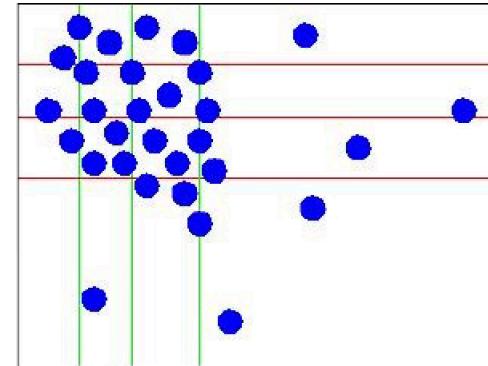


Granularity

- Granularity is determined by the decomposition level (number of task) on which we want divide the problem
- The degree to which task/data can be subdivided is limit to concurrency and parallel execution
- Parallelization has to become “topology aware”
 - coarse grain and fine grained parallelization has to be mapped to the topology to reduce memory and I/O contention
 - make your code modularized to enhance different levels of granularity and consequently to become more “platform adaptable”

Limitations of Parallel Computing

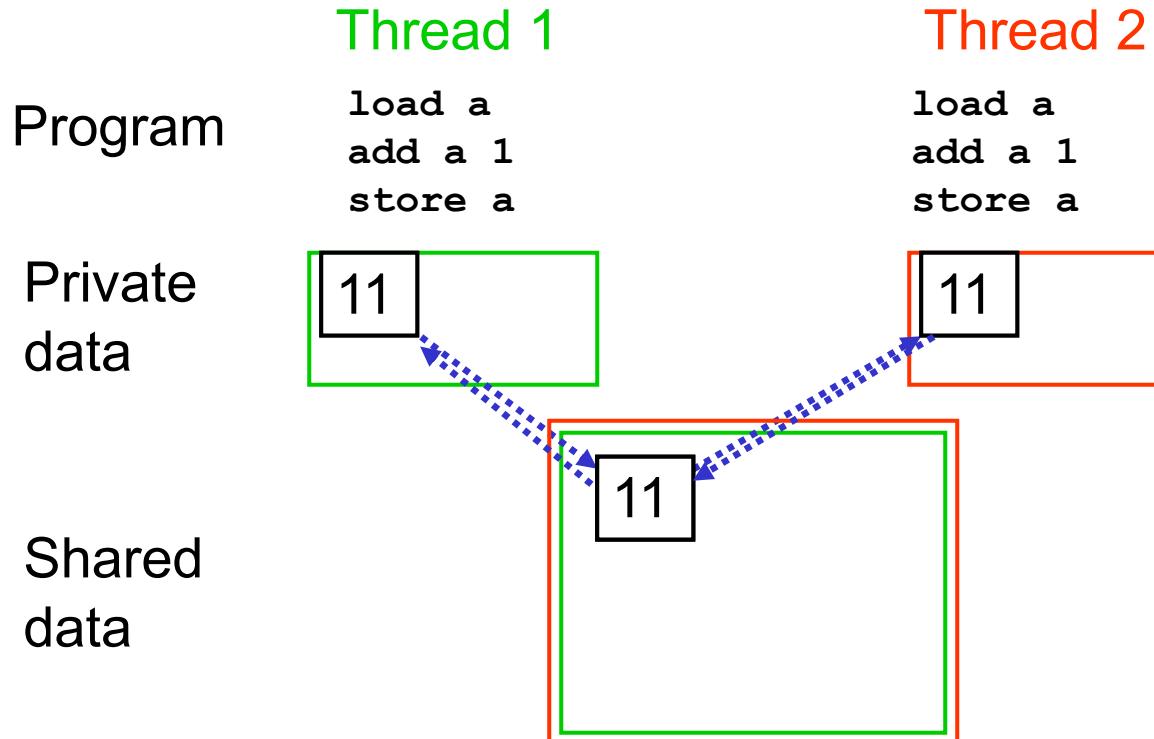
- Fraction of serial code limits parallel speedup
- Degree to which tasks/data can be subdivided is limit to concurrency and parallel execution
- Load imbalance:
 - parallel tasks have a different amount of work
 - CPUs are partially idle
 - redistributing work helps but has limitations
 - communication and synchronization overhead





Shared Resources

- In parallel programming, developers must manage exclusive access to shared resources
- Resources are in different forms:
 - concurrent read/write (including parallel write) to shared memory locations
 - concurrent read/write (including parallel write) to shared devices
 - a message that must be send and received





Fundamental Tools of Parallel Programming





MPI Program Design

- Multiple and separate processes (can be local and remote) concurrently that are coordinated and exchange data through “messages”
=> a “share nothing” parallelization
- Best for coarse grained parallelization
- Distribute large data sets; replicate small data
- Minimize communication or overlap communication and computing for efficiency
=> Amdahl's law: speedup is limited by the fraction of serial code plus communication



Phases of an MPI Program

1. Startup
 - Parse arguments (mpirun may add some!)
 - Identify parallel environment and rank of process
 - Read and distribute all data
2. Execution
 - Proceed to subroutine with parallel work (can be same or different for all parallel tasks)
3. Cleanup

CAUTION: this sequence may be run only once



```
program bcast
```

```
implicit none
```

```
include "mpif.h"
```

```
integer :: myrank, ncpus, imesg, ierr  
integer, parameter :: comm = MPI_COMM_WORLD
```

```
call MPI_INIT(ierr)
```

```
call MPI_COMM_RANK(comm, myrank, ierr)  
call MPI_COMM_SIZE(comm, ncpus, ierr)
```

```
imesg = myrank
```

```
print *, "Before Bcast operation I'm ", myrank, &  
        " and my message content is ", imesg
```

```
call MPI_BCAST(imesg, 1, MPI_INTEGER, 0, comm,  
ierr)
```

```
print *, "After Bcast operation I'm ", myrank, &  
        " and my message content is ", imesg
```

```
call MPI_FINALIZE(ierr)
```

```
end program bcast
```

```
program bcast
```

```
implicit none
```

```
include "mpif.h"
```

```
integer :: myrank, ncpus, imesg, ierr  
integer, parameter :: comm = MPI_COMM_WORLD
```

P₀

```
myrank = ??  
ncpus = ??  
imesg = ??  
ierr = ??  
comm =  
MPI_C...
```

P₁

```
myrank = ??  
ncpus = ??  
imesg = ??  
ierr = ??  
comm =  
MPI_C...
```

P₂

```
myrank = ??  
ncpus = ??  
imesg = ??  
ierr = ??  
comm =  
MPI_C...
```

P₃

```
myrank = ??  
ncpus = ??  
imesg = ??  
ierr = ??  
comm =  
MPI_C...
```

```
end program bcast
```

```
program bcast
```

```
implicit none
```

```
include "mpif.h"
```

```
integer :: myrank, ncpus, imesg, ierr
integer, parameter :: comm = MPI_COMM_WORLD
```

```
call MPI_INIT(ierr)
```

```
end program bcast
```

P₀

```
myrank = ??  

ncpus = ??  

imesg = ??  

ierr =  

MPI_SUC...  

comm =
```

P₂

```
myrank = ??  

ncpus = ??  

imesg = ??  

ierr =  

MPI_SUC...  

comm =  

MPI_C...
```

P₁

```
myrank = ??  

ncpus = ??  

imesg = ??  

ierr =  

MPI_SUC...  

comm =
```

P₃

```
myrank = ??  

ncpus = ??  

imesg = ??  

ierr =  

MPI_SUC...  

comm =  

MPI_C...
```

```
program bcast
    implicit none
    include "mpif.h"
    integer :: myrank, ncpus, imesg, ierr
    integer, parameter :: comm = MPI_COMM_WORLD
    call MPI_INIT(ierr)
    call MPI_COMM_SIZE(comm, ncpus, ierr)
    call MPI_COMM_RANK(comm, myrank, ierr)
```

P₀

```
myrank = ??  

ncpus = 4  

imesg = ??  

ierr =  

MPI_SUC...  

comm =
```

P₂

```
myrank = ??  

ncpus = 4  

imesg = ??  

ierr =  

MPI_SUC...  

comm =  

MPI_C...
```

P₁

```
myrank = ??  

ncpus = 4  

imesg = ??  

ierr =  

MPI_SUC...  

comm =
```

P₃

```
myrank = ??  

ncpus = 4  

imesg = ??  

ierr =  

MPI_SUC...  

comm =  

MPI_C...
```

```
program bcast
    implicit none
    include "mpif.h"
    integer :: myrank, ncpus, imesg, ierr
    integer, parameter :: comm = MPI_COMM_WORLD
    call MPI_INIT(ierr)
    call MPI_COMM_SIZE(comm, ncpus, ierr)
    call MPI_COMM_RANK(comm, myrank, ierr)
```

P₀

```
myrank = 0
ncpus = 4
imesg = ???
ierr =
MPI_SUC...
comm =
```

P₂

```
myrank = 2
ncpus = 4
imesg = ???
ierr =
MPI_SUC...
comm =
MPI_C...
```

P₁

```
myrank = 1
ncpus = 4
imesg = ???
ierr =
MPI_SUC...
comm =
```

P₃

```
myrank = 3
ncpus = 4
imesg = ???
ierr =
MPI_SUC...
comm =
MPI_C...
```

```

program bcast
implicit none
include "mpif.h"
integer :: myrank, ncpus, imesg, ierr
integer, parameter :: comm = MPI_COMM_WORLD
call MPI_INIT(ierr)
call MPI_COMM_RANK(comm, myrank, ierr)
call MPI_COMM_SIZE(comm, ncpus, ierr)
imesg = myrank
print *, "Before Bcast operation I'm ", myrank, &
        " and my message content is ", imesg
end program bcast
    
```

P₀

```

myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
    
```

P₂

```

myrank = 2
ncpus = 4
imesg = 2
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

P₁

```

myrank = 1
ncpus = 4
imesg = 1
ierr =
MPI_SUC...
comm =
    
```

P₃

```

myrank = 3
ncpus = 4
imesg = 3
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

```
program bcast
    implicit none
    include "mpif.h"
    integer :: myrank, ncpus, imesg, ierr
    integer, parameter :: comm = MPI_COMM_WORLD
    call MPI_INIT(ierr)
    call MPI_COMM_RANK(comm, myrank, ierr)
    call MPI_COMM_SIZE(comm, ncpus, ierr)
    imesg = myrank
    print *, "Before Bcast operation I'm ", myrank, &
              " and my message content is ", imesg
    call MPI_BCAST(imesg, 1, MPI_INTEGER, 0, comm, ierr)
```

P₀

```
myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
```

P₂

```
myrank = 2
ncpus = 4
imesg = 2
ierr =
MPI_SUC...
comm =
MPI_C...
```

P₁

```
myrank = 1
ncpus = 4
imesg = 1
ierr =
MPI_SUC...
comm =
```

P₃

```
myrank = 3
ncpus = 4
imesg = 3
ierr =
MPI_SUC...
comm =
MPI_C...
```

```
call MPI_BCAST( imesg, 1, MPI_INTEGER, 0, comm, ierr )
```

P₀

```
myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
```

P₁

```
myrank = 1
ncpus = 4
imesg = 1
ierr =
MPI_SUC...
comm =
```

P₂

```
myrank = 2
ncpus = 4
imesg = 2
ierr =
MPI_SUC...
comm =
```

P₃

```
myrank = 3
ncpus = 4
imesg = 3
ierr =
MPI_SUC...
comm =
```





```
call MPI_BCAST( imesg, 1, MPI_INTEGER, 0, comm, ierr )
```

P₀

```
myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
```

P₁

```
myrank = 1
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
```

P₂

```
myrank = 2
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
```

P₃

```
myrank = 3
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
```

```

program bcast
implicit none
include "mpif.h"
integer :: myrank, ncpus, imesg, ierr
integer, parameter :: comm = MPI_COMM_WORLD
call MPI_INIT(ierr)
call MPI_COMM_RANK(comm, myrank, ierr)
call MPI_COMM_SIZE(comm, ncpus, ierr)
imesg = myrank
print *, "Before Bcast operation I'm ", myrank, &
        " and my message content is ", imesg
call MPI_BCAST(imesg, 1, MPI_INTEGER, 0, comm,
               ierr)
print *, "After Bcast operation I'm ", myrank, &
        " and my message content is ", imesg
end program bcast
    
```

P₀

```

myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

P₂

```

myrank = 2
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

P₁

```

myrank = 1
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

P₃

```

myrank = 3
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

```

program bcast

implicit none

include "mpif.h"

integer :: myrank, ncpus, imesg, ierr
integer, parameter :: comm = MPI_COMM_WORLD

call MPI_INIT(ierr)
call MPI_COMM_RANK(comm, myrank, ierr)
call MPI_COMM_SIZE(comm, ncpus, ierr)

imesg = myrank
print *, "Before Bcast operation I'm ", myrank, &
        " and my message content is ", imesg

call MPI_BCAST(imesg, 1, MPI_INTEGER, 0, comm,
               ierr)

print *, "After Bcast operation I'm ", myrank, &
        " and my message content is ", imesg

call MPI_FINALIZE(ierr)

end program bcast

```

P₀

```

myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =

```

P₂

```

myrank = 2
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =

```

P₁

```

myrank = 1
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =

```

P₃

```

myrank = 3
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =

```

```

program bcast
implicit none
include "mpif.h"
integer :: myrank, ncpus, imesg, ierr
integer, parameter :: comm = MPI_COMM_WORLD
call MPI_INIT(ierr)
call MPI_COMM_RANK(comm, myrank, ierr)
call MPI_COMM_SIZE(comm, ncpus, ierr)
imesg = myrank
print *, "Before Bcast operation I'm ", myrank, &
        " and my message content is ", imesg
call MPI_BCAST(imesg, 1, MPI_INTEGER, 0, comm,
ierr)
print *, "After Bcast operation I'm ", myrank, &
        " and my message content is ", imesg
call MPI_FINALIZE(ierr)
end program bcast
    
```

P₀

```

myrank = 0
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
    
```

P₂

```

myrank = 2
ncpus = 4
imesg = 0
ierr = MPI_SUCC
comm =
MPI_C...
    
```

P₁

```

myrank = 1
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
    
```

P₃

```

myrank = 3
ncpus = 4
imesg = 0
ierr =
MPI_SUC...
comm =
MPI_C...
    
```

Programming Parallel Paradigms

- Are the tools we use to express the parallelism for on a given architecture (see also SPMD, SIMD, etc...)
- They differ in how programmers can manage and define key features like:
 - parallel regions
 - concurrency
 - process communication
 - synchronism





Parallelism - 101

- there are two main reasons to write a parallel program:
 - access to larger amount of memory (aggregated, going bigger)
 - reduce time to solution (going faster)

Static Data Partitioning

**The simplest data decomposition schemes for dense matrices are
1-D block distribution schemes.**

row-wise distribution

P_0
P_1
P_2
P_3
P_4
P_5
P_6
P_7

column-wise distribution

P_0	P_1	P_2	P_3	P_4	P_5	P_6	P_7
-------	-------	-------	-------	-------	-------	-------	-------

Distributed Data Vs Replicated Data

- Replicated data distribution is useful if it helps to reduce the communication among process at the cost of bounding scalability
- Distributed data is the ideal data distribution but not always applicable for all data-sets
- Usually complex application are a mix of those techniques => distribute large data sets; replicate small data

Global Vs Local Indexes

- In sequential code you always refer to global indexes
- With distributed data you must handle the distinction between global and local indexes (and possibly implementing utilities for transparent conversion)

Local Idx	<table border="1"><tr><td>1</td><td>2</td><td>3</td></tr></table>	1	2	3	<table border="1"><tr><td>1</td><td>2</td><td>3</td></tr></table>	1	2	3	<table border="1"><tr><td>1</td><td>2</td><td>3</td></tr></table>	1	2	3
1	2	3										
1	2	3										
1	2	3										
<hr/>												
Global Idx	<table border="1"><tr><td>1</td><td>2</td><td>3</td></tr></table>	1	2	3	<table border="1"><tr><td>4</td><td>5</td><td>6</td></tr></table>	4	5	6	<table border="1"><tr><td>7</td><td>8</td><td>9</td></tr></table>	7	8	9
1	2	3										
4	5	6										
7	8	9										

Block Array Distribution Schemes

Block distribution schemes can be generalized to higher dimensions as well.

P_0	P_1	P_2	P_3
P_4	P_5	P_6	P_7
P_8	P_9	P_{10}	P_{11}
P_{12}	P_{13}	P_{14}	P_{15}

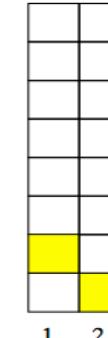
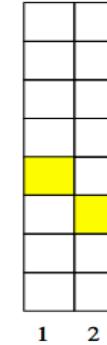
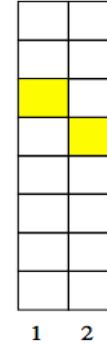
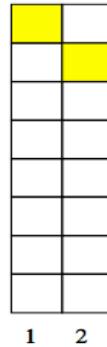
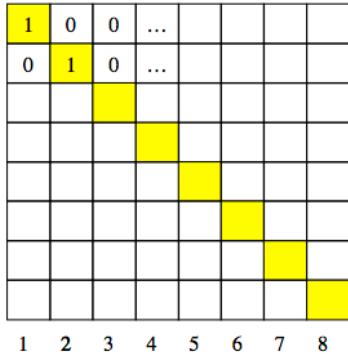
(a)

P_0	P_1	P_2	P_3	P_4	P_5	P_6	P_7
P_8	P_9	P_{10}	P_{11}	P_{12}	P_{13}	P_{14}	P_{15}

(b)

Degree to which tasks/data can be subdivided is limit to concurrency and parallel execution!!

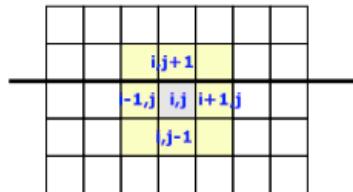
Collaterals to Domain Decomposition /1



**Are all the domain's dimensions
always multiple of the number
of tasks/processes we are
willing to use?**

Again on Domain Decomposition

sub-domain boundaries

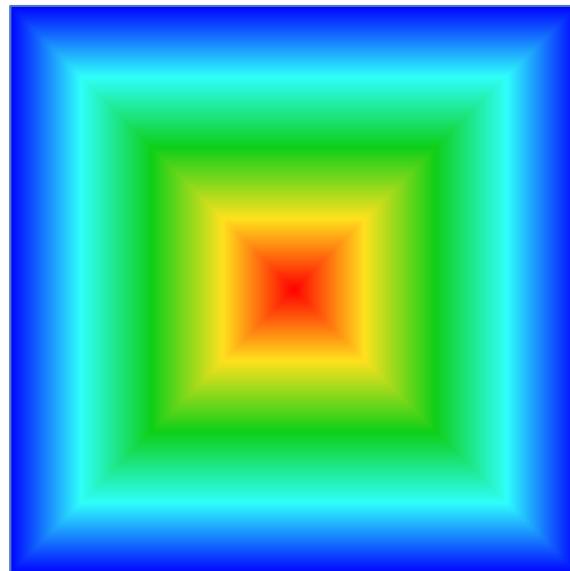




The Abdus Salam
**International Centre
for Theoretical Physics**



P_0



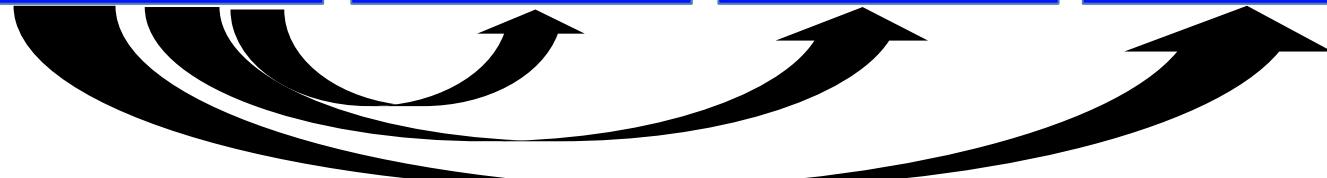
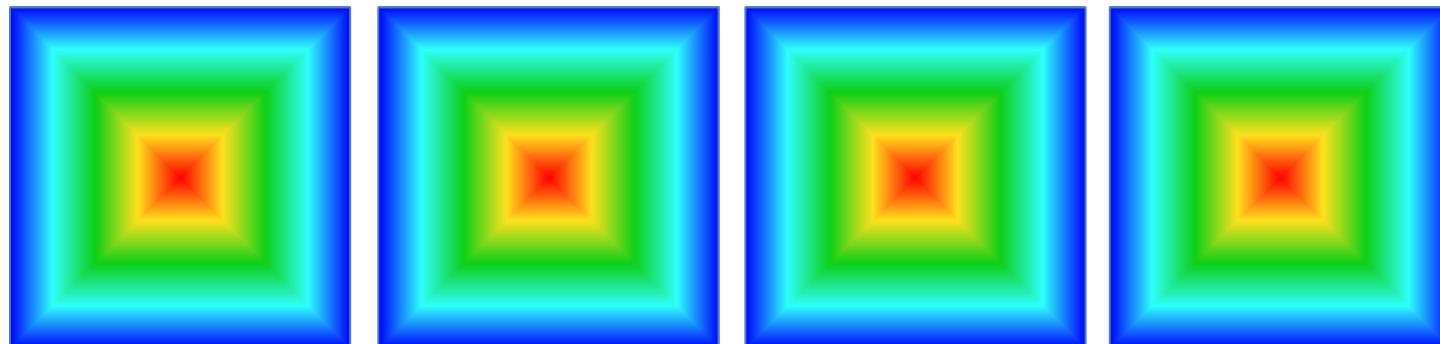
call MPI_BCAST(...)

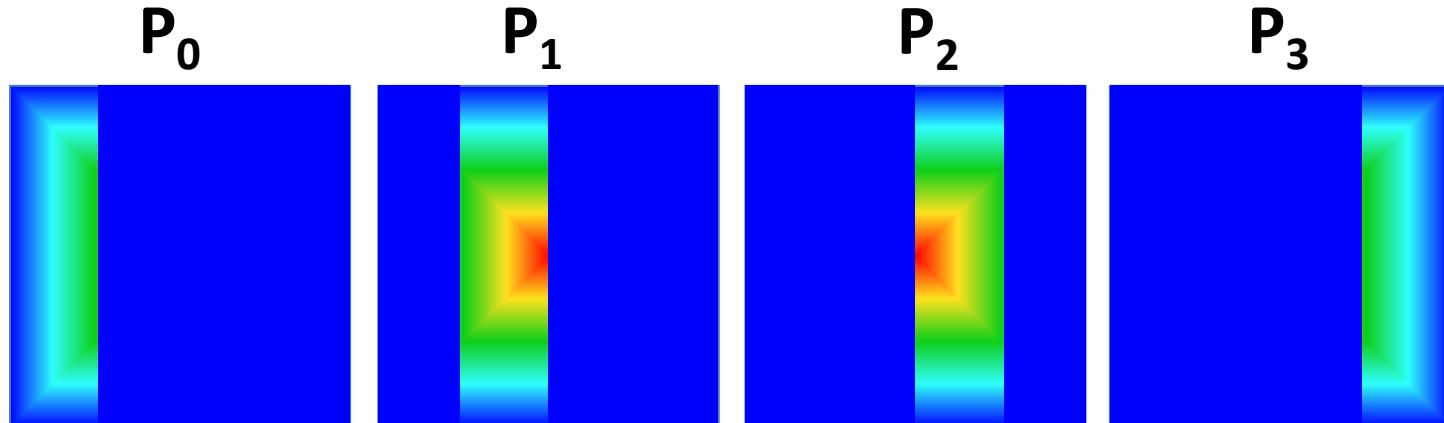
P_0 (root)

P_1

P_2

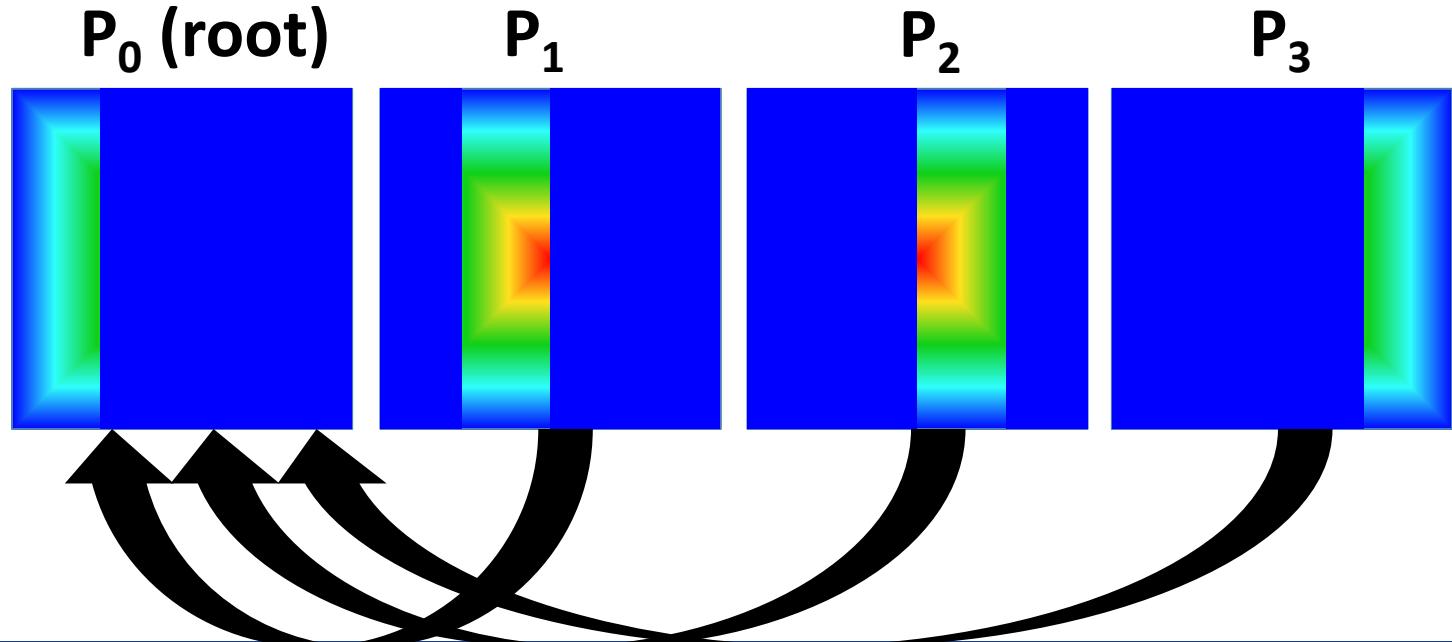
P_3





call evolve(dtfact)

call MPI_Gather(..., ..., ...)

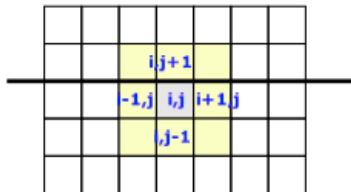


Replicated data

- Compute domain (and workload) distribution among processes
- Master-slaves: P_0 drives all processes
- Large amount of data communication
 - at each step P_0 distribute data to all processes and collect the contribution of each process
- Problem size scaling limited in memory capacity

Collaterals to Domain Decomposition /2

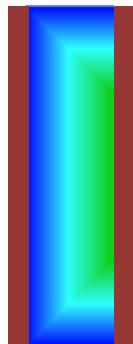
sub-domain boundaries



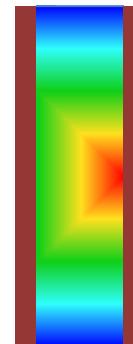


The Transport Code - Parallel Version

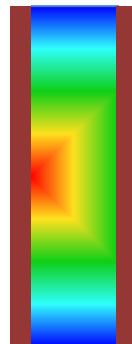
P_0



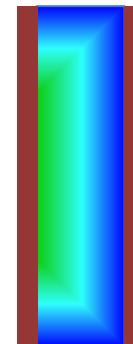
P_1



P_2

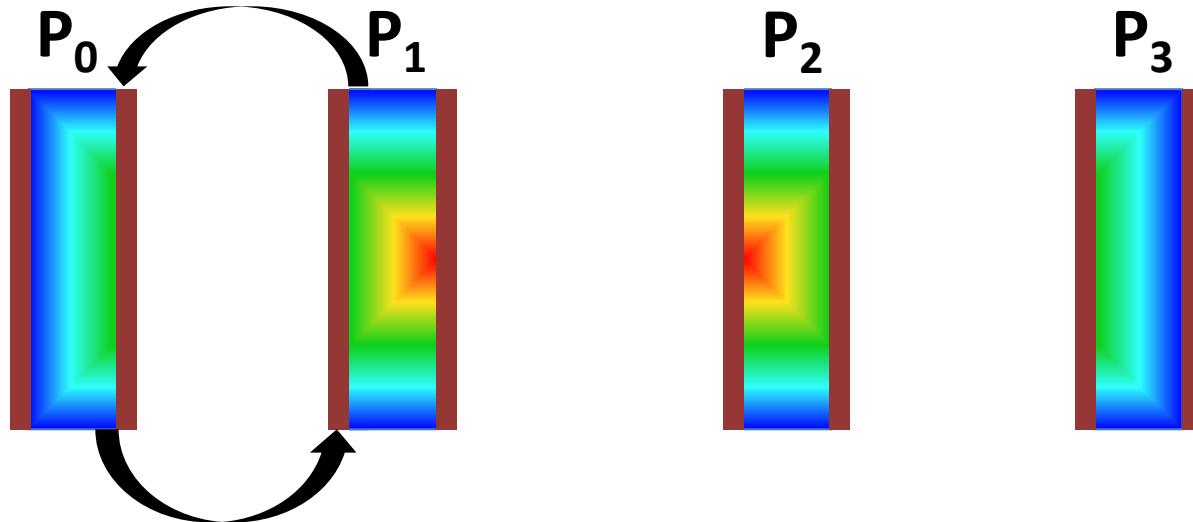


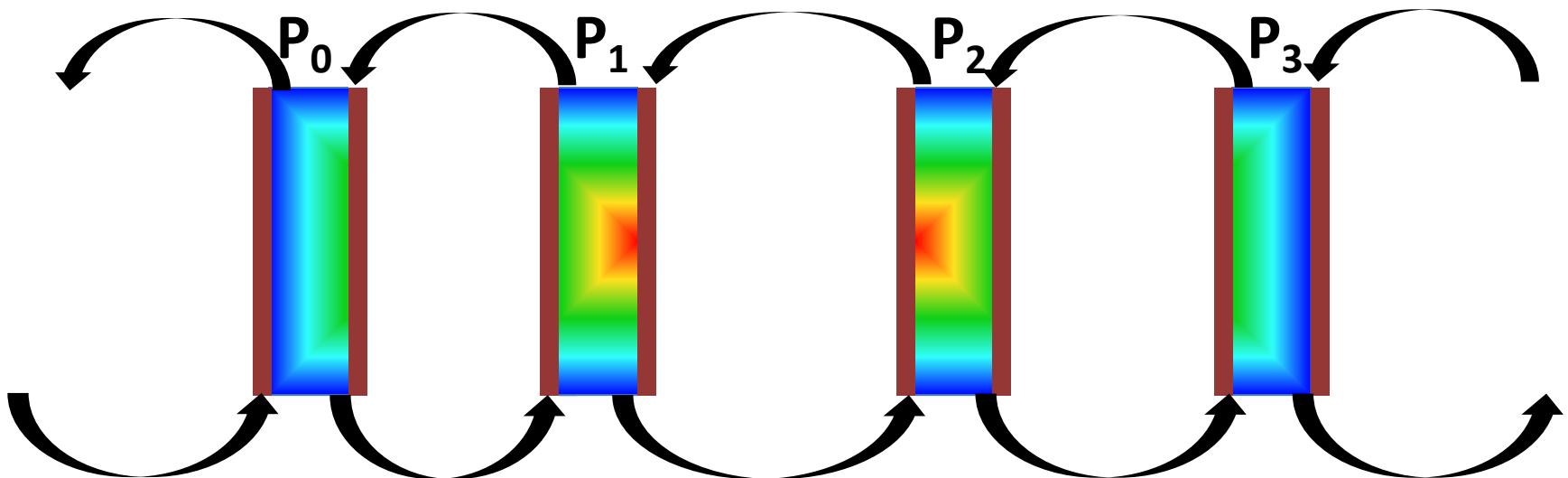
P_3



call evolve(dtfact)

Data exchange among processes



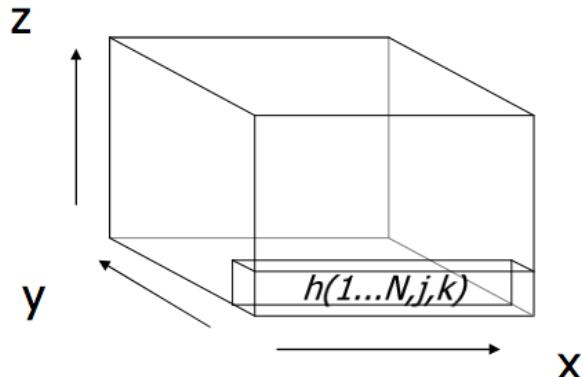
$$\text{proc_down} = \text{mod}(\text{proc_me} - 1 + \text{nprocs}, \text{nprocs})$$

$$\text{proc_up} = \text{mod}(\text{proc_me} + 1, \text{nprocs})$$



Distributed Data

- Global and Local Indexes
- Ghost Cells Exchange Between Processes
 - Compute Neighbor Processes
- Parallel Output

Multidimensional FFT



1) For any value of j and k transform the column $(1\dots N, j, k)$

2) For any value of i and k transform the column $(i, 1\dots N, k)$

3) For any value of i and j transform the column $(i, j, 1\dots N)$

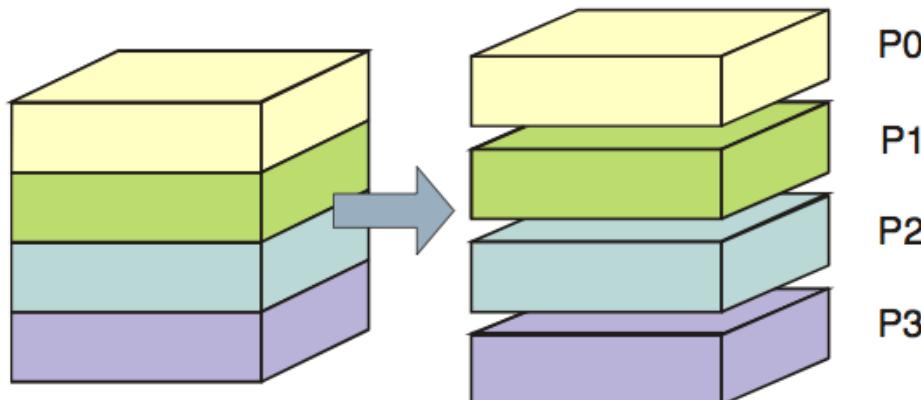
$$f(x, y, z) = \frac{1}{N_z N_y N_x} \sum_{z=0}^{N_z-1} \left(\sum_{y=0}^{N_y-1} \left(\sum_{x=0}^{N_x-1} \underbrace{F(u, v, w)}_{\text{DFT long x-dimension}} e^{-2\pi i \frac{xu}{N_x}} e^{-2\pi i \frac{yu}{N_y}} e^{-2\pi i \frac{zu}{N_z}} \right) \right)$$

 DFT long x-dimension

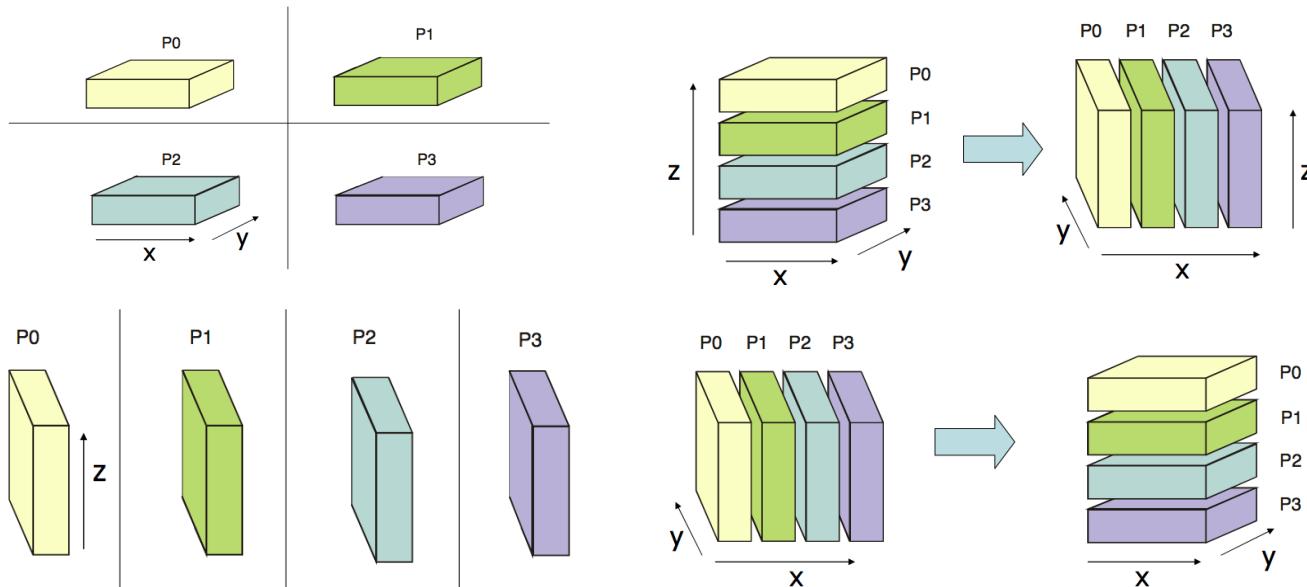
 DFT long y-dimension

 DFT long z-dimension

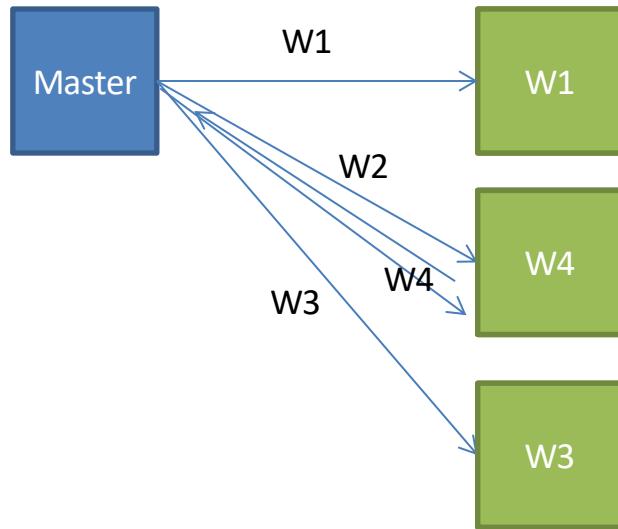
Parallel 3DFFT / 1



Parallel 3DFFT / 2



Master/Slave



Task Farming

- Many independent programs (tasks) running at once
 - each task can be serial or parallel
 - “independent” means they don’t communicate directly
 - Processes possibly driven by the mpirun framework

```
[igirotto@localhost]$ more my_shell_wrapper.sh
#!/bin/bash
#example for the OpenMPI implementation
./prog.x --input input_${OMPI_COMM_WORLD_RANK}.dat

[igirotto@localhost]$ mpirun -np 400 ./my_shell_wrapper.sh
```



Easy Parallel Computing

- Farming, embarrassingly parallel
 - Executing multiple instances on the same program with different inputs/initial cond.
 - Reading large binary files by splitting the workload among processes
 - Searching elements on large data-sets
 - Other parallel execution of embarrassingly parallel problem (no communication among tasks)
- Ensemble simulations (weather forecast)
- Parameter space (find the best wing shape)

Three nice reason to use python in scientific programming

NumPy:

<http://www.numpy.org/>

SciPy:

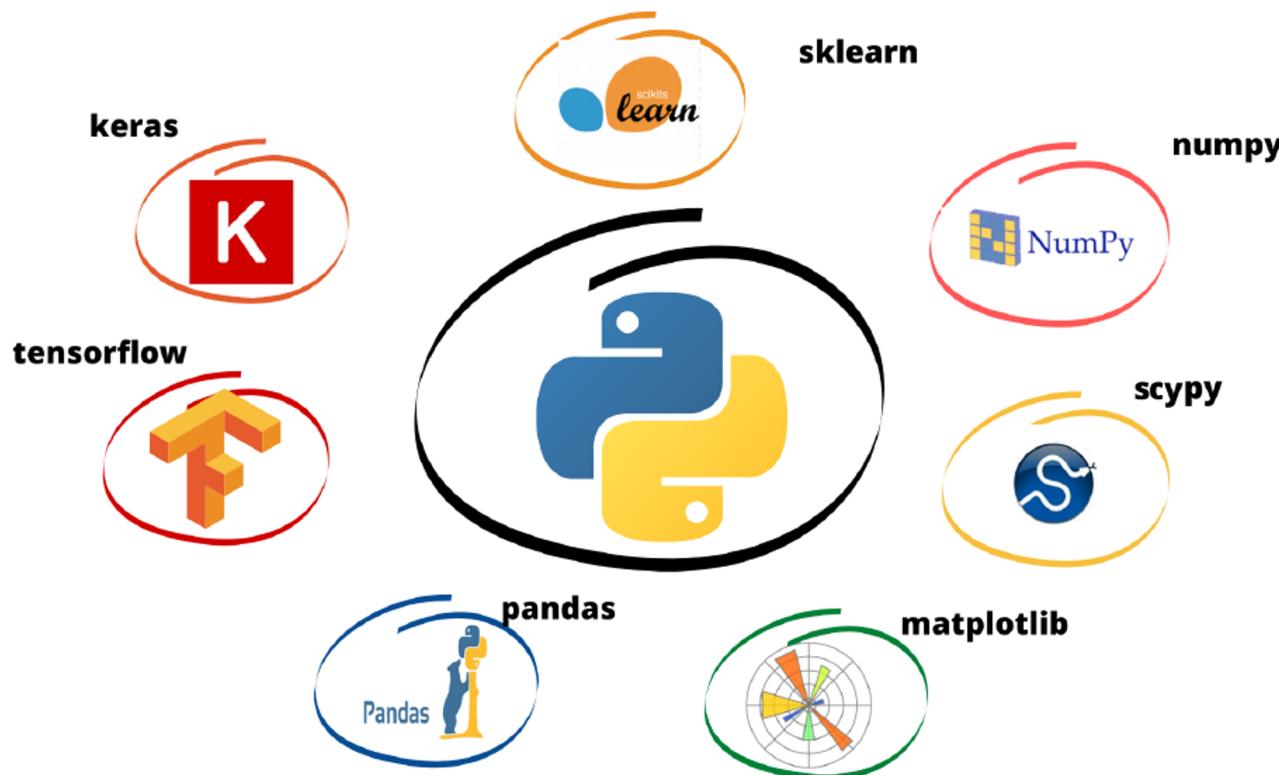
<http://www.scipy.org/>

Matplotlib:

<http://matplotlib.org/>

All are open source!

And more



Nice reason to use python in scientific programming

NumPy provides functionality to create, delete, manage and operate on large arrays of type raw" data (like Fortran and C/C++ arrays).

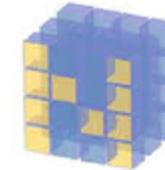
SciPy extends NumPy with a collection of useful algorithms like minimization, Fourier transforms, regression and many other applied mathematical techniques.

Both packages are add-on packages (not part of the Python standard library) containing Python code and compiled with (fftpack, BLAS).

MatPlotLib is a nice library to plot (but there are others as Seaborn or Bokeh).

How to import libraries?

```
import numpy as np  
import scipy as sp  
import matplotlib.pyplot as pp  
import matplotlib.pyplot as plt
```

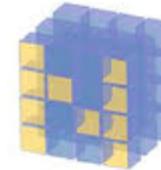


NumPy

Numpy

NumPy is the fundamental package for scientific computing in Python.

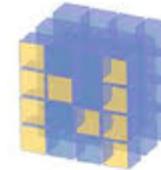
- A powerful N-dimensional array object.
- Sophisticated (broadcasting) functions.
- Tools for integrating C/C++ and Fortran code.
- Useful linear algebra, Fourier transform, and random number capabilities.
- Operations on matrices and vectors in NumPy are very efficient because they are linked to compiled in BLAS/LAPACK code.



NumPy

NumPy functionality:

- Polynomial mathematics.
- Statistical computations.
- Pseudo random number generators.
- Discrete Fourier transforms.
- Size / shape / type testing of arrays.

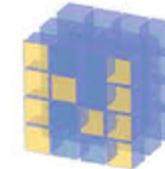


NumPy

Fundamental thing from Numpy: np.array

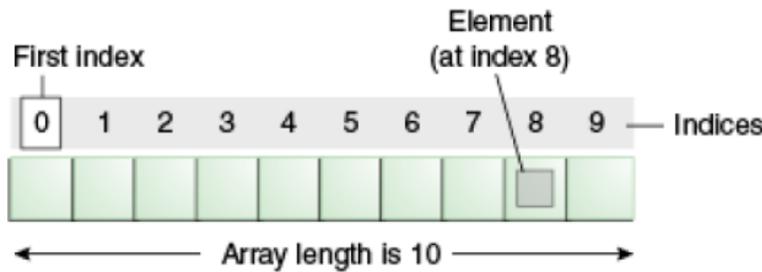
There are 5 general mechanisms for creating arrays:

- Conversion from other Python structures (e.g., lists, tuples).
- Intrinsic numpy array creation objects (e.g., arange, ones, zeros, etc.)
- Reading arrays from disk, either from standard or custom formats.
- Creating arrays from raw bytes through the use of strings or buffers.
- Use of special library functions (e.g., random).

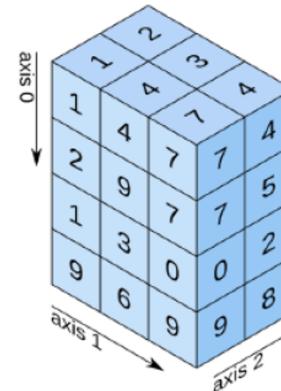


NumPy

How is an Array?



3D array



2D array

5.2	3.0	4.5
9.1	0.1	0.3

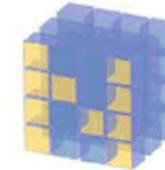
axis 0

axis 1

1D array

7	2	9	10
---	---	---	----

axis 0



NumPy

Examples of how to create arrays

```
x = np.array([2, 3, 1, 0])
print(x)
[2 3 1 0]
```

```
np.arange(10)
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

```
np.arange(2, 10, dtype=np.float)
array([2., 3., 4., 5., 6., 7., 8., 9.])
```

```
np.arange(2, 3, 0.1)
array([2., 2.1, 2.2, 2.3, 2.4, 2.5, 2.6, 2.7, 2.8, 2.9])
```

Special Functions

```
a = np.zeros((5, 2))  
print(a)  
[[0. 0.]  
 [0. 0.]  
 [0. 0.]  
 [0. 0.]  
 [0. 0.]]  
  
b = a.reshape((2, 5))  
print(b)  
[[0. 0. 0. 0. 0.]  
 [0. 0. 0. 0. 0.]]
```

Linear Algebra Operations

```
x = np.array([[1., 2., 3., 4.], [1., 2., 3., 4.]])  
print(x)  
[[1. 2. 3. 4.]  
 [1. 2. 3. 4.]]  
  
b= x.T  
print(b)  
[[1. 1.]  
 [2. 2.]  
 [3. 3.]  
 [4. 4.]]
```



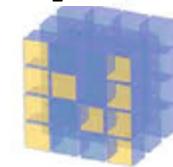
Simple Linear Algebra Operations

```
v=np.array([-1, 2,5])
w=np.array([1, 2,3])
print(v.dot(w))
18
print(np.dot(v, w))
18
#-----
```

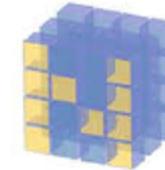
```
v1=np.array([[ -1, 2,5],[1, 3,5]])
w1=np.array([1,10,3])
print(v1.dot(w1))
[34 46]
print(np.cross(v,w))
[-4 8 -4]
```

```
x=np.array([1, 2,3])
y=np.array([11, -2,5])
print(x + y)
[12 0 8]
print(np.add(x, y))
[12 0 8]
#-----
```

```
print(x - y)
[-10 4 -2]
print(np.subtract(x, y))
[-10 4 -2]
```



NumPy



Opening a txt or cvs file

```
import numpy as np

file_name_you_want = np.loadtxt(fname,delimiter=" ")

print "First column element: ",file_name_you_want[0]
#to get the full column:

Transpose_your_file = file_name_you_want.T

print "First column: ", Transpose_your_file[0]
```

Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science **Interactively** at www.DataCamp.com



NumPy

The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```

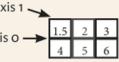


NumPy Arrays

1D array

```
[1 2 3]
```

axis 0



2D array

axis 1

3D array

axis 1



Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1,5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1,5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]),
      dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3,4))          Create an array of zeros
>>> np.ones((2,3,4),dtype=np.int16) Create an array of ones
>>> d = np.arange(10,25,5)    Create an array of evenly spaced values (step value)
>>> np.linspace(0,2,9)       Create an array of evenly spaced values (number of samples)
>>> e = np.full((2,2),7)     Create a constant array
>>> f = np.eye(2)            Create a 2x2 identity matrix
>>> np.random.random((2,2))  Create an array with random values
>>> np.empty((3,2))         Create an empty array
```

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savetxt('array.npz', a, b)
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("myfile.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter="")
```

Data Types

```
>>> np.int64
>>> np.float32
>>> np.complex
>>> np.bool_
>>> np.object_
>>> np.string_
>>> np.unicode_
```

Signed 64-bit integer types
Standard double-precision floating point
Complex numbers represented by 128 floats
Boolean type storing TRUE and FALSE values
Python object type
Fixed-length string type
Fixed-length unicode type

Inspecting Your Array

```
>>> a.shape
>>> len(a)
>>> b.ndim
>>> b.size
>>> b.dtype
>>> b.dtype.name
>>> b.astype(int)
```

Array dimensions
Length of array
Number of array dimensions
Number of array elements
Data type of array elements
Name of data type
Convert an array to a different type

Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

Array Mathematics

Arithmetic Operations

```
>>> g = a - b
>>> np.subtract(a,b)
>>> b + a
>>> np.add(b,a)
>>> a / b
>>> np.divide(a,b)
>>> a * b
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> e.dot(f)
>>> np.array([[ 1.,  2.,  3.], [ 4.,  5.,  6.], [ 7.,  8.,  9.]]))
```

Subtraction
Addition
Addition
Division
Multiplication
Multiplication
Exponentiation
Square root
Print sines of an array
Element-wise cosine
Element-wise natural logarithm
Dot product

Comparison

```
>>> a == b
>>> np.array_equal([True, False, True], [False, True, False], dtype=bool)
>>> a < 2
>>> np.array_equal([True, False, False], dtype=bool)
>>> np.array_equal(a, b)
```

Element-wise comparison
Element-wise comparison
Array-wise comparison

Aggregate Functions

```
>>> a.sum()
>>> a.min()
>>> b.max(axis=0)
>>> b.cumsum(axis=1)
>>> a.mean()
>>> b.median()
>>> a.corrcoef()
>>> np.std(b)
```

Array-wise sum
Array-wise minimum value
Maximum value of an array row
Cumulative sum of the elements
Mean
Median
Correlation coefficient
Standard deviation

Copying Arrays

```
>>> h = a.view()
>>> np.copy(a)
>>> h = a.copy()
```

Create a view of the array with the same data
Create a copy of the array
Create a deep copy of the array

Sorting Arrays

```
>>> a.sort()
>>> c.sort(axis=0)
```

Sort an array
Sort the elements of an array's axis

Subsetting, Slicing, Indexing

Subsetting

```
>>> a[2]
3
>>> b[1,2]
1, 2, 3
4, 5, 6
```

Slicing

```
>>> a[0:2]
array([1, 2])
>>> b[0:2,1]
array([ 1,  2,  5.], [ 4,  5,  6])
```

```
>>> b[:,1]
array([1.5, 2., 3.])
>>> c[1,:,:]
array([[ 3.,  2.,  1.], [ 4.,  5.,  6.]])
```

Reversed array

```
>>> a[-1]
array([1, 2, 3])
>>> a[a<2]
array([1])
```

Fancy Indexing

```
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]
array([ 1.,  2.,  6.,  1.5])
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]
array([[ 1.,  2.,  6.,  1.5], [ 1.5,  2.,  3.,  1.5], [ 1.5,  2.,  3.,  1.5], [ 1.5,  2.,  3.,  1.5]])
```

Select elements from a less than 2
Select elements (1,0),(0,1),(1,2) and (0,0)
Select a subset of the matrix's rows and columns

Array Manipulation

Transposing Array

```
>>> i = np.transpose(b)
>>> i.T
```

Permute array dimensions
Permute array dimensions

Changing Array Shape

```
>>> b.ravel()
>>> g.reshape(3,-2)
```

Flatten the array
Reshape, but don't change data

Adding/Removing Elements

```
>>> h.reshape(2,6)
>>> np.append(h,g)
>>> np.insert(a, 1, 5)
>>> np.delete(a, [1])
```

Return a new array with shape (2,6)
Append items to an array
Insert items in an array
Delete items from an array

Combining Arrays

```
>>> np.concatenate((a,d),axis=0)
array([ 1,  2,  3, 10, 15, 20])
>>> np.vstack((a,b))
array([[ 1.,  2.,  3.], [ 1.5,  2.,  3.], [ 4.,  5.,  6.]])
>>> np.r_[e,f]
>>> np.hstack((e,f))
array([[ 7.,  8.,  0.,  1.], [ 7.,  8.,  0.,  1.]])
>>> np.column_stack((a,d))
array([[ 1,  2,  3, 10, 15, 20], [ 3,  4,  5, 16, 21, 20]])
>>> np.c_[a,d]
```

Concatenate arrays

Stack arrays vertically (row-wise)

Stack arrays vertically (row-wise)
Stack arrays horizontally (column-wise)

Create stacked column-wise arrays

Create stacked column-wise arrays

Split the array horizontally at the 3rd index
Split the array vertically at the 2nd index





SciPy

One of the core packages. SciPy is built on top of NumPy and implements many specialized scientific computation tools:

- Mainly user-friendly.
- Efficient numerical routines such as routines for numerical integration and optimization.
- Clustering.
- Fourier transforms.
- Numerical integration, interpolations, data I/O, LAPACK.

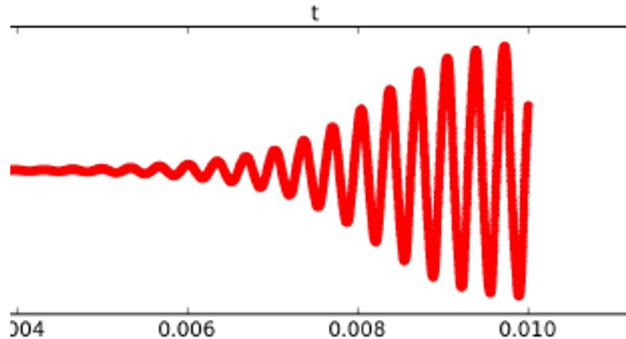
Also:

- sparse matrices, linear solvers, optimization.
- signal processing.
- statistical functions



Nice for numerical integration of equations

```
from scipy.integrate import odeint  
dy/dt = func(y, t0, ...)  
sol = odeint(func, x0, t)
```



<https://docs.scipy.org/doc/scipy/reference/generated/scipy.integrate.odeint.html>



Or for example for a simple fit

```
import numpy as np
import matplotlib.pyplot as pp
from scipy.optimize import curve_fit

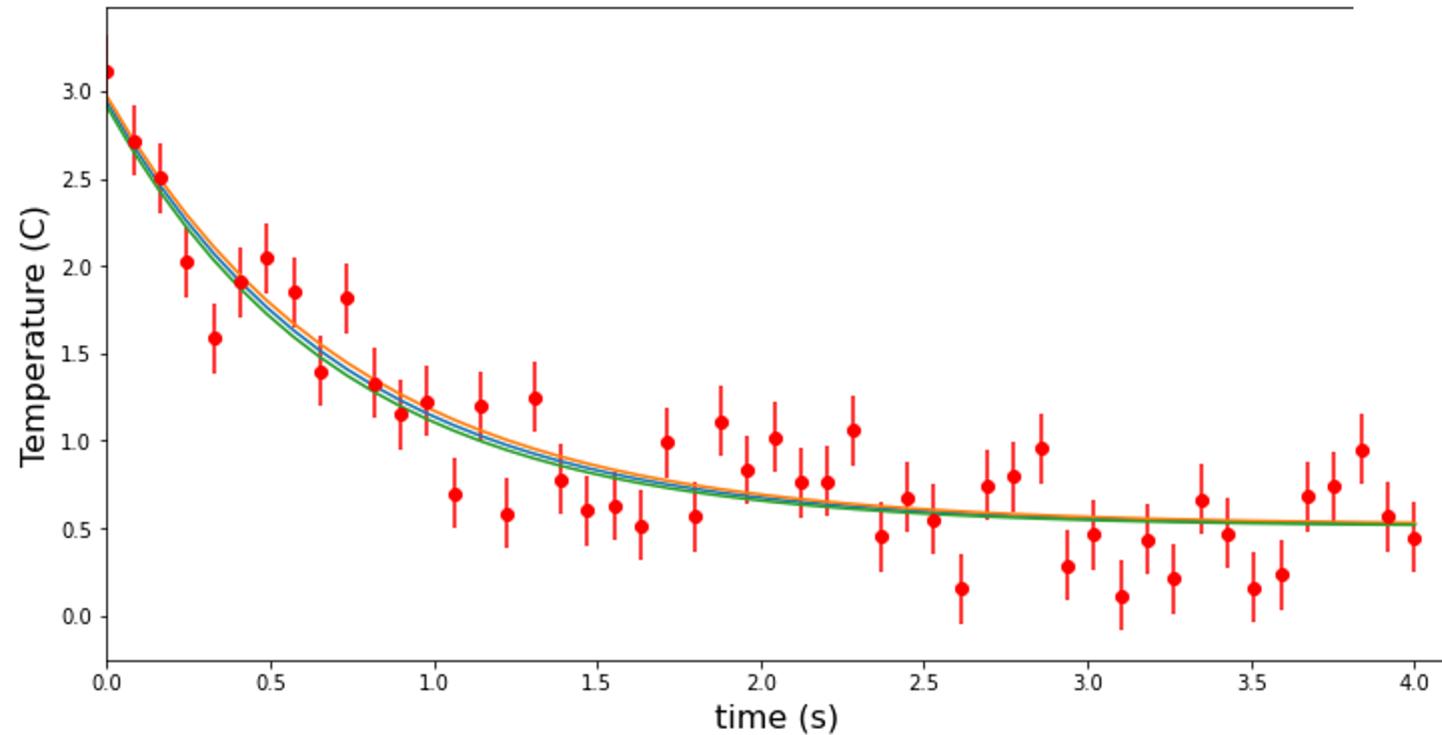
def fitFunc(t, a, b, c):
    return a*np.exp(-b*t) + c

t = np.linspace(0,4,50)
temp = fitFunc(t, 2.5, 1.3, 0.5)
noisy = temp + 0.25*np.random.normal(size=len(temp))
fitParams, fitCovariances = curve_fit(fitFunc, t, noisy)

pp.figure(figsize=(12, 6))
pp.ylabel('Temperature (C)', fontsize = 16)
pp.xlabel('time (s)', fontsize = 16)
pp.xlim(0,4.1)
pp.errorbar(t, noisy, fmt = 'ro', yerr = 0.2)
sigma = [fitCovariances[0,0], fitCovariances[1,1], fitCovariances[2,2] ]
pp.plot(t, fitFunc(t, fitParams[0], fitParams[1], fitParams[2]))
pp.plot(t, fitFunc(t, fitParams[0] + sigma[0], fitParams[1] - sigma[1], fitParams[2] + sigma[2]))
pp.plot(t, fitFunc(t, fitParams[0] - sigma[0], fitParams[1] + sigma[1], fitParams[2] - sigma[2]))
pp.savefig('dataFitted.pdf', bbox_inches=0, dpi=600)
pp.show()
```



And the result!





With the Fit parameters

[2.595658 1.74438726 0.69809511]

And covariance matrix:

[[0.02506636 0.01490486 -0.00068609]

2 [0.01490486 0.04178044 0.00641246]

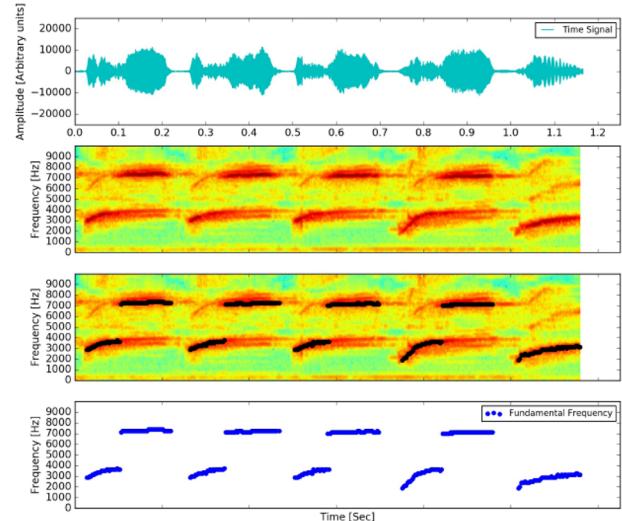
3 [-0.00068609 0.00641246 0.00257799]]

Also possible to apply for signal analysis and time series

DOI:<https://doi.org/10.31527/analesafa.2018.29.2.51>

<https://doi.org/10.1016/j.mex.2018.12.011>

DOI:<https://doi.org/10.31527/analesafa.2019.30.3.68>



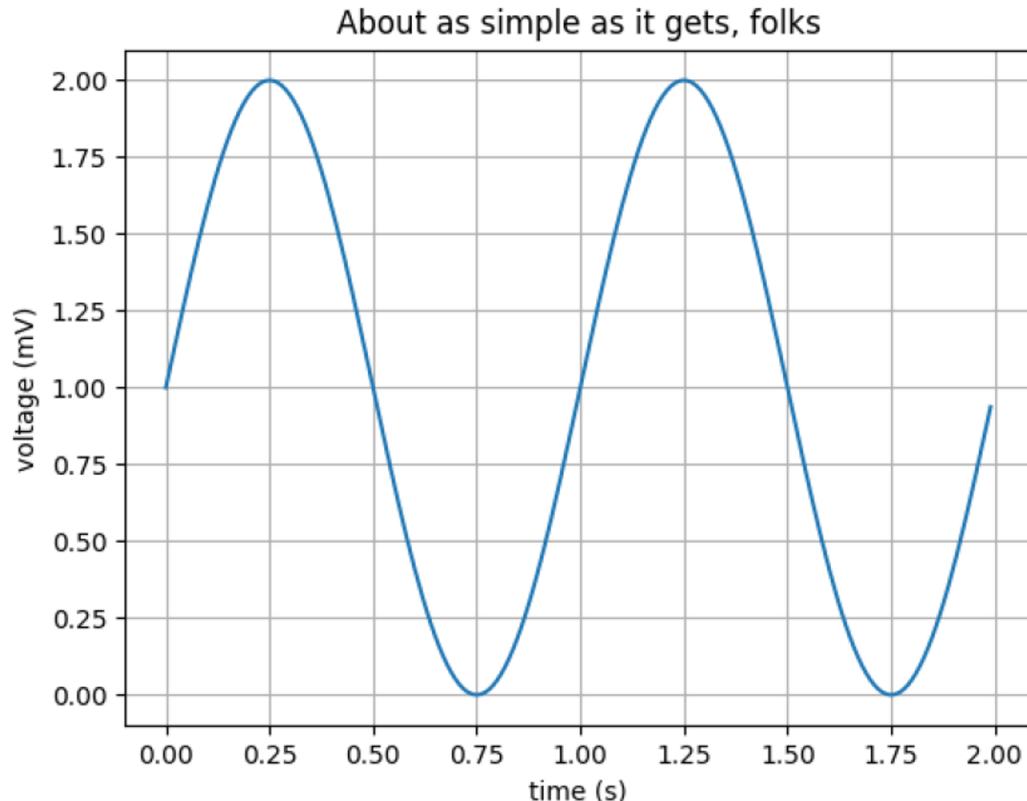


Matplotlib

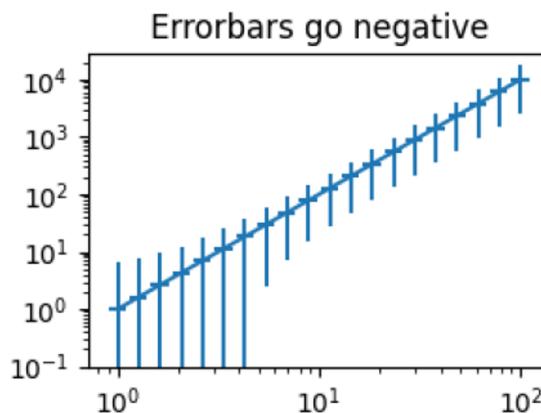
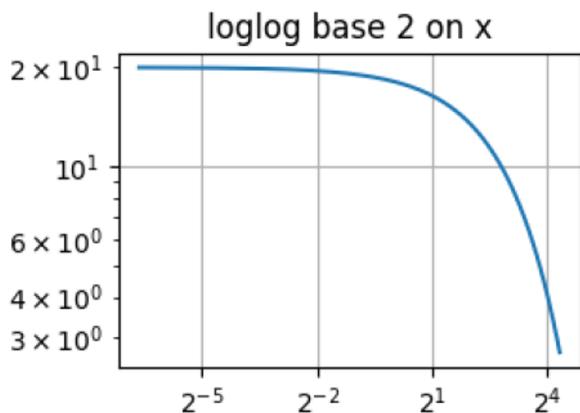
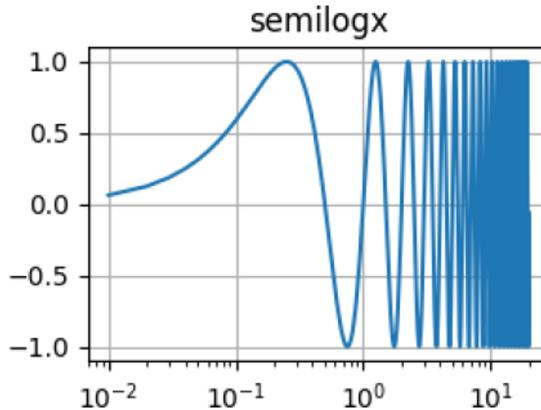
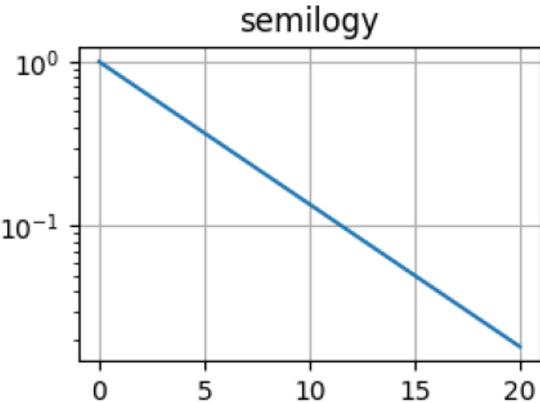
Is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms. Can be used in:

- Python scripts.
- The Python and IPython shell.
- The jupyter notebook.
- Web application servers.
- Graphical user interface toolkits.

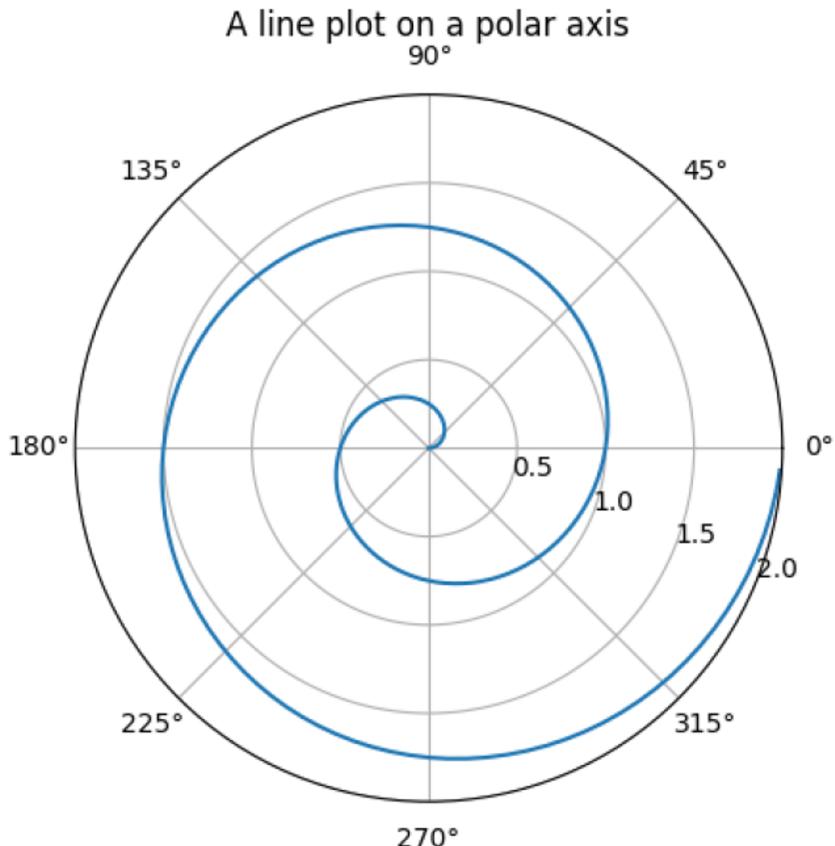
Different kind of plots: y vs. x



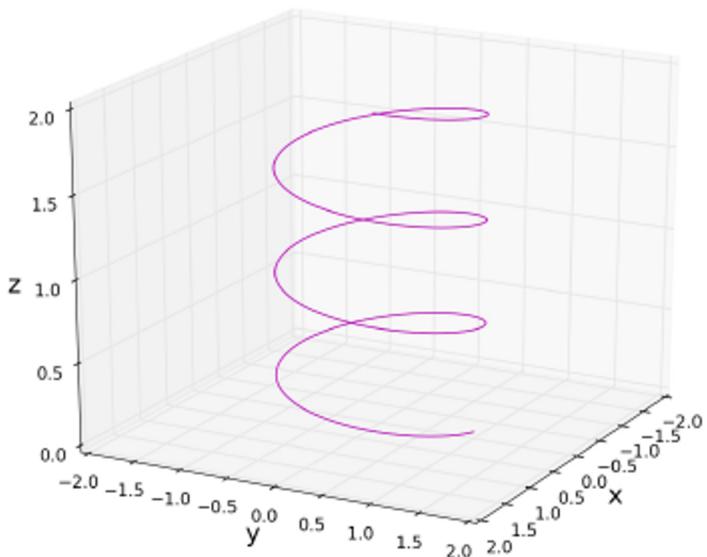
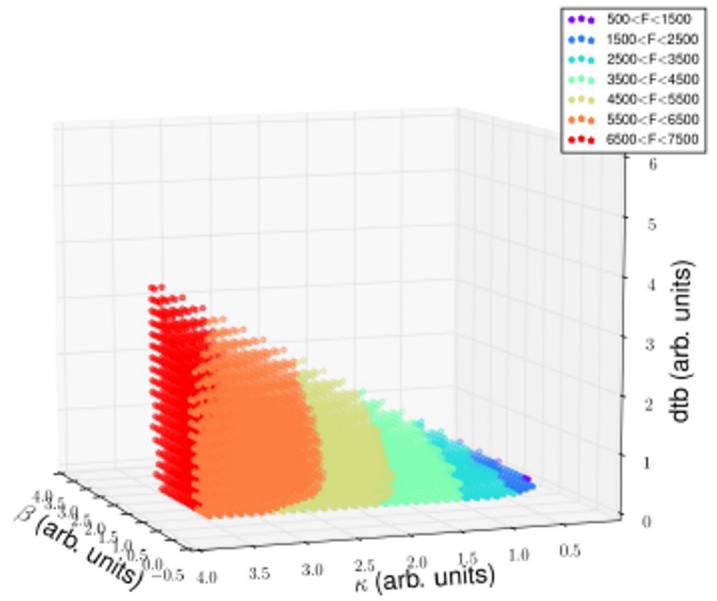
Different scales:



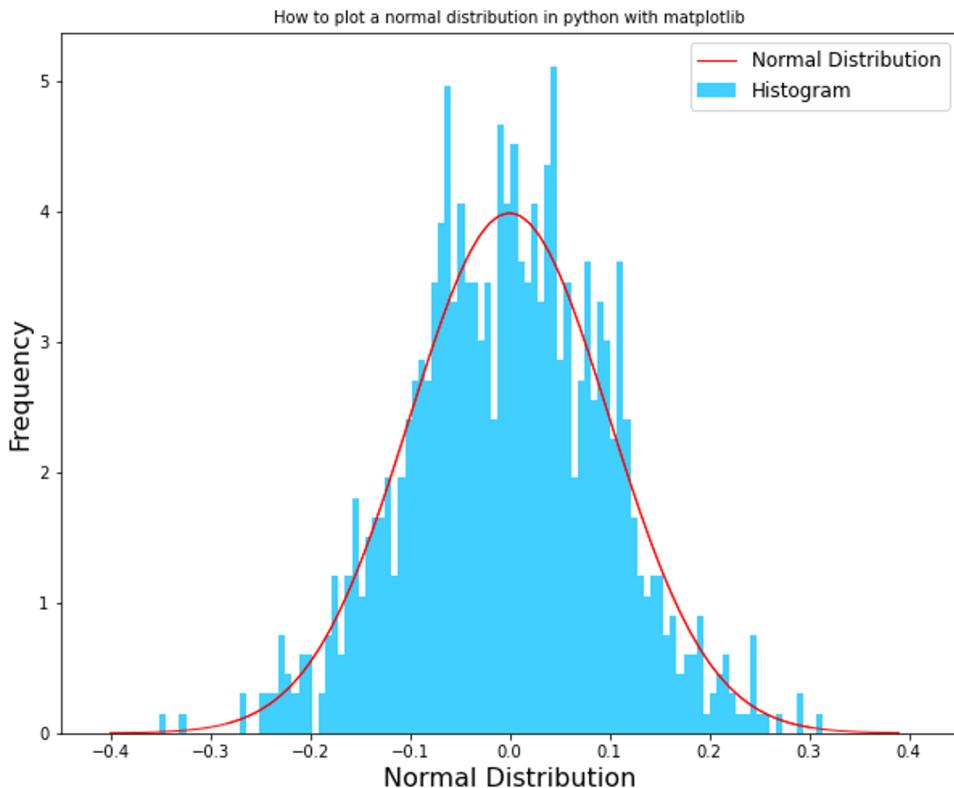
Different coordinates



3d plots:



Histograms and distributions





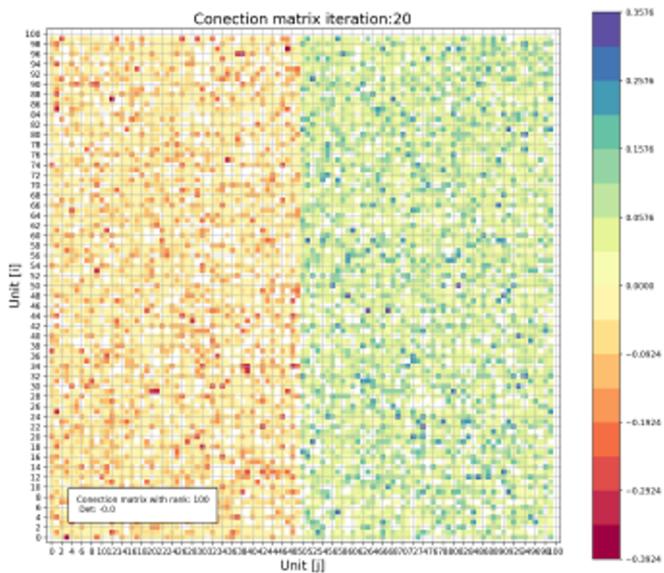
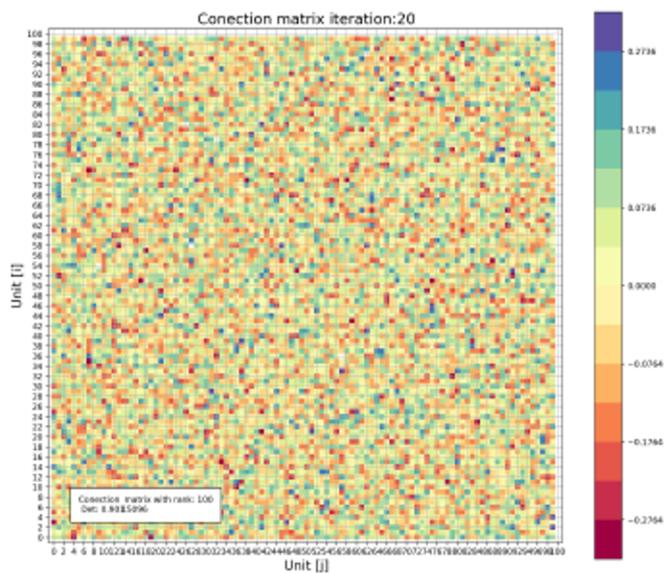
The code for the histogram:

```
import numpy as np
from scipy.stats import norm
import matplotlib.pyplot as plt
mu, sigma = 0, 0.1 # mean and standard deviation
s = np.random.normal(mu, sigma, 1000)
x_n = np.arange(-0.4,0.4,0.01)
y_n = norm.pdf( x_n, mu, sigma)

fig      = plt.figure(figsize=(10,8))
plt.title('Histogram Normal Distribution', fontsize = 18)
plt.hist(s, 100, density=1, facecolor='deepskyblue',label='Histogram', alpha=0.75)
plt.plot(x_n, y_n, 'r-', linewidth=1,label='Normal Distribution')
plt.show()
```

Matrices:

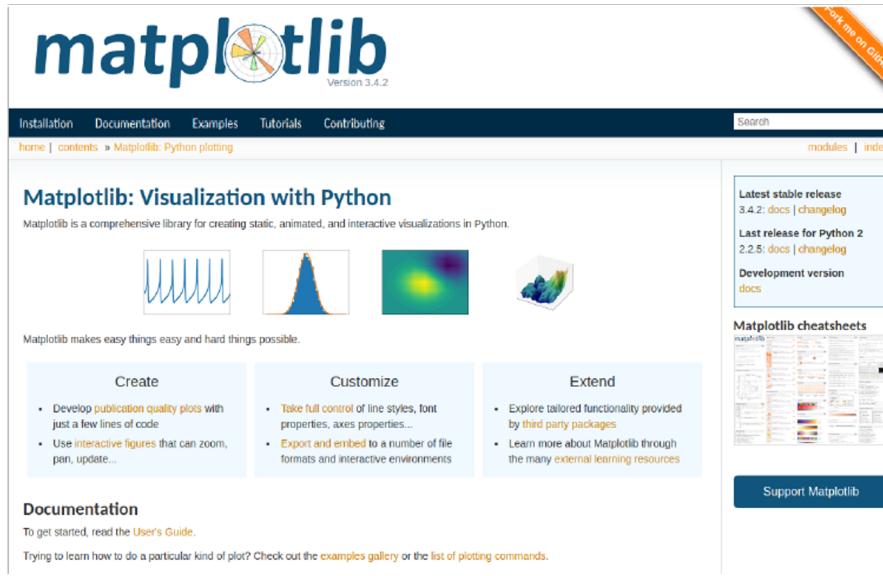
Numpy Scipy Matplotlib





Please follow the examples in Documentation
there are a lot!

<https://matplotlib.org/stable/gallery/index.html>



The screenshot shows the official matplotlib documentation website. At the top, there's a navigation bar with links for Installation, Documentation, Examples, Tutorials, and Contributing. Below the navigation bar, a main heading reads "Matplotlib: Visualization with Python". A sub-headline states: "Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python." To the right of this text are four small images illustrating different types of plots: a line plot with oscillations, a histogram-like plot with a peak, a heatmap, and a 3D surface plot. Below the main heading, there are three sections: "Create", "Customize", and "Extend", each with a bulleted list of features. At the bottom left, there's a "Documentation" section with a link to the User's Guide. On the right side, there's a sidebar with links for the latest stable release (3.4.2), the last release for Python 2 (2.2.5), and development version links. It also includes a "Matplotlib cheatsheets" section with two small preview images of the cheatsheets and a "Support Matplotlib" button at the bottom.

More on visualization with python

<https://clauswilke.com/dataviz/>

<https://www.python-graph-gallery.com/>

Python For Data Science Cheat Sheet

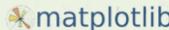
Matplotlib

Learn Python Interactively at www.DataCamp.com



Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.



1 Prepare The Data

Also see [Lists & NumPy](#)

1D Data

```
>>> import numpy as np  
>>> x = np.linspace(0, 10, 100)  
>>> y = np.cos(x)  
>>> z = np.sin(x)
```

2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))  
>>> data2 = 3 * np.random.random((10, 10))  
>>> I, X, Y = np.indices([3:100j, -3:100j])  
>>> V = -1 + X**2 + Y**2  
>>> from matplotlib.cbook import get_sample_data  
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

2 Create Plot

```
>>> import matplotlib.pyplot as plt
```

Figure

```
>>> fig = plt.figure()  
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

Axes

All plotting is done with respect to an `Axes`. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()  
>>> ax1 = fig.add_subplot(221) # row-col-num  
>>> ax3 = fig.add_subplot(212)  
>>> fig3, axes = plt.subplots(rows=2,ncols=2)  
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3 Plotting Routines

1D Data

```
>>> fig, ax = plt.subplots()  
>>> lines = ax.plot(x,y)  
>>> ax.scatter(x,y)  
>>> axes[0,0].bar([1,2,3],[4,4,5])  
>>> axes[1,0].barh([0.5,1.25],[0,1,2])  
>>> axes[1,1].axhline(0.45)  
>>> axes[0,1].axvline(0.65)  
>>> ax.fill(x,y,color='blue')  
>>> ax.fill_between(x,y,color='yellow')
```

Draw points with lines or markers connecting them
Draw unconnected points, scaled or colored
Plot vertical rectangles (constant width)
Plot horizontal rectangles (constant height)
Draw a horizontal line across axes
Draw a vertical line across axes
Draw filled polygons
Fill between y-values and 0

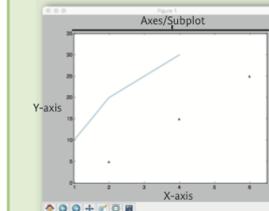
2D Data or Images

```
>>> fig, ax = plt.subplots()  
>>> im = ax.imshow(img,  
>>> cmap='gist_earth',  
>>> interpolation='nearest',  
>>> vmin=-2,  
>>> vmax=2)
```

Colormapped or RGB arrays

Plot Anatomy & Workflow

Plot Anatomy



Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
 - 2 Create plot
 - 3 Plot
 - 4 Customize plot
 - 5 Save plot
 - 6 Show plot
- ```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4] Step 1
>>> y = [10,20,25,30]
>>> fig = plt.figure() Step 2
>>> ax = fig.add_subplot(111) Step 3
>>> ax.plot(x, y, color='lightblue', linewidth=3) Step 3.4
>>> ax.scatter([2,4,6],
>>> [5,15,25],
>>> color='darkgreen',
>>> marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig('foo.png')
>>> plt.show() Step 6
```

## 4 Customize Plot

### Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha= 0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img,
>>> cmap='seismic')
```

### Markers

```
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker=".")
>>> ax.plot(x,y,marker="o")
```

### LineStyles

```
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,ls='solid')
>>> plt.plot(x,y,ls="--")
>>> plt.plot(x,y,'-.',x**2,y**2,'-.')
>>> plt.setp(lines,color='r',linewidth=4.0)
```

### Text & Annotations

```
>>> ax.text(1,-2.1,
>>> "Example Graph",
>>> style='italic')
>>> ax.annotate('Simpler',
>>> xy=(18, 0),
>>> xycoords='data',
>>> xytext=(10.5, 0),
>>> textcoords='data',
>>> arrowprops=dict(arrowstyle="->",
>>> connectionstyle="arc3"),)
```

### Mathtext

```
>>> plt.title(r'$\sigma_{\mathrm{i}}=15$', fontsize=20)
```

### Limits, Legends & Layouts

#### Limits & Autoscaling

```
>>> ax.margins(x=0.0,y=0.1)
>>> ax.axis('equal')
>>> ax.set_xlim(0,10.5)
>>> ax.set_xlim(0,10.5)
```

#### Legends

```
>>> ax.set(title='An Example Axes',
>>> xlabel='Y-Axis',
>>> ylabel='X-Axis')
```

#### Ticks

```
>>> ax.xaxis.set(ticks=range(1,5),
>>> ticklabels=[3,100,-12,'foo'])
>>> ax.tick_params(axis='y',
>>> direction='inout',
>>> length=10)
```

#### Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,
>>> hspace=0.3,
>>> left=0.125,
>>> right=0.9,
>>> top=0.9,
>>> bottom=0.1)
```

#### Axis Spines

```
>>> ax1.spines['top'].set_visible(False)
>>> ax1.spines['bottom'].set_position(('outward',10))
```

Add padding to a plot  
Set the aspect ratio of the plot to 1  
Set limits for x-and y-axis  
Set limits for x-axis

Set a title and x-and y-axis labels

No overlapping plot elements

Manually set x-ticks

Make y-ticks longer and go in and out

Adjust the spacing between subplots

Fit subplot(s) in to the figure area

Make the top axis line for a plot invisible  
Move the bottom axis line outward

## 5 Save Plot

### Save figures

```
>>> plt.savefig('foo.png')
>>> Save transparent figures
>>> plt.savefig('foo.png', transparent=True)
```

## 6 Show Plot

```
>>> plt.show()
```

## Close & Clear

```
>>> plt.clf()
>>> plt.cla()
>>> plt.close()
```

Clear an axis  
Clear the entire figure  
Close a window



<http://pandas.pydata.org/>



# To import library

```
import pandas as pd
```

A unidimensional array with labels:

```
s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

```
print(s)
```

```
a 3
```

```
b -5
```

```
c 7
```

```
d 4
```

```
dtype: int64
```

# A bi-dimensional array

```
data = {'Country': ['Belgium', 'India', 'Brazil'], 'Capital': ['Brussels',
'New Delhi', 'Brasilia'], 'Population': [11190846, 1303171035, 207847528]}
```

```
df = pd.DataFrame(data,columns=['Country', 'Capital', 'Population'])
```

```
print(df)
```

|   | Country | Capital   | Population |
|---|---------|-----------|------------|
| 0 | Belgium | Brussels  | 11190846   |
| 1 | India   | New Delhi | 1303171035 |
| 2 | Brazil  | Brasilia  | 207847528  |

## To read a file and write it with pandas

```
pd.read_csv('file.csv', header=None, nrows=5)

df.to_csv('myDataFrame.csv')
```

## To read multiple sheets of the same file

```
xlsx = pd.ExcelFile('file.xls')
df = pd.read_excel(xlsx, 'Sheet1')
```

# To request help!

```
help(pd.Series.loc)
```

```
mark ii 1 4
viper mark ii 7 1
 mark iii 16 36
```

Single label. Note this returns a DataFrame with a single index.

```
>>> df.loc['cobra']
 max_speed shield
mark i 12 2
mark ii 0 4
```

Single index tuple. Note this returns a Series.

```
>>> df.loc[('cobra', 'mark ii')]
max_speed 0
shield 4
Name: (cobra, mark ii), dtype: int64
```

Single label for row and column. Similar to passing in a tuple, this returns a Series.

```
>>> df.loc['cobra', 'mark i']....
```

# To select one element

```
s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
s['b']
-5
```

# To select a subset

```
data = {'Country': ['Belgium', 'India', 'Brazil'], 'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
'Population': [11190846, 1303171035, 207847528]}
```

```
df = pd.DataFrame(data,columns=['Country', 'Capital', 'Population'])
```

```
df[1:]
```

| Country | Capital | Population |                |
|---------|---------|------------|----------------|
| 1       | India   | New Delhi  | 130317103<br>5 |
| 2       | Brazil  | Brasilia   | 207847528      |

# Retrieving Series/DataFrame Information

```
#(rows, columns)
df.shape

#Describe index
df.index

#Describe DataFrame columns
df.columns
```

```
#Info on DataFrame

df.info()

#Number of non-NA values

df.count()
```

```
<class
'pandas.core.frame.DataFrame'>
RangeIndex: 3 entries, 0 to 2
Data columns (total 3 columns):
 # Column Non-Null Count
 Dtype
--- ---

 0 Country 3 non-null
 object
 1 Capital 3 non-null
 object
 2 Population 3 non-null
 int64
 dtypes: int64(1), object(2)
 memory usage: 200.0+ bytes
Country 3
Capital 3
Population 3
dtype: int64
```

# Applying Functions

```
f = lambda x: x*2
#Apply function
df.apply(f)

#Apply function element-wise
df.applymap(f)
```