







Modelling Distributed Computing Infrastructures for High Energy Physics Applications

26th International Conference on Computing in High Energy & Nuclear Physics

Maximilian Horzela, Henri Casanova, Manuel Giffels, Artur Gottmann, Robin Hofsaess, Günter Quast, Simone Rossi Tisbeni, Achim Streit, Frederic Suter | 08. May 2023



www.kit.edu

Planning a Computing Infrastructure for HEP



- We only have one grid to play with
 - Many connections, (managed) storages and compute resources
 - Dynamically changing workloads and data access patterns
 - Resource availability
 - Complex scheduling systems with many parameters
- Many competing ideas for infrastructure designs, enough resources only for one

Most efficient realization of grid infrastructure is not obvious

How to study such systems?





Simulating Infrastructure was Appealing 20 Years Ago

• Proposed & Used 20 years ago: MONARC [1, 2] \rightarrow original WLCG



Simulating Infrastructure is Appealing Today





[6]

- Proposed & Used 20 years ago: MONARC $[1, 2] \rightarrow$ original WLCG
- There have been many attempts on modelling complex computing systems → many failed, some succeeded
- SIMGRID [3]: Library of low-level simulation abstractions for distributed systems in general (Actors using Resources through Activities)
 - Fluid-model for usage of network, storage, and CPU resources
 - Activities defined by a total and remaining amount of work to accomplish
 - Manager of a queue of activities, which start new and stop other activities triggered by events
- WRENCH [4]: High-level building blocks
 - Services specifying Activities
- Addition of (HEP-)specific adaptions (e.g. job-, dataset- & workflow-model, data streaming, ...) and logic starting and using services



DCSim: Implementation of (HEP) Extensions and Simulator

Define workloads of jobs with:

- Number of operations to execute
- Memory
- Size of in- and output files
- Submission time
- Oefine platform (network & hosts) with
 - Properties: N_{core}, CPU-speed, RAM, disk, bandwidth
 - roles: worker, storage, data cache, scheduler, ...
- Instantiate initial deployment of files on storage systems
- Start the simulation!
 - Jobs are scheduled and run
 - Input-files are streamed and cached
 - Caches evict files if necessary
 - → Job dynamics are monitored



Karlsruhe Institute of Technology

Calibration and Validation

- \blacksquare Model not based on first principles \rightarrow free parameters in simulation need to be tuned
- Validate tuned simulator by comparing with measured data e.g. job-execution times over fraction of cached input files



After calibration simulation accurately reproduces general behaviour of real systems, differences in details

Mitigating Computational Complexity of Simulation



- Superlinear increase of runtime with number of simultaneous activities
 - \rightarrow Detailed simulation of big platforms is problematic
- ightarrow Scale size of platform to reduce number of simultaneously active entities (links, cores, hosts, \ldots)



- Reduce number of active slots by decreasing number of cores on worker node hosts
- Reduce RAM on hosts to conserve the number of job slots
- Adjust link and storage bandwidths to conserve bandwidth per core
- Decrease the number of jobs per workload accordingly



Mitigating Computational Complexity of Simulation

- Superlinear increase of runtime with number of simultaneous activities
 - \rightarrow Detailed simulation of big platforms is problematic
- ightarrow Scale size of platform to reduce number of simultaneously active entities (links, cores, hosts, \ldots)



Scaling decreases runtime significantly (up to several orders of magnitude), while conserving realism



HEP Job Mix

- Job characteristics gathered from job monitoring data over a representative time frame
- Group jobs into representative classes
- Main job types
 - Analysis: All types of jobs submitted by users, mainly MC production and data reduction
 - ReadoutSim: Simulation of detector readout electronics
 - Processing: Collective chunks of multiple data and MC processing steps from reconstruction to analysis objects data files
 - Merge: Reduction of number of files by merging multiple into a single



Karlsruhe Institute of Technology

HEP Job Mix

- Job characteristics gathered from job monitoring data over a representative time frame
- Group jobs into representative classes
- Main job types
 - Analysis: All types of jobs submitted by users, mainly MC production and data reduction
 - ReadoutSim: Simulation of detector readout electronics
 - Processing: Collective chunks of multiple data and MC processing steps from reconstruction to analysis objects data files
 - Merge: Reduction of number of files by merging multiple into a single

Job types differ in their characteristics



HEP Workload Sampling



Random sampling of job characteristics

- CPU-time of jobs
- Memory required or used
- Number of bytes read/written by the job → input/output size

Keeping proportions relative to the monitoring data

8000 8000 7000 6000 6000 ting 5000 4000 tuno 4000 3000 2000 2000 1000 50 60 70 80 20 40 80 100 0 10 20 30 40 0 60 lob CPU-time / h Data read / GB

Example for **Processing** jobs

Use-case Study:

Interplay Between Tier 1 and Tier 2 with Unmanaged Storage

- Study hypothetical interplay of a Tier 1 with one Tier 2 with grid storage replaced by cache (Tier 2')
- Model both sites as local networks of worker nodes and storages with network in between

Tier 1	Tier 2'	
$\mathcal{O}(40k)$ cores	$\mathcal{O}(20k)$ cores	
80 Gbps bandwidth storage	80 Gbps bandwidth cache	
2 \otimes 100 Gbps local network 40 Gbps local network		
100 Gbps network between sites		

- Scan through fraction of prefetched files at cache (prefetch-rate)
 - $p = 0 \Leftrightarrow$ single grid storage serving both sites
 - $p = 1 \Leftrightarrow$ separate grid storage for Tier 2







Interplay Tier 1 \leftrightarrow Tier 2': Job Runtime & Transfer Duration



- Increasing presence of input-files at cache reduces overall wall-times of jobs because of faster I/O
- Benefit flattens out at Tier 2'(jobs throttled by cache and compute speed), further reduces at Tier 1
- → Both sites benefit from reduced I/O on grid storage

Interplay Tier 1 \leftrightarrow Tier 2': CPU Efficiency







- Job efficiency drops significantly with lower fraction of input-files in cache due to higher load on Tier 1 grid storage
- ⇒ Cache as replacement only reasonable with performance improvements of Tier 1 grid storage

What if...? – Tier 2 with Cache Connected with Lower Bandwidth (e.g. Cloud/HPC)



Tier 1

Tier 2'



Job efficiency significantly boosted with cache running at \Rightarrow "HPC" without efficiency loss at Tier 1

0.9 10



Conclusions and Future Work

- Developed flexible, performant and accurate tool
 - Calibrated and validated using measured data
 - Enabling practicable research on interesting and realistic computing infrastructures without deploying them
- More studies on scalability, realism and more complex scenarios in the pipeline
 - Improve speed of the simulation
 - Influence of network load introduced by other sites
 - Datasets instead of job-specific input-files
 - Automated calibration and validation

On behalf of the collaboration, contact maximilian.horzela@kit.edu, if you are interested in using or in helping us develop the tool!



References I



- [1] Iosif C. Legrand and Harvey B. Newman. "The MONARC Toolset for Simulating Large Network-Distributed Processing Systems". <u>Proceedings of the 32nd Conference on Winter Simulation</u>. WSC '00. Orlando, Florida: Society for Computer Simulation International, 2000, pp. 1794–1801. ISBN: 0780365828.
- [2] Harvey B. Newman and Iosif C. Legrand. "Simulating distributed systems". <u>AIP Conference Proceedings</u> 583.1 (Aug. 2001), pp. 164–166. ISSN: 0094-243X. DOI: 10.1063/1.1405293. eprint: https://pubs.aip.org/aip/acp/article-pdf/583/1/164/11368715/164_1_online.pdf.
- [3] Henri Casanova et al. "Versatile, Scalable, and Accurate Simulation of Distributed Applications and Platforms". Journal of Parallel and Distributed Computing 74.10 (Oct. 2014), pp. 2899–2917. DOI: 10.1016/j.jpdc.2014.06.008.
- [4] Henri Casanova et al. "Developing Accurate and Scalable Simulators of Production Workflow Management Systems with WRENCH". <u>Future Generation Computer Systems</u> 112 (2020), pp. 162–175. DOI: 10.1016/j.future.2020.05.030.
- [5] SimGrid Team. "simgrid.org". 2022. URL: https://simgrid.org/.



References II

- [6] WRENCH Team. "wrench-project.org". 2022. URL: https://wrench-project.org/.
- [7] CMS Collaboration. "The Higgs Boson turns 10: Results from the CMS experiment". 2022. URL: https://cms.cern/news/higgs-boson-turns-10-results-cms-experiment.
- [8] Pedro Velho et al. "On the Validity of Flow-Level Tcp Network Models for Grid and Cloud Simulations". <u>ACM Trans. Model. Comput. Simul.</u> 23.4 (Dec. 2013). ISSN: 1049-3301. DOI: 10.1145/2517448. URL: https://doi.org/10.1145/2517448.

Backup

Static Diskless Tier 2 with Cache



- Study interplay of a Tier 1 with one (or more) Tier 2 with grid storage replaced by cache
- Investigate "static" collection of production and analysis jobs

→ Scan prefetch-fraction of relevant input-files at cache (0: no cache, 1: grid storage)



Static Diskless Tier 2 with Cache

- Study interplay of a Tier 1 with one (or more) Tier 2 with grid storage replaced by cache
- Investigate "static" collection of production and analysis jobs
- Approximate dispensable sites' effects with randomized links/bws
- → Scan prefetch-fraction of relevant input-files at cache (0: no cache, 1: grid storage)





HEP Job Characteristics



- Job characteristics gathered from job monitoring data of time frame 24.02.-07.03.23
- CMS jobs executed at sites T1_DE_KIT and T2_DE_DESY analysed
- Job types
 - Analysis: All types of jobs submitted by users, mainly MC production and data reduction
 - ReadoutSim: Simulation of detector readout electronics
 - Processing: Collective chunks of multiple data and MC processing steps from reconstruction to analysis objects data files
 - Merge: Reduction of number of files by merging multiple into a single
 - Others including
 - Generation of MC events and simulation of Detector effects
 - Reconstruction of digital detector signals
 - Analysis object data derivation
 - Collection of job logs
 - Cleanup of logs and data
 - Test jobs



CMS Job Characteristics – T1_DE_KIT and T2_DE_DESY





CMS Job Characteristics – T1_DE_KIT





CMS Job Characteristics – T2_DE_DESY



Underlying Models

SIMGRID and WRENCH

SIMGRID [3]

- Library of exposed functions (low-level simulation abstractions) written in C++
- Framework for building own simulator of distributed computer systems in C/C++, Python or Java
- Accurate (validated), scalable (low ratio of simulated versus real time) and expressive (able to simulate arbitrary platforms, applications and execution scenarios)
- Large user-base

WRENCH [4]

- High-level simulation abstractions based on SIMGRID
- C++ API to define services and execution controllers for the implementation of a simulator









SIMGRID Engine

- "Actors" using "Platform" of resources through corresponding "Activities"
- "Activities" have both qualitative (synchronization between actors / locking) and quantitative (consumption
 of resource capacities) components
- "Actor" management by a central actor, called "Maestro", in scheduling round:
 - Assign control flow to each actor not blocked on a simcall, start sub-scheduling:
 - Actors execute activities and return intermediate simcalls that take no time to execute (e.g. spawn new actor)
 - Repeat until all actors return a blocking simcall or terminate
 - Advance time to that point at which first next activity terminates
 - Start new scheduling round

Context switching between actor and maestro is highly optimized

SIMGRID Resource-Activity Model



- Same analytical flow-model for simulation of network, storage and CPU
- Activities defined by a total and remaining amount of work to accomplish



- Resource capacity of resource C_r is assigned to a set of concurrent activities A at time t₀ is determined by solving max [min_{a∈A} (x_a)] under constraints ∑_{a∈A} x_a ≤ C_r
- Simulation time is advanced to time at which first activity is completed t₁



SIMGRID CPU and Storage Model

- Max-min of n resources sharing a single resource r leads to fair share Cr/n assigned to each activity
- For CPU-shares optional scaling by normalized priorities possible
- Fair sharing of storage I/O bandwidth plus additional fixed initial delay at simulation time advance due to seek time

TCP Network Model

- TCP doesn't show Max-Min fairness, two options:
- Packet level network simulator ns-3
- Improved flow-level network model with
 - modified constraints accounting for RTT unfairness of TCP and throughput degradation due to reverse traffic [8]
 - and improved execution time

$$T = \alpha l_f + \frac{V}{\beta x_f}$$

with TCP version specific parameters α and β tuned by packet-level simulation

valid for transfers of data of size \geq 100KiB

WRENCH

- Layer of high level building blocks on top of SimGrid
- Collection of Services using the simulated platform to define activities and events, i.e.:
 - Computing services implement usage of CPU resources
 - Storage services provide management of storage resources and read/write of files
 - Execution controllers schedule and start actions on other services
- Provides an API to define services and execution controllers for the implementation of a simulator







WRENCH Services

- Adds high level abstractions ("services") on top of SIMGRID
 - Compute services knows how and where to compute tasks, e.g. bare-metal, cloud, virtualized cluster, batch-scheduled cluster platforms and HTCondor
 - Storage services know how to store and give access to files
 - File-registry services know where files reside
 - Network proximity services monitor network and maintain database of host-to-host distances
 - Energy-meter services periodically measure energy-consumption of all resources
- All services introduce their own messages (activities) and according payloads



HTCondor Compute-Service

- Starting HTCondor-compute-service spawns a central-manager and a negotiator service
 - Compute service: entry point for job submission to WMS
 - Central manager: management of available resources (pool of bare-metal or batch-compute-services), manages job submission, initializes negotiation cycles
 - Negotiator: matches jobs to resources (based on #CPU and memory requirements)
- Actual task execution and resource allocation on host managed by matched compute-service
- Main focus so far on simulation of grid-universe jobs, need to be reviewed in the future to support more complex scheduling scenarios



Batch and Streaming Jobs

Batch jobs



Sequentially read a batch of data, compute and write

Streaming jobs



- Can concurrently read and compute blocks of data (enabled by XRootD in HEP)
- \rightarrow More compact pattern and reduced duration



Caching Job

- XRootD inspired caching of input-files
- Information about reachable caches definite only once job started
- \rightarrow Active job
- Once job is scheduled
 - Determine reachable (local, same local network, same network tier) caches
 - Is For each file, determine source storage
 - If cache, read from cache
 - If remote, read from remote, cache (instantly create) & evict (LRU), when applicable
 - Perform I/O and computation (copy & compute, stream blocks)

Validation & Calibration

Calibration/Validation Strategy

Karlsruhe Institute of Technology

Real-world system:

- Assemble real test systems with a control on parameters
- Start collections of jobs with known/steerable characteristics
- Measure job dynamics

Simulation:

- Define platforms as close to the test-system as reasonable
- Start workloads of jobs with similar job characteristics
- Read-out job dynamics

Define a set of observables, which are sensitive to simulation parameters

Tune the simulation parameters until simulation fits measurements Combine differently configured measurements to learn something about the system

 \Rightarrow Use chain of slightly varied scenarios to calibrate in steps





WN 1	WN 2	WN 3
24 job slots	12 job slots	12 job slots
HDD data cache each		
10 Gbps link to switch each		





WN 1	WN 2	WN 3
24 job slots	12 job slots	12 job slots
HDD data cache each		
10 Gbps link to switch each		





WN 1	WN 2	WN 3
24 job slots	12 job slots	12 job slots
HDD data cache each		
10 Gbps link to switch each		





WN 1	WN 2	WN 3
24 job slots	12 job slots	12 job slots
HDD data cache each		
10 Gbps link to switch each		

Karlsruhe Institute of Technology

Storage Effectively control remote-bandwidth Switch Service CPU slots 3 Input-file Worker Node 1 Worker Node 3 Worker Node 2 Cache

Benchmark Setup

WN 1	WN 2	WN 3
24 job slots	12 job slots	12 job slots
HDD data cache each		
10 Gbps link to switch each		



Benchmark Jobs

- Exact replications of the same data analysis job
 - Same executable, executing data transformation and reduction
 - Read (copies of) same input-files via XRootD in the same order
 - Fraction of input-files read from local cache
 - Number of jobs matches number of available slots
- Input-files on remote storage and prefetched on caches
- Only metadata transfer at stage-out
- Job monitoring part of the executable



Calibration Step 1: 10G Remote with Memory-Cache

- Setting gateway machine's network interface to 10 Gbit s⁻¹
- Each job reading same input-files → Memory cache operative



Measurement

 \blacksquare Clearly CPU limited workflow, no I/O throttling \rightarrow Tune CPU speeds in simulation

Calibration Step 1: 10G Remote with Memory-Cache

Setting gateway machine's network interface to 10 Gbit s⁻¹

36/13

08.05.2023

• Each job reading same input-files \rightarrow Memory cache operative



 \blacksquare Clearly CPU limited workflow, no I/O throttling \rightarrow Tune CPU speeds in simulation



Calibration Step 2: 1G Remote with Memory-Cache



- Lower gateway machine's network interface to 1 Gbit $s^{-1} \rightarrow$ expect throttling by network
- ightarrow Reuse pre-calibration information from previous measurement



Measurement

Limited by I/O via remote network

Calibration Step 2: 1G Remote with Memory-Cache



Reuse pre-calibration information from previous measurement \rightarrow

Measurement

Worker Node 1 60 60 Worker Node 1 Worker Node 2 Worker Node 2 50 Worker Node 3 50 Worker Node 3 jobtime / min 30 50 jobtime / min 40 30 20 10 10 . 0 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0 hitrate hitrate

• Limited by I/O via remote network \rightarrow 15% higher network bandwidth than expected





Calibration Step 3: 10G Remote with HDD-Cache



- Setting gateway machine's network interface to 10 Gbit s⁻¹
- \blacksquare Each job reading copies of same input-files \rightarrow data cache operative



Measurement

- Limited I/O by cache bandwidth
- ⇒ Now all simulation parameters are set!

Calibration Step 3: 10G Remote with HDD-Cache



 \blacksquare Each job reading copies of same input-files \rightarrow data cache operative



Measurement

- Limited I/O by cache bandwidth \rightarrow hard drive bandwidth $O(10 \, \text{MB s}^{-1})$
- ⇒ Now all simulation parameters are set!



Simulation

Validation: 1G Remote with HDD-Cache

- Setting gateway machine's network interface to 1 Gbit s⁻¹
- Each job reading copies of same input-files \rightarrow data cache operative

300 Worker Node 1 250 Worker Node 2 Worker Node 3 * jobtime / min 200 150 100 50 97.5 percentile 1.5 percentile 0 0.4 0.5 0.6 0.7 0.8 0.9 1.0 0.0 0.1 0.2 0.3 hitrate

Measurement



Simulation

39/13 08. 05. 2023 Maximilian Horzela: Modelling Distributed Computing Infrastructures for HEP Applications

Simulation reproduces measurement without further adjustments

• Limited I/O by cache bandwidth & network to remote \rightarrow Simulation parameters obtained previously

The simulator is able to reproduce reality!





Simulation Scaling

Scaling of $\mathcal{O}(10)$ -core system



Performant enough to simulate (parts-of) the grid with O (10⁵) jobs?



• Linear scaling of memory and simulator run-time \rightarrow 15 h runtime, 20 GB memory

Scaling of $\mathcal{O}(1000)$ -core system



• Performant enough to simulate (parts-of) the grid with $\mathcal{O}(10^5)$ jobs?



• Exponential scaling of simulator run-time

Scaling of $\mathcal{O}(1000)$ -core system



• Performant enough to simulate (parts-of) the grid with $\mathcal{O}(10^5)$ jobs?

- Simplifying platform (condense workers, simplify network) → Simpler assignment of resource shares
- $ightarrow\,$ Factor pprox 4 saving in runtime



Still exponential scaling of simulator run-time

Scaling of $\mathcal{O}(1000)$ -core system



• Performant enough to simulate (parts-of) the grid with $\mathcal{O}(10^5)$ jobs?

- Scaling I/O \rightarrow less total number of activities \propto SIMGRID solver calls
- ightarrow Factor pprox 20 saving in runtime



- Exponential Linear scaling of simulator run-time driven by exponential!
- → Exponential increase with platform size

Scaling to $\mathcal{O}(100k)$ -core systems



- Hard to win against exponential (this project's no-free-lunch)
- Modelling tricks: platform simplification, I/O scaling, buffer- & block-size tweaking, ...
 - Proof-of-concept (show that tricks retain realism) ongoing
- $\rightarrow\,$ Scale down the exponential, but still very much exponential
- How to proceed?



Simulation-Scaling Challenge

- Exponential increase of simulation time with number of concurrently active job slots
- $\rightarrow\,$ No point in running test simulations with \approx 75k concurrent jobs without changes
- Scaling I/O bandwidths and file sizes or XRootD-block and buffer-sizes reduces simulation time by constant factors O(10)
- ightarrow Maybe run it once everything is set





Improve Scaling by Exponential

- Exponential scaling with number of slots is a problem
- Why not scale along this exponentially manifesting dimension?
- $\rightarrow\,$ Scale the size of the platform, number of slots and bandwidths
 - Since jobs anyway drawn from distribution, decreasing the statistics is fine
 - \rightarrow Statistics can be recovered by repeating simulation with different random seed
 - Must assure that situation on platform is still the same





Improve Scaling by Exponential

- Exponential scaling with number of slots is a problem
- Why not scale along this exponentially manifesting dimension?
- $\rightarrow\,$ Scale the size of the platform, number of slots and bandwidths
 - Since jobs anyway drawn from distribution, decreasing the statistics is fine
 - \rightarrow Statistics can be recovered by repeating simulation with different random seed
 - Must assure that situation on platform is still the same
 - \Rightarrow Works for validation platform!

