

# Parallel Computing Survey

- Basic Definitions
- Illustrative Examples
- Computer models
- Speedup
- Engineering principles
- Multiple processors
- Parallel programming models
- Evolution of Performances

# Parallel Computing

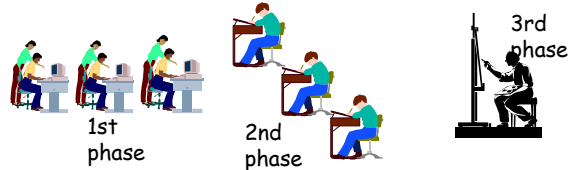
- Parallel computer + parallel program = parallel computing
- A parallel computer is a set of processors that are able to work cooperatively to solve a common computational problem.
  - parallel supercomputers with hundreds or thousands of processors,
  - networks of workstations,
  - multiple-processor workstations, and
  - embedded systems.
- A parallel program is a code being able to be executed by parallel computers.

## Parallel Terminology

- There are many different expressions used in the literature for similar principles.
- Processor, Computer, Node, etc.
- Parallel execution, Concurrent execution, Distributed execution
- Multicomputer, Multiprocessor, Parallel processor, Distributed processor
- Communication channel, Link,
- We will use parallel computer when communication by messages is used and parallel processor if shared memory is used for communication.

## Illustrative Examp.1 - Writing Projects

- Writing projects (some communication in the first phase, no communication during work - all writers are independent, and no communication on grading).



- Finally, several projects have been done in the time needed to finish a single project plus some overhead because of initial personalised communication.

## Illustrative Examp.2 - Rowing a Boat

- Synchronisation -> helmsman & oarsmen,
- Communication -> oarsman- oarsman, On-line corrections -> oarsmen and helmsman



- The boat is moving faster if the crew is larger but not proportionally and only for some extent. Too much people could even sink the boat.

## Fundamental Problems to be Solved

- Problem decomposition,
- Functional and data dependencies,
- Synchronisation, Communication,
- Partitioning,
- Visualisation of results,
- Performance measurements.

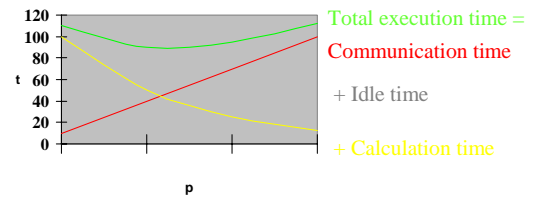
## Parallel Applications

- Numerical simulations of complex systems such as weather, climate, mechanical devices, electronic circuits,
- Control of manufacturing processes, and chemical reactions,
- Commercial applications for large amounts of data,
- Video servers (real-time video).
- The need for faster computers is driven by the demands of both data-intensive applications in commerce and computation-intensive applications in science and engineering.

University of Salzburg, Department of Scientific Computing, HPSC SS 2005

7

## Computation + Communication

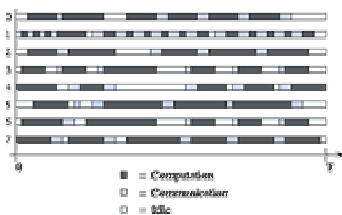


- execution time = computation + communication + idle time
- some optimal number of processors  $p$  exist for a given problem size  $N$ ,
- increase processor performances (technology),
- shorten communication and idle time (algorithm).

University of Salzburg, Department of Scientific Computing, HPSC SS 2005

8

## Computation + Communication (cont.)



Activity plot during execution of a parallel program on eight processors. Each processor spends its time computing, communicating, or idling.

University of Salzburg, Department of Scientific Computing, HPSC SS 2005

9

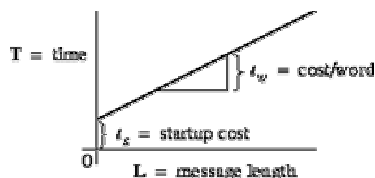
## Overlapping Computation with Communication

- A worker that has sent a request to the input task waits for the parameters to arrive.
- In some cases, efficiency can be improved by pre-fetching, that is, requesting the next parameter before it is needed.
- The worker can then perform computation while its request is being processed by the input task.
- A part of Computation and Communication can be executed in parallel - overlapped in time.
- In such a way idle time can be reduced.

University of Salzburg, Department of Scientific Computing, HPSC SS 2005

10

## Communication time - Model

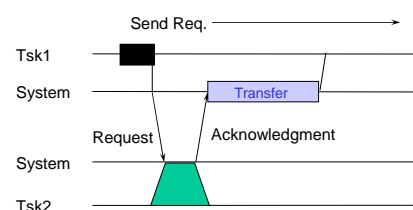


-Time  $T$  versus message length  $L$  is given, the slope of the line corresponds to the cost per word transferred  $t_w$  ( $1/t_w$  is bandwidth).  $T(0)$  is the message startup cost  $t_s$ .  
 -This idealised model of communication performance is adequate for many purposes but does break down in some situations (congestion of communication channels).

University of Salzburg, Department of Scientific Computing, HPSC SS 2005

11

## Communication time (exm.)

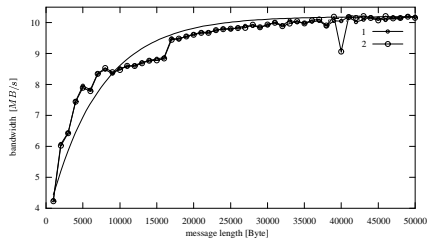


- After a send request a fixed period is spend to service ( $t_s$ ). Then data start to transfer ( $t_w$ ).

University of Salzburg, Department of Scientific Computing, HPSC SS 2005

12

### Communication time (cont.)



Measured bandwidth by 100Mb/s FE (two runs & theory).  
 - Short messages have lower bandwidth because of startup time  $t_s$ .  
 - By longer messages the bandwidth is approaching to the theoretical value (11 MB/s).

### Speedup and Efficiency

- Execution time of a program on a single processor is  $T_1$ .
- Execution time of the same parallelized program on  $p$  (possibly equally loaded) processors is  $T_p$ .
- $T_p$  is composed of probably shorter execution time, some communication time, and some idle time.
- Speedup  $S = T_1/T_p > 1$ ,
- Efficiency  $E = S/p < 1$ .
- Ideally,  $S=p$  and  $E=1$ , typically, speedup increases slower than  $p$  and efficiency decreases.
- Isoefficiency function tells us how the amount of computation must scale with  $P$  to keep  $E$  constant. Algorithm with isoeff. fun. of  $O(P)$  is highly scalable.

### Performance Analysis

- We estimate the total computation time for a calculation domain of dimension  $f(N)$  by analysing the code and the number of floating point operations.

$$T_{comp} = t_c * f(N)$$

- The communication time depends on the number of messages  $M(N,P)$  per calculation step and on the amount of data  $D(N,P)$  to be transferred. Clearly both parameters depend on the number of processors.

$$T_{comm} = P * M(N,P) * D(N,P)$$

- If idle time is 0, because of well balanced computation the execution time of a parallelized program is:

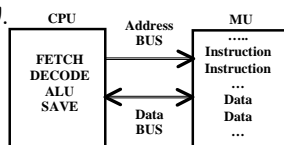
$$T_p = (T_{comp} + T_{comm}) / P$$

### How to increase computer performances ?

- By increasing the data parallelism in a processor (parallel busses, etc.)
- By using multiple functional units and pipelining,
- By using vectorized processors able to operate on several vectorized items at the same time,
- By using multiple computing and memory units,
- By increasing the system clock.
- ???

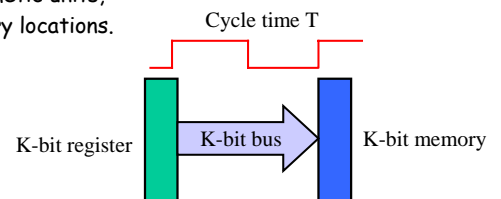
### Single Computer Model (Von Neuman)

- Control unit (CU) - fetch and decode instructions,
- Processing Unit (PU) - execute instructions,
- $CU + PU = CPU$ ,
- Random Access Memory (RAM) Unit (MU) - save data and instructions,
- CPU read instructions and data from MU and write results back to the MU.



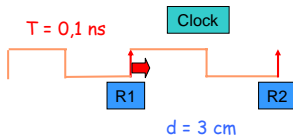
### Data Parallelism

- Parallel bits  $\Rightarrow$  bytes  $\Rightarrow$  words  $\Rightarrow$  etc.
- Registers, data buses, address busses,
- Arithmetic units,
- Memory locations.



## Physical Limits

- Speed of light and **electromagnetic field**  
 $v = 3 \cdot 10^8 \text{ m/s}$ ,
- $f = 10^9 \text{ Hz} = 1 \text{ GHz} \Rightarrow T = 10^{-9} \text{ s} = 1 \text{ ns} \Rightarrow \text{Distance}$   
 $d = v \cdot T = 3 \cdot 10^8 \text{ [m/s]} \cdot 10^{-9} \text{ [s]} = 30 \text{ cm}$ ,
- For **10 GHz** we get the distance of **3cm**,



## Pipelining

- Suppose that for finishing a complex repetitive task (e.g., car production) time  $T$  is needed for one person.
- If the task can be divided into  $k$  equally time-demanding subtasks,  $k$  workers can work in parallel, each its own subtask.
- The resources needed for the implementation of our task can be exploited totally with  $k$  workers.
- The first car will come from the pipe in time  $T$ , but the second car will be out already after  $T + T/k$ , and third after an additional delay of  $T/k$ .

## Pipelining (cont.)

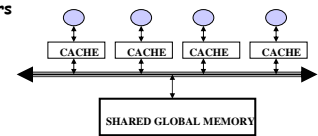
Time  $T$  needed for a single person to complete the task

1/5 of task	1/5 of task	1/5 of task	1/5 of task	1/5 of task
1/5 of task	1/5 of task	1/5 of task	1/5 of task	1/5 of task
1/5 of task	1/5 of task	1/5 of task	1/5 of task	1/5 of task
1/5 of task	1/5 of task	1/5 of task	1/5 of task	1/5 of task
1/5 of task	1/5 of task	1/5 of task	1/5 of task	1/5 of task
1/5 of task	1/5 of task	1/5 of task	1/5 of task	1/5 of task

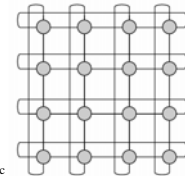
- What will happen if subtasks have diff. time complexity?

## Increasing number of processors

- Shared memory processors (SMP)**  
 (number of processors is limited ( $< 6$ ) because of shared memory access).



- Distributed memory computers or Message Passing Processors (MPP)** (Network Of Workstations (NOW), Clusters (COW), GRIDs, Constellations (clusters of SMP or MPP)  
 (communication channels are needed).



## Properties of SMP

- Advantages**
  - no need for explicit communication mechanisms
  - the data distribution is completely transparent to the user
  - simpler and more compact program code
- Disadvantages**
  - the connection of the CPUs to each other and to the memory (caches can alleviate the bandwidth problem)
  - speed of the memory
  - in practice the number of processors limited to 10-20

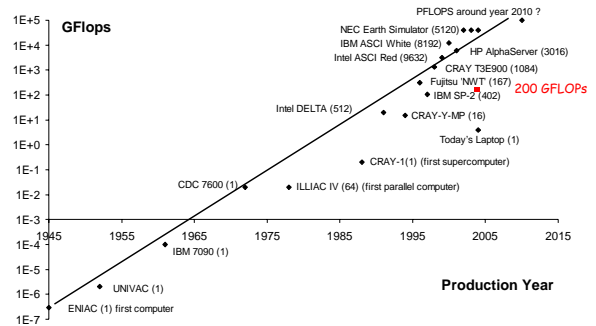
## Properties of MPP

- Advantages**
  - communication **bandwidth scales** up automatically with the number of processors
  - no need** for very fast memory
  - unlimited** number of processors
- Disadvantages**
  - global communication** is much slower (synchronization slower too)
  - user has to distribute** the data over the processors
  - data **exchange** between processors has to be performed explicitly by calling a communication routine.
  - more **complex program code** (actual speed-up  $\ll$  theoretical speed-up)

## Performance Measures

- Floating Point Operation is an arithmetic operation on 64 bit wide data.
- Number of Floating Point Operations per second (R):  
Mflop/s =  $10^6$ , Gflop/s =  $10^9$ , Tflop/s =  $10^{12}$
- Clock of the Central Processing Unit (MHz =  $10^6$ /s, GHz =  $10^9$ /s),
- Memory (Cache, Bulk, Back-up) (Mbit=Mb, Mbyte=MB, Mword=Mw).
- A problem has order of magnitude  $f(N)=O(g(N))$  if there is a positive constant  $c$  such that  $|f(N)| \leq c \cdot |g(N)|$  (e.g., for  $2/3 \cdot N^2 + N$  operations the order of magnitude is  $O(N^2)$ , because terms of lower order than  $N^2$  become insignificant if  $N$  becomes large.

## R<sub>max</sub> in the period 1945 - 2004



## Moore's Law

- We can learn from the previous graphs that the computer performance increased **exponentially**.
- From  $3 \cdot 10^{-7}$  GFlops in 1945 to  $4.1 \cdot 10^4$  GFlops in 2004 (**factor of  $1.4 \cdot 10^{11}$  in 59 years**).
- $(R_2/R_1) = 2^{(t_2-t_1)/T_d}$   
( $T_d$ =time for double performance;  $\log_2=0.3$ )
- $T_d = 59 \cdot \log(2) / \log(2 \cdot 10^{11}) = 1.57$
- Moore's law** states that the computer performance doubles each **1.5 year** (long-term agreement OK).

## Top10 from Top500 list June 2004

Rank	Site	Country / Year	Computer / Processors	Computer Family	Inst. Type	Peak	Memory	EXPLANATION OF THE PEAK
1	Earth Simulator Center	Japan/2002	Earth Simulator / 3120 SPS	NEC Vector	Research	35840	1,0732e+0	268240
2	Lafayette Laboratory National Laboratory	United States/2004	Think Intel Xeon/D, Intel L, Xeon, Quadrics / 4096 Dellware Digital Corporation	NOW - Intel Itanium Dynamic Tige4 Cluster - Quadrics	Research	19940	979000	110000
3	Los Alamos National Laboratory	United States/2002	ASCI Q, AlphaServer SC43, L3000 / 8192	HP AlphaServer Alpha Server Cluster	Research	13980	630000	20480
4	IBM - Rochester	United States/2004	BlueGene/L D01 Prototype D, Data PowerPC 980 W4000 / 4192	IBM BlueGene/L BlueGene/L	Vendor	11880	331778	
5	NCSA	United States/2003	Fujitsu PowerEdge 3730, P6 Super L3000 Site, System / 2500	Dell Cluster PowerEdge 1750, Myraid	Academic	9819	630000	
6	SCOPF	United Kingdom/2004	server adobe 380 (L3000) Intel	IBM SP PowerPC / 2112	Research	8985	180000	
7	Institute of Physics and Astronomy	United Kingdom/2004	RIXM Super Combined Cluster / 2048	Fujitsu Cluster Fujitsu Cluster	Research	8728	474000	150000
8	IBM - Thomas Watson Research Center	United States/2004	BlueGene/L D02 Prototype (P6) Data PowerPC 980 / 4096 IBM L3000	IBM BlueGene/L BlueGene/L	Research	8625	294911	
9	Lafayette Laboratory National Laboratory	United States/2003	Intel Xeon/D, Intel L, Xeon, Quadrics / 18432	HP Cluster Integrity v2000 Integrity Cluster	Research	8620	630000	140000
10	International Supercomputer Center	China/2004	Compaq 8600, SuperM 2.2 SPS, System / 2345	NOW AMD NOW Cluster - AMD - Genomg	Research	8581	728400	180000

**R<sub>max</sub>** and **R<sub>peak</sub>** values are in GFlops. **R<sub>max</sub>** of the cluster in Itzling is about 200 GFlops (200 times smaller than No.1).

