



## SMPI status

SUD, June 13, 2012

# Motivations: why use simulation?

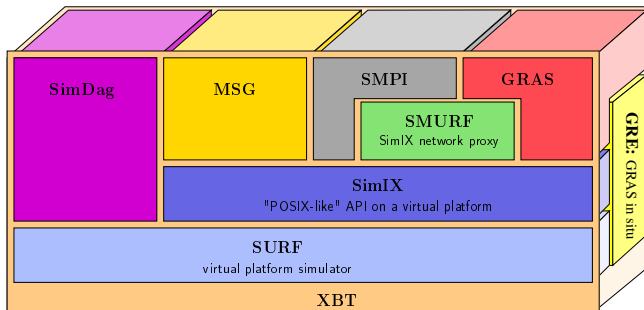
- ▶ Performance Prediction (“what-if?” scenarios)
  - ▶ Platform dimensioning
  - ▶ Tuning of application parameters
- ▶ Teaching (parallel programming, HPC)
  - ▶ No need for real hardware
  - ▶ Handy environment

# Challenges

- ▶ **Accuracy**: How well does the simulation reflect reality?
- ▶ **Scalability**: Which problem size can we simulate? On which platforms?
- ▶ **Speed**: How fast is the simulation as compared to the real execution?
- ▶ **Reproducibility**: Are the simulation results stable and reusable?

# SMPI overview

- ▶ Partial implementation of MPI on top of SimGrid

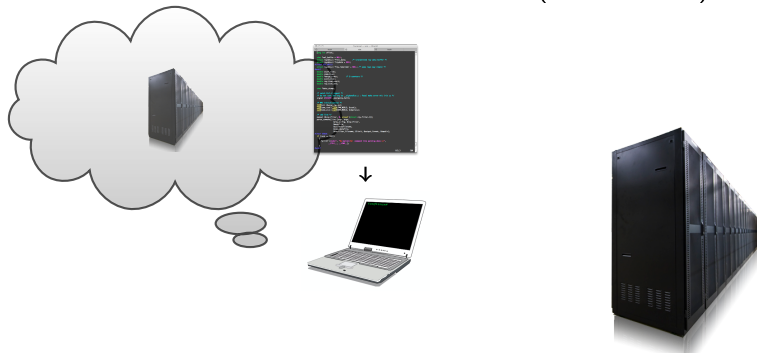


# SMPI overview

- ▶ Partial implementation of MPI on top of SimGrid
- ▶ No or few modifications to the source code (C or Fortran)

# SMPI overview

- ▶ Partial implementation of MPI on top of SimGrid
- ▶ No or few modifications to the source code (C or Fortran)



# Howto

```
[laptop]% smpicc mympiprogram.c -o mympiprogram
```

```
[laptop]% smpirun -np 32 -hostfile machines.txt  
-platform griffon_cluster.xml ./mympiprogram
```

## Howto (with tracing)

```
[laptop]% smpicc mympiprogram.c -o mympiprogram -laky
```

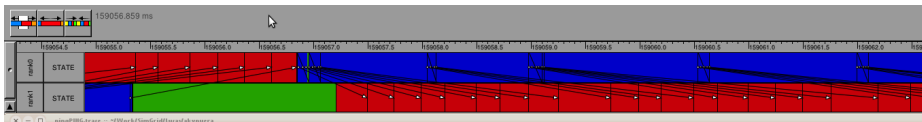
```
[laptop]% smpirun -np 32 -hostfile machines.txt  
-platform griffon_cluster.xml -trace ./mympiprogram
```



# Howto (with tracing)

```
[laptop]% smpicc mympiprogram.c -o mympiprogram -laky
```

```
[laptop]% smpirun -np 32 -hostfile machines.txt  
-platform griffon_cluster.xml -trace ./mympiprogram
```



# SMPI overview

- ▶ Computations: real execution on the host computer
  - ▶ CPU bursts are benched
  - ▶ Scale linearly CPU time according to power ratios

# SMPI overview

- ▶ Computations: real execution on the host computer
  - ▶ CPU bursts are benched
  - ▶ Scale linearly CPU time according to power ratios
- ▶ Communications: simulated
  - ▶ Network models are flow-based models (TCP)
  - ▶ **Validity of these models for MPI applications**

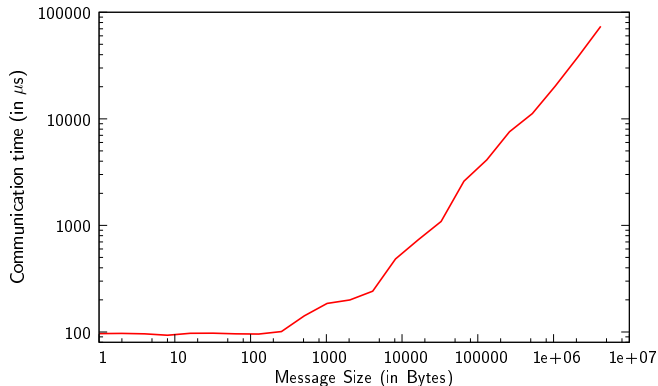
# SMPI overview

- ▶ Computations: real execution on the host computer
  - ▶ CPU bursts are benched
  - ▶ Scale linearly CPU time according to power ratios
- ▶ Communications: simulated
  - ▶ Network models are flow-based models (TCP)
  - ▶ **Validity of these models for MPI applications**
- ▶ Folding of the parallel program processes onto a single node
  - ▶ Serialization of computations
  - ▶ Single address space
  - ▶ **Requires to reduce**
    - ▶ **Memory footprint (scalability)**
    - ▶ **Simulation time (speed)**

# Existing Network Models in SimGrid

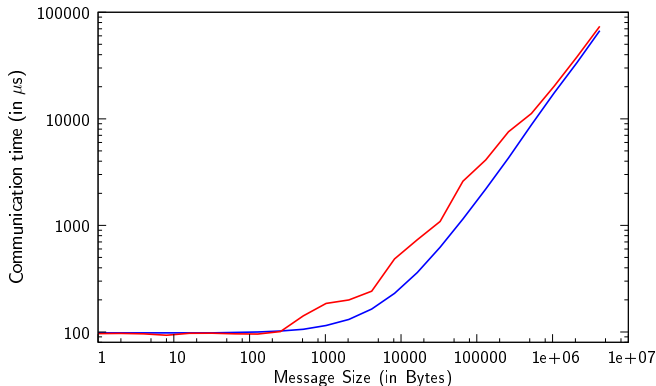
- ▶ Arbitrary topology, endpoints connected through multi-hop paths
- ▶ Network link characteristics: latency (L) and bandwidth (B)
- ▶ Simulation using flows
  - ▶ Simulation is fast ( $\neq$  packet-level simulation)
  - ▶ Contention evaluation is simple
- ▶ Simple Model:  $T(S) = L + \frac{S}{B}$ 
  - ▶ Shown to be valid for  $S \geq 10$  MB
- ▶ Improved model:  $T(S) = \alpha \cdot L + \frac{S}{\min(\beta \cdot B, \frac{\gamma}{2 \cdot L})}$ 
  - ▶  $\alpha$  accounts for TCP slow-start
  - ▶  $\beta$  accounts for the overhead induced by TCP/IP headers (e.g 92%)
  - ▶  $\gamma$  enables the modeling of the TCP window induced behavior
  - ▶ Model valid for  $S \geq 100$  KiB, does not address a lot of message sizes found in MPI applications
- ▶ Need for a new, accurate network model when  $S < 100$  KiB

# Point-to-point Communication



Experimental measurement using SKaMPI

# Point-to-point Communication

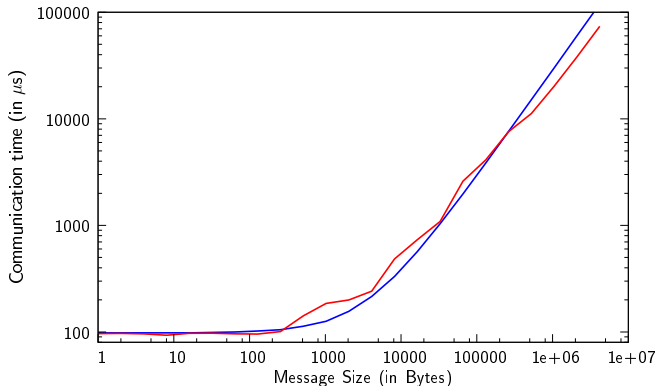


Experimental measurement using SKaMPI

Default linear model, error: 32.1%

Ok with asymptotic message sizes,  
but wrong for 1KiB-1MiB messages

# Point-to-point Communication



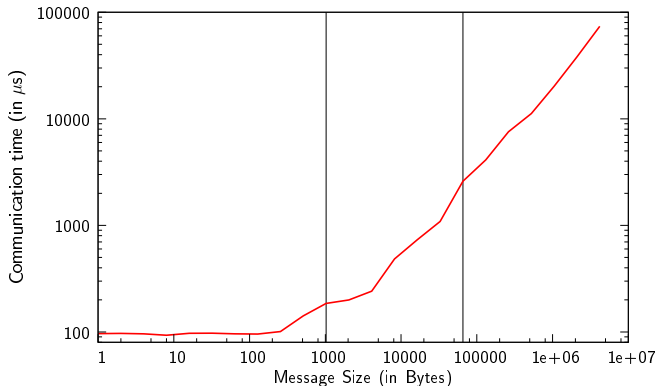
Experimental measurement using SKaMPI

Best-fitted linear model  $(\alpha, \beta, \gamma)$ , error: 18.5%

Better for a lot of sizes,  
but cannot fit all real values



# Point-to-point Communication

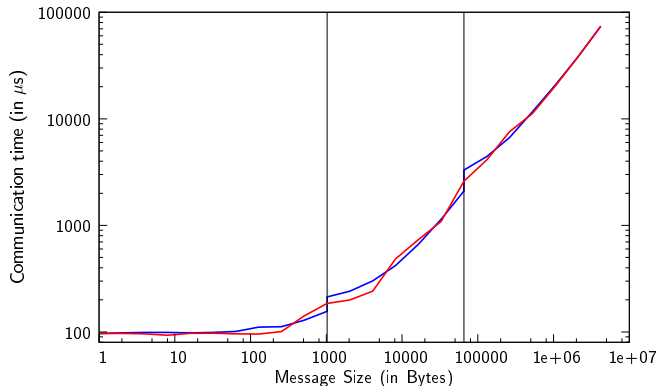


Experimental measurement using SKaMPI

Breakdown depending on message size

- packet size < MTU,
- eager/rendezvous switch limit

# Point-to-point Communication



Experimental measurement using SKaMPI

New piece-wise linear model, error: 8.63%

Correctly adjust linear segments

# Calibration of the piece-wise linear model

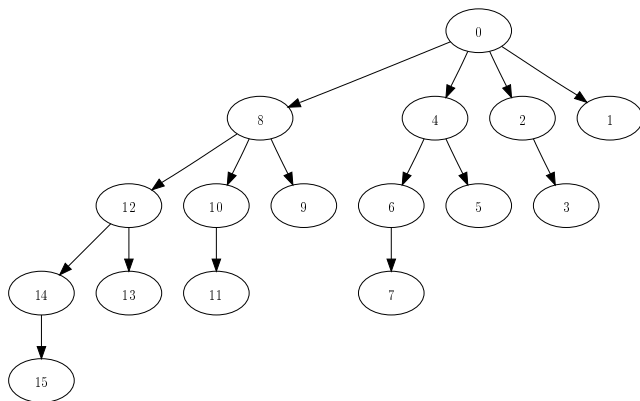
- ▶ Instantiate 9 parameters instead of 3
  - ▶ 2 segment frontiers
  - ▶ 2 factors  $\alpha$  and  $\beta$  per segment
  - ▶ 1 global factor  $\gamma$
- ▶ A calibration script comes with SMPI. Computes parameters given:
  - ▶ 1 SKaMPI-formatted datafile of a ping-pong performance measurement
  - ▶ The number of physical links crossed by packets in the ping-pong
  - ▶  $L$  and  $B$  values for the links
  - ▶ segment bounds (computed by another script)

# Collectives

## Assess contention

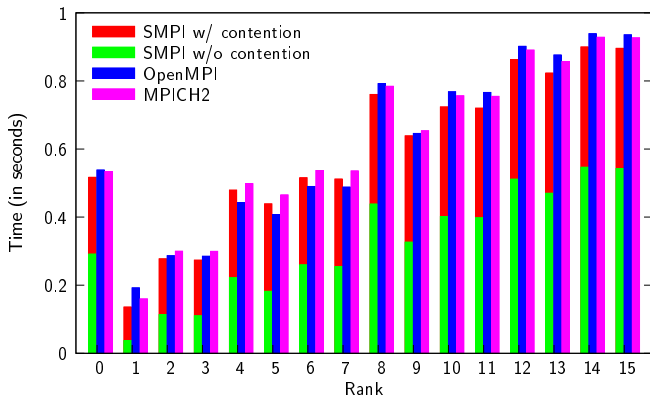
- ▶ Real world (OpenMPI, MPICH2)
  - ▶ Dynamic selection of tuned algorithms
  - ▶ Depends on the number of processes and message size
- ▶ Simulated world (SMPI)
  - ▶ Smaller variety of algorithms
- ▶ For a sake of comparison: use a manual implementation for real and in simulation

## One-to-many: MPI\_Scatter



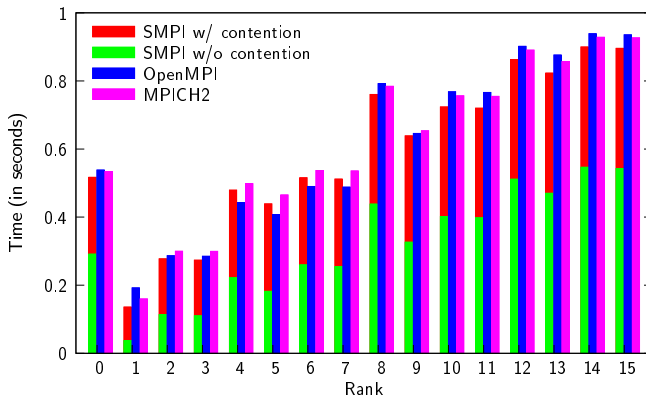
- ▶ Algorithm: A binomial tree
- ▶ 64 MiB at the root, 4 MiB per process

# Scatter: 16-processes test



- Comparison SMPI/MPICH2  $\Leftrightarrow$  OpenMPI/MPICH2: error 5.3%

# Scatter: 16-processes test



- ▶ Comparison SMPI/MPICH2  $\Leftrightarrow$  OpenMPI/MPICH2: error 5.3%
- ▶ Taking contention into account is important

# Reducing the Memory Footprint

- ▶ Idea: Share arrays between processes
  - ▶ Pros: Allocate once, use plenty
  - ▶ Pros: Simulated times stay valid
  - ▶ **Cons:** Computed results become erroneous
- ▶ Implemented as (optional) macros

```
double* data = (double*)SMPI_SHARED_MALLOC(...);  
...  
SMPI_SHARED_FREE(data);
```



# Reducing the Simulation Time

- ▶ Idea: Do not execute all the iterations
- ▶ Use sampling instead
  - ▶ LOCAL: each process executes a specified number of iterations
  - ▶ GLOBAL: a specified number of samples is produced by all processors
- ▶ Remaining iterations are replaced by average of measured values
- ▶ Implemented as (optional) macros

```
for(i = 0; i < n; i++) SMPI_SAMPLE_LOCAL( 0.75*n , 0.01 ) {  
    ...  
}  
...  
for(j = 0; j < k; j++) SMPI_SAMPLE_GLOBAL(0.5*k,0.01) {  
    ...  
}
```

# Reducing the Simulation Time

- ▶ Idea: Do not execute all the iterations
- ▶ Use sampling instead
  - ▶ LOCAL: each process executes a specified number of iterations
  - ▶ GLOBAL: a specified number of samples is produced by all processors
- ▶ Remaining iterations are replaced by average of measured values
- ▶ Implemented as (optional) macros

```
for(i = 0; i < n; i++) SMPI_SAMPLE_LOCAL( 0.75*n , 0.01 ){  
    ...  
}  
...  
for(j = 0; j < k; j++) SMPI_SAMPLE_GLOBAL(0.5*k,0.01) {  
    ...  
}
```

max part of iterations performed

threshold average variability

# Requirements

- ▶ C or F77 fortran code (f2c is used)
- ▶ No global variable
  - ▶ f77: automatic privatization of common (smpif2c)
  - ▶ C: script `patch_source.sh` in `src/smpi`  
`int v`; becomes  
`SMPI_VARINIT_GLOBAL(int, v)`; and then  
`SMPI_VARGET_GLOBAL(v)`;
- ▶ TCP network only, IB (hopefully) in the short term.  
Therefore, we only consider one process per host mappings (no intra-host comms).
- ▶ Not all the MPI API is covered (about 70 primitives)

# Future Work

- ▶ Mid term
  - ▶ Automate the privatization process that enables unmodified code to compile
  - ▶ Handle burst of point-to-point communications
  - ▶ Handle *receiver*-side gap for collectives.
  - ▶ Model other network interconnects: Myrinet, Infiniband
- ▶ Long term
  - ▶ I/O simulation
  - ▶ Automatic memory factoring and loop sampling
  - ▶ Simulation of a full implementation (OpenMPI or MPICH2)