# Introduction to Dimemas
## Performance Analysis Tool for Parallel Programs and Platforms

**dimemas@cepba.upc.es**

European Center for Parallelism of Barcelona
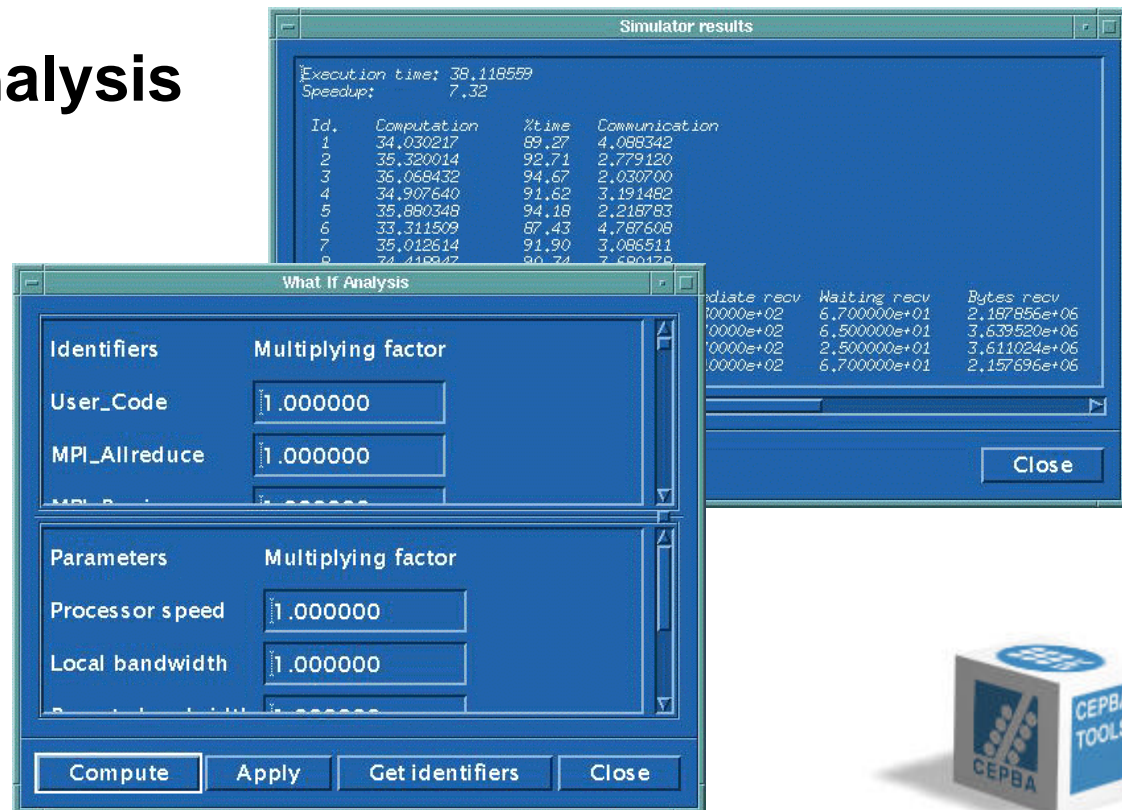
Technical University of Catalonia
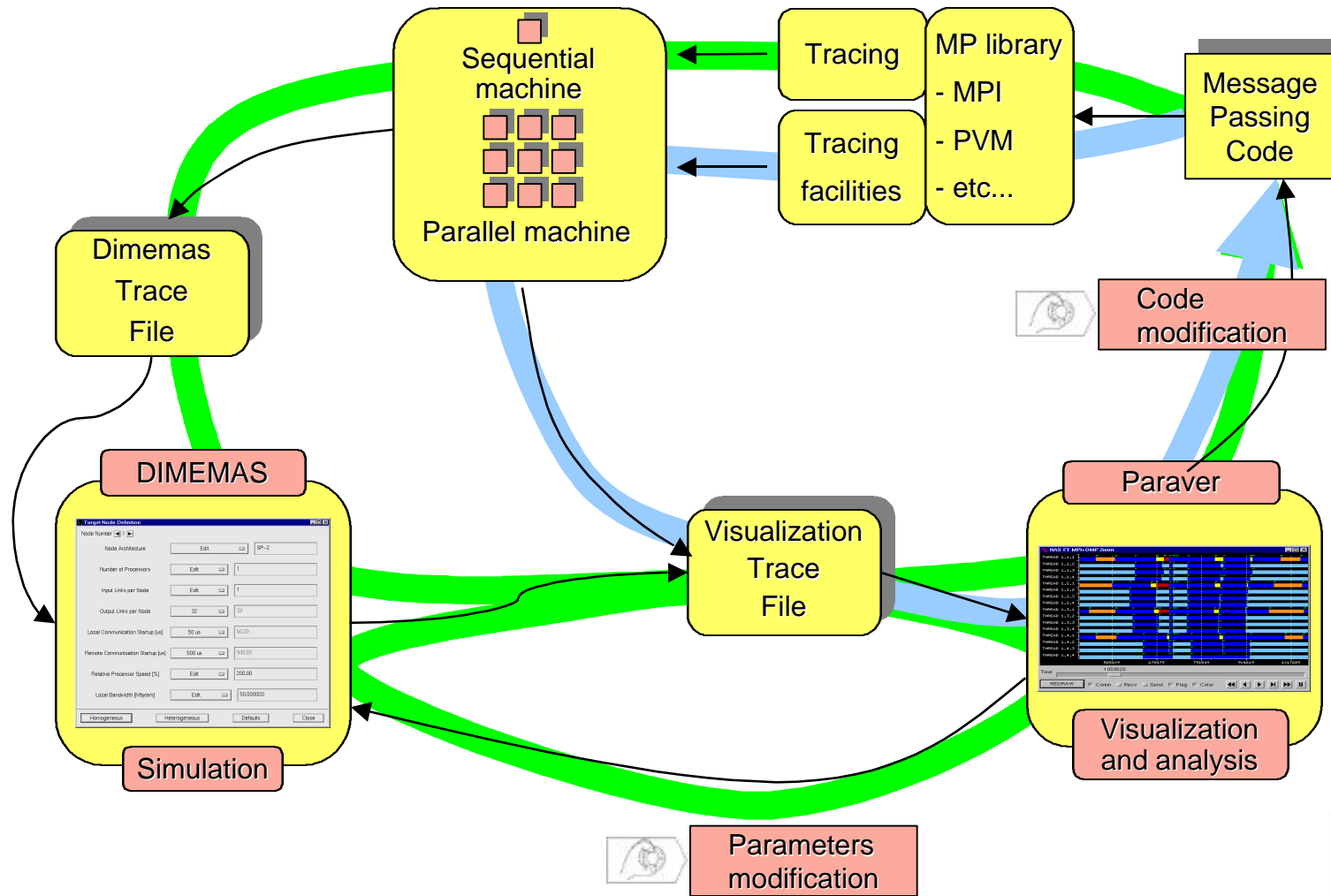
Barcelona, Spain

# Dimemas

- **Application performance analysis tool for message passing programs**

- **In development since 1992**

- **Perform all the analysis on a workstation**
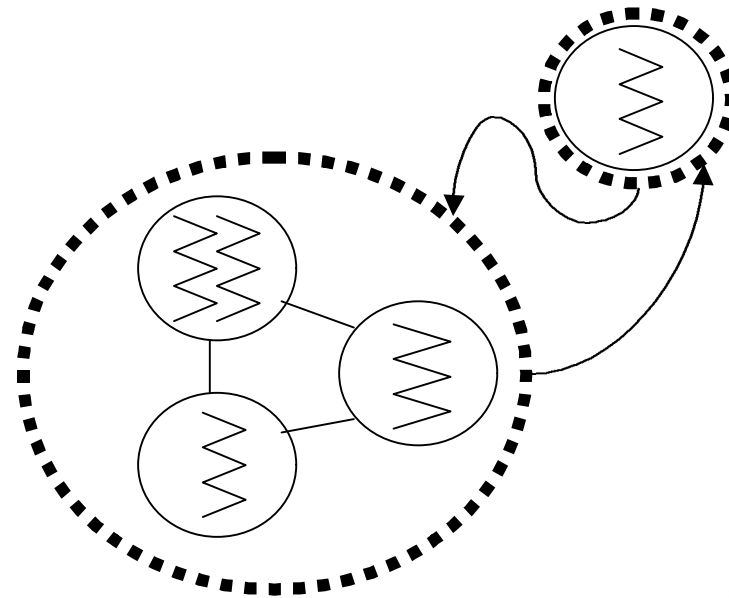
# Tuning Methodology



dimemas@cepba.upc.es

## Characterizes application

- Sequence of resource demands for each task
  - ✓ **CPU bursts**
  - ✓ **Communication**
- Sequence of events:
  - ✓ **Block/routine entry and exit flags**
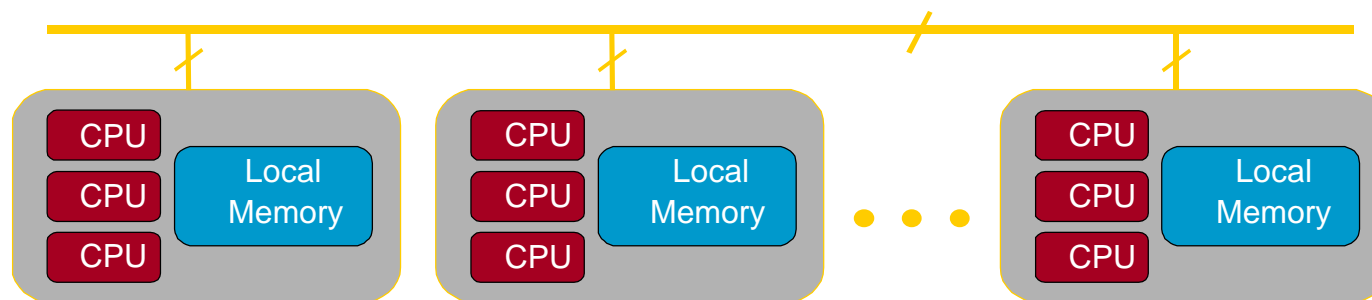
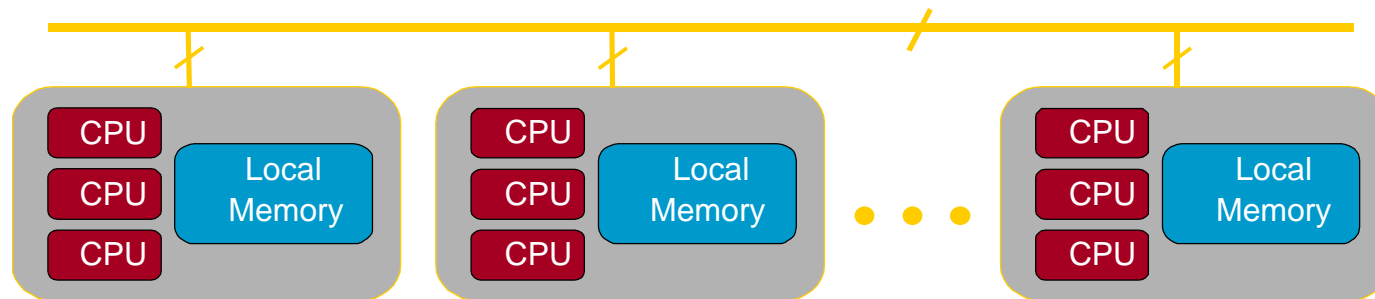## Application model

# Simulated Architecture

- **"Abstract" architecture**
  - Simple and general
    - ✓ **Network of SMPs**
  - Fast simulation
  - Key factors influencing performance
  - Abstract interconnect
    - ✓ **Local/remote latency/BW**
    - ✓ **Injection mechanism (#links, half/full duplex)**
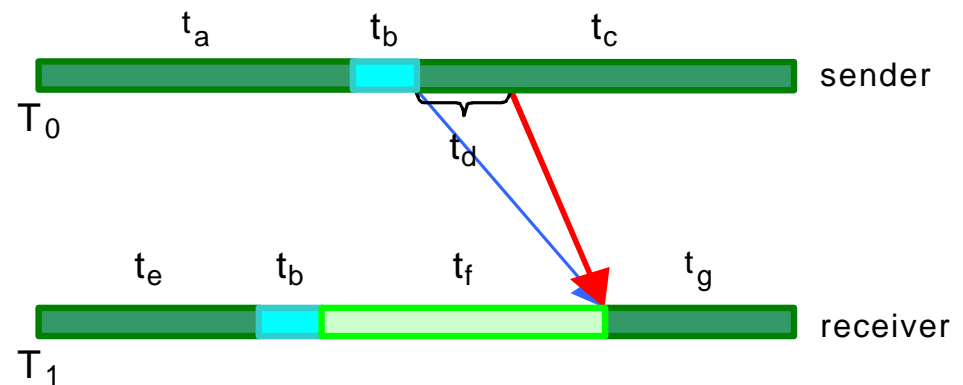    - ✓ **Bisection BW, contention**

# System

- **Process to processor mapping**

- **Multiprogramming**
  - Tasks sharing node
  - Different applications

$$T = \text{Latency} + \frac{\text{Size}}{\text{Bandwidth}}$$
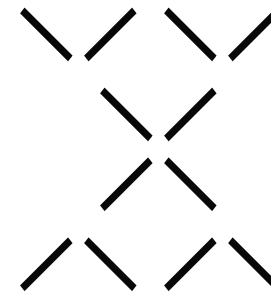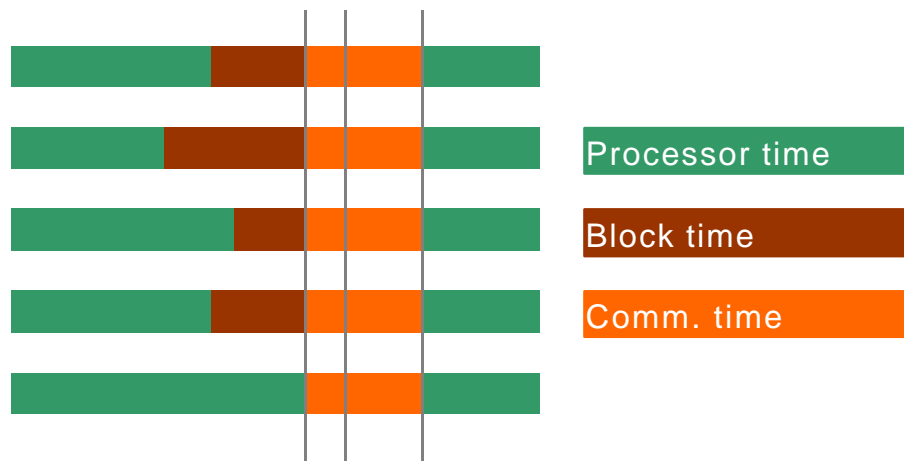
- **Latency**

- **Bandwidth**

- **Resource contention**

# Collective Communication Model

■ **Model hypotheses**

- Implicit barrier
- Fan-in/fan-out phases
  - ✓ **Null/Constant/Linear/Logarithmic model for each phase**



Processor time

Block time

Comm. time

■ **Communication time**

$$Time = \left( Latency + \frac{Size}{Bandwidth} \right) * MODEL\_FACT \quad OR$$

■ **Model factor**

| Model | Factor |
|---|---|
| Null | 0 |
| Constant | 1 |
| Linear | P |
| Logarithmic | $Nsteps = \sum_{i=1}^{\lceil \log_2 P \rceil} steps_i, steps_i = \left\lceil \frac{C}{B} \right\rceil$ |

# Starting Dimemas

■ **Invoke Dimemas gui:** dimemas

■ **Dimemas main menu**



```
                            dimemas
  Configuration  Simulator  Database        Information
  Current configuration file:
```

■ **Configuration**
  ● First time for each user
    **mkdir $HOME/.DIMEMAS_defaults**
    **cp $DIMEMAS_HOME/DIMEMAS_defaults/* $HOME/.DIMEMAS_defaults**
  ● Every session
    **setenv DIMEMAS_HOME /aplic/DIMEMAS**
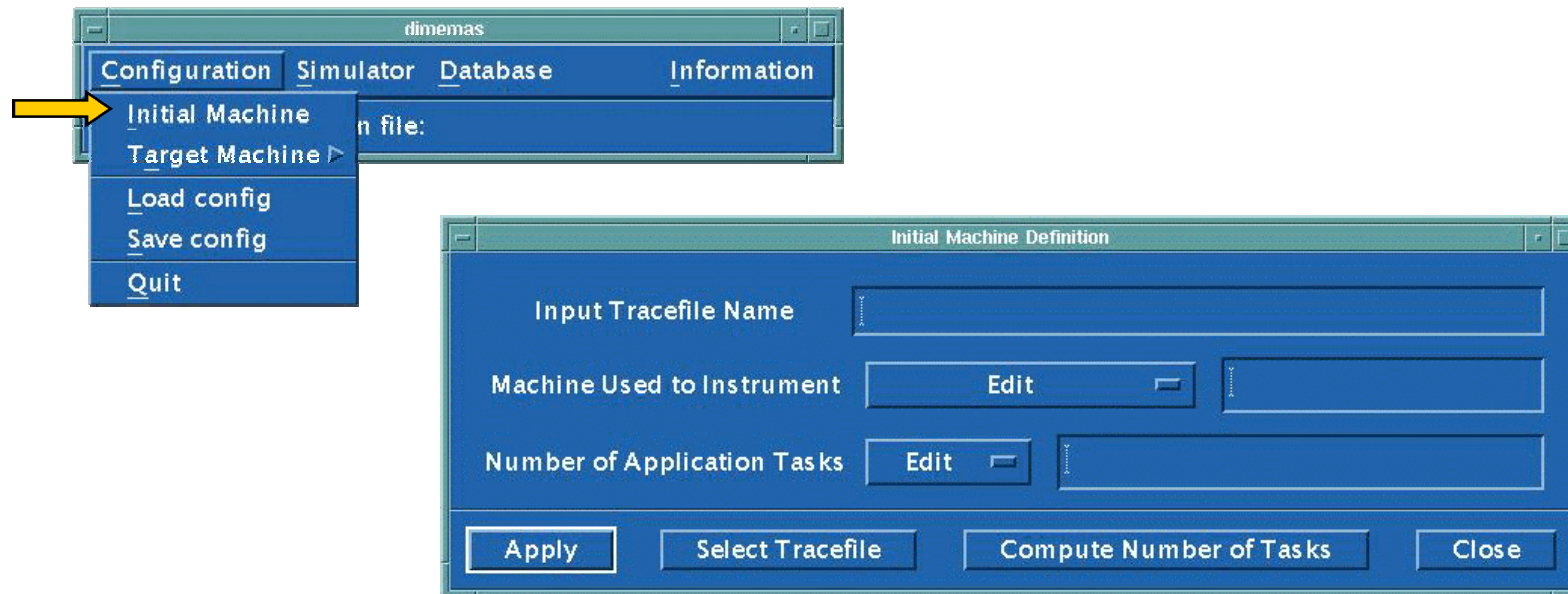
# Dimemas files

- **Configuration files, located in $HOME/.DIMEMAS_defaults**
  - Dimemas
    - ✓ **Configuration file for X-Windows**
  - Machines.db
    - ✓ **Information of machines database**
    - ✓ **Some predefined values**
  - Network.db
    - ✓ **Information of network database**
    - ✓ **Some predefined values**

- **Executables, located in $DIMEMAS_HOME/bin**

- **License file, $DIMEMAS_HOME/etc/license.dat**

# Initial machine definition

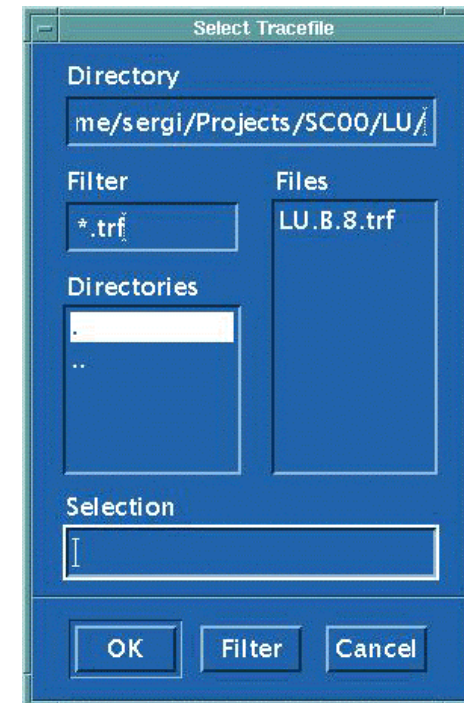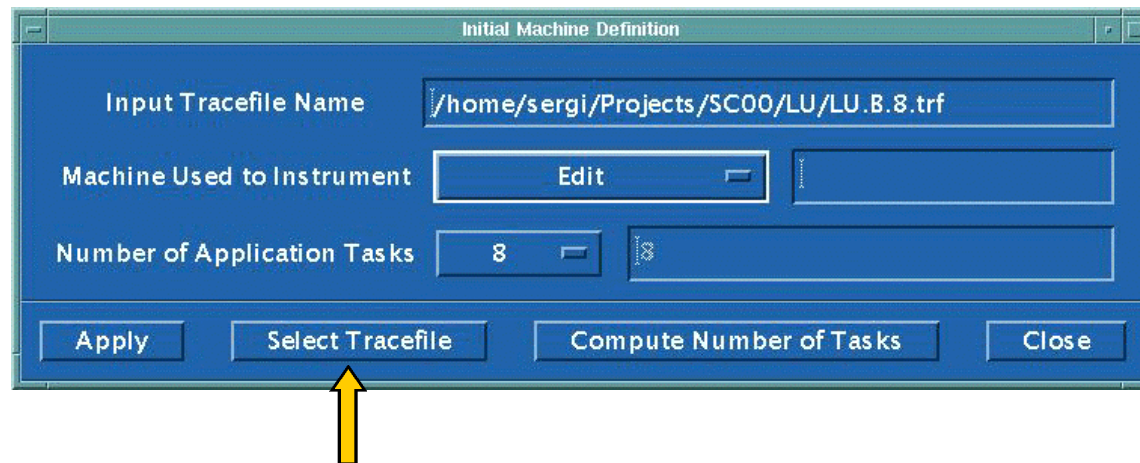■ **To define where the tracefile comes from**



■ **Select an input tracefile and define the initial machine**

# Initial machine definition
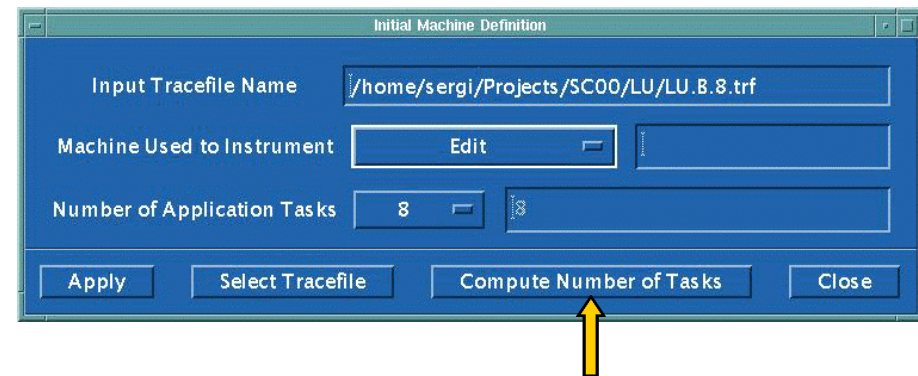
■ **Select a tracefile**

# Initial machine definition

■ **Specify trace file**

- Type or use Select browser

■ **Compute the number of tasks**

✓ **Scans the trace file to compute the number of tasks**
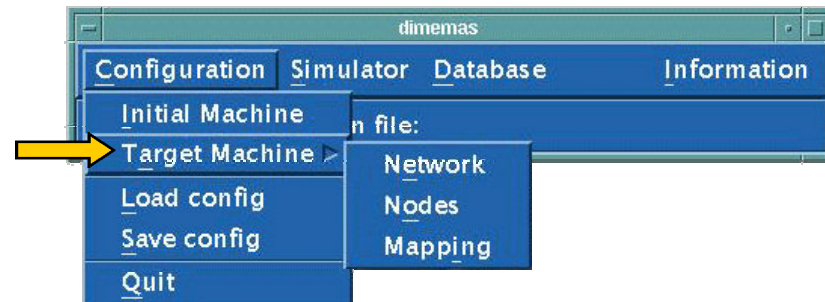


■ **Enter the initial machine (optional)**

- Specifies the machine where the trace was obtained
- If used, the relative speed between initial and target machine can be obtained from the machines.db file

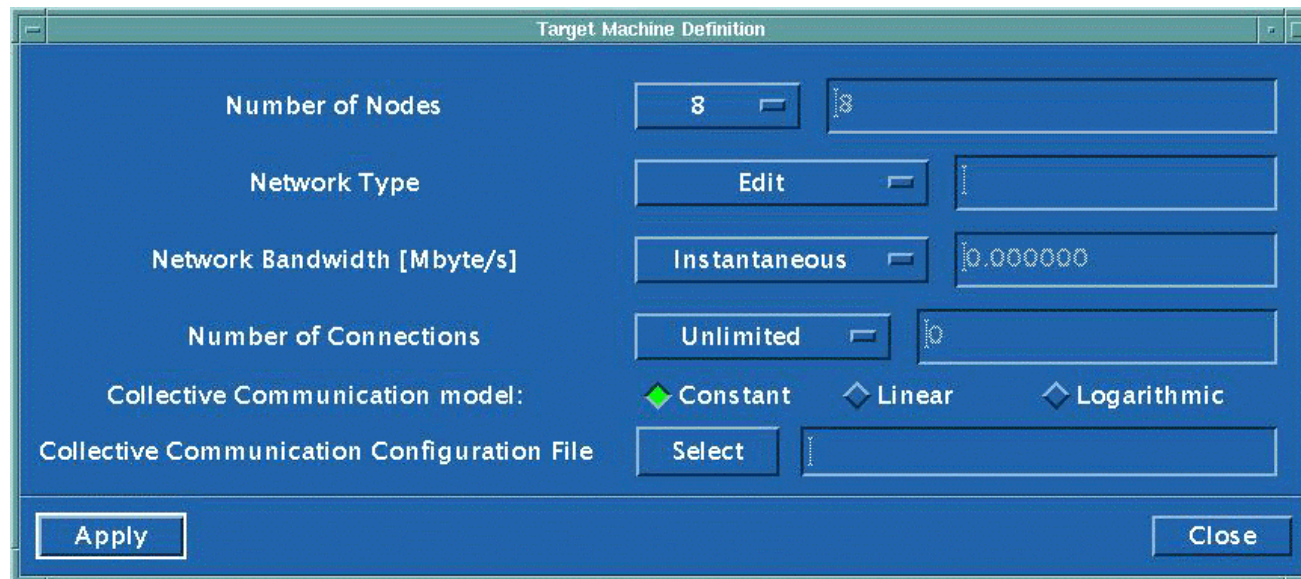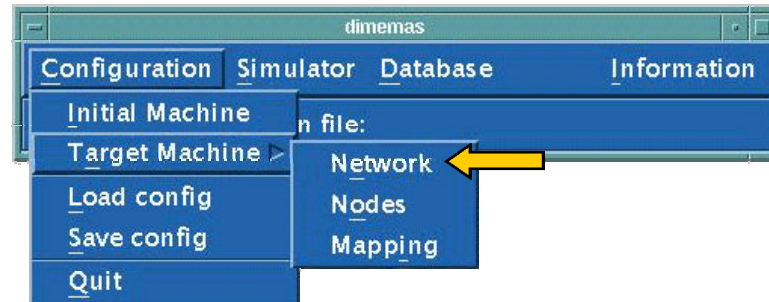■ **Click the Apply button to confirm the input**

# Target machine

■ **Machine simulated by Dimemas**

■ **Description composed by**



- target network definition
- target node definition
- process to node mapping
- file system parameters

# Network definition

# Network definition

- **Number of nodes in the system**

- **Network type (optional)**
  - If used, reads parameters from network database

- **Network bandwidth**
  - Bandwidth for inter-node messages

- **Number of connections**
  - Maximum number of simultaneous messages in transit (simple model of network contention)

- **Collective communication model**
  - Default model for all collective operations (constant/linear/log)

- **Collective communication configuration file (optional)**
  - If used, specifies a detailed model for each collective operation

# Network definition

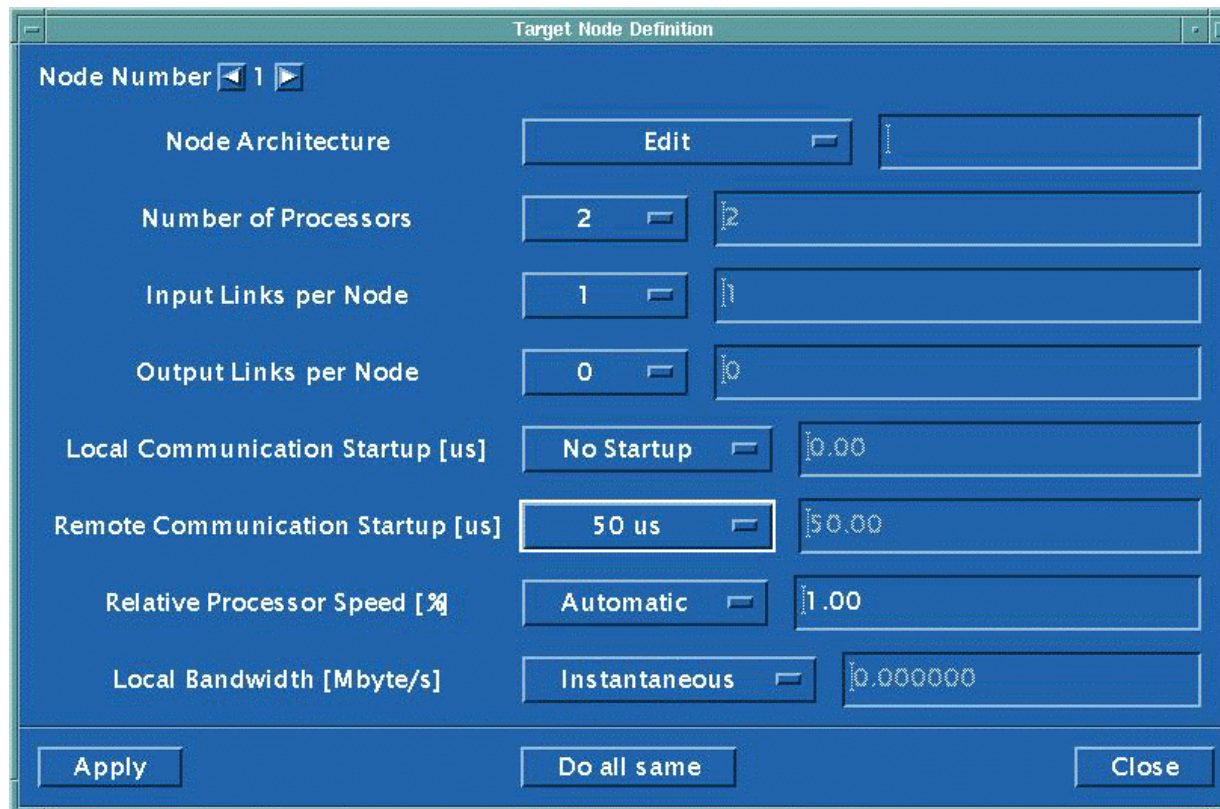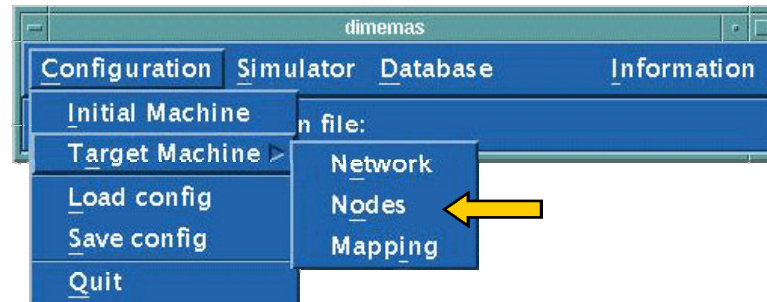■ **Collective communication configuration file format**

- Model and message length to use for each phase
- Example

| Policy: FIFO | | | | | | This line is mandatory |
|---|---|---|---|---|---|---|
| 0 | LIN | MAX | 0 | MAX | 1.0 | MPI_Barrier |
| 1 | LOG | MAX | 0 | MAX | 1.0 | MPI_Bcast |
| 2 | LOG | MEAN | 0 | MAX | 1.0 | MPI_ Gather |
| 3 | CTE | MAX | CTE | MAX | 1.0 | MPI_Gatherv |
| 4 | LOG | MEAN | 0 | MAX | 1.0 | MPI_Scatter |
| 5 | CTE | MAX | CTE | MAX | 1.0 | MPI_Scatterv |
| 6 | LIN | MIN | LIN | MIN | 1.0 | MPI_Allgather |
| 7 | LIN | MIN | LIN | MIN | 1.0 | MPI_Allgatherv |
| 8 | LIN | MIN | LIN | MIN | 1.0 | MPI_Alltoall |
| 9 | CTE | MAX | CTE | MAX | 1.0 | MPI_Alltoallv |
| 10 | LOG | MAX | 0 | MAX | 1.0 | MPI_Reduce |
| 11 | LOG | MAX | LOG | MAX | 1.0 | MPI_Allreduce |
| 12 | LOG | MAX | LIN | MIN | 1.0 | MPI_Reduce_Scatter |
| 13 | LOG | MAX | LOG | MAX | 1.0 | MPI_Scan |

Funct. ID    Fan-in Model    Fan-in length    Fan-out model    Fan-out length

# Node definition

# Node definition

- **Node Number**
  - Node to which the parameters apply

- **Node architecture (optional)**
  - If used, loads parameters from machine database

- **Number of processors**
  - Size of SMP node

- **Input/output links**
  - Maximum number of simultaneously incoming/outgoing messages to/from the node. Simple model of node injection mechanisms.
  - To specify half duplex links, input or output must be zero

- **Local communication startup**
  - Latency used for messages within the node

# Node definition

- **Remote communication startup**
  - Latency used for messages through the network

- **Relative processor speed**
  - Processor speed ratio between the target processor and the processor where the trace was obtained

- **Local bandwidth**
  - Bandwidth for communication within SMP

- **Click the Apply button to confirm the input for the current node**

- **Click the Do all same button to copy the confirmed values of the current node to remaining nodes**

# Process mapping

- **Define the process– to– node mapping**



- **Click the button close to Node number, to open the node definition dialog**

# Configuration file

■ **File storing the initial and target model parameters**

- Save the initial and target machine definitions with the menu function
- Load configuration files with the menu function

# Simulator execution



- **Once the initial machine and the target machine have been defined, the simulator can be used**

- **Click on the Call Simulator button to run the simulator**

# Simulator results

- **The simulator will display the performance results**



- **Click on the Save button to save the simulation results to disk (the configuration file is also saved)**

# Simulator results

- **Global statistics:**
  - execution time: Modeled elapsed time in seconds
  - speedup: Total CPU time / execution time

- **Per– process statistics:**
  - computation time: absolute time in seconds and percentage
  - communication time: absolute time in seconds

  - Sent: number of messages and number of bytes sent
  - Received: number of messages received without blocking (the message had arrived before the reception request), number of messages received with blocking and number of bytes received
  - Group operations: number and total size

# Simulation parameters



- **Load trace into memory**
  - If set, speeds up the simulation. May cause swapping with large trace files.

- **Break time**
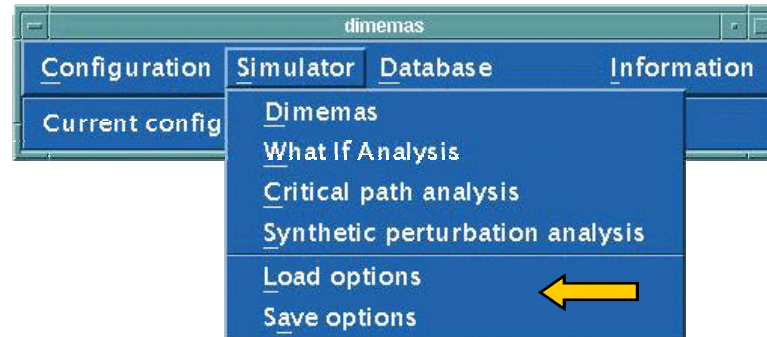  - Define the final simulation time

# Simulation parameters



- **Output tracefile (visualization tracefile)**

  - None, Paraver and Vampir (if license available)

  - Tracefile type: ASCII or Binary

  - Start time and stop time
    - ✓ Specify the initial time and final time for the visualization tracefile
    - ✓ Messages starting before "start time" are delivered at time 0
    - ✓ Messages ending after "stop time" are delivered at stop time

# Simulator options



- **Save execution options for Dimemas (command line options to the simulator)**

- **Load previously saved execution options**

■ **Initial studies that may help understanding the application characteristics**

- Does the application have load balanced and dependence problems?
  - ✓ **BW = ∞,   L = 0**
- Would we benefit from grouping messages?
  - ✓ **L = ... ,   BW = ∞**
- Is bandwidth the problem?
  - ✓ **BW = ...  , L = 0**
- Is network contention the problem?
  - ✓ **BW = target , L = target , Buses = 1, 2, …**

# Methodology

- Perform some of the above extreme simulations generating Paraver traces. Visualize the traces to perceive the general behavior and on which parts of the application the different parameters are more relevant.

- Perform parametric studies (many simulations without generating traces, just reporting the execution time) to quantify the influence of the different parameters.

- Iterate the process going from very large dynamic range of the parameters to points closer to the estimated operation range of the target machine. If the visualization identifies that problems are only on a specific trace section, you can generate with Paraver a tracefile of only that section. In this way the simulation will be faster and the parametric studies will reflect the effect of the parameters on the section of interest.

- **The following tools within Dimemas may be helpful in this analysis**

# What if analysis



- **Analyze the behavior of the application under relative parameter modifications**

# What if analysis

- **Upper area: subroutines/block names**
  - The multiplying factor applies to de CPU demand within the subroutine or block.
  - Examples
    - ✓ a value of 0.75 for a given routine could be used to estimate the effect of sequential code optimizations that speed up the computation within that routine by 25%.
    - ✓ A value of 0 for a given routine can be used to estimate the effect on the global performance of totally eliminating that routine.

- **Lower area: architectural parameters**
  - The multiplying factor applies to the parameter as defined in the target mode definition section.
  - Examples:
    - ✓ a value of 2 in relative processor speed could be used to estimate the effect on the execution time of a CPU twice as fast as the one defined by the value specified in the target node definition section
    - ✓ a value of 0 in the network latency could be used to estimate the effect of an ideal (null) network latency

# What if analysis



- **Click on the Compute button to predict effects**

- **Click on the Apply button to save this changes to the current configuration**

- **Click on the Get Identifiers button to analyze the tracefile to get the application identifiers**

# Critical Path analysis



- **Percent computing and computation time on critical path**

- **Click on the Save button to save the critical path to disk**
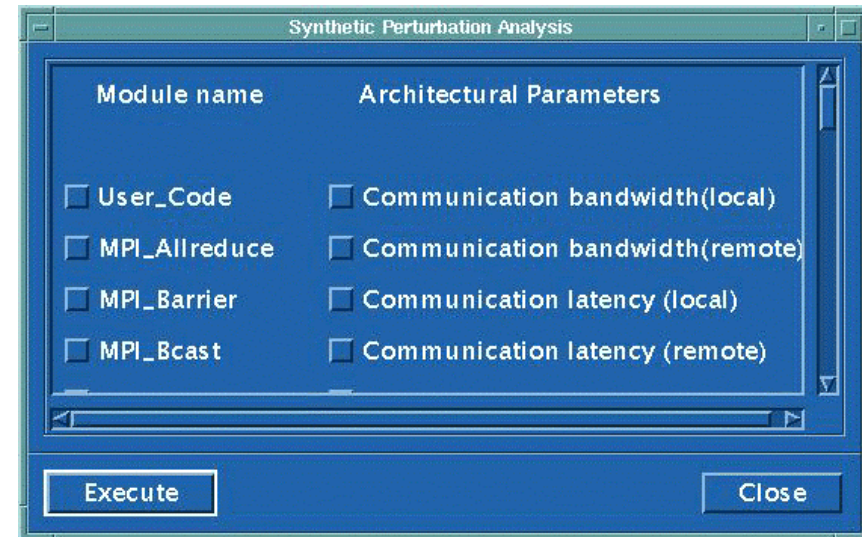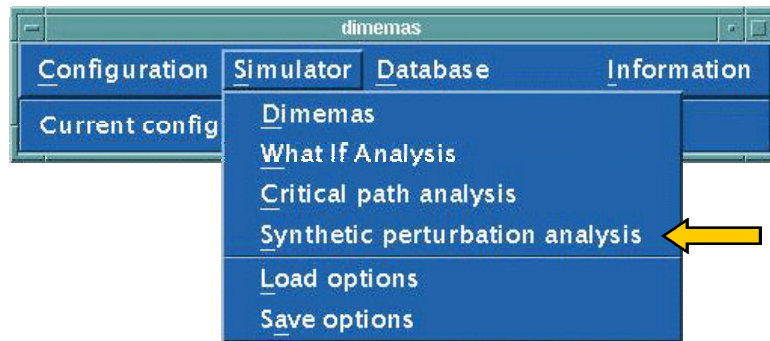
- **Paraver is able to show you the Critical Path**

# Critical Path analysis

- **Visualization of CP using Paraver for LU, class B, 2 iterations, 8 tasks, contention network**

# Synthetic Perturbation Analysis

■ **Full factorial analysis**



■ **Select "interesting" factors**
  ● routine sequential computation demands
  ● architectural parameters

■ **Click on Execute button to compute model**

# Synthetic Perturbation Analysis

- **Model**

$$Time = T_0 + T_1 * F_1 + T_2 * F_2 + T_{12} * F_1 * F_2$$

- **$T_0$ is the model execution time for the nominal parameter set up.**

- **$F_i$ represents the percentage variation of the $i^{th}$ parameter value.**
  - F=1 represents a +10% variation of the parameter
  - F=-1 represents a -10%

- **$V_i$ is the coefficient for factor i**

- **$V_{ij}$ are coupled coefficient for factors i and j**

- **SPA displays**
  - Influence of each factor: Percentage of total variability explained by each factor
  - Model coefficients

# Synthetic Perturbation Analysis

- **Example output**



- **ssor** is the most influent routine (more than 85%)

- Comparing the coefficients with $T_0$ exposes the actual magnitude of the variation (2 orders of magnitude difference in this example)

# Machine database



- **To maintain the information of the machines**

- **To define the ratio speed between to machines in DB**

- **Information to automatically complete fields in the GUI**

- **Load and Save database to disk**

# Machine database
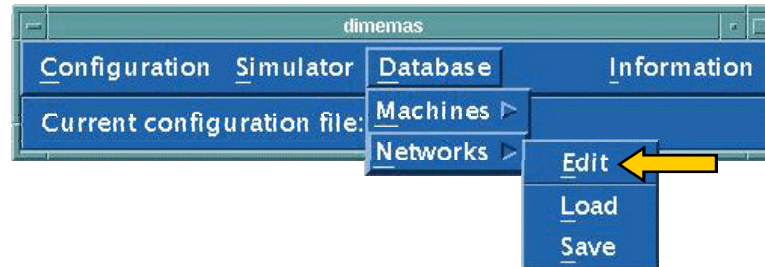
**New machine**

# Machine database

■ **Modify relative processor speed**
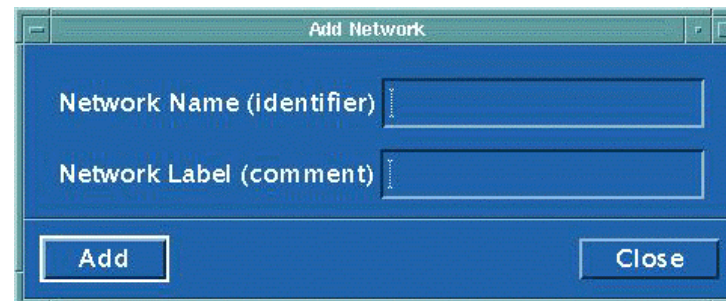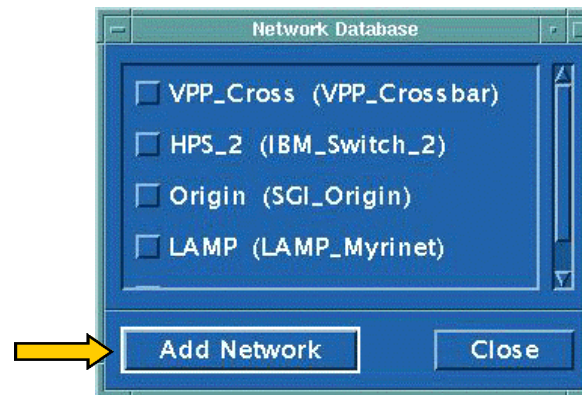
# Network database



- **To maintain the information of the networks**

- **Information to automatically complete fields in the GUI**
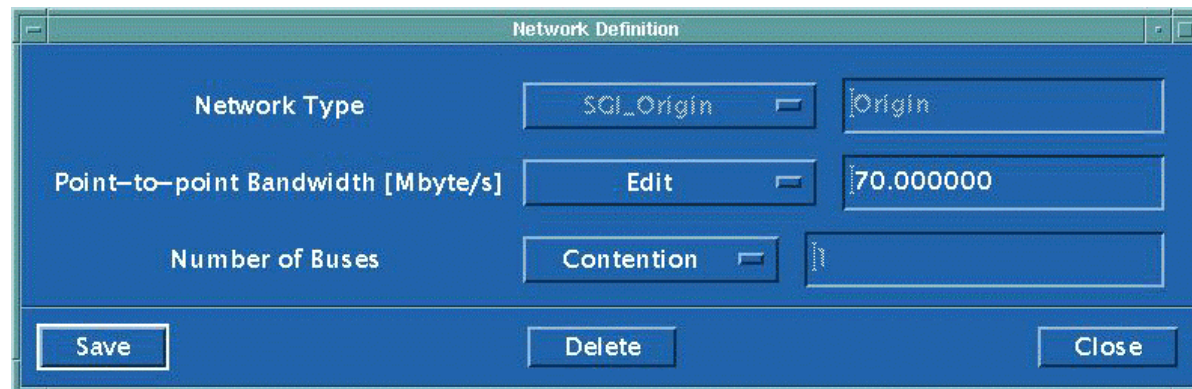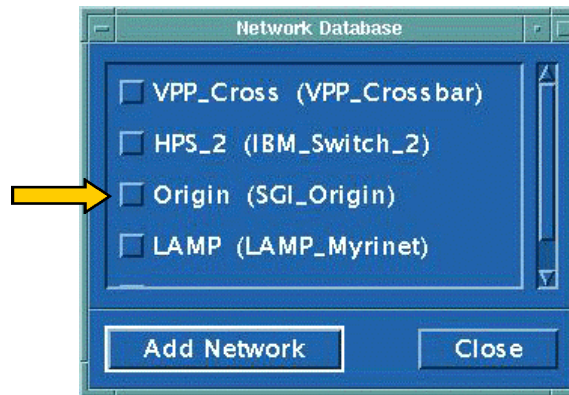
- **Load and Save database to disk**

# Network database

- **New network**

# Network database

■ **Network parameters modification**

# More information

## http://www.cepba.upc.es/dimemas

## dimemas@cepba.upc.es