

# **PROTEUS:** A SCIENTIFIC WORKFLOW GATEWAY ENRICHED BY PROVENANCE FOR LARGE-SCALE EXPERIMENTS

Felipe F. Horta

M.Sc. Student at COPPE/UFRJ

Associate Researcher at NACAD - *High  
Performance Computing Center (COPPE-UFRJ)*

fhorta@cos.ufrj.br

Advisor: Marta Mattoso

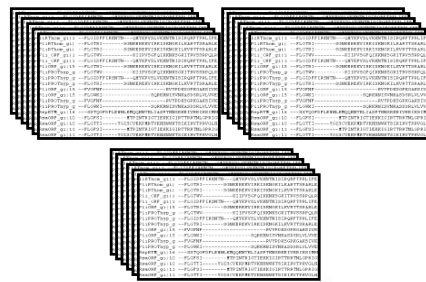
Co-advisor: Renato N. Elias



# Scientific Workflow Scenario

The analysis uses a chain of programs that may require HPC, e.g. clouds

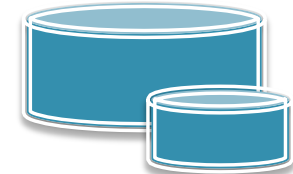
1. Scientific experiments are analyzed...



2. And modeled as scientific workflows



3. Large volume of data produced ...



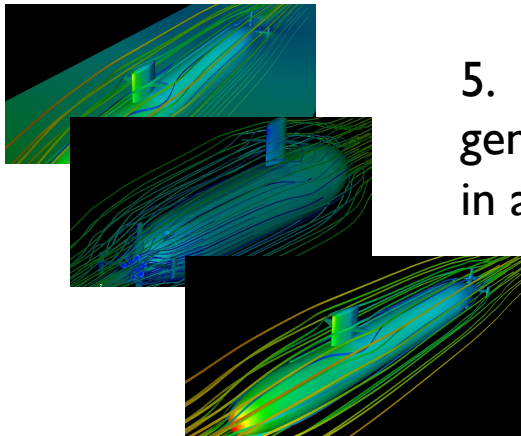
*Provenance Data*



5. Final results generated in a feasible time



4. ...which need to be processed by a computing-intensive environment



CFD

# Scientific Workflows in HPC

- The same problems of parallel programs
  - ▣ Volatility of computational resources
  - ▣ Failure occurrences at runtime
  - ▣ Difficulties in debugging failures
- Large-Scale Experiments Scenario
  - ▣ Black Box execution
    - Use of different computer programs and scripts
    - Heterogeneity/granularity
  - ▣ Hard to visualize partial results
    - Traceability
    - Fragmented Data
    - Laborious data transfers

# SWfMS

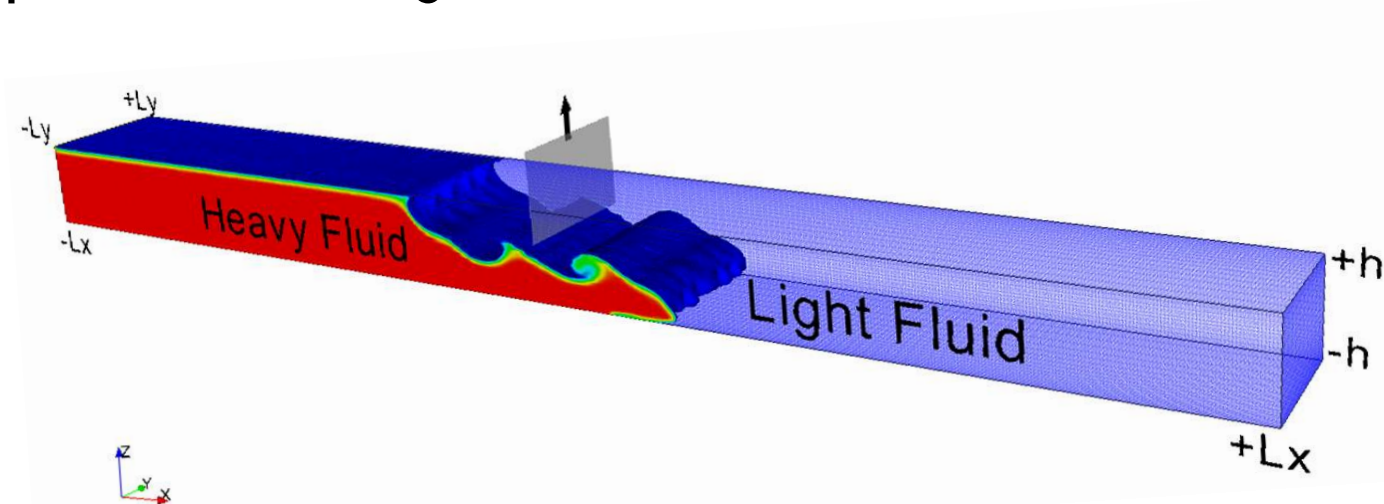


- ❑ SWfMS can take over debugging (app layer) by taking advantage from “knowing” what is behind workflow tasks and data-flow
  - Queries about activities that presented failures
  - Adjustments in parameters or workflow modeling
- ❑ Improves experiment management
- ❑ May provide data uniformity
- ❑ **May generate Provenance Data**



# Provenance, 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?



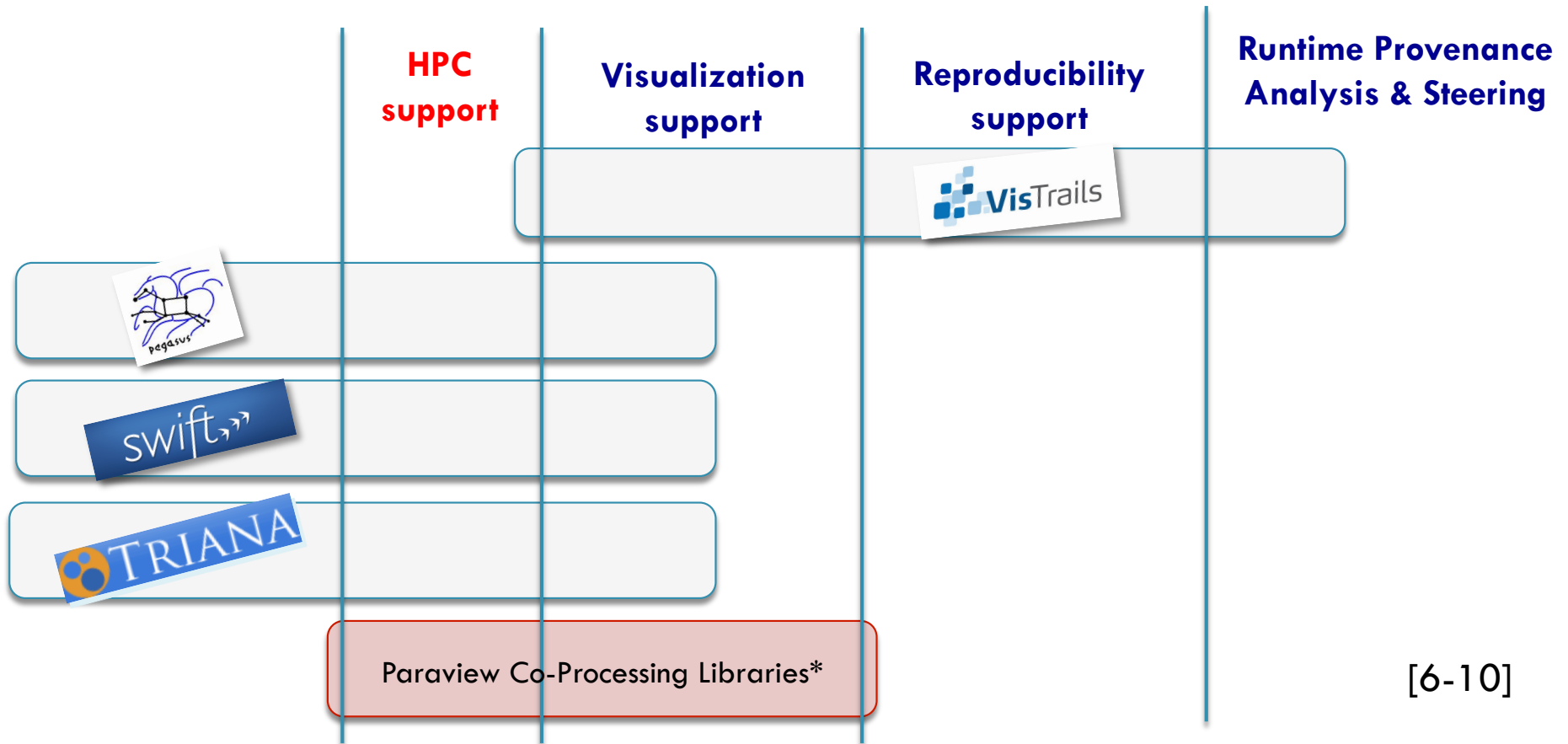
# Provenance Visualization

- 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

✓ Provenance enriched visualization may support **runtime** and systematic analysis on results from large-scale experiments.

# Related works

No solution that integrates HPC execution with visualization enriching analysis with runtime provenance data

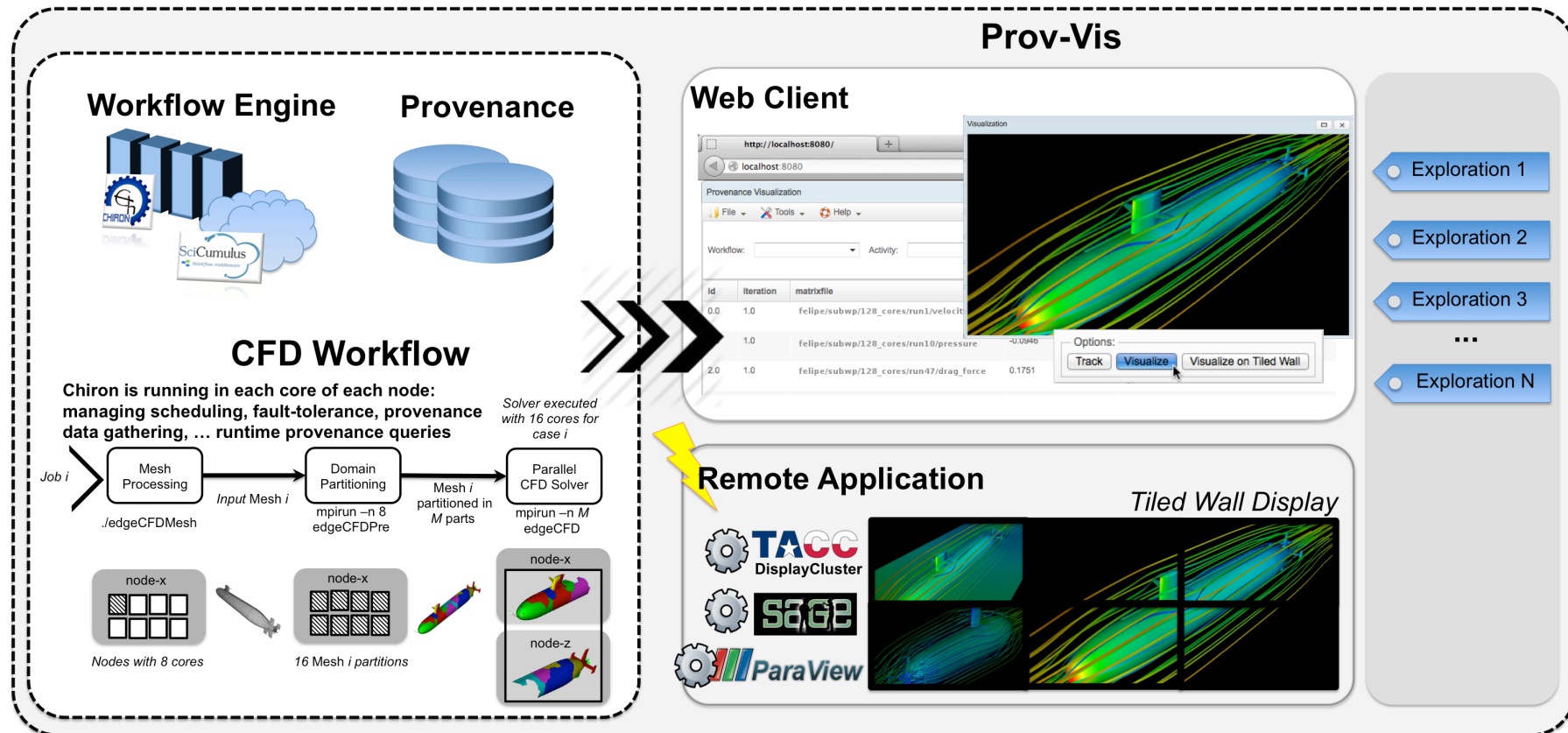


\* Not related with SWfMS

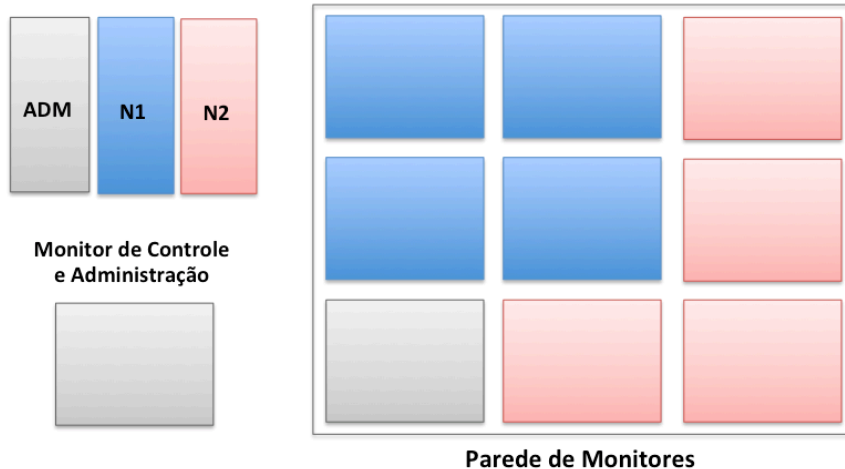
Uncoupled  
Supporters



# Provenance Visualization

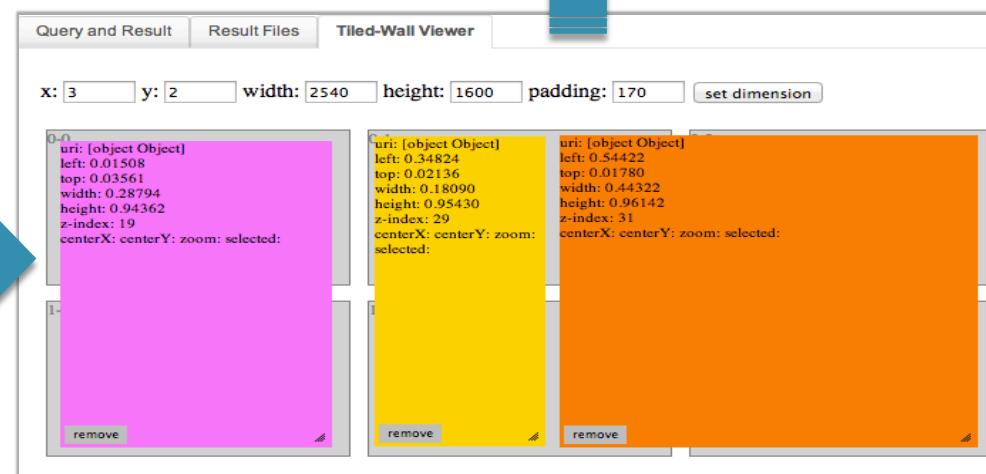


# Visualization Environment



- 3 DELL workstations
  - ▣ 24 cores
    - Intel Xeon E5506 @ 2.13GHz
  - ▣ 36 GB RAM memory
  - ▣ 5 NVIDIA Quadro 6000
  - ▣ 2,5 TB hard disk
- 10 DELL displays
  - ▣ 31" w/ 2560x1600 resolution
  - ▣ 9 displays for visualization
  - ▣ 1 for administration

[5]



# Issues

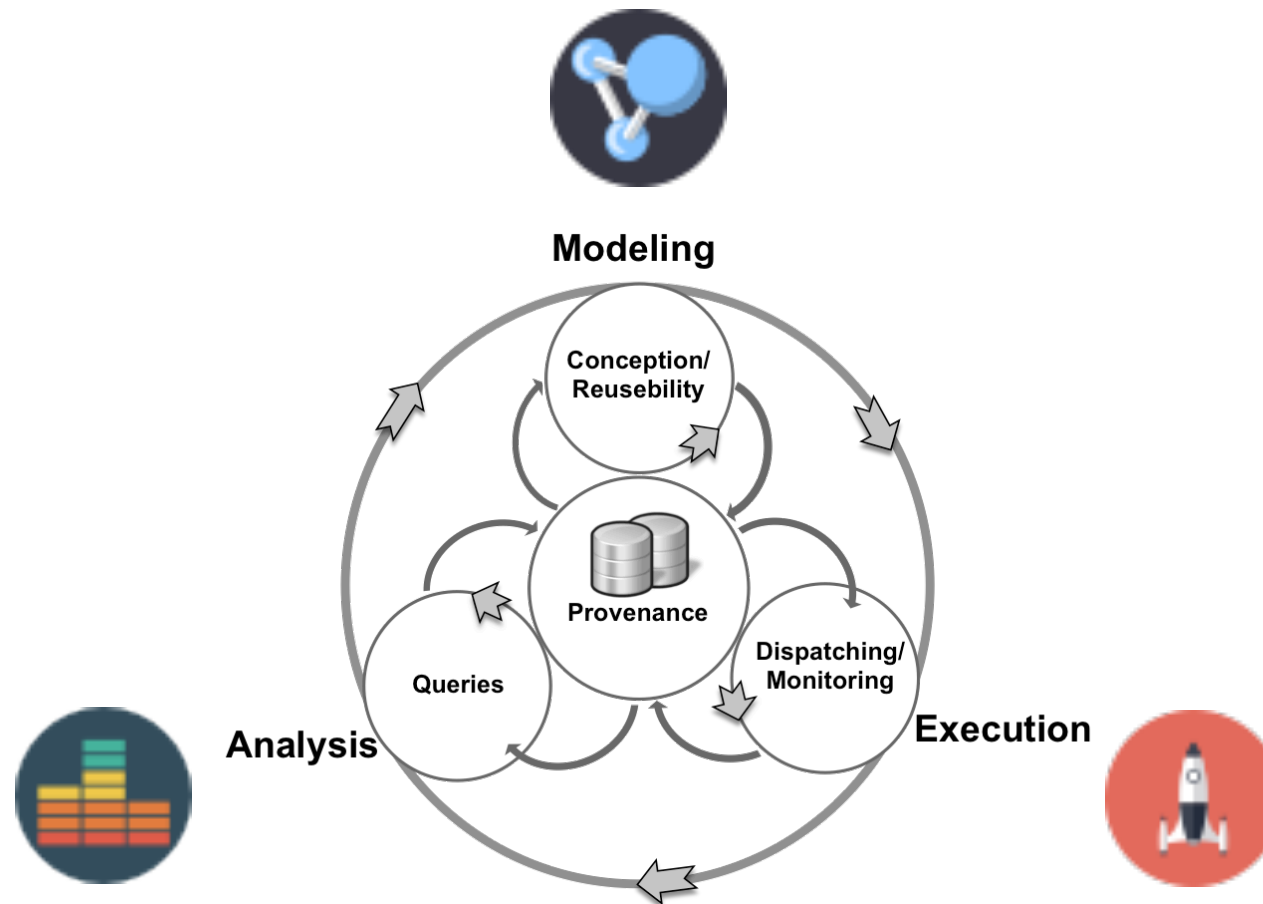
- Provenance Visualization scope
  - ▣ Workflow engine (Chiron)
    - Soon, version would be deprecated by a new one
  - ▣ Cluster still “hardcoded” attached to solution
    - No credential management
- Chiron Engine scope
  - ▣ Handle XML submissions may be hard to manage
  - ▣ Sharing workflow (model/results) is not easy
- Research Group scope
  - ▣ No integration with another projects on our research group
    - GExp-Line; Wf Modeling; User-Steering; or Distributed Provenance;

# Proteus Scientific Gateway: Goals

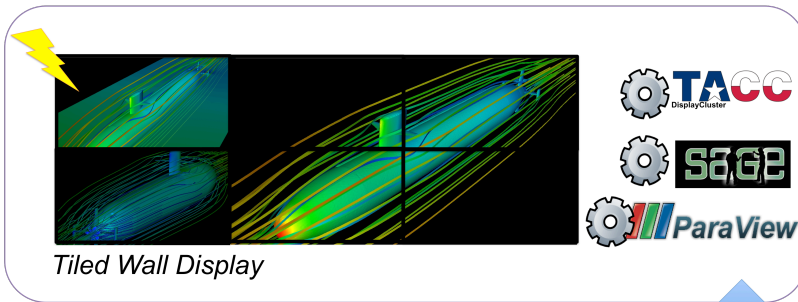
- Consolidation
  - ▣ Keep control of Chiron engine's change release (**backward compatibility**)
  - ▣ Expose current research solutions as features
- Platform for applications
  - ▣ Improve interoperability and reusability
  - ▣ Support workflow data integration
- Why not support full scientific workflow life-cycle?
- Support Reproducibility
  - ▣ Uncoupled environments, data, workflow models
  - ▣ Access control & user management



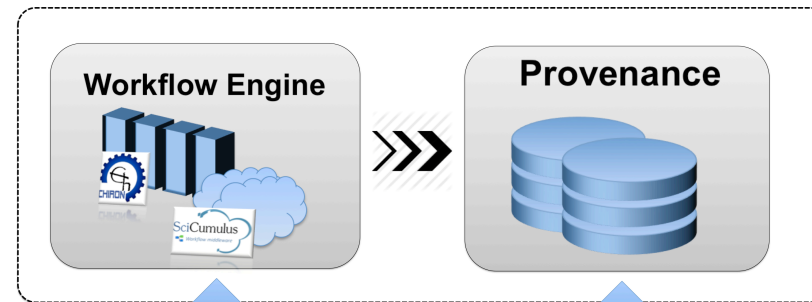
# Scientific Workflow Life Cycle



## Remote Visualization

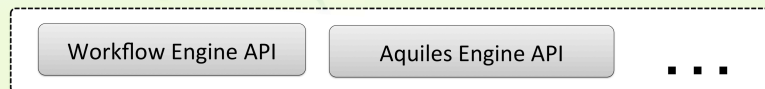
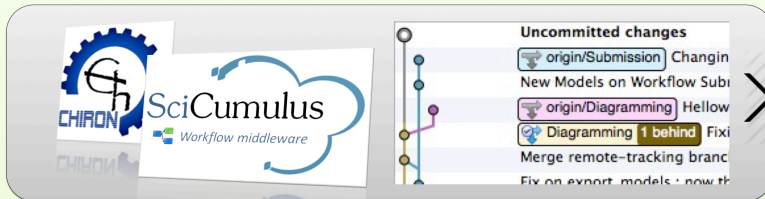


## Remote HPC and/or Provenance

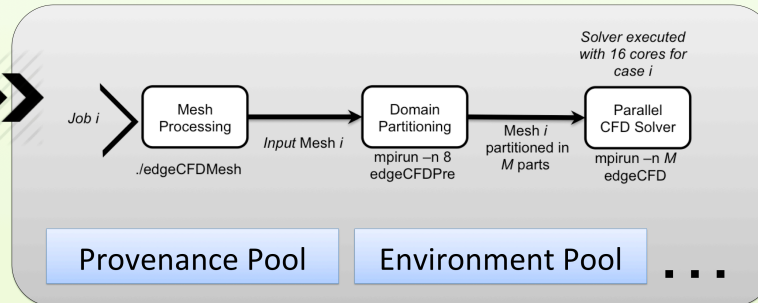


## Proteus Scientific Gateway

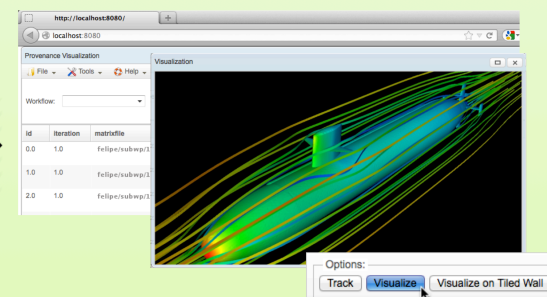
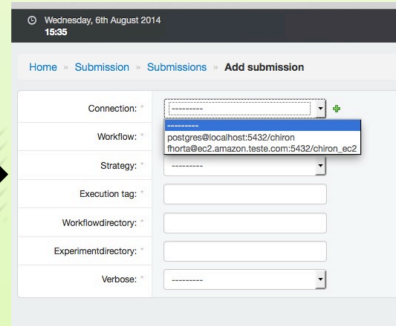
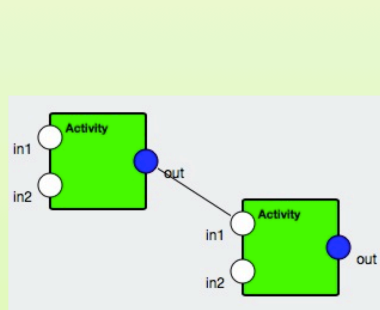
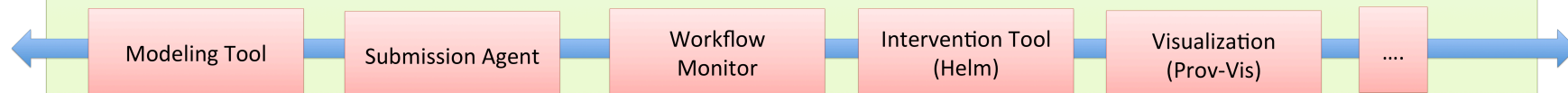
### Workflow Engine Management



### Proteus Database

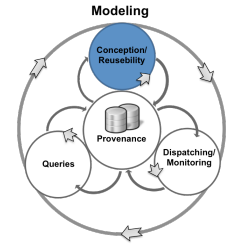


## Backend Interface

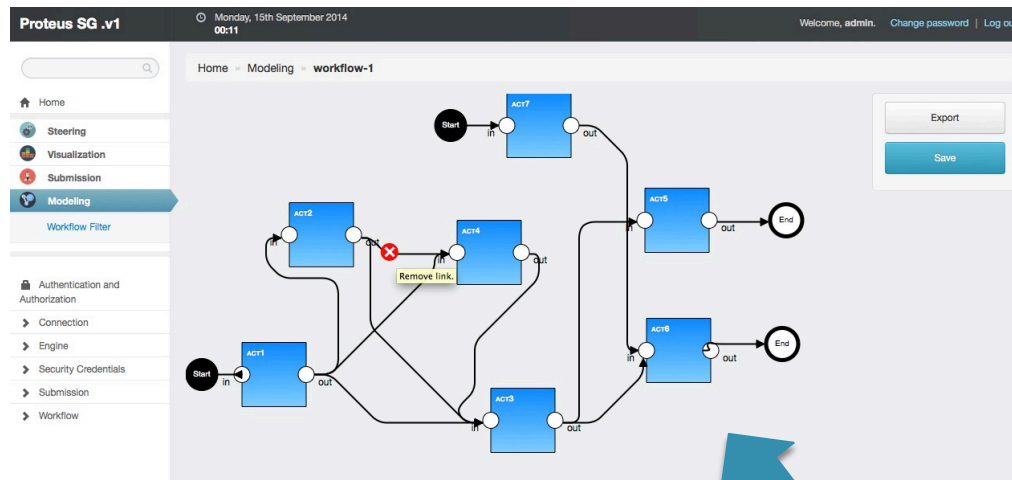


## Access Management

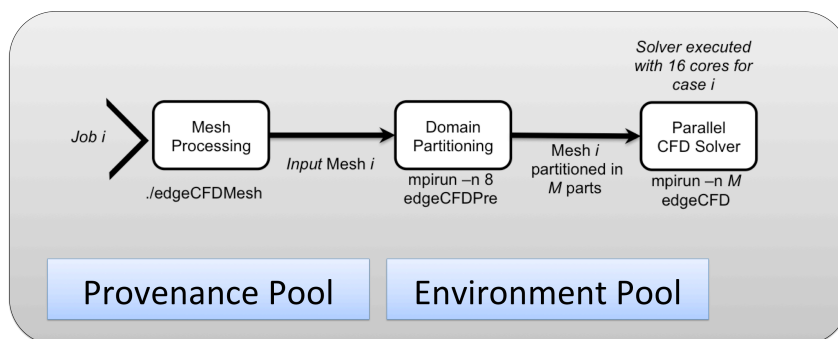
# State of art - Modeling



## Create/edit workflow models



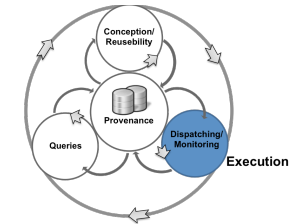
## Persist/Update models in Proteus DB



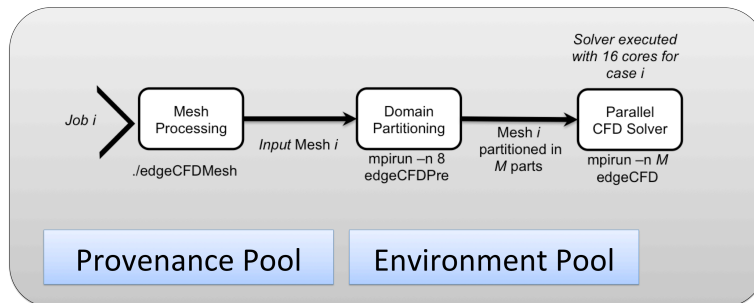
## Import/Export workflow models



# State of art - Submission



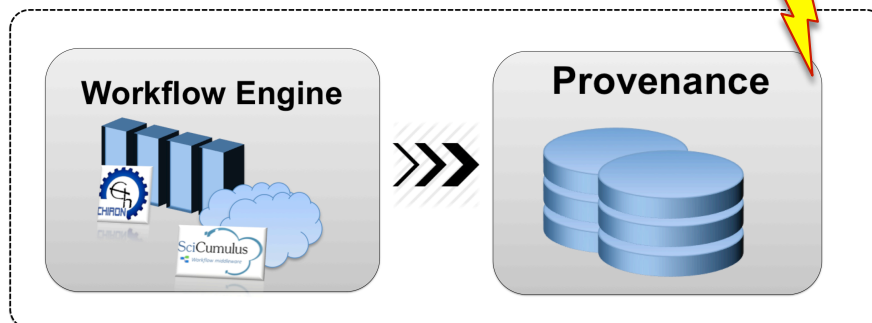
## Proteus Database



## (Re) Parameter workflow submission

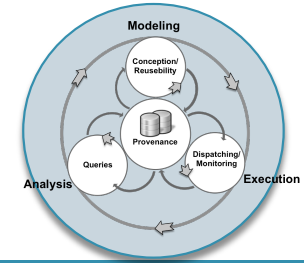
A screenshot of the Proteus SG .v1 web interface. The 'Add submission' form is displayed, showing fields for 'Created by', 'Name', 'Credential', 'Environment', 'Binary', 'Machine', 'Constraint', 'Workspace', 'Database', 'Query', 'Conceptual workflow', 'Execution workflow', and 'Status'. The 'Status' field is set to 'workflow-1'. The form includes 'Save', 'Save and continue editing', and 'Save and add another' buttons.

## Remote HPC and/or Provenance



```
<tabbox id="tb" width="100%" height="100%">
  <tabs>
    <tab label="Query and Result" />
    <tab label="Result Files" />
    <tab label="Tiled-Wall Viewer" />
  </tabs>
  <tabpanel>
    <table id="tableTitle" multiline="true" style="font-size:13px;">
      <tr>
        <td>dbresults</td>
      </tr>
    </table>
  </tabpanel>
  <tabpanel>
    <table id="dbfiles" width="100%" height="100%">
      <tr>
        <td>dbfiles</td>
      </tr>
    </table>
  </tabpanel>
  <tabpanel>
    <a style="display:none;" href="#" id="send_xml" class="send_xml" >Refresh!</a>
    <zscript>
      send_xml.setWidgetListener(Events.ON_CLICK, "sendXmlToServer()");
    </zscript>
    <include src="viewer.html">
      <custom-attributes org.zkoss.zul.include.html.defer="true"/>
    </include>
  </tabpanel>
</tabpanel>
</tabbox>
```

# State of art - Reproducibility



## Manage Permission Groups

## Provide permission credentials on your publication

### Prov-Vis: Large-Scale Scientific Data Visualization

Felipe Horta<sup>1</sup>, Jonas Dias<sup>1</sup>, Renato Elias<sup>1</sup>, Daniel de Oliveira<sup>2</sup>, Alvaro L. G.

<sup>1</sup> COPPE- Federal University of Rio de Janeiro, <sup>2</sup> IC- Fluminense  
{fhorta, jonasdias, marta}@cos.ufrj.br, {renato, alvaro}@nacad.ufrj.br

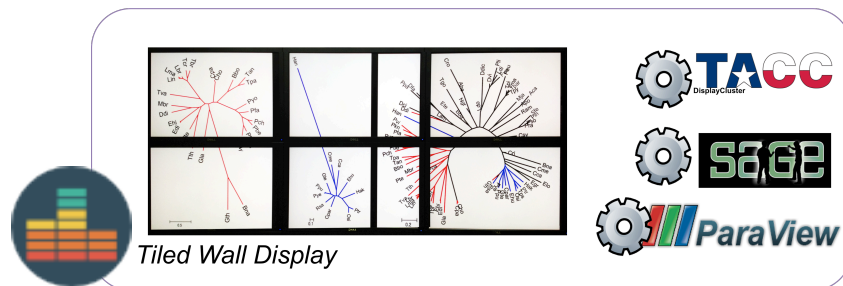
**Abstract** — Large-scale scientific computing often rely on intensive tasks chained through a workflow. Scientists need to check the status of the execution at particular points, to discover if anything odd has happened and take actions. To achieve that, they need to track partial result files, which is usually complex and laborious. When using a scientific workflow system, provenance data keeps track of every step of the execution. If traversing provenance data is allowed at runtime, it is easier to monitor and analyze partial results. However, visualization of partial results is necessary to be done in sync to the workflow provenance. Prov-Vis is a scientific data visualization tool for large-scale workflows that is based on runtime prove-

hard to identify their subset of the outputs. Scientific Work may improve data. They make it easier to through provenance SWMS since it also happened during two high-level and don't "What are the maximum a given simulation e-

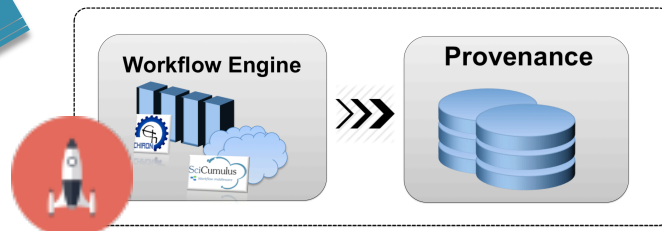
## Guest sign in

## Guest (re)parameter or (re)run workflow

Guest may change provenance, env., input data etc..



## Remote HPC and/or Provenance



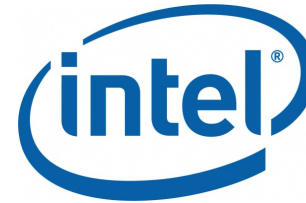
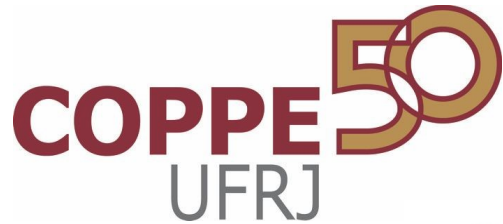
# Contributions

- Backward compatibility
  - ▣ Keep control of Chiron engine's change release
- Support workflow life-cycle integration
  - ▣ Improve interoperability and reusability
    - Chiron Engine [1]
    - Modeling App – Lucas Carneiro (IC student, JICTAC 2014)
    - Submission App– Kaique Rodrigues (IC student, JICTAC 2014)
    - Steering Engine- Jonas Furtado Dias (Ph.D. project) {reference} [2]
    - Steering Interface App – F. Pinheiro (undergraduate final project, B.E Computer UFRJ)
    - Prov-Vis App – Felipe Horta (me) {reference}
  - ▣ Support workflow data integration
- Support experiment reproducibility
  - ▣ Uncoupled environments, data, workflow models
  - ▣ Access control & user management

# References

- [1] E. Ogasawara, J. Dias, V. Silva, F. Chirigati, D. de Oliveira, F. Porto, P. Valduriez, and M. Mattoso, “Chiron: a parallel engine for algebraic scientific workflows,” *Concurrency and Computation: Practice and Experience*, vol. 25, no. 16, pp. 2327–2341, Nov. 2013.
- [2] M. Mattoso, K. Ocaña, F. Horta, J. Dias, E. Ogasawara, V. Silva, D. de Oliveira, F. Costa, and I. Araújo, “User-steering of HPC workflows: state-of-the-art and future directions,” in *Proceedings of the 2nd ACM SIGMOD Workshop on Scalable Workflow Execution Engines and Technologies*, 2013, p. 4.
- [3] F. Horta, J. Dias, R. Elias, D. de Oliveira, A. Coutinho, and M. Mattoso, “Prov-Vis: Large-Scale Scientific Data Visualization using Provenance,” 2013.
- [4] F. Horta, J. Dias, K. A. C. S. Ocana, D. de Oliveira, E. Ogasawara, and M. Mattoso, “Abstract: Using Provenance to Visualize Data from Large-Scale Experiments,” 2012, pp. 1418–1419.
- [5] K. A. C. S. Ocaña, D. de Oliveira, F. Horta, J. Dias, E. Ogasawara, and M. Mattoso, “Exploring Molecular Evolution Reconstruction Using a Parallel Cloud-based Scientific Workflow,” in *Advances in Bioinformatics and Computational Biology*, vol. 7409, Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 179–191.
- [6] Y. Zhao, M. Hategan, B. Clifford, I. Foster, G. von Laszewski, V. Nefedova, I. Raicu, T. Stef-Praun, and M. Wilde, “Swift: Fast, Reliable, Loosely Coupled Parallel Computation,” in *3rd IEEE World Congress on Services*, Salt Lake City, USA, 2007, pp. 206, 199.
- [7] E. Deelman, G. Mehta, G. Singh, M.-H. Su, and K. Vahi, “Pegasus: Mapping Large-Scale Workflows to Distributed Resources,” in *Workflows for e-Science*, Springer, 2007, pp. 376–394.
- [8] S. P. Callahan, J. Freire, E. Santos, C. E. Scheidegger, C. T. Silva, and H. T. Vo, “VisTrails: visualization meets data management,” in *SIGMOD International Conference on Management of Data*, Chicago, Illinois, USA, 2006, pp. 745–747.





# THANK YOU!

PROTEUS: A SCIENTIFIC WORKFLOW GATEWAY  
ENRICHED BY PROVENANCE FOR LARGE-SCALE  
EXPERIMENTS

Felipe Horta