

# Validation of Dimemas communication model for MPI collective operations

Sergi Girona, Jesús Labarta, Rosa M. Badia

European Center for Parallelism of Barcelona

Departament d'Arquitectura de Computadors

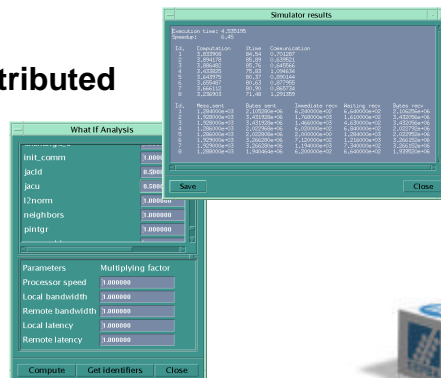
Technical University of Catalonia

Barcelona, Spain



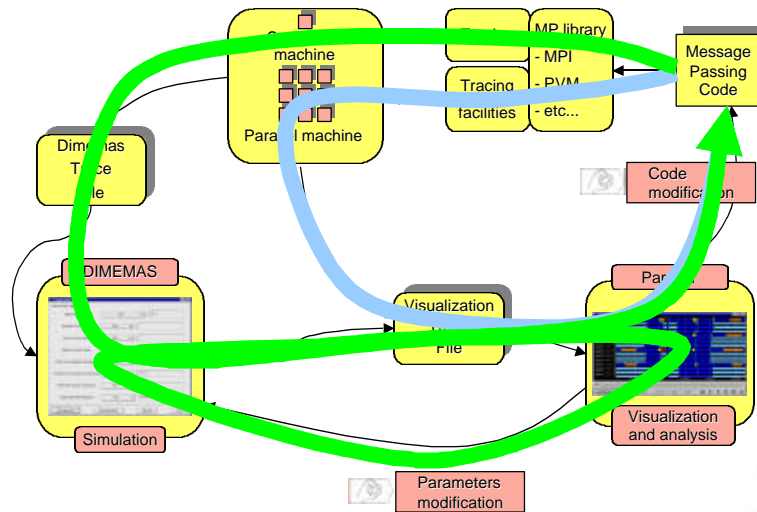
## Dimemas

- Application performance analysis tool for message passing programs
- In development since 1992
- On a workstation
- Dimemas currently distributed by CEPBA



Sergi Girona, EuroPVM/MPI'2000

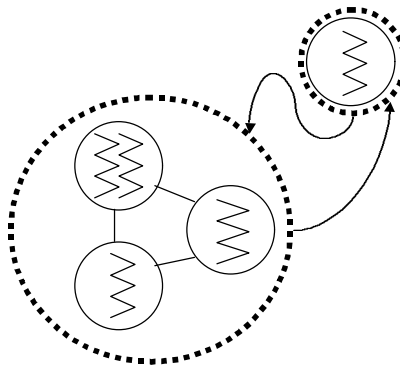
## Tuning Methodology



Sergi Girona, EuroPVM/MPI'2000

## Tracefile

- **Characterizes application**
  - Sequence of resource demands for each task
  - Sequence of events: communication
- **Application model**

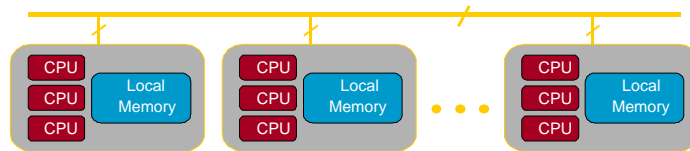


Sergi Girona, EuroPVM/MPI'2000

## Simulated Architecture

### ■ “Abstract” architecture

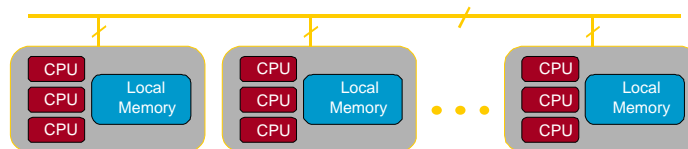
- Simple/General
  - ✓ Network of SMPs
- Fast simulation
- Key factors influencing performance
- Abstract interconnect
  - ✓ Local/remote latency/BW
  - ✓ Injection mechanism (#links, half/full duplex)
  - ✓ Bisection BW, contention



Sergi Girona, EuroPVM/MPI'2000

## System

- Process to processor mapping
- Multiprogramming
  - Tasks sharing node
  - Different applications

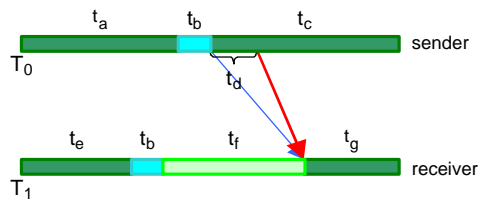


Sergi Girona, EuroPVM/MPI'2000

## Point to Point Communication

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

- Latency
- Bandwidth
- Resource contention

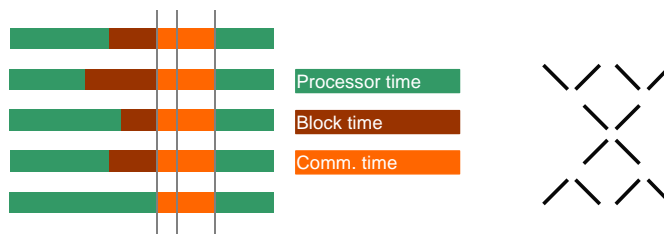


Sergi Girona, EuroPVM/MPI'2000



## Collective Communication Model

- Barrier
- Fan-in/fan-out phases
  - Size of message
  - Null/Const/Lin/Log



Sergi Girona, EuroPVM/MPI'2000



## Collective Communication Model

### ■ Communication time

$$\text{Time} = \left( \text{Latency} + \frac{\text{Size}}{\text{Bandwidth}} \right) * \text{MODEL\_FACTOR}$$

### ■ Model factor

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

Sergi Girona, EuroPVM/MPI'2000



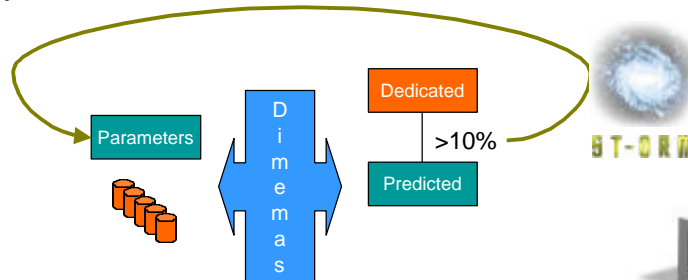
## Parameters Acquisition

### ■ Execution of PBM on SGI Origin

- Dedicated: execution time
- Shared: traces for Dimemas

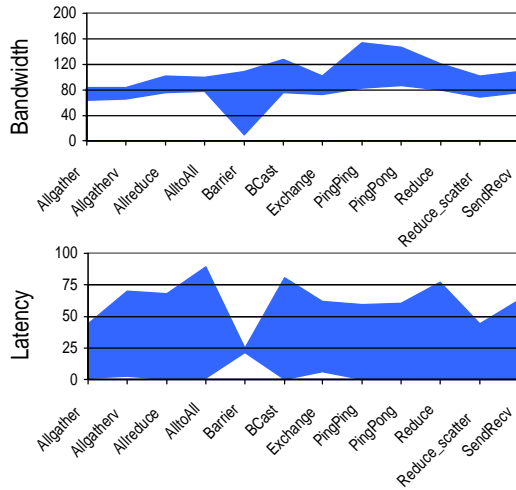
### ■ Compute latency, bandwidth, links, buses, phases, ...

- ST-ORM <http://www.cepba.upc.es/ST-ORM>
- Objective: Predicted time with less than 10% error



Sergi Girona, EuroPVM/MPI'2000

# System Characterization



Sergi Girona, EuroPVM/MPI'2000



# System Characterization

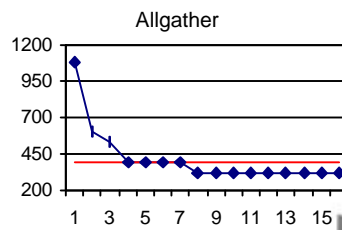
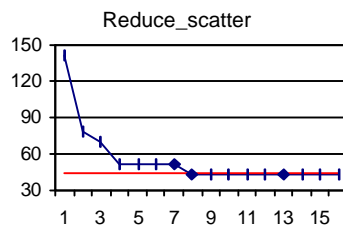
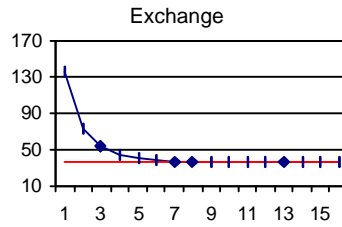
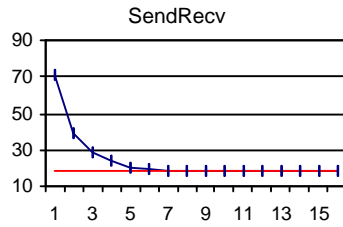
| Operation      | IN    |      | OUT   |      |
|----------------|-------|------|-------|------|
|                | Model | Size | Model | Size |
| Barrier        | LIN   | MAX  | LIN   | MAX  |
| Bcast          | LOG   | MAX  | NULL  |      |
| Gather         | LOG   | MEAN | NULL  |      |
| Gatherv        | LOG   | MEAN | NULL  |      |
| Scatter        | NULL  |      | LOG   | MEAN |
| Scatterv       | NULL  |      | LOG   | MEAN |
| Allgather      | LOG   | MEAN | LOG   | MEAN |
| Allgatherv     | LOG   | MEAN | LOG   | MEAN |
| Alltoall       | LOG   | MEAN | LOG   | MAX  |
| Alltoallv      | LOG   | MEAN | LOG   | MAX  |
| Reduce         | LOG   | 2MAX | NULL  |      |
| Allreduce      | LOG   | 2MAX | LOG   | MAX  |
| Reduce_Scatter | LOG   | 2MAX | LOG   | MIN  |
| Scan           | LOG   | MAX  | LOG   | MAX  |

- Latency = 25  $\mu$ seconds
- Bandwidth = 87.5 MB/s
- 1 HD link per node

Sergi Girona, EuroPVM/MPI'2000



## Influence of Buses

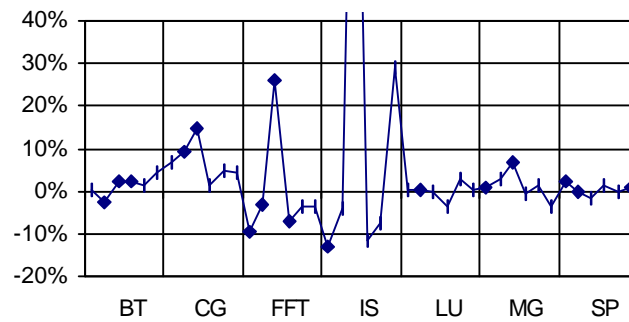


Sergi Girona, EuroPVM/MPI'2000



## Validation

- NAS benchmarks
- Classes W, A
- Size: 8/9, 16, 25/32



Sergi Girona, EuroPVM/MPI'2000



## Conclusions

- Simple but accurate formulation for collective communication
- Methodology for model validation
- Dimemas is a valid tool for performance analysis of message passing programs, parallel machines and message passing libraries
- Future: RMA and I/O operations pending for validation

Sergi Girona, EuroPVM/MPI'2000

