

Ultra Scalable Simulation with SimGrid USS SimGrid (ANR 08 SEGI 022)

<http://uss-simgrid.gforge.inria.fr>

Coordinated by Martin Quinson (Nancy University)

Lyon, January 4 2012



Our Scientific Objects: Distributed Systems

Grid Computing: Distributed infrastructure for Computational Science

- ▶ Massive systems federating numerous organization worldwide
- ▶ **Main issues:** Large production infrastructures, challenging to experiment with

P2P Systems

- ▶ Exploit resources at network edges (storage, CPU, human presence)
- ▶ **Main issues:** Intermittent connectivity (churn); Network locality; Anonymity

Cloud Computing

- ▶ Large infrastructures underlying commercial Internet (eBay, Amazon, Google)
- ▶ **Main issues:** Optimize costs; Keep up with the load (flash crowds)

Systems already in use, but characteristics hard to assess

- ▶ **Performance:** makespan, economics, energy, ... ← context of this project
- ▶ **Correction:** absence of crash, race conditions, deadlocks and other defects

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

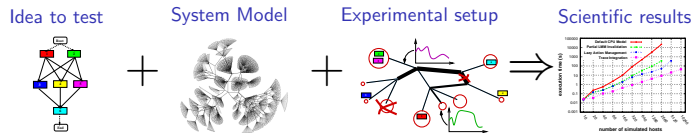
[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 1/14

Assessing Distributed Applications' Performance

Most Performance Studies are conducted through Experimentation

- ▶ **Experimental Facilities:** real applications on real platform (*in vivo*)
- ▶ **Emulation:** real applications on models of platforms (*in vitro*)
- ▶ **Simulation:** models (prototypes) of applications on system's models (*in silico*)

Simulating Distributed Systems ← context of this project



Simulation's Advantages

- ▶ Less simplistic than proposed **theoretical models** (which are useful too)
- ▶ Better XP control (~ reproducible) than **production systems** (+ not disruptive)
- ▶ Not as tedious, time/labor consuming than **experimental platforms**
- ▶ **Plus:** Lower technical burden; Quick and easy experiments; What if analysis

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 2/14

USS-SimGrid

Purpose of the SimGrid Project

- ▶ Allow a scientific approach of Large-Scale Distributed Systems simulation
- ▶ Propose ready to use tools enforcing methodological best practices

Main challenges

- ▶ **Validity:** Get realistic results (controlled experimental bias)
- ▶ **Scalability:** Simulate *fast enough* problems *big enough*
- ▶ **Usability:** Associated Tools; Ease of use; Applicability to context of interest

The USS-SimGrid project

- ▶ **Main Goal:** Make SimGrid usable in studies mandating extreme scaling
- ~ Perimeter increase from Grid Computing to Peer-to-peer
- ~ Improving simulation scalability: mandatory but not enough
- ~ Campaign data management pre- & post-processing not trivial anymore

Coming next: Some scientific achievements on these main challenges

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 3/14

Validity Challenge

SotA: Models in most simulators are either simplistic, wrong or not assessed

- ▶ **PeerSim:** discrete time, application as automaton; **GridSim:** naive packet level
- ▶ **OptorSim, GroudSim:** documented as wrong on heterogeneous platforms
- ▶ **Validity evaluation:** tricky, requires meticulous attention & sound methodology

Quality Levels of Validity

- ▶ Level -1: not validated (probably plainly wrong)
- ▶ Level 0 (visually ok): a few curves that look similar (generally hides a lot)
- ▶ Level 1 (ratios ok): $A < B$ in Simulation $\Leftrightarrow A < B$ in Reality
- ▶ Level 2 (prediction abilities): bounded distance between simulation and reality

SIMGRID validity before USS: Research focus in SimGrid since 2002

- ▶ **Several models:** NS3; **Fast precise model** (preferred); Fast simplistic model
- ▶ Hunting (and fixing) worst case scenarios since 10 years
- ▶ Error in percents if: TCP steady state (flows > 10Mb), latency-bound (WAN)

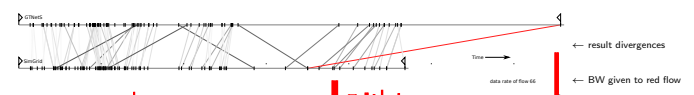
USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 4/14

Validity: SimGrid compared to Packet-Level Tools

Settings: *Synthetic App.* + *Synthetic WAN*. Compare against *GTNetS*

- ▶ Some **errors** were **hunted down** + unexpected **phenomenon** were **understood**
- ▶ Sharing mechanism from theoretical literature experimentally proved wrong
- ~ The model and its instantiation were considerably **improved**
- Widen validity range to flows > 100Kb and WAN with small latencies
- ▶ SimGrid and packet-level simulators now mostly diverge in **extreme** WAN cases



In this scenario, GTNetS and SG agree on termination date of most flows. The most diverging gets no bandwidth for a while although all others are done.

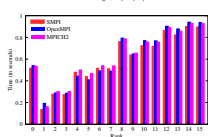
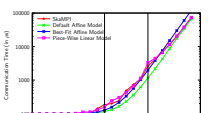
Going Further: developed **SMPI** ~ **Real App.** (NAS PB) + clusters (**LAN**)

- ▶ Good prediction for short messages is crucial; Numerical instabilities deadly
- ▶ Accurately modeling MPI semantic (asynchronous & collectives ops) is tricky
- ▶ Need to account for MPI overhead; what is Real with several MPI implems?

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 5/14

Accuracy of MPI simulations



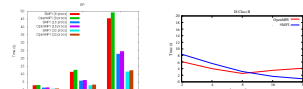
Timings of each communication

- ▶ $\lambda + \text{size} \times \tau$ not sufficient (TCP congestion)
- ▶ No affine function can match for all message sizes
- ▶ A 3-parts piecewise affine gives satisfying results

Taking resource sharing into account

- ▶ Rather good (visual) accuracy
- ▶ Our "error" \approx difference between runtimes
- ▶ This is only one collective

Still a work in progress for complete MPI applications



- ▶ Performance prediction not correct
- ▶ Trashing particularly challenging

- ▶ Although not perfect, accuracy comparable / better to other MPI simulators
- ▶ Ways better than the most precise existing P2P simulators

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 6/14

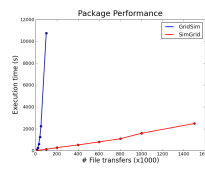
Scalability Challenge

Scalability constitutes the main objective of the USS SimGrid

- ▶ **Two aspects:** Big enough (large platforms) \oplus Fast enough (large workload)

Situation before the project

- ▶ Timings from CERN guys
- ▶ Maximal amount of user processes
 - ▶ **GridSim:** 10,922 (hard limit)
 - ▶ **SimGrid:** 200k (memory limit, 4Gb)
- ▶ But needs of the users:
 - ▶ **CERN:** 300 \times bigger than that (10 days/run)
 - ▶ **BOINC:** 600k volatile hosts over a year
- ▶ **PeerSim** simulates millions of processes
 - ▶ but with simplistic models only



USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022)

[Introduction](#) [Validity](#) [Scalability](#) [Usability](#) [Conclusion](#) 7/14

Approaches to Scalability in USS-SimGrid

Algorithmic optimization

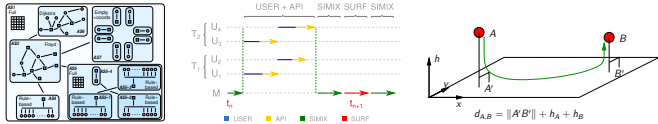
- **Compact Routing Representation:** From $O(n^2)$ to $O(n)$ memory consumption
- **Lazy Evaluation:** Arbitrary speedups on loosely coupled scenarios

Leverage several computing units

- **Parallel simulation:** P2P's grain so fine that classical // schema not applicable
- **Distributed simulation:** Still TBD, but not needed due to other optimization

Simpler models (but with potential loss of realism)

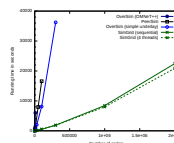
- **Coordinate-based:** extremely efficient, but only encodes latency
- **Last-mile models** (Manhattan distances): very efficient; controlled realism loss



USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 8/14

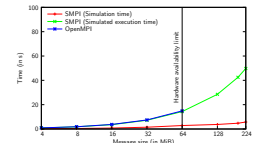
SimGrid Scalability Results

Millions of small processes



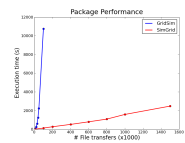
Chord simulation

Dozen of huge processes



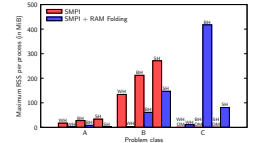
Binomial scatter with 16 processes

Large Workload



Simulation from CERN users

Hundreds of large processes



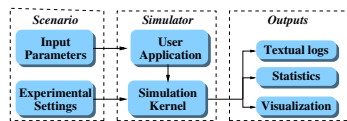
DT with up to 448 processes

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 9/14

Usability Challenge

Workflow to any Experiments through Simulation

1. Prepare the experimental scenarios
 2. Launch thousands of simulations
 3. Post-processing and result analysis
- ~ Each simulation is only a brick

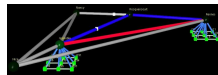
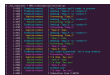


Situation before the project

- **Others simulators** come with *ad hoc* tools (but many *demowares*)
- **SimGrid:** nothing public/generic, but each user grow home-made scripts

Building a *demoware* is easy. Helping understanding is harder

- Often specific to a given simulator; often scalability issues
- Show only what the authors needed (platform/app. state, tracing/profiling)



USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 10/14

Approaches to Usability in USS-SimGrid

USS-SimGrid Proposal

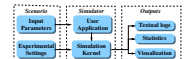
1. Workload generation:

- **Platforms:** Simulacrum (generation), PDA (archive) and MintCAR (mapping)
- **Applicative Workload:** Tau-based trace collection + replay
- **Background Workload:** Pilgrim (trace aggregation tool)

2. Campaign management: Workflow engine

3. Single simulation analysis: Visualization

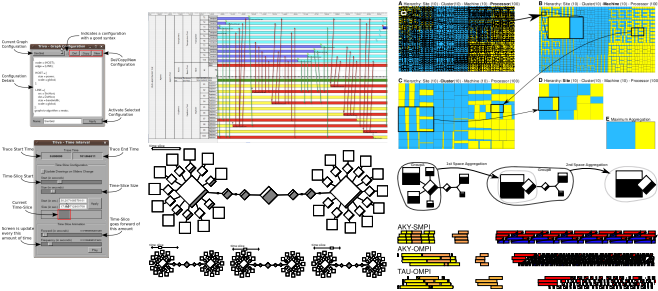
- Builds upon separate established projects: Triva and Paje
- **Generic and dedicated to visualization:** SimGrid only produces adapted traces (but SimGrid heavily modified to that extend by Triva author)



USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 11/14

Visualizing SimGrid Simulations

- **Visualization scriptable:** easy but powerful configuration; *Scalable tools*
- **Right Information:** both platform and applicative visualizations
- **Right Representation:** gantt charts, spatial representations, tree-graphs
- **Easy navigation in space and time:** selection, *aggregation*, animation
- **Easy trace comparison:** Trace *diffing* (still partial ATM)



USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 12/14

Conclusion: Project Outcomes

Scientific Production

► Publications

- 22 international publications (including 5 multi-site publications)
- 4 submitted articles (including 2 multi-site publications)

► Software: 11 releases of SimGrid (including 4 major releases)

- **Visualization:** 4 releases of Triva
- **Automatic Platform Mapping:** release of MintCar and UMCTool
- **Synthetic Platform Generation:** release of Simulacrum

Dissemination

- 2 Tutorials: HPCS'10, CLCAR'10; 2 Invited talks: P2P'09, RGE
- **SuperComputing presence** every year of project (@INRIA booth)
- 3-day Workshop: The SimGrid User Days (SUD'10)

USS-SimGrid as a Flagship (collaborative projects associated to this)

- Collaboration with ANR CIP (that use SimGrid to assess P2P HPC middleware)
- PHC Tournesol with the University of Antwerp (on scalable simulation)
- PICS CNRS Hawai'i/Villeurbanne (on MPI simulation)
- INRIA ADT (engineering forces devoted to SimGrid)

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 13/14

Conclusion and Open questions

Answers to good questions lead to new questions

- The work planned in this project were be done on time
- But these developments gave us new ideas about going even further
- These new ideas are paving our future work

SONGS (Simulation Of Next Generation Systems) ANR project

► Making SimGrid usable in 2 more domains

- Task 1: **[Data]Grid** (distributed Data mgnt for LHC; Hierarchical Storage System)
- Task 2: **Peer-to-Peer and Volunteer Computing** (Replica Placement in VOD;...)
- Task 3: **IaaS Clouds** (study from client or provider POV; energy metrics, EC2 APIs)
- Task 4: **High Performance Computing** (exascale; memory & energy models)

► Further improve our Simulation Fundamentals

- Task 5: **Simulation Kernel** (Efficient Simulation Kernel; DEVS Standard)
- Task 6: **Concepts and Models** (energy, storage, memory, networks, volatility)
- Task 7: **Analysis and Visualization** (Scalable Visualization, Causes Inference)
- Task 8: **Support to Experimental Methodology** (Open Science, DoE)

- Project funded as platform project on INFRA call for 4 years (2012-2016)

- The USS adventure revealed to be the first step of the campaign...

USS-SimGrid (Ultra Scalable Simulation with SimGrid – ANR 08 SEGI 022) Introduction Validity Scalability Usability Conclusion 14/14

Outline

• Context and Motivation

• Validity of Simulation Results

Context, Challenge and State of the Art
Validity: SimGrid compared to Packet-Level Simulators
Accuracy of MPI simulations

• Scalable Simulations

Context, Challenges, and State of the Art
Approaches to Scalability in USS-SimGrid
Results on Scalability in SimGrid

• Usability Challenge

Context, Challenges, and State of the Art
Approaches to Usability in USS-SimGrid
Visualizing SimGrid Simulations

• Conclusion

Project Outcomes
Conclusion and Open questions