Hoscar Bordeaux Workshop'13



INTERACTING WITH SCIENTIFIC WORKFLOWS FOR COMPUTATIONAL ANALYSES

Marta Mattoso, Patrick Valduriez (INRIA – Montpellier), Jonas Dias, Daniel de Oliveira, Kary Ocana, Eduardo Ogasawara, Flavio Costa, Felipe Horta, Vitor Silva and Igor Araujo.

Marta Mattoso - COPPE/Federal University of Rio de Janeiro



HOSCAR - Objectives



• "...taking full benefits of the processing capabilities of future high performance massively parallel architectures in the framework of very large-scale datasets and numerical simulations..."

Computer Science contributions to

- Group 3 (HPC software systems): basis for everyone
- Intergroup collaboration Group 1 (numerical simulations) and Group 2 (large-scale datasets)

Managing a Scientific Experiment PPE

- How to re-execute only part of a workflow ?
- What is the difference in execution time of the same program with different parameters
- What results are generated for particular sets of input values ?
- How to find previous results ?
- How can a user can trace the "process" that led to the aggregation of executions producing a particular output ?
- How to extract knowledge from previous experiments ?

SWfMS – support analysis AFTER execution

Several SWfMS are available, being used on different apps

Kepler

Your Science, Enabled.



Vis Trails

Making Science Flow For You





Triden

SWfMS keeps track of the execution through provenance support

SWfMS Workflow execution



- SWfMS manages parameter exploration
- N algorithms,
 - K parameters each,
 - L values in each parameter \rightarrow
 - SWfMS manages the algorithms, parallel processing, provenance
- End workflow execution
- Results & provenace are analyzed

SWfMS execution life cycle



- SWfMS manages parameter exploration
- N algorithms,
 - K parameters each,
 - L values in each parameter ->
 - SWfMS manages the algorithms, parallel processing, provenance
- End workflow execution
- Results & provenace are analyzed

New parameters are defined

Workflow Execution Life Cycle



CFD analyses take several factors into COP consideration

- geometry,
- material properties, e.g. viscosity,
- mesh partitioning,
- time step size,
- wall time
- the frequency at which the results are stored

Based on the produced outcome: explore the simulation differently



- refine the mesh,
- change time step size
- store more or less results during specific simulation time intervals

Matisse's painting Life Cycle



Matisse's painting Life Cycle



Workflow execution Off-line (black-box) X On-line (steering) to Luiz Combra de Engenharia

- Only after the whole workflow execution :
 - Change parameters
 - Replace some algorithms
- Interrupt the execution

During workflow execution

- Partial results & provenance are analyzed
- Fine tuning of parameters
- Skipping/replace some algorithms
- Skipping some parameter combinations





Workflow execution Off-line (black-box) X On-line (steering) to Luiz Coimbra de U

✓ Hard to visualize partial results "Off-line"

- A lot of information to be filtered
- Traceability of which input produces output
- Some results may not be useful
 - Data and meta-data association
- Fragmented Data
- Relevance of data
 - Each analysis may focus on a subset of the obtained outputs

Time consuming data transfers

Hoscar Workshop 2013

Some real examples that need "online"

Bioinformatics

- query and traverse partial results to change activities or parameters
- Oil & Gas
 - engineers need to explore the parameter space in slices skipping input data dynamically
- Numerical simulations, Uncertainty Quantification
 - snapshots of current simulation results to refine the model (iterations) during runtime









Provenance is a key feature



- Keeps track of everything that happens during experiment execution
- A log that can be queried
- Allows for high-level and domain-specific queries What are the maximum values for velocity and pressure on a given CFD simulation?
- A powerful association: Experiment meta-data with strategic experimental results
 - Runtime analysis
 - Good for time-consuming experiments
- User-steering for convergence analysis

User interacts with data of interest at runtime Interfere in the workflow execution to adjust

Open issues



Workflow execution monitoring

- verifying the status of the execution at specific points
- tools to handle data staging, consolidation, statistics and visualization
- Data analysis at runtime
 - discover if anything odd has happened

Dynamic interference in the execution

- change parameter values or even programs of the workflow during its execution
- truncate the iteration based on the error

Current approaches: Data Analysis at Runtime



"access to preliminary results to analyze them empowered with data from provenance repositories and statistical and visualization tools" \rightarrow N/A

- VisTrails- visualize results and compare them on a single screen or on tiled wall displays – no HPC
- Swift/Turbine or Pegasus
 - HPC in very large scale with provenance
 - No **runtime** provenance support
- ParaView- instruments the code and select produced results to be exported for visualization in their sw, no SWfMS

Dynamic Workflow Interference



- Needs monitoring and runtime analysis
- Simple interventions:
 - stopping a workflow, stopping an activity, or reexecuting activities
- Complex interventions:
 - changing an activity of the workflow, interfering in the iterative process, changing the structure of the workflow and changing parameters or complete parameter sets

Current approaches: Dynamic Workflow Interference



- Pegasus failure handling in scientific workflows
 - Regression trees to classify activity executions as failed or successful. Based on provenance data the framework learns activity execution behavior and predicts future failures.
- Dias *et al.* present dynamic parameter sweep in HPC workflows
 - Based on partial result data from runtime provenance scientists can decide when to stop a time dependent iteration (loop) of a scientific workflow.

Our experience with Chiron workflow engine



- Designed for HPC environments
- Algebraic engine based on relational algebra *
 - Runs legacy and third party applications
- SciCumulus in the Cloud
- Model to represent provenance data
 - PROV extension
 - Prospective and retrospective provenance
 - Available for querying at runtime

 Ogasawara et al., An Algebraic Approach for Data-Centric Workflows. VLDB, 2011

Algebraic engine and provenance oriented

Hoscar Workshop 2013

Trigger – Query



- Submit a query to the cloud (or cluster) to :
 - monitor the execution time of a specific activity execution and check if it is under or above the average time of previous executions of the same activity
 - obtain results from one specific parameter and do a visualization preview before the end of execution

Turbulence UQ analysis on a submarine

COPPE Instituto Alberto Luiz Coimbra de Pós-Graduação e Pesquisa de Engenharia



on Uncertainty Quantification and Stochastic Modeling 2012

Hoscar Workshop 2013



Uncertainty Quantification using Adaptive Sparse Grid Collocation





Large Eddy Simulation Workflow Parallel Execution in Chiron

Instituto Alberto Luiz Coimbra de Pós-Graduação e Pesquisa de Engenharia



- Provenance gathers data to be used to evaluate statistical moments
- We can track the error between MC and ASGC and determine an interpolation level to satisfy a minimum error
- Provenance capabilities in Chiron can automatically increase the interpolation level and resubmit the experiment to obtain statistical moments with a prefixed error



- Go fast to the results you need to analyze
 - ✓ Navigation over data results (partial) added by provenance (log)
 - ✓ Filter results and stage out only the data you need to analyze
 - ✓ Interact with a web-based user interface
- Export and analyze partial results to take decisions
 - ✓ Efficient data consolidation and data staging features to aggregate and stage out only relevant information
 - Use statistical and visualization tools empowered by provenance data

Draw preliminary conclusions



- Monitor and Analyze workflow execution
 - ✓ Awareness of the current workflow status
 - $\checkmark\,$ Identify and solve problems during execution
 - ✓ Recover from failures and misconceptions
 - ✓ Reduce time spent in processing low-quality data
 - $\checkmark\,$ Cut execution time and reduce financial costs
 - Stage data out can be costly!
 - Remote clusters



- Use high-resolution display environments
 - Visualize multiple results to establish comparisons
 - On parameters exploration scenarios
 - ✓ Analyze detailed images and simulations
 - Can take advantage of available tiled display technologies such as TACC DisplayCluster, SAGE, CGLX, and Paraview

Large-Scale Scientific Data COPPE Visualization using Provenance UFRJ

Web interface connected to Provenance

Database and visualization environment

- ✓ Interactive interface
- $\checkmark\,$ Navigate through its activities and i/o data
- $\checkmark\,$ Filtering options to select the desired results
- ✓ Stage data out of the execution environment
- Text, images, PDF, video and Paraview visualization

Large-Scale Scientific Data COPPE Visualization using Provenance UFRJ

- Web service in the visualization cluster to display staged out data in the Tiled Wall Display
 - ✓ Implemented as a Web Service
 - \checkmark Interface accessed by the web module in runtime
 - Can also stage data out of the execution environment
 - Pluggable to support different tiled wall displays middlewares.



Hoscar Workshop 2013

Conclusions



Traces of a real-time provenance database

- It is possible to take more advantage of experiment's provenance data
- Real-time provenance offers valuable information to help scientists
 - Monitor, analysis, steering, interference ...
- Effort to follow PROV recommendation
 - How to define stereotypes for HPC in PROV
 - Considering hybrid workflows
 - Running in the cluster or Cloud with distinct systems



Acknowledgements







Conselho Nacional de Desenvolvimento Científico e Tecnológico





à Pesquisa do Estado do Rio de Janeiro



Hoscar Workshop 2013

Hoscar Bordeaux Workshop'13





Chiron workflow engine

 http://sourceforge.net/projects/ chironengine/



SciCumulus in the Cloud

 http://sourceforge.net/projects/ scicumulus/

Marta Mattoso Federal Univ Rio de Janeiro



Hoscar Workshop 2013