

The SimGrid Project Simulation and Deployment of Distributed Applications



Arnaud Legrand CNRS, MESCAL INRIA project Laboratoire Informatique et Distribution

Martin Quinson Université Henri Poincaré, Nancy 1 LORIA

Henri Casanova Kayo Fujiwara Dept. of Information and Computer Sciences University of Hawaii at Manoa

Developing efficient large-scale concurrent applications poses many challenges:

- Understanding the performance behavior of the code is non-trivial
- Conducting experiments in real-world large-scale platforms is non-trivial
 - Requires a fully functional implementation
 - Limited to a few particular platform configurations
 - In many cases, non-repeatable
- Accurate and validated simulations results are elusive
 - Potentially more accurate emulation is extremely time consuming
- Simulation code is often "throw-away" and may differ from the real code

The SimGrid project addresses all the above challenges via a multi-component software infrastructure for application prototyping, development, and deployment.

> Available at http://simgrid.gforge.inria.fr/

Examples of target applications:

- A parallel linear system solver on a commodity cluster
- A parallel rendering application running on a network of workstations
- A scientific simulation running on a multi-site high-end grid platform
- A network monitoring application running on a wide-area network
- A peer-to-peer file-sharing application running on volatile Internet hosts

Features

- Fast and accurate simulation capabilities (SURF)
- Ability to run the same code in full or partial simulation mode or in real-world mode (GRAS, SMPI)
- An API for rapid application prototyping to test and evaluate distributed algorithms (MSG)
 - Only in simulation mode
- An API for application development to obtain fast, robust and portable application code (GRAS)

17.9ms

- Either in simulation or in real-world mode
- An API for MPI application simulation to study the effect of platform heterogenity (SMPI)
 - In partial simulation mode

Application and algorithm prototyping

- Enables the the easy prototyping of distributed algorithms
 - No need to realize a complete implementation
 - Just focus on the fundamentals of distributed computing

m_task_t local, remote, ack; m host t destination;

/* simulated data transfer */

/* simulated task execution */

MSG_task_execute(local);

return O;

return O;

/* simulate data reception */

MSG_task_get(&ack, PORT_23);

destination = MSG_get_host_by_name(server_host_name)

remote = MSG_task_create("Remote", 30.0, 3.2);

/* 30.0 MFlop, 3.2 MB */

MSG_task_put(remote, destination, PORT_22);

local = MSG_task_create("Loca", 10.50, 3.2);

/* 10.50 MFlop, 3.2 MB */

Application development

Convenience

x86

- API for rapid development of real-world distributed applications
- Simple and cross-architecture communication of complex data structures

Portability (Linux, Mac OSX, Solaris, AIX, IRIX; 12 CPU architectures)

• Performance via efficient communication (computation unchanged)

• Unmodified code run in simulation mode or in real-world mode

• Automatic benchmarking of application code for simulation (CPU)

• Automatic computation of communication volume for simulation (network)

• Resulting application is production, not prototype

Application testing and evaluation



int client(int argc, char * * argv) {

gras_socket_t peer, from;

int ping=1234, pong;

gras_init(&argc, argv);



- Applications consist of processes
- Processes can be created, suspended, resumed and terminated dynamically
- Processes can synchronize by exchanging tasks
- Tasks have a communication payload and an execution payload
- All processes are in the same address space
- Enables convenient communication via global data structure







N Matrix A

Gantt chart for an execution of the above code for **2 servers** and 3 clients

Dark portions denote computations, light portions denote communications

Concurrent communications interfere with each other as the TCP flows share network links





Average time to exchange one Pastry message on a LAN (in seconds) for MPICH, OmniORB, PBIO, and XML-based communication, between PowerPC, Sparc, and x86 architectures



Average time to exchange one Pastry message on a WAN (in seconds) for MPICH, OmniORB, PBIO, and XML-based communication,

Simulation Program Simulation Program With GRAS Without GRAS

gras_os_sleep(1); /* Wait for the server startup */

gras_msgtype_declare("ping", gras_datadesc_by_name("int"));

/* dest, msgtype, payload */

/* timeout, wanted msgtype, &source, &payload */

/ * name, payload */

Code



gras_init(&argc, argv); gras_msgtype_declare("ping", gras_datadesc_by_name("int")); gras_msgtype_declare("pong", gras_datadesc_by_name("int")); gras_cb_register(gras_msg_type_by_name("ping"), ping_callback); gras_socket_server(4000); /* wait for next message (up to 600s) and handle it */gras_msg_handle(600.0);

gras_exit(); return O;

> int ping_callback(gras_socket_t experditor, void *payload_data) { int msg = *(int *)payload_data;

GRAS_BENCH_ALWAYS_BEGIN(); /* Some computation whose duration should be simulated */

GRAS_BENCH_ALWAYS_END() /* Send data back as payload of pong message to the ping's source */ gras_msg_send(source, gras_msgtype_by_name("pong"), &msg);

Work in Progress

- Parallel simulation for better scalability
- Native multi-threading support Port to Windows





between PowerPC, Sparc, and x86 architectures (WAN: California - France)

Grid Application Toolbox

- Platform monitoring (CPU and network)
- Network topology discovery

Why three interfaces?

- MSG: Rapid prototyping of distributed applications / algorithms
- Development of production distributed applications • GRAS:
- Study how an existing MPI application reacts to platform heterogenity • SMPI:

Simulation and Concurrency

- MSG: All simulated application processes run within a single process
- Subsets of simulated application processes run within multiple • GRAS: processes on multiple hosts, for increased scalability
- Each application process runs as a separate process • SMPI:

Work in Progress

Simulation of an existing MPI application

- Automatic (but directed) benchmarking of communication and computation costs during an application execution on an homogeneous platform
- Easy simulation of the application on a heterogeneous platform
- No code modification required beyond inserting benchmarking commands so that the simulation can be instanciated

Example: 1-D Matrix Multiplication in MPI

• Matrices are distributed among processors using a vertical strip decomposition • Column blocs are broadcasted at every step





• Computation

S

R

Point-to-point Communication

Features and Capabilities

- Simulation of complex communications (multi-hop routing)
- Simulation of resource sharing
- Simulation LAN and WAN links
- Topology can be imported from topology generators (such as BRITE)
- Trace-based simulation of performance variations due to external load
- Trace-based simulation of dynamic resource failures



Simulation of resource sharing

- Consider a set of resources, R.
- Consider a set of "tasks", T
- Each task is defined as the subset of R it uses
- SURF uses the unifying **MaxMin Fairness** model



- Used for computation and communication resources
 - Multiple TCP flows sharing links
 - Multiple CPU-bound processes sharing a CPU
 - Interference of communication and computation • Parallel tasks



• Can be implemented efficiently

GTNetS

SSFNet SURF

NS

• Is very accurate in many scenarios

Experiment comparing SURF, GTNetS, SSFNet, and NS-2 for a range of message sizes • 10 Mbs link, 10ms latency • Message size: 1KB to 1GB • Effective data transfer rate is measured

SURF is very accurate for data > 1MB But it is less accurate for <1MB because of transient TCP behavior (e.g., slow start)

Experiment comparing SURF, GTNetS, SSFNet, and NS-2 for a range of bandwidths • 10ms latency, 10MB message size • Physical bandwidth: 100KBps to 1GBps • Effective data transfer rate is measured

SURF is very accurate for the whole range of physical bandwidths