## Research on Performance Modeling and Evaluation at TU Delft (2004—)



#### **Alexandru Iosup**

Parallel and Distributed Systems Group Delft University of Technology The Netherlands



**Our team: Undergrad** Gargi Prasad, Arnoud Bakker, Nassos Antoniou, Thomas de Ruiter, ... **Grad** Siqi Shen, Nezih Yigitbasi, Ozan Sonmez **Staff** Henk Sips, Dick Epema, Alexandru Iosup **Collaborators** Ion Stoica and the Mesos team (UC Berkeley), Thomas Fahringer, Radu Prodan (U. Innsbruck), Nicolae Tapus, Mihaela Balint, Vlad Posea (UPB), Derrick Kondo, Emmanuel Jeannot (INRIA), Assaf Schuster, Mark Silberstein, Orna Ben-Yehuda (Technion), ...



December 14, 2011

SPEC RG Cloud Meeting

#### Main Goal:

# Understand and Control the Performance of Large-Scale Distributed Systems

- Systems
  - Cluster Computing
  - Grid Computing
  - Cloud Computing
  - Peer-to-Peer (P2P) Systems
- Applications
  - E-Science
  - File-Sharing
  - Multi-media, esp. online gaming (MMOG, MSG)



### Approach: Real Traces, Models, Real Tools, Real-World Experimentation (+ Simulation)

- Formalize real-world scenarios
- Exchange real traces
- Model relevant operational elements
- Scalable tools for meaningful and repeatable experiments
- Comparative studies, almost like benchmarking





1. Main Goal and Approach

#### 2. Performance Modeling and Characterization

- 1. Real Traces
- 2. Models
- 3. Performance Evaluation
- 4. Conclusion



#### 2.1. Performance Characterization and Modeling, Data The Grid Workloads Archive

- Motivation: little is known about real grid use
  - No grid workloads (except "my grid")
  - No standard way to share them
- The Grid Workloads Archive: easy to share grid workload traces and research associated with them
  - Understand how real grids are used
    - Address the challenges facing grid resource management (both research and practice) Develop and test grid resource management solutions
  - **Perform** realistic simulations

http://gwa.ewi.tudelft.nl

A. Iosup, H. Li, M. Jan, S. Anoep, C. Dumitrescu, L. Wolters, D. Epema, The Grid Workloads Archive, FGCS 24, 672-686, 2008.

**Delft University of Technology** 

## 2.1. The Grid Workloads Archive Approach

- Standard data format (GWF)
  - Share traces with the community
  - Use extensions for specific modeling aspects
  - Text-based, easy to parse for custom tasks
  - Additional SQL-compatible data format (GWF-SQLite)
- Automated trace analysis
  - Provide ready-to-use tools to the community
  - Promote results availability and comparability
- Automated trace ranking
  - Help non-experts with their trace selection process

Delft University of Technology

A. Iosup, H. Li, M. Jan, S. Anoep, C. Dumitrescu, L. Wolters,D. Epema, The Grid Workloads Archive, FGCS 24, 672-686, 2008.

## 2.1. The Grid Workloads Archive Content

			Number of observed					
ID	System	Period	Sites	CPUs	Jobs	Groups	Users	
GWA-T-1	DAS-2	02/05-03/06	5	400	602K	12	332	
GWA-T-2	Grid'5000	05/04-11/06	15	$\sim 2500$	951K	10	473	
GWA-T-3	NorduGrid	05/04-02/06	$\sim 75$	$\sim 2000$	781K	106	387	
GWA-T-4	AuverGrid	01/06-01/07	5	475	404K	9	405	THE GRID WORKLOADS ARCHIVE
$GWA-T-5^{\diamond}$	NGS	02/03-02/07	4	${\sim}400$	632K	1	379	
$GWA-T-6^{\diamond}$							206	
$GWA-T-7^{\ddagger}$	http	)://gwa	a.ev	vi.tuo	lelft	t.nl	18	
$GWA-T-8^{\ddagger}$							19	<b>o</b> traces
GWA-T-9 <sup>‡</sup>	TeraGrid	08/05-03/06	1*	96	1.1M	26	121	online
	Total	13.51 yrs	136	>10000	>7M	191	2340	
	Average	1.5 yrs	15	1151	>750	<b>K</b> 21	>250	
A. Los D. Epe	A. Iosup, H. Li, M. Jan, S. Anoep, C. Dumitrescu, L. Wolters, D. Epema, The Grid Workloads Archive, FGCS 24, 672-686, 2008.							



December 14, 2011

## 2.1. The Grid Workloads Archive Approach: Automated Trace Analysis

- General information
- System-wide characteristics
  - Utilization



## Auto-reporting useful for benchmarking

oser and group characteristics

- Analysis for Top10 users
- Analysis for Top10 groups
- Performance
  - # running/waiting jobs
  - Throughput, # completed jobs





#### 2.1. Performance Characterization and Modeling, Data The Failure Trace Archive

System	Туре	# of Nodes	Target Component	Period	Year
SETI@home	Desktop Grid	226,208	CPU	1.5 years	2007-2009
Overnet	P2P	3,000	host	2 weeks	2003
Microsoft	Desktop	51,663	host	35 days	1999
LANL	SMP, HPC Clusters	4750	host	9 years	1996-2005
HPC2	HPC Clusters	256	IO	2.5 years	1996-2005
NERSC Skype	htt	o://f	ta.inria.	fr	008
Web sites	Web servers	129	host	8 months	2001-2002
DNS	DNS servers	62,201	host	2 weeks	2004
PlanetLab	P2P	200-400	host	1.5 year	2004-2005
Grenouilleog	DSL	4800	host	1 year	2003
Grenouilleo5	DSL	4800	host	1 year	2005
EGEE	Grid	2500 queues	CE queue	1 month	2007
Grid'5000	Grid	1288	host	1.5 years	2005-2006





D. Kondo, B. Javadi, A. Iosup, D. Epema, The Failure Trace Archive: Enabling Comparative Analysis of Failures in Diverse Distributed Systems, CCGrid 2010 (Best Paper Award)



#### 2.1. Performance Characterization and Modeling, Data The P2P Trace Archive

Trace ID	Community	Measurement Period	Sampling Rate	No. Files	No. Sessions	Traffic	Contributor	
<u>T1'03</u> [153MB]	SuprNova, (general)	06 Dec 2003 ~ 17 Jan 2004	2.5 min	12	28,423,470	n/a	PDS, TU Delft	
<u>T2'05</u> [212MB]	ThePirateBay, (general)	06 May 2005 ~ 11 May 2005	2.5 min	4800	35,881,338	12 PB/year	PDS, TU Delft	
<u>T3'05</u> [9.5GB]	Filelist.org, (general)	14 Dec 2005 ~ 04 Apr 2006	6 min	3000	2,172,738	n/a		
<u>T4'05</u> [6.5MB]	LegalTorrents.com, (general)	22 Mar 2005 ~ 19 Jul 2005	5 min	41	n/a	698 GB/		
<u>T4'09</u> [44.3]493	LegalTorrents.com, (general)	24 Sep 2009 ~ Feb 2009	5 min	183	n/a	1.1 ТВ/		
<u>T5</u> [8.:						9 GB/y		
[19.	attn://n2n	ta awi t		lft s		n/a		
<u>T5</u> [27	inth'\hat		uue	11 601		143 GB/		
<u>T6</u> [72M						735 GB/year	PDS, TU Delft	
<u>T6'09</u> [28MB]	tlm-project.org, (Linux OS)	24 Sep 2009 ~ Feb 2009	10 min	74	21,529	15 GB/		
T7'05 [69MB]	transamrit.net, (Slackware OS)	22 Mar 2005 ~ 19 Jul 2005	5 min	14	130,253	258 GB/	16 +r	2000
T7'09 [78MB]	transamrit.net, (Slackware OS)	24 Sep 2009 ~ Feb 2009	5 min	60	61,011	840 GB/		aces
T8'05 [183MB]	unix-ag.uni-kl.de, (Knoppix OS)	22 Mar 2005 ~ 19 Jul 2005	5 min	11	279,323	<mark>493 GB</mark> /	on	line
T8'00			- ·			249.00	<b>~</b> 111	
[163MB]	unix-ag.uni-kl.de, (Knoppix OS)	24 Sep 2009 ~ Feb 2009	5 min	12	160,522	348 GB/		

B. Zhang, A. Iosup, J. Pouwelse, and D. Epema (2010). The peer-to-peer trace archive: design and comparative trace analysis. CoNEXT Workshops.



#### 2.1. Performance Characterization and Modeling, Data The Cloud Workloads Archive

- Looking for invariants
  - Wr [%] ~40% Total IO, but absolute values vary

Trace ID	Total IO [MB]	Rd. [MB]	Wr [%]	HDFS Wr[MB]
CWA-01	10,934	6,805	38%	1,538
CWA-02	75,546	47,539	37%	8,563

- # Tasks/Job, ratio M:(M+R) Tasks, vary
- Understanding workload evolution





1. Main Goal and Approach

#### 2. Performance Modeling and Characterization

- 1. Real Traces
- 2. Models
- 3. Performance Evaluation
- 4. Conclusion



## 2.2. Characterization: Grid Workloads Single-Node Jobs

 Average job size is 1 (that is, there are no [!] tightlycoupled, only conveniently parallel jobs)



# 2.2. Characterization: Grid Workloads VO, Group, and User Characteristics

- Top 2-5 groups/users dominate the workload
- Top groups/users are constant submitters
- The week's top group/user is not always the same



#### 2.2. Characterization: Grid Workloads Analysis Summary: Grids vs. Parallel Production Systems



A. Iosup, D.H.J. Epema, C. Franke, A. Papaspyrou, L. Schley,
B. Song, R. Yahyapour, On Grid Performance Evaluation using Synthetic Workloads, JSSPP'06.

#### 2.2. Characterization and Modeling: Grid Workloads More Analysis: Special Workload Components **Bags-of-Tasks** (**BoTs**) Workflows (WFs)





WF = set of jobs with precedence(think Direct Acyclic Graph)

## 2.2. Characterization: Grid Workloads BoTs are predominant in grids

- Selected Findings
  - Batches predominant in grid workloads; up to 98% CPUTime
  - Average batch size (∆≤120s) is 15-30 (500 max)
  - 75% of the batches are sized
     20 jobs or less

Trace	Observed	Percer	ntage From Total
D	BoTs	Jobs	CPUTime
GWA-T-1	57k	92%	78%
GWA-T-2	26k	85%	30%
GWA-T-3	50k	94%	90%
GWA-T-6	43k	95%	95%
GWA-T-7	13k	95%	96%
GWA-T-8	302k	94%	98%
GWA-T-10	16k	93%	92%
GWA-T-11	5k	96%	97%
GWA-T-12	135K	94%	96%
GWA-T-13	68K	96%	86%

A. Iosup, M. Jan, O. Sonmez, and D.H.J. Epema, The Characteristics and Performance of Groups of Jobs in Grids, Euro-Par, LNCS, vol.4641, pp. 382-393, 2007.

A. Iosup and D.H.J. Epema, Grid Computing Workloads, IEEE Internet Computing 15(2): 19-26 (2011)

### 2.2. Grid Workloads Workflows exist, but they seem small

Traces



75% + WFs are sized 40 jobs or less, 95% are sized 200 jobs or less



#### 3.1. Grid Workloads Modeling Grid Workloads: Feitelson adapted



- Adapted to grids: percentage parallel jobs, other values.
- Validated with 4 grid and 7 parallel production env. traces

A. Iosup, D.H.J. Epema, T. Tannenbaum, M. Farrellee, and M. Livny. Inter-Operating Grids Through Delegated MatchMaking, ACM/IEEE Conference on High Performance Networking and Computing (SC), pp. 13-21, 2007.

#### 2.2. Grid Workloads Modeling Grid Workloads: adding users, BoTs



- Single arrival process for both BoTs and parallel jobs
- Reduce over-fitting and complexity of "Feitelson adapted" by removing the RunTime-Parallelism correlated model
- Validated with 7 grid workloads

A. Iosup, O. Sonmez, S. Anoep, and D.H.J. Epema. The Performance of Bags-of-Tasks in Large-Scale Distributed Systems, HPDC, pp. 97-108, 2008.

#### 2.2. Grid Infrastructure Resource dynamics in cluster-based grids

- Environment: Grid'5000 traces
  - jobs 05/2004-11/2006 (30 mo., 950K jobs)
  - resource availability traces 05/2005-11/2006 (18 mo., 600K events)
- Resource availability model for multi-cluster grids





A. Iosup, M. Jan, O. Sonmez, and D.H.J. Epema, On the Dynamic Resource Availability in Grids, Grid 2007, Sep 2007.



December 14, 2011

#### 2.2. Grid Infrastructure Correlated Failures

• Correlated failure



Maximal set of failures (ordered according to increasing event time), of time parameter  $\Delta$  in which for any two successive failures E and F,  $TS(F) \leq TS(E) + \Delta$  where  $TS(\cdot)$  returns the timestamp of the event;  $\Delta = 1$ -3600s.



## 2.2. Grid Infrastructure Dynamics Model



- Assume no correlation of failure occurrence between clusters
- Which site/cluster?
  - f<sub>s</sub>, fraction of failures at cluster *s*
- Weibull distribution for IAT
  - Shape parameter > 1: increasing hazard rate the longer a node is online, the higher the chances that it will fail

A. Iosup, M. Jan, O. Sonmez, and D.H.J. Epema, On the Dynamic Resource Availability in Grids, Grid 2007, Sep 2007.



December 14, 2011

## 2.2. Grid Infrastructure **Evolution Model**



**Delft University of Technology** 

## 2.2. Cloud Workload Model MapReduce Workload Model

	Job	Task	Value
Inter arrival time	Х	Х	Seconds
Executable ID	Х		Integer ID
Run time		Х	Seconds
Number of tasks	Х		Count
Map/Reduce ratio	Х		Fraction: Maps/Reduces
Forced Quit Time	Х		Seconds
CPU's		Х	Count
Disk I+O		Х	Bytes
I/O Ratio		Х	Fraction: Input/Output
Memory		Х	Bytes
Network		Х	Bytes
Exit State	Х	Х	Integer Coded State

- Statistical model
- Traces (10s of millions of tasks) from:
  - Leading Social Networking company
  - 2 x Leading Search company





- 1. Main Goal and Approach
- 2. Performance Modeling and Characterization
  - 1. Real Traces
  - 2. Models

#### 3. Performance Evaluation

- 1. Observation Tools
- 2. Measurement Tools
- 3. Experiments

#### 4. Conclusion



## 3.1. MultiProbe Observing P2P Systems at Large



- Environment:
  - 700 BitTorrent swarms actively observed
  - 300 computer nodes coordinated around the world

#### Largest measurement from 2005 to 2010

A. Iosup et al., Correlating Topology and Path Characteristics of Overlay Networks and the Internet, CCGrid Workshops 2006.

27

## 3.1. BTWorld Observing the Global Public BitTorrent

Metric	Value
Tracing Period	Jan 03-09, 2010
Number of trackers	912
Number of alive trackers	769
Number of swarms	10,329,950
Number of hashes	6,314,318
Number of hashes with known size	1,024,573 (16.2%)
Number of swarm samples	899,537,250

#### • Environment:

- Over 10M BitTorrent swarms actively observed
- Only 4 computer nodes coordinated around the world

#### • Largest measurement from 2011-

M. Wojciechowski, M. Capota, J. Pouwelse, and A. Iosup: BTWorld: towards observing the global BitTorrent filesharing network. HPDC Workshops 2010.



- 1. Main Goal and Approach
- 2. Performance Modeling and Characterization
  - 1. Real Traces
  - 2. Models

#### 3. Performance Evaluation

- 1. Observation Tools
- 2. Measurement Tools
- 3. Experiments
- 4. Conclusion



#### 3.2. GrenchMark: Testing in LSDCSs GrenchMark: a Framework for Analyzing, Testing, and Comparing Grids

 What's in a name? grid benchmark → working towards a generic tool for the whole community: help standardizing the testing procedures,

> GrenchMark evolved from a grid-specific testing framework to a framework for testing large-scale distributed computing systems

- Easy-to-use tools to create synthetic grid workloads
- Flexible, extensible framework

Alexandru Iosup, Dick H. J. Epema: GRENCHMARK: A Framework for Analyzing, Testing, and Comparing Grids. CCGRID 2006: 313-320

#### 3.2. GrenchMark: Testing in LSDCSs Architecture Overview





December 14, 2011

#### 3.2. GrenchMark: Testing in LSDCSs

... but More Complicated Than You Think

#### Workload structure

- User-defined and statistical models
- Dynamic jobs arrival
- Burstiness and self-similarity
- Feedback, background load
- Machine usage assumptions
- Users, VOs

#### Metrics

- A(W) Run/Wait/Resp. Time
- Efficiency, MakeSpan
- Failure rate [!]

#### Notions

• Co-allocation, interactive jobs, malleable, moldable, ...

#### Measurement methods

- Long workloads
- Saturated / non-saturated system
- Start-up, production, and cool-down scenarios
- Scaling workload to system
- Applications
  - Synthetic
  - Real
- Workload definition language
  - Base language layer
  - Extended language layer
- Other
  - Can use the same workload for both simulations and real environments



#### 3.2. GrenchMark: Testing in LSDCSs ServMark, a Distributed GrenchMark



GrenchMark



- Tackles two orthogonal issues:
  - Multi-sourced testing (multi-user scenarios, scalability)
  - Generate and run dynamic test workloads with complex structure (real-world scenarios, flexibility)

Adds

- **Coordination and automation layers** 
  - Fault tolerance module

December 14, 2011

## 3.2. GrenchMark: Testing in LSDCSs SkyMark, GrenchMark for LaaS Clouds



- Provisioning and allocation queues + policies
- Short-, Many-Task workloads

D. Villegas, A. Antoniou, S. Sadjadi, and A. Iosup: An Analysis of Provisioning and Allocation Policies for LaaS Clouds. Submitted to CCGRID 2012.

## 3.2. RTSenv: Testing Online Games RTSenv: A Tool for testing RTS games



- Abstractions for RTS games: units, map structure, etc.
- Metrics for RTS game performance and experience
- Replayability

Siqi Shen, Otto Visser, Alexandru Iosup: RTSenv: An experimental environment for real-time strategy games. NETGAMES 2011: 1-6



- 1. Main Goal and Approach
- 2. Performance Modeling and Characterization
  - 1. Real Traces
  - 2. Models

#### 3. Performance Evaluation

- 1. Observation Tools
- 2. Measurement Tools
- 3. Experiments
- 4. Conclusion



## 3.3. GrenchMark: Testing in LSDCSs Testing a Large-Scale Environment

• Performance metrics

system-, job-, operational-, application-, and service-level



#### 3.3. GrenchMark: Testing in LSDCSs Experiments: Testing Performance

- Testing application performance: test the performance of an application (for sequential, MPI, Ibis applications)
  - Report runtimes, waiting times, grid middleware overhead
  - Automatic results analysis

Table 2: A summary of time and run/success percentages for different job types.

	Job	Job	Turnaround [s]		Runt	time [s]		$\operatorname{Run}+$
	name	type	Avg	Range	Avg	Range	$\operatorname{Run}$	Success
	sser	seq	129	16 - 926	44	1-588	100%	97%
	smpi1	MPI	332	21 - 1078	110	1 - 332	80%	81%
_	NQueens	Ibis	99	15 - 1835	31	1 - 201	70%	85%

- What-if analysis: evaluate potential situations
  - System change
  - Grid inter-operability
  - Special situations: spikes in demand





December 14, 2011

## 3.3. GrenchMark: Testing in LSDCSs Testing a Large-Scale Environment

- Testing a 1500-processors Condor environment
  - Workloads of 1000 jobs, grouped by 2, 10, 20, 50, 100, 200
  - Test finishes 1h after the last submission
  - Results
    - >150,000 jobs submitted
    - >100,000 jobs successfully run, >2 yr CPU time in 1 week
    - 5% jobs failed (much less than other grids' average)
    - 25% jobs did not start in time and where cancelled





#### 3.3 Performance Evaluation in Real-World Environments Raw Perf.: Performance vs. Res. Consumption

Middleware	MS [s]
DAGMan	$1,327 \pm 138$
Karajan	$1,111 \pm 154$

#### Karajan performs better than DAGMan, but runs quickly out of resources.



**Delft University of Technology** 



<sup>1-</sup> losup et al., Performance Analysis of Cloud Computing Services for Many Tasks Scientific Computing, IEEE TPDS, 2011,

c1.xlarge

http://www.st.ewi.tudelft.nl/~iosup/cloud-perf10tpds\_in-print.pdf

2- losup et al., On the Performance Variability of Production Cloud Services, CCGrid 2011, pds.twi.tudelft.nl/reports/2010/PDS-2010-002.pdf

## 3.3. IaaS Cloud Results Provisioning and Allocation



- Tested in 3 real envs., including Amazon EC2
- Performance-cost trade-off
- New metrics: utility, cost efficiency.

D. Villegas, A. Antoniou, S. Sadjadi, and A. Iosup: An Analysis of Provisioning and Allocation Policies for LaaS Clouds. Submitted to CCGRLD 2012.

# 3.3. RTSenv: Testing Online Games Testing with Many Players



• Assess performance

GRENCHMARK

- Assess gameplay experience
- Main finding: performance and gameplay scale differently





## Agenda

- 1. Main Goal and Approach
- 2. Performance Modeling and Characterization
  - 1. Real Traces
  - 2. Models
- 3. Performance Evaluation
  - 1. Observation Tools
  - 2. Measurement Tools
  - 3. Experiments

#### 4. Conclusion



## A Take-Home Message

- Understanding how real large-scale distributed systems work
  - Real traces
  - Models for workload and infrastructure
- Building tools for performance observation and evaluation
  - Observation: MultiProve; BTWorld; ...
  - Evaluation: GrenchMark + ServMark + SkyMark; RTSenv; ...

#### Performance evaluation in real environments

- Grids
- IaaS clouds
- Online Gaming
- ... and many others
- Much to be done, esp. in clouds



http://www.flickr.com/photos/dimitrisotiropoulos/4204766418/

December 14, 2011



**Delft University of Technology** 

## Thank you for your attention! Questions? Suggestions? Observations?

#### More Info:

- http://www.st.ewi.tudelft.nl/~iosup/research.html
- <u>http://www.st.ewi.tudelft.nl/~iosup/research\_gaming.html</u>
- <u>http://www.st.ewi.tudelft.nl/~iosup/research\_cloud.html</u>

#### **Alexandru Iosup**

#### A.Iosup@tudelft.nl

http://www.pds.ewi.tudelft.nl/~iosup/ (or google "iosup") Parallel and Distributed Systems Group Delft University of Technology

Do not hesitate to

contact me



**Delft University of Technology**