

Runtime Data Analytics in Computational Science

Alvaro Coutinho
alvaro@nacad.ufrj.br

High Performance Computing Center
COPPE/Federal University of Rio de Janeiro
www.nacad.ufrj.br



*With the invaluable help of Marta Mattoso, Fernando Rochinha (COPPE) and
Jose Camata (UFJF)*

Contents

- Computational Science and Engineering
- High Performance Computing
- In-situ Visualization
- Predictive Computing
- Predictive Data Analytics
- Provenance in Data Analytics
- Data Science and Machine Learning
- Conclusions and Discussion



COMPUTATIONAL SCIENCE AND ENGINEERING

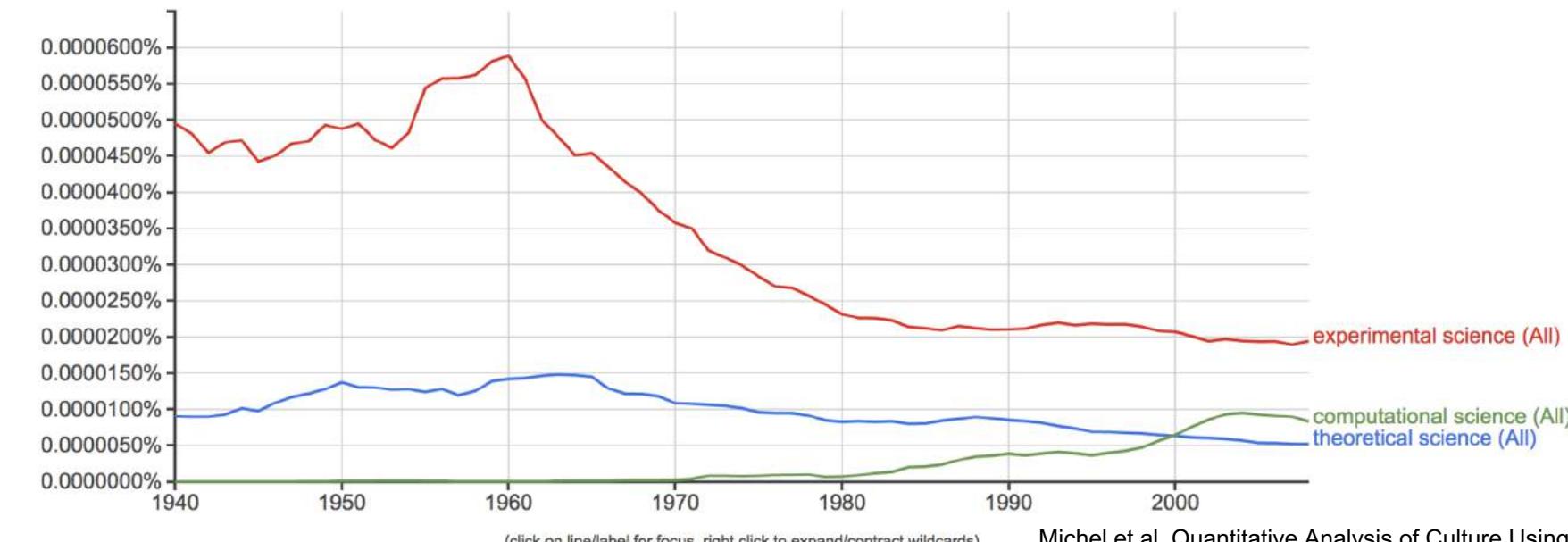
Computational Science and Engineering

Google Books Ngram Viewer

Graph these comma-separated phrases: theoretical science,experimental science,computational science case-insensitive
 between and from the corpus English with smoothing of 3

[G+ Share](#) 0
[Tweet](#)

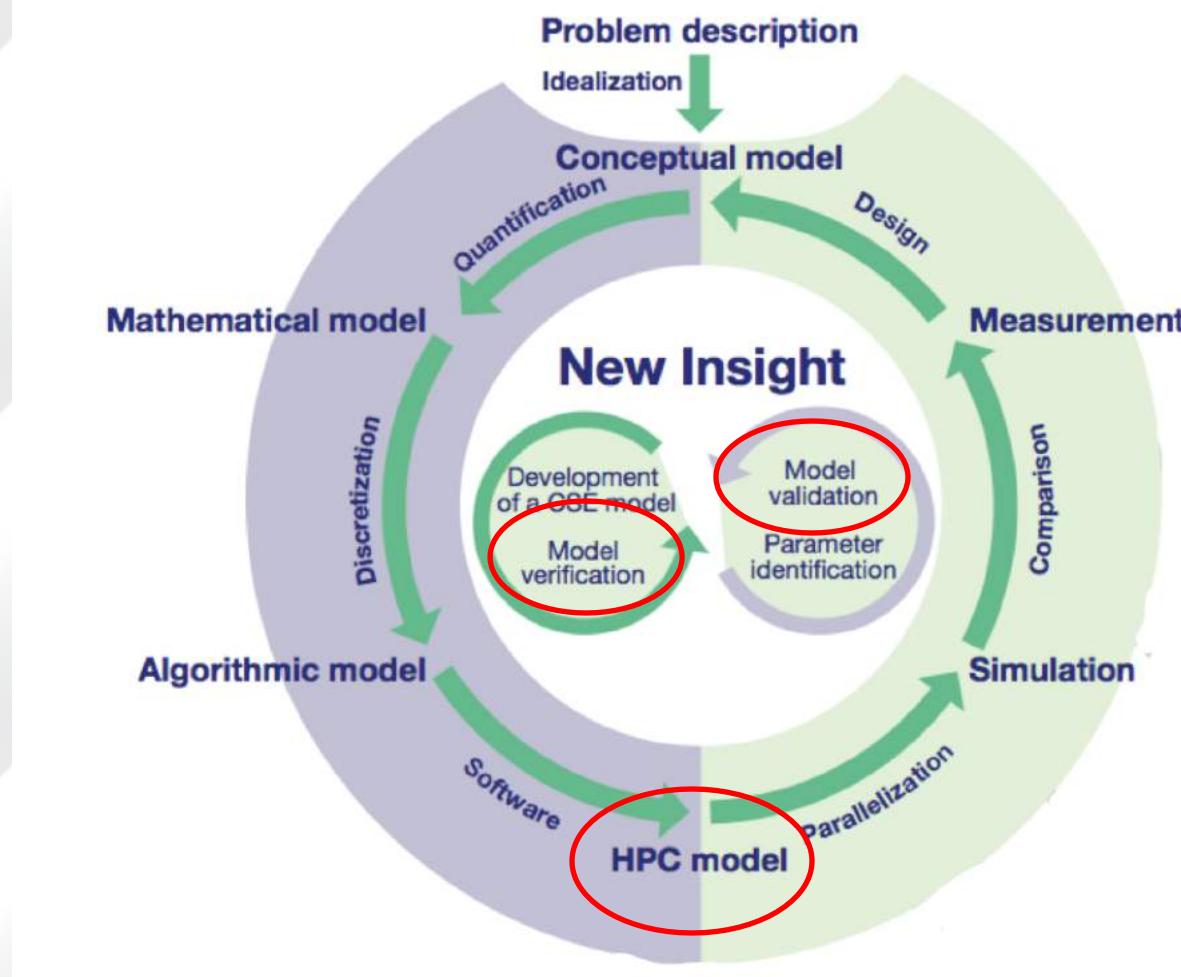
[Embed Chart](#)



Michel et al, Quantitative Analysis of Culture Using Millions of Digitized Books, Science, 2011

- CSE encompasses methods of HPC and it is central in Data Sciences

CSE Pipeline¹



HIGH PERFORMANCE COMPUTING

High Performance Computers or Supercomputers

Supercomputers are the fastest and most powerful general purpose scientific computing systems available at any given time.

Dongarra et al, "Numerical Linear Algebra for High-Performance Computers", SIAM, 1998



Turing's Bombe, UK, 1941

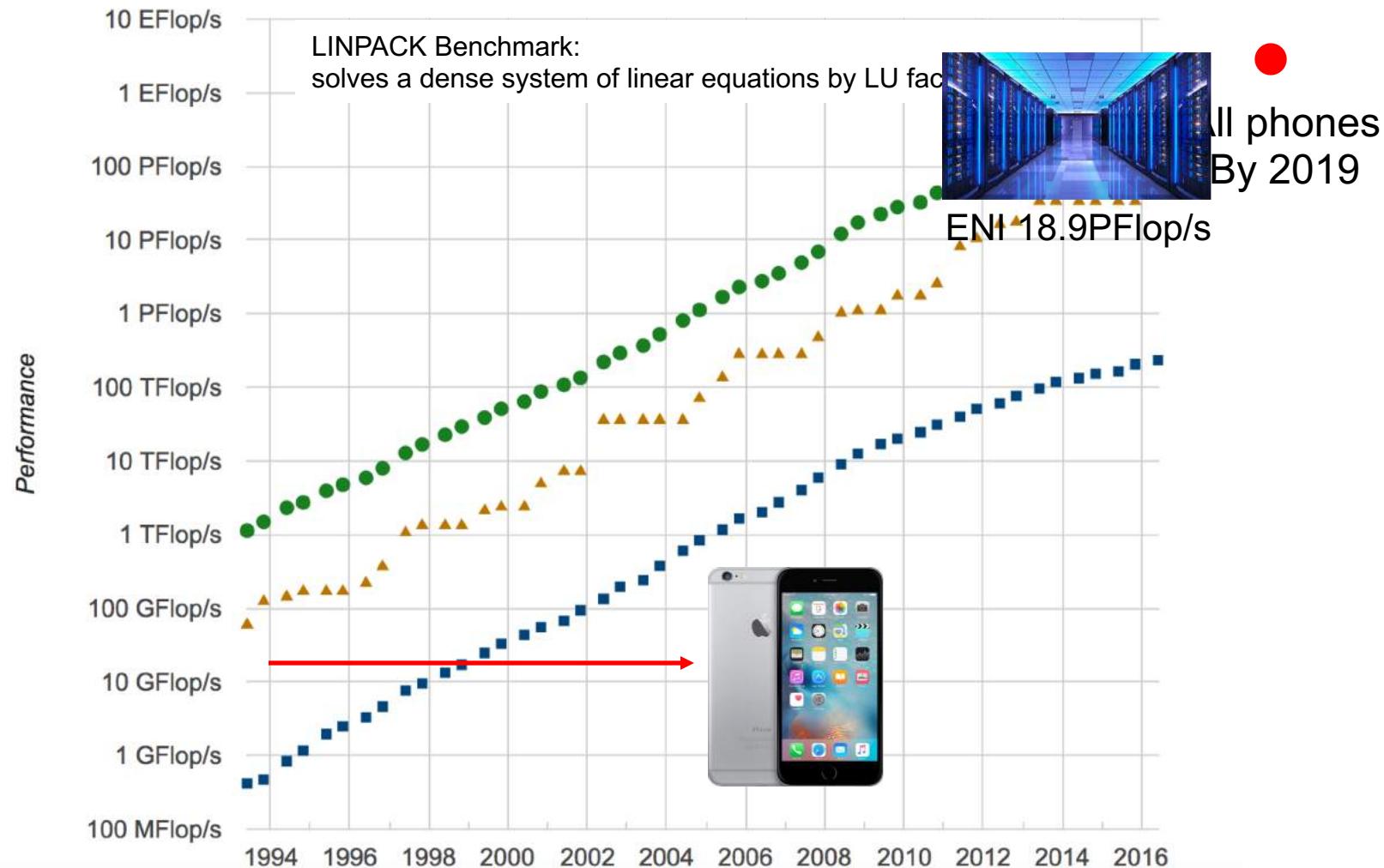
Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, Cores: 2,414,592; Memory: 2,801,664 GB, Rmax: 148,600 TFlop/s; Rpeak: 200,795 TFlop/s; Power: 10,096.00 kW



<http://www.olcf.ornl.gov/olcf-resources/compute-systems/summit/>

Historical Trends TOP500 List

Performance Development



Brazil in the TOP500 List (Jun 2019)

Rank	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
142	Fênix - SYS-1029GQ-TRT, Xeon Gold 5122 4C 3.6GHz, Infiniband EDR, NVIDIA Tesla V100 , Bull, Atos Group Petróleo Brasileiro S.A Brazil	48,384	1,836.0	4,297.4	287
419	BC1 - Lenovo C1040, Xeon E5-2673v4 20C 2.3GHz, 40G Ethernet , Lenovo Cloud Provider Brazil	38,400	1,123.2	1,413.1	
431	BC2 - Lenovo C1040, Xeon E5-2673v4 20C 2.3GHz, 40G Ethernet , Lenovo Software Company (M) Brazil	38,400	1,123.2	1,413.1	

- High Performance Science -
- To provide
- Develop a especially
 - Energy
 - Civil,
 - Environ.
 - Comp.
 - Biolog.
- Member of

Intel® Parallel Computing Centers

Apply Today for a Grant

Financial analyses, Climate modeling & weather prediction, Computer Aided Design & Manufacturing, Energy - Seismic, Medical

ABOUT THE PROGRAM CURRENT CENTERS BECOME A CENTER NEWS

Click on the logos to learn more about what each of these Intel® Parallel Computing Centers is doing.

QUICK LINKS

[Intel® Software Academic Program](#)

[Academic Courseware](#)

[Intel® Many Integrated Core Architecture Forum](#)

[Intel® Xeon Phi™ Coprocessor Developer Starter Kit](#)

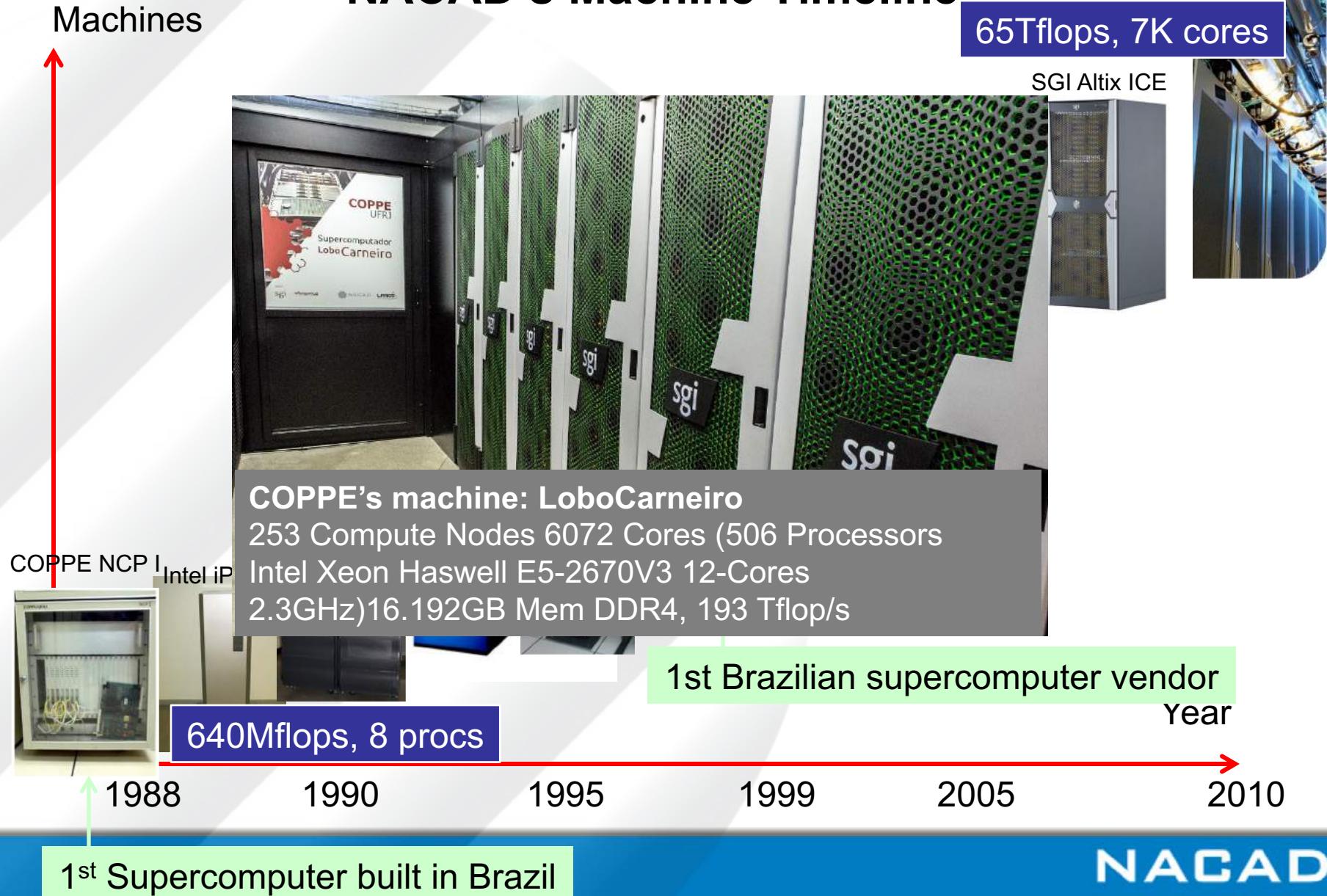
[FAQs](#)

[Send us your comments](#)

tional
88
puting
evance,
and Gas

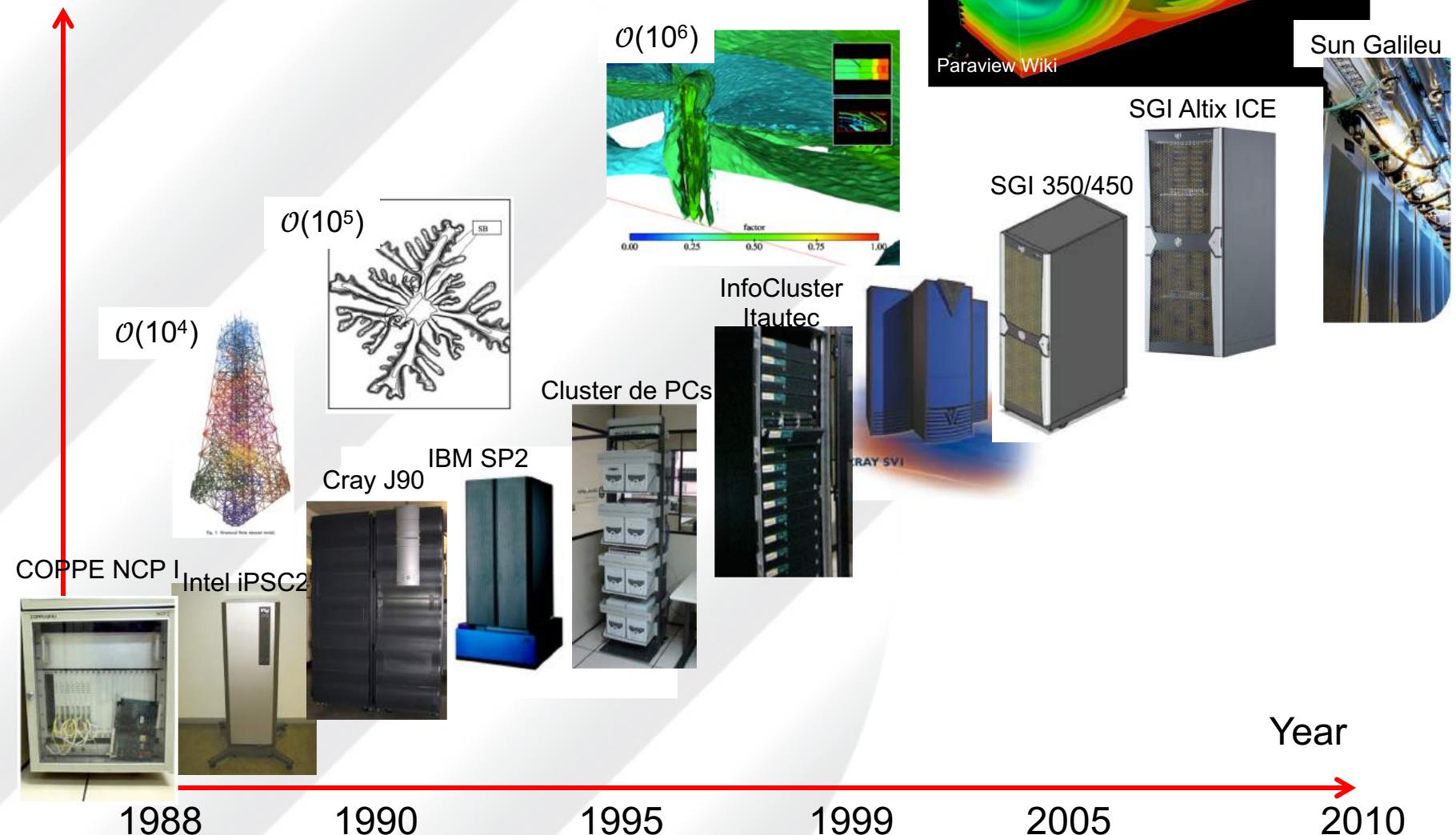
NACAD's Machine Timeline

Sun Galileu

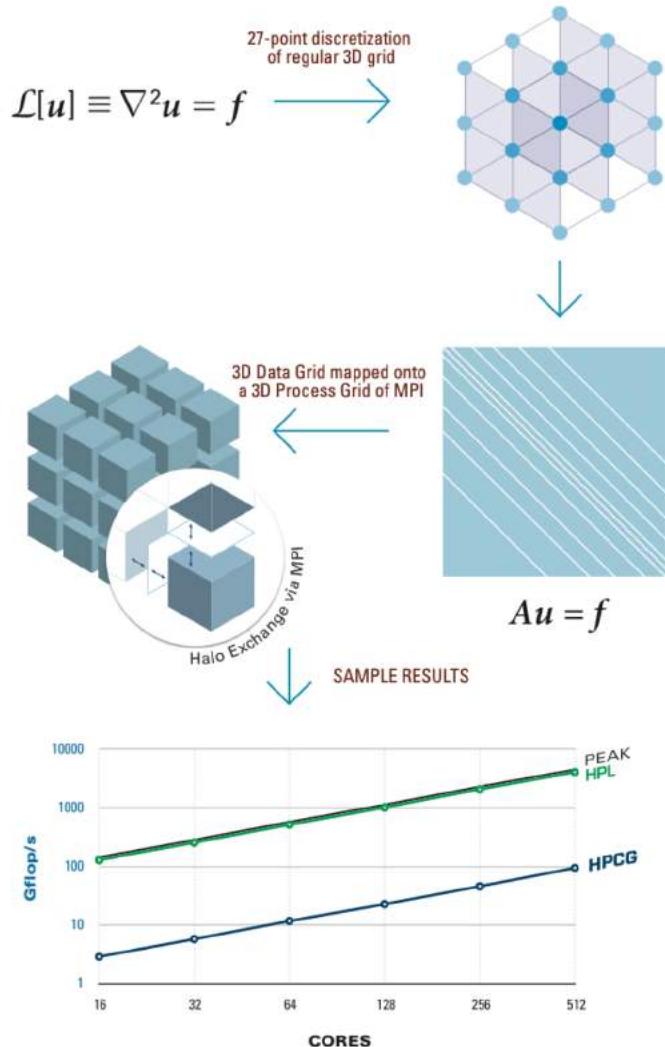


NACAD's AppsTimeline

App Complexity



The HPCG Benchmark¹



HPCG on LoboCarneiro

Distributed Processes: 6048

Global Problem Dimensions:

nx: 1040 ny: 1040 nz: 1456

Number of Equations: 1,574,809,600

Number of Nonzero Terms: 42,445,920,184

GFLOP/s rating of: 4520.67 ~2.34% peak

HPL: 193.09 Tflop/s

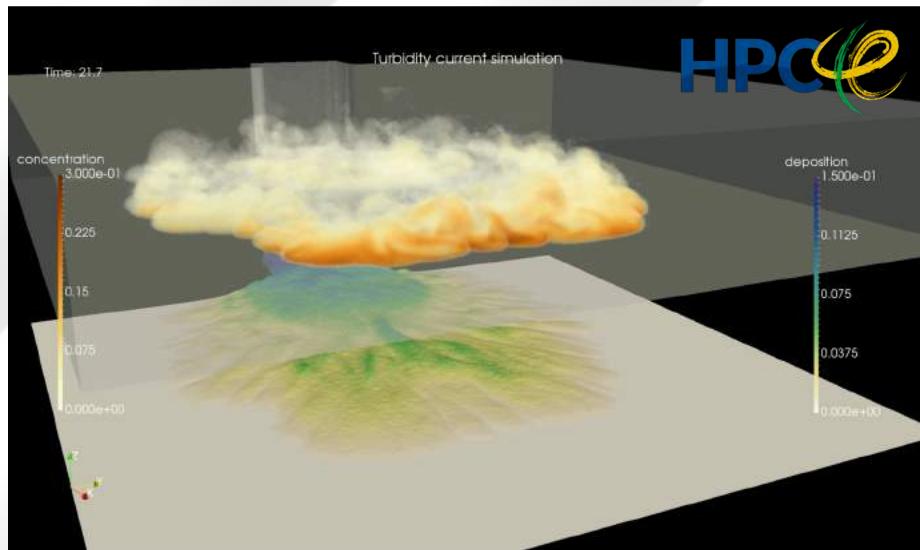
HPCG List for June 2019

TOP500			Rmax [TFlop/s]	HPCG [TFlop/s]
Rank	Rank	System	Cores	
1	1	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband , IBM DOE/SC/Oak Ridge National Laboratory United States	2,414,592	148,600.0 2925.75
2	2	Sierra - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband , IBM / NVIDIA / Mellanox DOE/NNSA/LLNL United States	1,572,480	94,640.0 1795.67
3	20	K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect , Fujitsu RIKEN Advanced Institute for Computational Science [AICS] Japan	705,024	10,510.0 602.74
4	7	Trinity - Cray XC40, Xeon E5-2690v3 16C 2.3GHz, Intel Xeon Phi 7250 68C 1.4GHz, Aries interconnect , Cray Inc. DOE/NNSA/LANL/SNL United States	979,072	20,158.7 546.12
5	8	AI Bridging Cloud Infrastructure [ABCi] - PRIMERGY CX2570 M4, Xeon Gold 6148 20C 2.4GHz, NVIDIA Tesla V100 5XM2, Infiniband EDR , Fujitsu National Institute of Advanced Industrial Science and Technology [AIST] Japan	391,680	19,880.0 508.85
6	6	Piz Daint - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100 , Cray Inc. Swiss National Supercomputing Centre (CSCS) Switzerland	387,872	21,230.0 496.98
7	3	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway , NRPCP National Supercomputing Center in Wuxi China	10,649,600	93,014.6 480.85

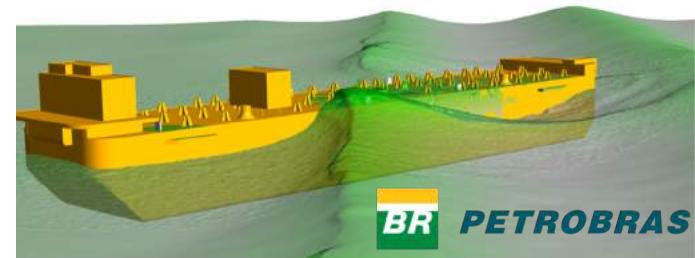
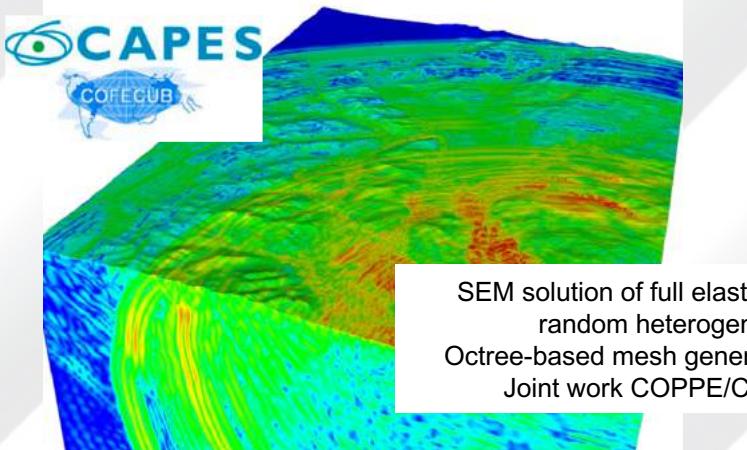
~2.0% peak

¹from: //software.sandia.gov/hpcg/

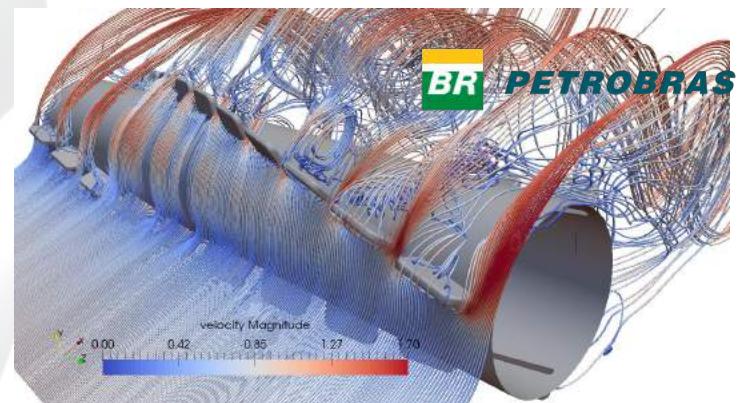
Real life computations



Turbulent turbidity current simulation
FE-VMS, EdgeCFD, 30M tets, LoboCarneiro

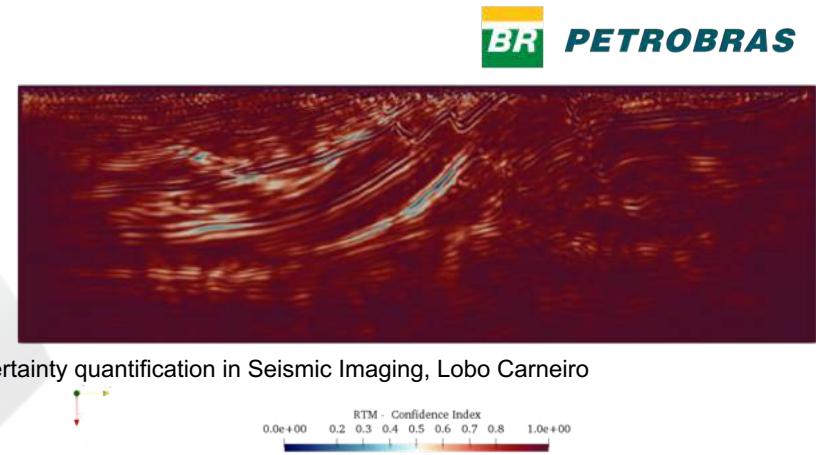
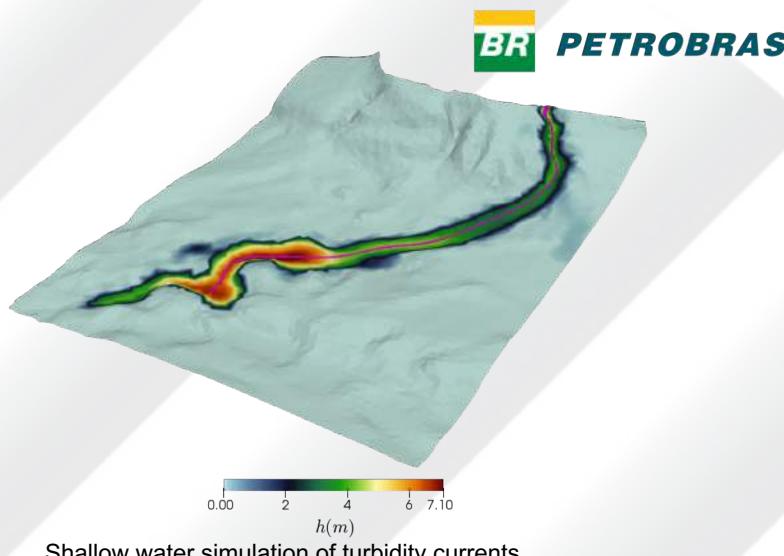


Fluid-Structure Interaction: Wave reaching midship region, simulation and experiments visualization.
Computed with FLUENT, Cluster Petrobras

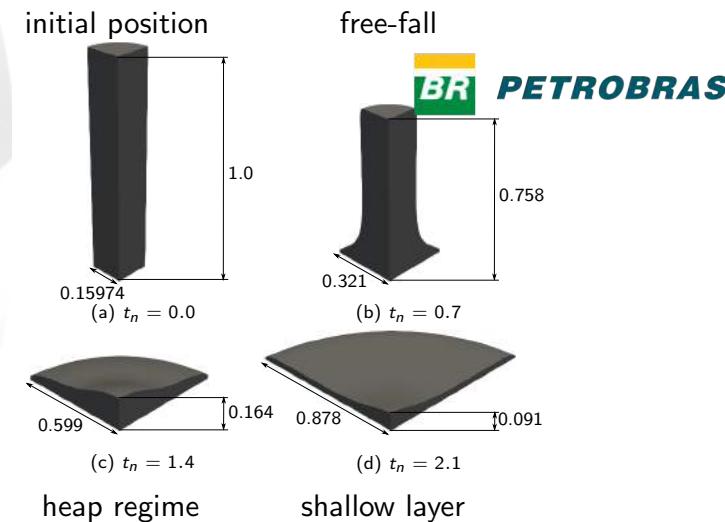


VIV on a rigid riser with strakes, Re=10K,
Computed with EdgeCFD 7M tet4, LoboCarneiro

Real life computations



Material science: Finite element simulation of Cahn-Hilliard equation in Lobo-Carneiro, FEniCS, 4.3M equations, 6000 time steps, wall time 360h



Finite element simulation of the collapse of a column of dense granular material in Lobo Carneiro libMesh, 54M Hex8

A Simplified CSE Workflow

Pre-processing and mesh generation: select parameters (Δt_{max} , tolerances, solvers, AMR/C fractions, etc.)

Time stepping:

$t \leftarrow 0$

while $t < t_{max}$ **do**

- a. Solve nonlinear equations
- b. Adapt mesh/time step (Δt)
- c. Update solution
- d. Save data on disk when required
- e. $t \leftarrow t + \Delta t$

end while

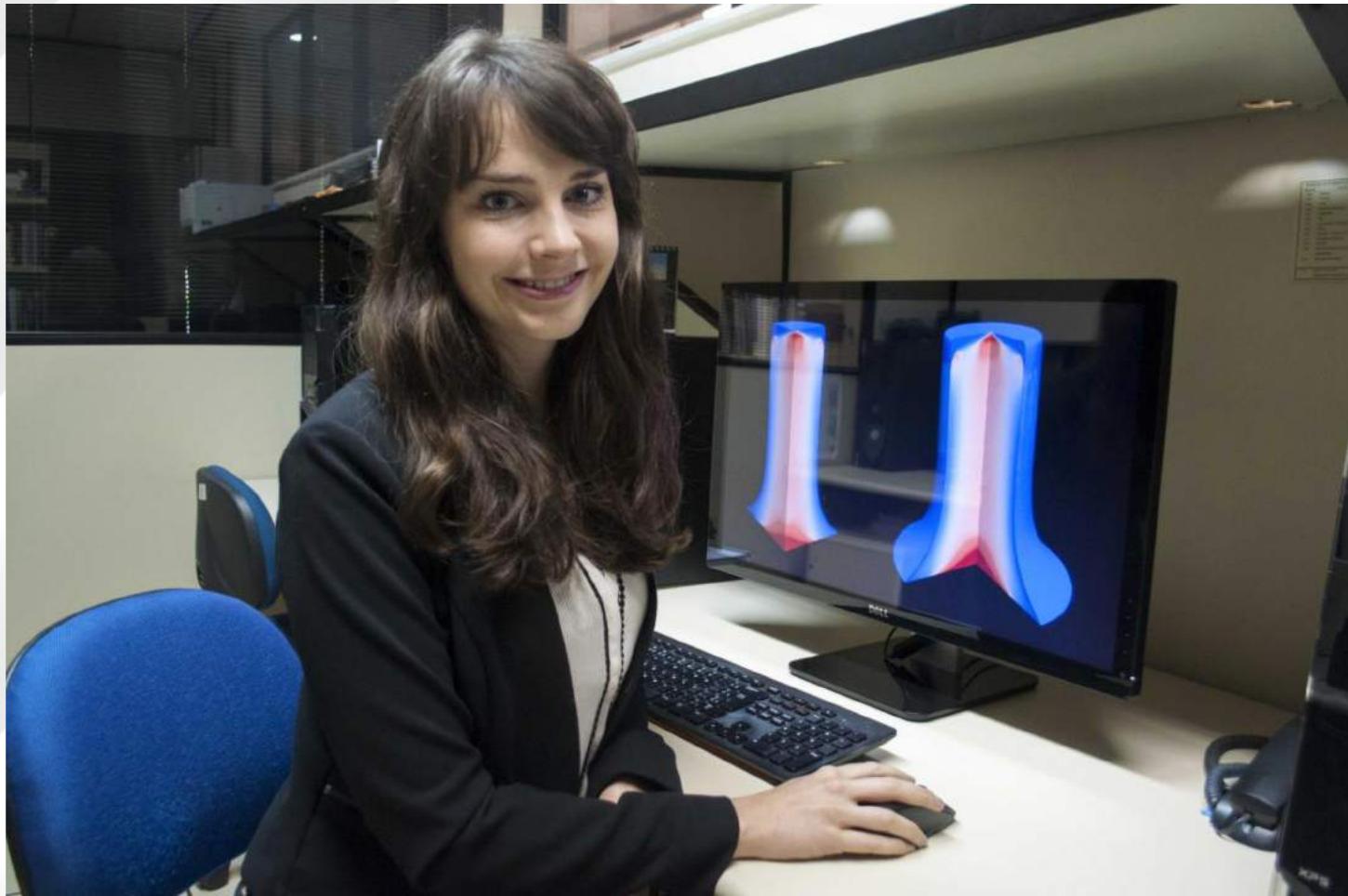
Post-processing, visualizing the data and extracting relevant QoIs

- ▶ **Step a:** computer time
- ▶ **Step b:** accuracy
- ▶ **Step d:** I/O in persistent storage
- ▶ **Optimal computational complexity¹:** $\mathcal{O}(n_{eq}^{\frac{4}{3}})$

CSE Software

- ▶ Translate complex mathematical models into predictive tools
- ▶ Usually coded in Fortran, C/C++, Python, Java or a mix of all of them
- ▶ Often invokes components of CSE frameworks and libraries
- ▶ Components are invoked to provide for:
 - ▶ support for PDE discretization methods like libMesh, FEniCS, MOOSE, deal.II, GRINS, OpenFOAM, PetIGA;
 - ▶ mesh generation, like Gmsh;
 - ▶ building blocks for solving numerical problems with parallel computations, like PETSc, LAPACK, SLEPc;
 - ▶ visualizations, like ParaView, VisIt;
 - ▶ in-situ, like ParaView Catalyst, SENSEI
 - ▶ I/O data management, like ADIOS.
- ▶ Modification of parameters at runtime allowed in several of these, e.g., PETSc.
- ▶ **Log files have to be manually inspected to take decisions; no query support at runtime**

Linda Gesenhues Named Recipient of 2018 ACM-IEEE CS George Michael Memorial HPC Fellowships



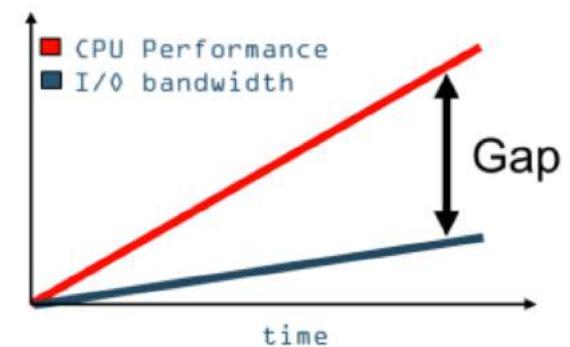
IN-SITU VISUALIZATION

WHY DO WE NEED IN-SITU VISUALIZATION?

- The old way: post-processing visualization

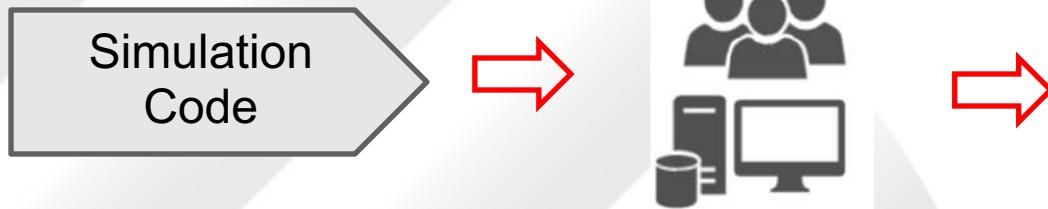


- Large-scale simulations can produce TB of data
- Moving data between computer centers has became a major bottleneck



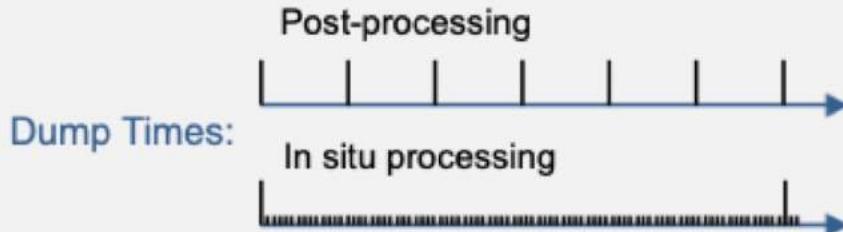
WHAT IS IN-SITU VISUALIZATION?

- ❑ Basics:



Visualizing results of a running simulation
without the need to copy the data to a storage device.

- ❑ Data Compression
 - Reduce traffic I/O and volume
- ❑ Access to more data
 - ❑ Enable real-time simulation exploration
 - ❑ Reduce latency to first result
 - ❑ Avoid large scale post processing



WHAT IS IN-SITU VISUALIZATION?

- ❑ Basics:



Visualizing results of a running simulation
without the need to copy the data to a storage device.

But...

- ❑ Simulation and Visualization code...
 - ... might have different internal data structures
 - ...might scale best with different parallelization strategies

Coupling simulation and visualization codes is a challenge.

HOW TO START WITH IN-SITU VISUALIZATION?

- **ParaView Catalyst**
 - In situ analysis and visualization library that enables using ParaView's visualization capabilities in-situ workflows.
 - Allows to execute complex analysis pipelines synchronized with the simulation, and connecting with the ParaView GUI for live, interactive visualization.
- **Visit Libsim**
 - Library that makes available the full complement of features from VisIt
 - Enables VisIt to connect interactively to running simulations for live exploration
- **ADIOS**
 - adaptive I/O service that is designed to allow applications to easily change between different I/O service providers
- **SENSEI**
 - generic data interface that allows access to multiple in situ infrastructures

In-Situ Visualization with ParaView Catalyst



- Typically 3 calls between simulation code and adaptor
 - Initialize()
 - *MPI communicator (optional)*
 - *Add analysis scripts*
 - CoProcess()
 - *Does the work (potentially)*
 - Finalize()
- Information provided by solver to adaptor
 - Time, time step, force output
 - Grids and fields
- Information provided by adaptor
 - Pipelines to execute
 - Time, time step, force output
 - Grid and fields when needed
 - MPI communicator
- Information provided by Catalyst
 - If co-processing needs to be done
 - What grids and fields are needed
- User data can be shared both ways

The new paradigm

PREDICTIVE COMPUTATIONAL SCIENCE

Predictive Science¹

Definition: Predictive science is the scientific discipline concerned with accessing the predictability of mathematical and computational models of physical events in the presence of uncertainties. It embraces the process of model selection, calibration, validation, verification, and their use in forecasting features of physical events with quantified uncertainty.

Models: mathematical constructions based on physical principles or empirical relations-generally based on inductive theories which attempt to characterize abstractions of physical reality

¹J. Tinsley Oden, Ivo Babuska, and Danial Faghihi, Predictive Computational Science: Computer Predictions in the Presence of Uncertainty, Encyclopedia of Computational Mechanics, Wiley

The Imperfect Paths to Knowledge

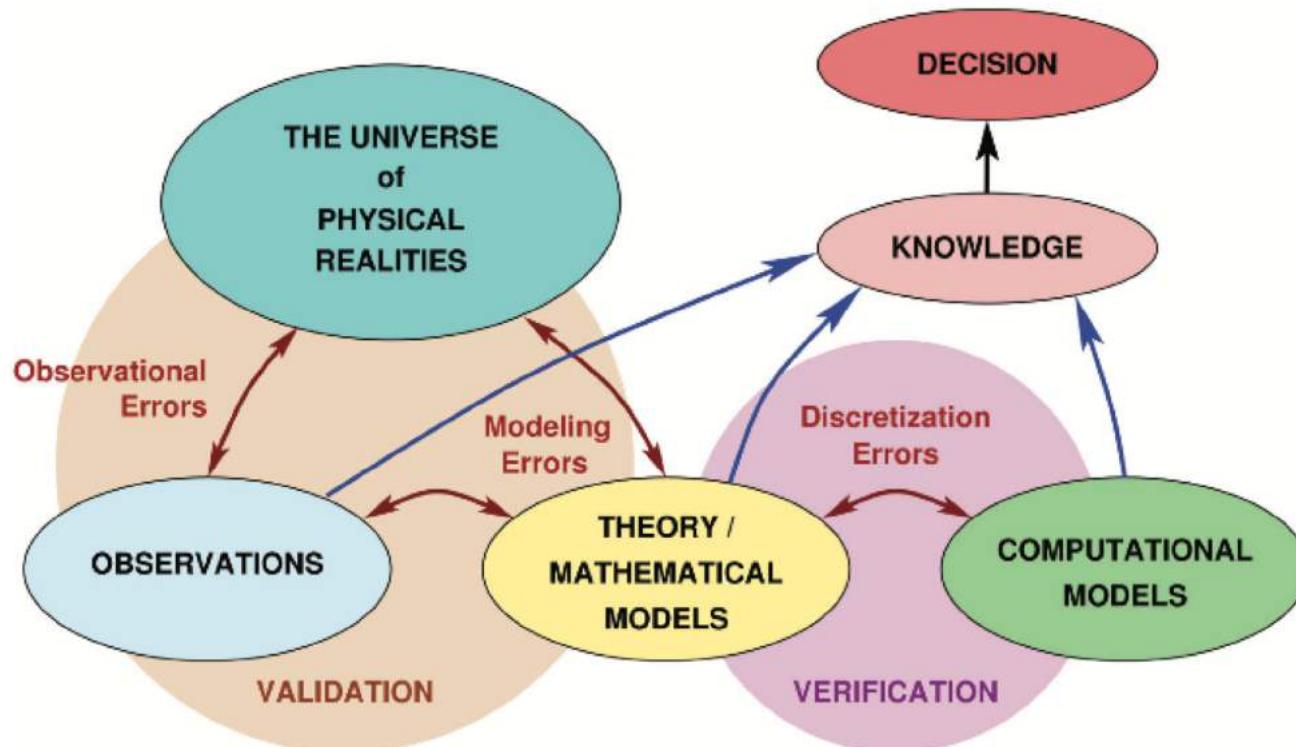
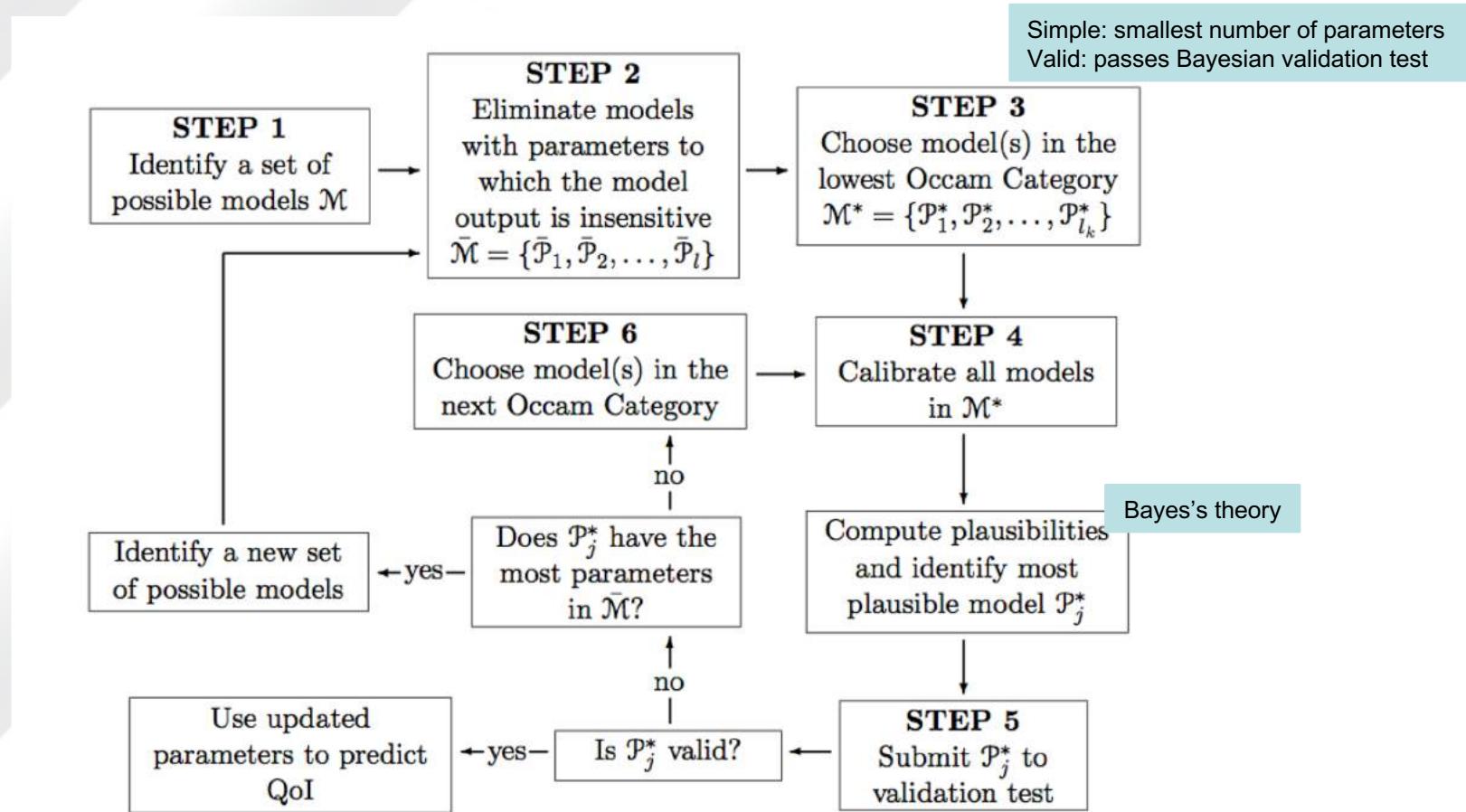


Figure 1. Imperfect computational modeling: Imperfections in the mathematical models, incomplete observational data, observations delivered by imperfect instruments, and corruption of the model itself in the discretization needed for computation all lead to imperfect paths to knowledge. Reproduced from J.T. Oden, "A Brief View of V & V & UQ," a presentation to the Board on Mathematical Sciences and Their Applications, National Research Council, October 2009.

The OCCAM Plausibility Algorithm¹



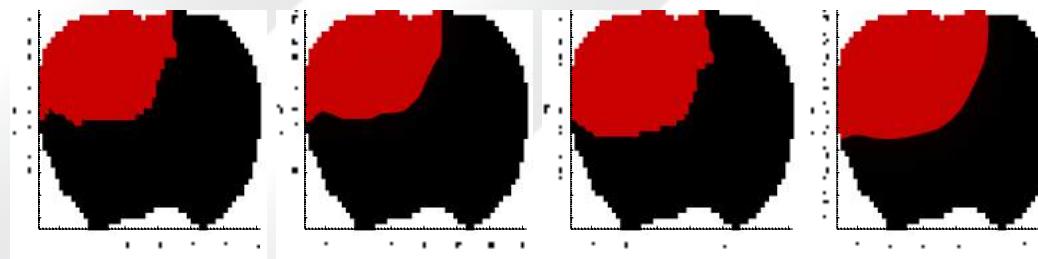
¹Farrell, Oden, Faghihi, A Bayesian framework for adaptive selection, calibration, and validation of coarse-grained models of atomistic systems, JCP, 2015

OPAL Results: Selection, Calibration, and Validation of Models of Tumor Growth

Model	Occam Category	Plausibility	$D_{KL} (\sigma_{16}^2)$	$D_{KL} (\sigma_{18}^2)$
RD01	1	n/a	0.92	0.70
PF01	2	n/a	0.65	0.74
RD02	3	0.22		

- General framework for model selection, calibration, validation in the presence of uncertainties in the data, model or parameters
- OPAL is an adaptive paradigm for resolving model inadequacy, for designing validation experiments, and for determining valid models for specified Qols

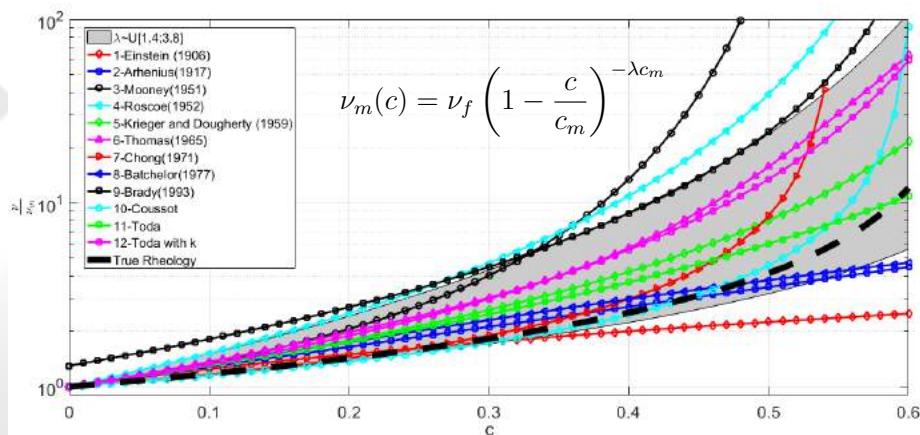
MD04	5	0.02		
PF03	5	0.34		
PF04	5	0.64	0.70	0.66
PF05	6	n/a		



(a) 32.73 mm^2 (data) (b) 32.45 mm^2 (PF04) (c) 37.92 mm^2 (data) (d) 41.96 mm^2 (PF04)

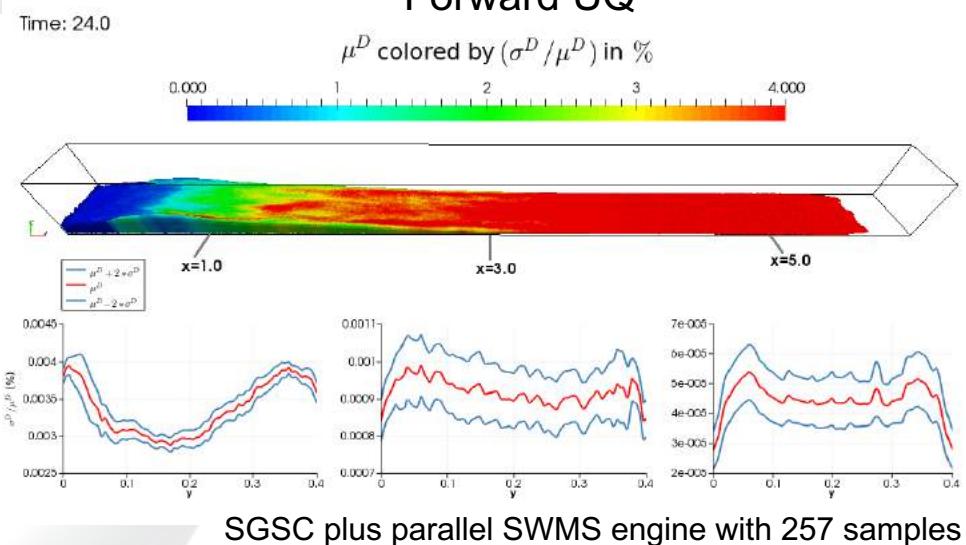
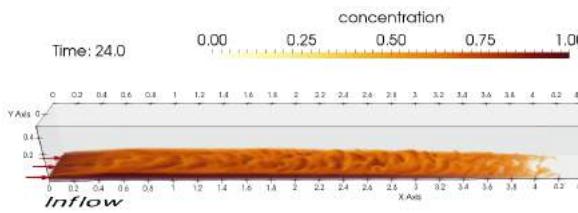
Uncertain Rheology of Non-Dilute Currents

Author	Equation
Einstein (1906) [25]	$\nu_m = \nu_f (1 + 2.5c)$
Mooney (1951) [26]	$\nu_m = \nu_f \left[\exp \left(\frac{2.5c}{1 - c/c_m} \right) \right]$
Krieger and Dougherty (1959) [7]	$\nu_m = \nu_f \left[1 - \frac{c}{c_m} \right]^{-2.5c_m}; c_m = 0.74$
Batchelor (1977) [27]	$\nu_m = \nu_f [1 + 2.5c + 6.2c^2]$
Brady (1993) [28]	$\nu_m = 1.3\nu_f \left[1 - \frac{c}{c_m} \right]^{-2.0}$
Toda and Hisamoto (2006) [29]	$\nu_m = \nu_f \left[\frac{1-0.5c}{(1-c)^3} \right]$
Toda and Hisamoto with k (2006) [29]	$\nu_m = \nu_f \left[\frac{1+0.5kc-c^2}{(1-kc)^2(1-c)} \right]; k = 1 + 0.6c$

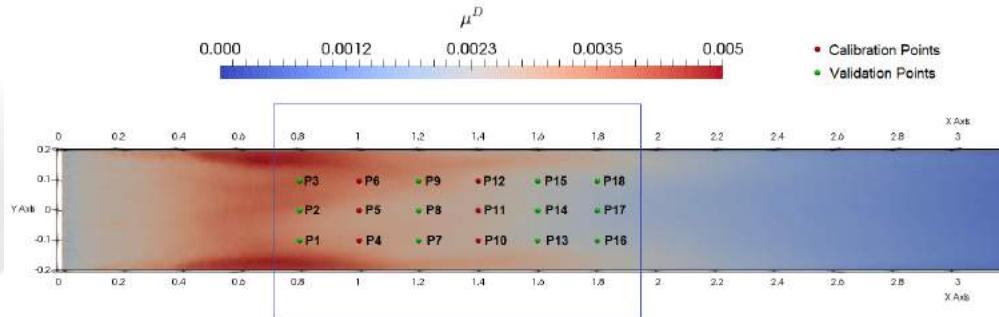


Computational Setup: closed channel with sustained current

- Channel dimensions, $xc = 6$, $yc = 0.4$, $zc = 0.5$, inlet windows $yw = 0.4$ $zw = 0.04$. Computational setup inspired on a experimental one (calibration and validation)
- Initial relative concentration = 0.11 (normalization constant)
- Reynolds number $Re = 1.5 \times 10^4$, used to allow the formation of turbulent structures. Transient flow features.
- No-slip and no-penetration in all walls with inflow velocity = 0.5



Uncertain Rheology of Non-Dilute Currents: Bayesian Calibration, Validation and Prediction



MCMC ran for 104521 samples with 50% burn-in, and for each step of the chain, 104522 samples for the Monte-Carlo integration.

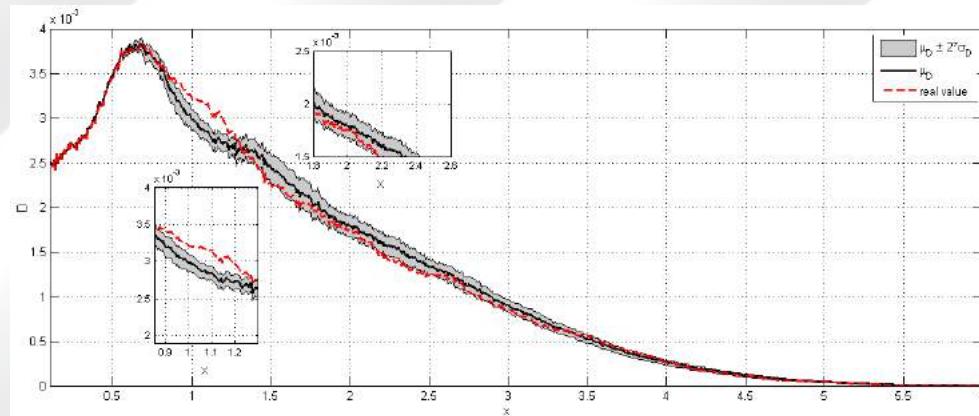
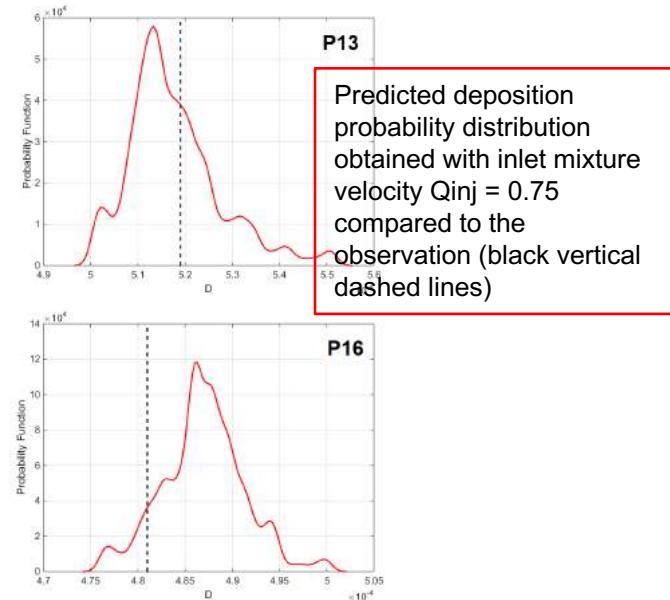
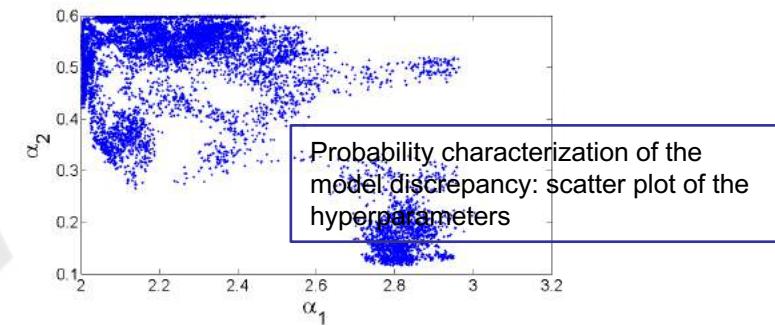


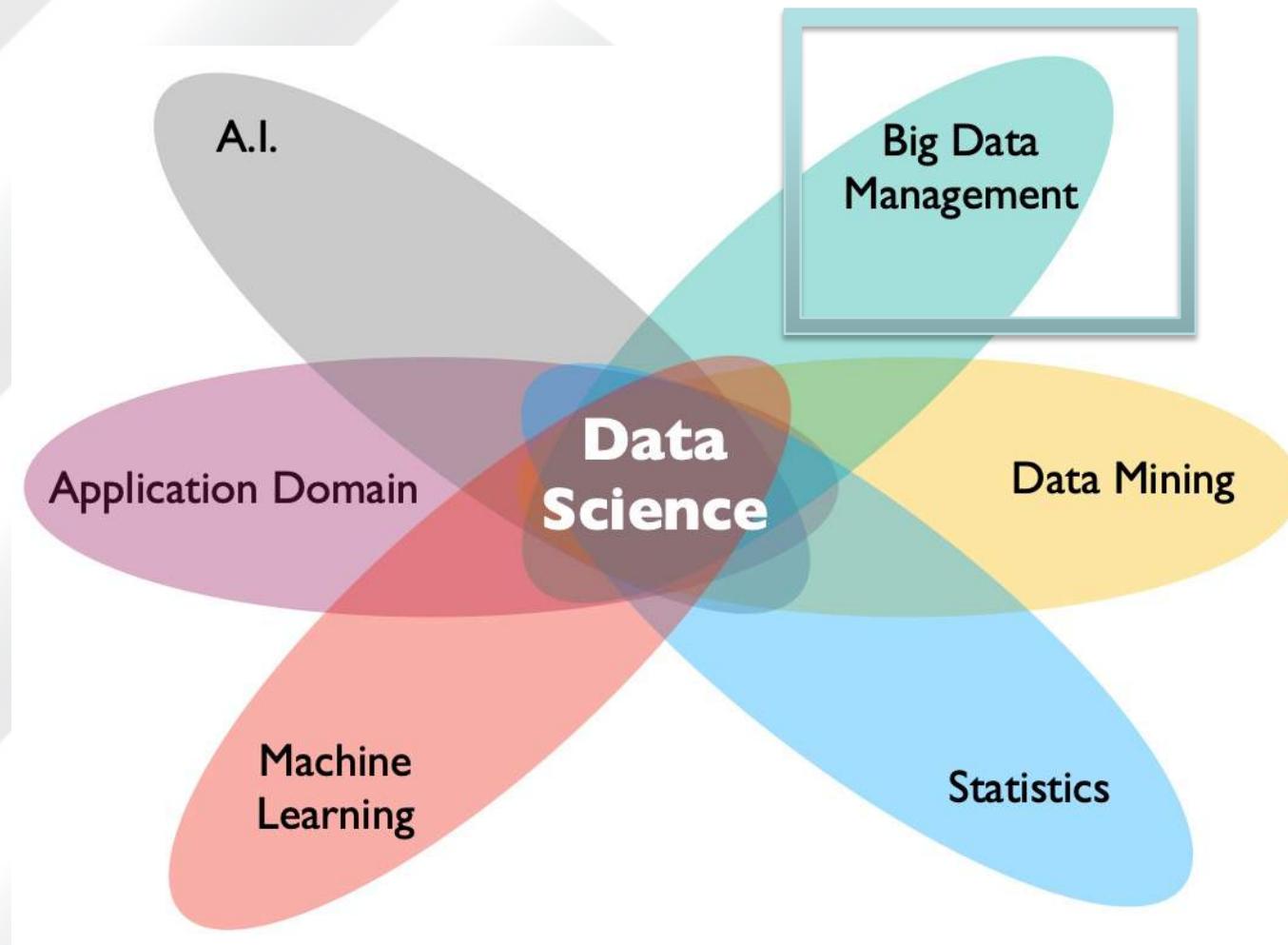
Figure 10: Predictive mean and confidence interval of deposition along the centerline in x direction compared to observation (red dashed line).



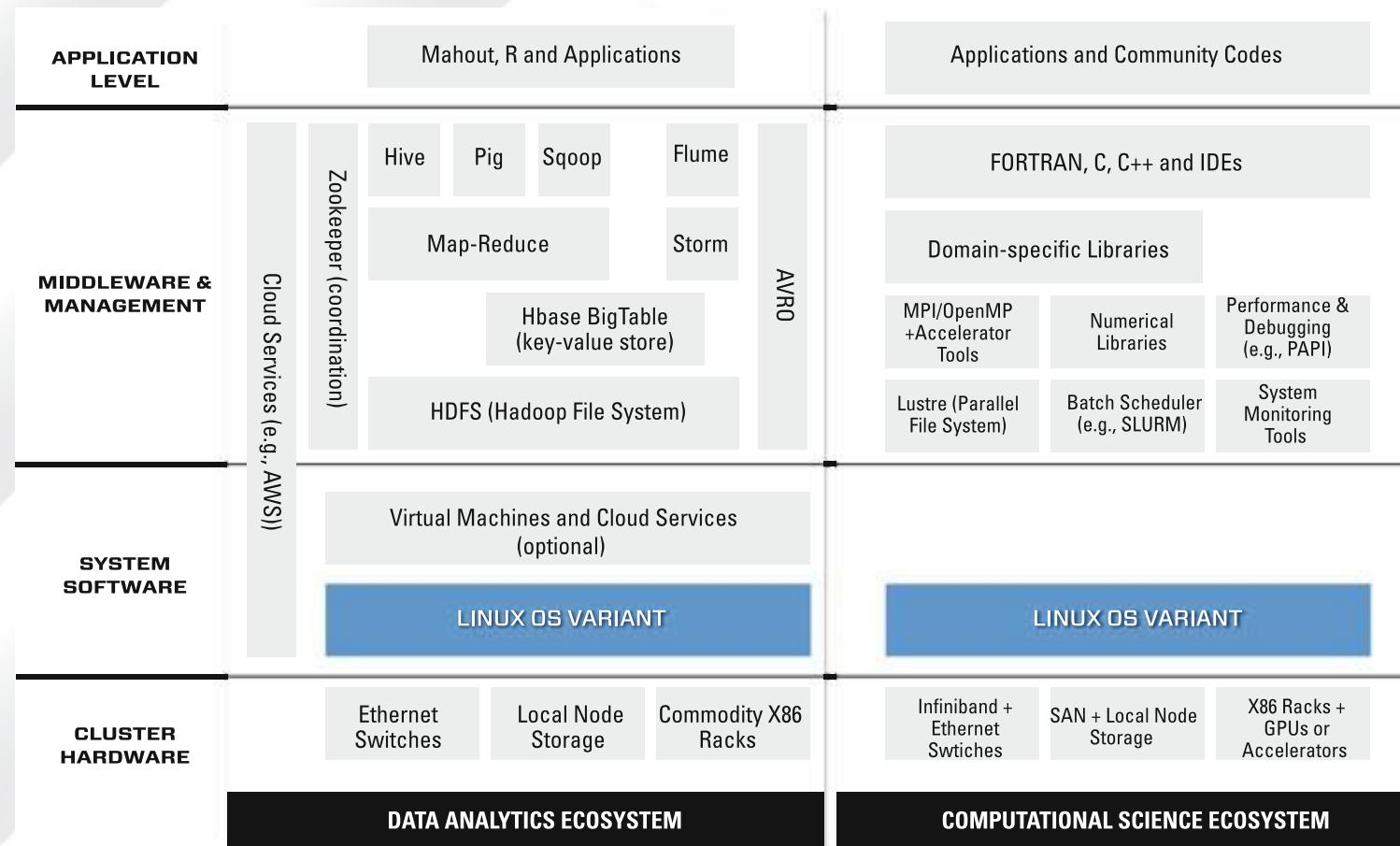
CHALLENGES IN DATA ANALYSIS



Data Science ≠ Machine Learning



Software Stack for HPC and Big Data¹



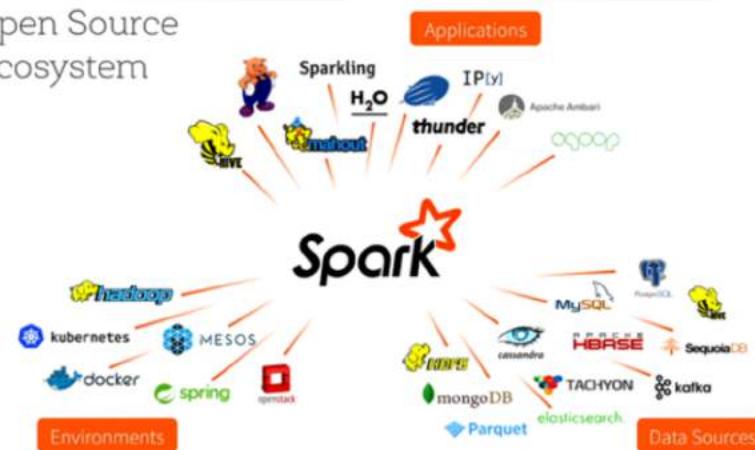
¹H. Anzt, J. Dongarra, M. Gates, J. Kurzak, P. Luszczek, S. Tomov and I. Yamazakio, Bringing High Performance Computing to Big Data Algorithms, A.Y. Zomaya and S. Sakr (eds.), Handbook of Big Data Technologies, 2017

Why Data Management differs

Business data

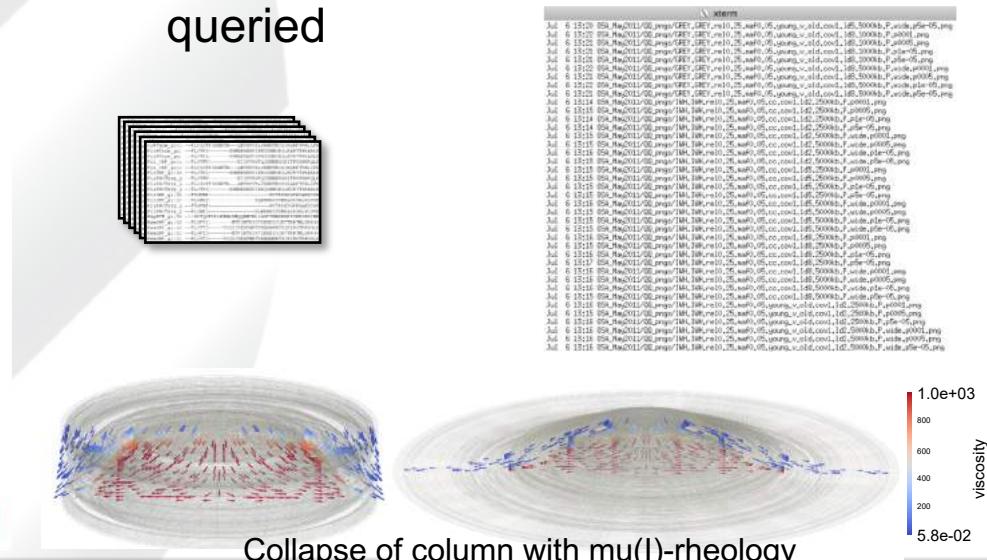
- ❑ easy to understand
- ❑ text format
- ❑ all manipulation in SQL
- ❑ most of data stored is traversed for queries

Open Source Ecosystem



Scientific data

- ❑ complex math / domain
- ❑ binary shared formats
- ❑ specific programs access
- ❑ only a small fraction of data is queried



How to relate data ? Manually? Annotation?

- Who produced this image?
- What were the programs used?
- Which were the input data?
- Which were the data transformations?
- Where are the data?

- Do we need a new model for each analysis?

Tracking Data Transformations

- ❑ data transformations – *ad-hoc*
- ❑ files generated independently
- ❑ parallel processing unaware of data-flow
- ❑ analysts need to manually manage the larger life cycle of big data flow analysis

```
# Output file
file = File("output.pvd", "compressed")
# Step in time
t = 0.0
T = 50*dt
prev = t3;
while (t < T):
    t += dt
# Solver execution
u0.vector()[:] = u.vector()
solver.solve(problem, u.vector())
# Visualization
file << (u.split()[0], t)
```

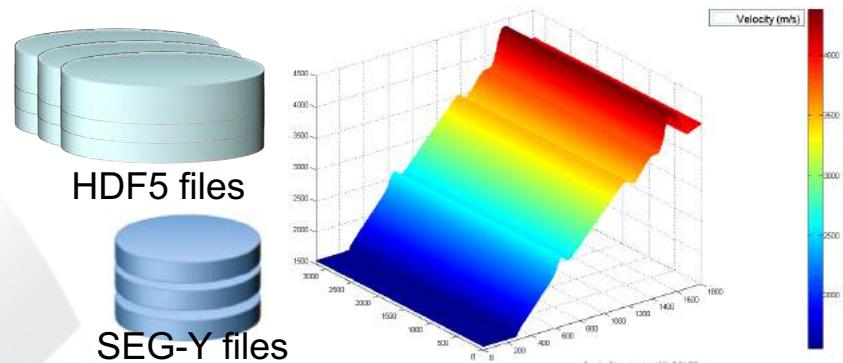
parameters



performance data



Log files



```
xterm
Jul 6 13:20 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld5,5000kb,P,wide,p5e-05.png
Jul 6 13:21 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,1000kb,P,p0001.png
Jul 6 13:21 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,1000kb,P,p0005.png
Jul 6 13:21 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,1000kb,P,p1e-05.png
Jul 6 13:21 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,1000kb,P,p5e-05.png
Jul 6 13:22 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,wide,p0001.png
Jul 6 13:22 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,wide,p0005.png
Jul 6 13:22 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,wide,p1e-05.png
Jul 6 13:22 OSA_May2011/00_pnrgs/GREY.GREY,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,wide,p5e-05.png
Jul 6 13:23 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,2500kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,2500kb,P,p0005.png
Jul 6 13:14 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,2500kb,P,p1e-05.png
Jul 6 13:14 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,2500kb,P,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,5000kb,P,wide,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,5000kb,P,wide,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,5000kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld2,5000kb,P,wide,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld5,2500kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld5,2500kb,P,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld5,2500kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld5,2500kb,P,wide,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,2500kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,2500kb,P,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,2500kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,2500kb,P,wide,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,5000kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,5000kb,P,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,5000kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,cc,cov1,ld8,5000kb,P,wide,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld2,2500kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld2,2500kb,P,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld2,2500kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld2,2500kb,P,wide,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,2500kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,2500kb,P,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,2500kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,2500kb,P,wide,p5e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,p0001.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,p0005.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,wide,p1e-05.png
Jul 6 13:15 OSA_May2011/00_pnrgs/IMH.IMH,re10,25,naf0,05,young_v_old,cov1,ld8,5000kb,P,wide,p5e-05.png
```

BLOG@CACM

Data Science Workflow: Overview and Challenges, by Philip Guo

NACAD

Tracking Data Transformations

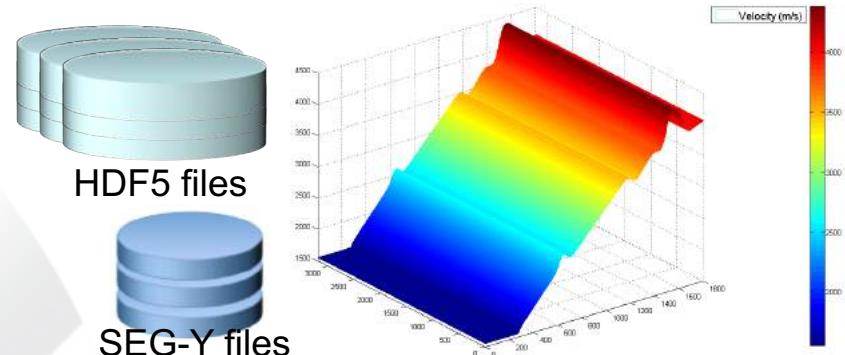
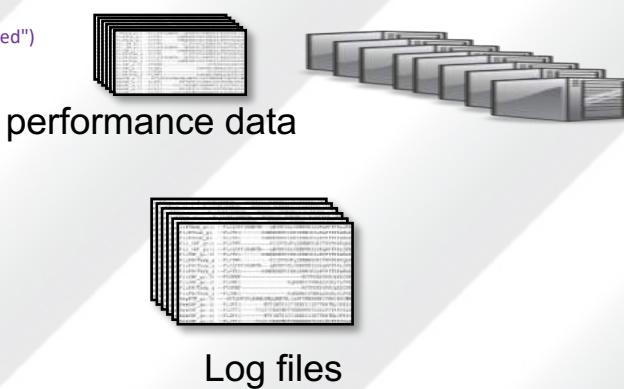
- ❑ **Steering is done manually**
 - Fine tuning configurations?
 - Debugging?
 - ❑ **Monitoring file directories**
 - Parsing files of intermediate data results
 - Trace back dataflow;
 - Data extraction through programming
 - Programming basic stats functions
 - Queries

```

# Output file
file = File("output.pvd", "compressed")
# Step in time
t = 0.0
T = 50*dt
prev = t3;
while (t < T):
    t += dt
# Solver execution
    u0.vector()[:] = u.vector()
    solver.solve(problem, u.vector())
# Visualization
    file << (u.split()[0], t)

```

parameters



Data Science Workflow: Overview and Challenges, by Philip Guo

NACAD

Analyzing Scientific Data Challenges

□ Preserving autonomy of scientific data

- Impossible to download all binary data to a DBMS
- Most data will not be queried, e.g. meshes
- Domain specific libraries (will not access a DBMS like SciDB)

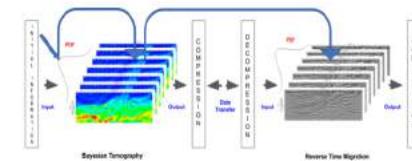


□ Data analysis at runtime to adapt the execution (HIL)

- Data needs to be extracted for analysis during the execution
 - *standard formats like HDF5, netCDF, ...*
- Complementing visualizations at runtime

□ Representing data relationships

- Files derivation paths
- Relating elements from data inside files



Putting the human in the loop

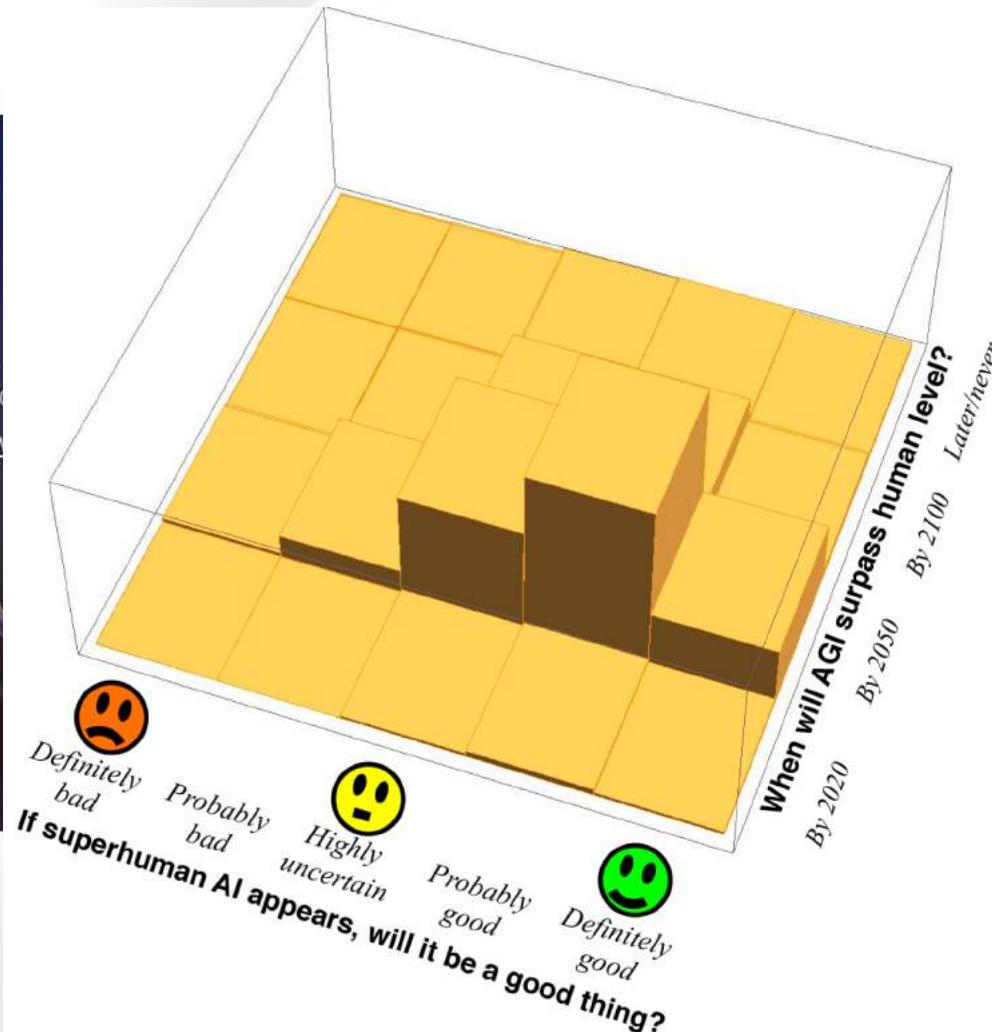
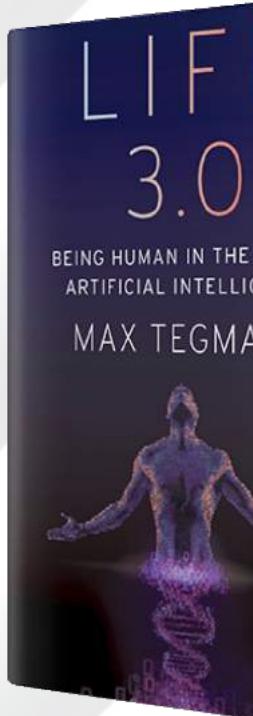
“In spite of the tremendous advances made in computational analysis, there remain many patterns that humans can easily detect but computer algorithms have a difficult time finding.”

Exploring the inherent technical challenges in realizing the potential of Big Data.

BY H.V. JAGADISH, JOHANNES GEHRKE, ALEXANDROS LABRINIDIS, YANNIS PAPAKONSTANTINOU, JIGNESH M. PATEL, RAGHU RAMAKRISHNAN, AND CYRUS SHAHABI

Big Data and Its Technical Challenges

What AI has to say about HIL?



idation aren't enough to
od control: ability for a
change its behavior if
ystems to work well, it's
ation be effective. In this
veniently alert you if
open."

rely building good user
figuring out how to
outer teams—for
ol should be transferred,
y to the highest-value
ntrollers with a flood of

Human In-the-Loop Challenges

Monitoring

- Define what data to monitor
- Provide data monitoring
- Give context (provenance) to data
- Show data at run time

Run Time Analysis

- Interact with data monitoring
- Define queries and visualization at real time

Fine-tuning

- Prepare execution to change
- Change the configuration at run time
- Keep data consistency

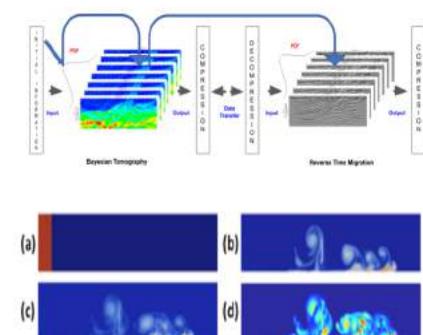
DfANALYZER

DfADAPTER

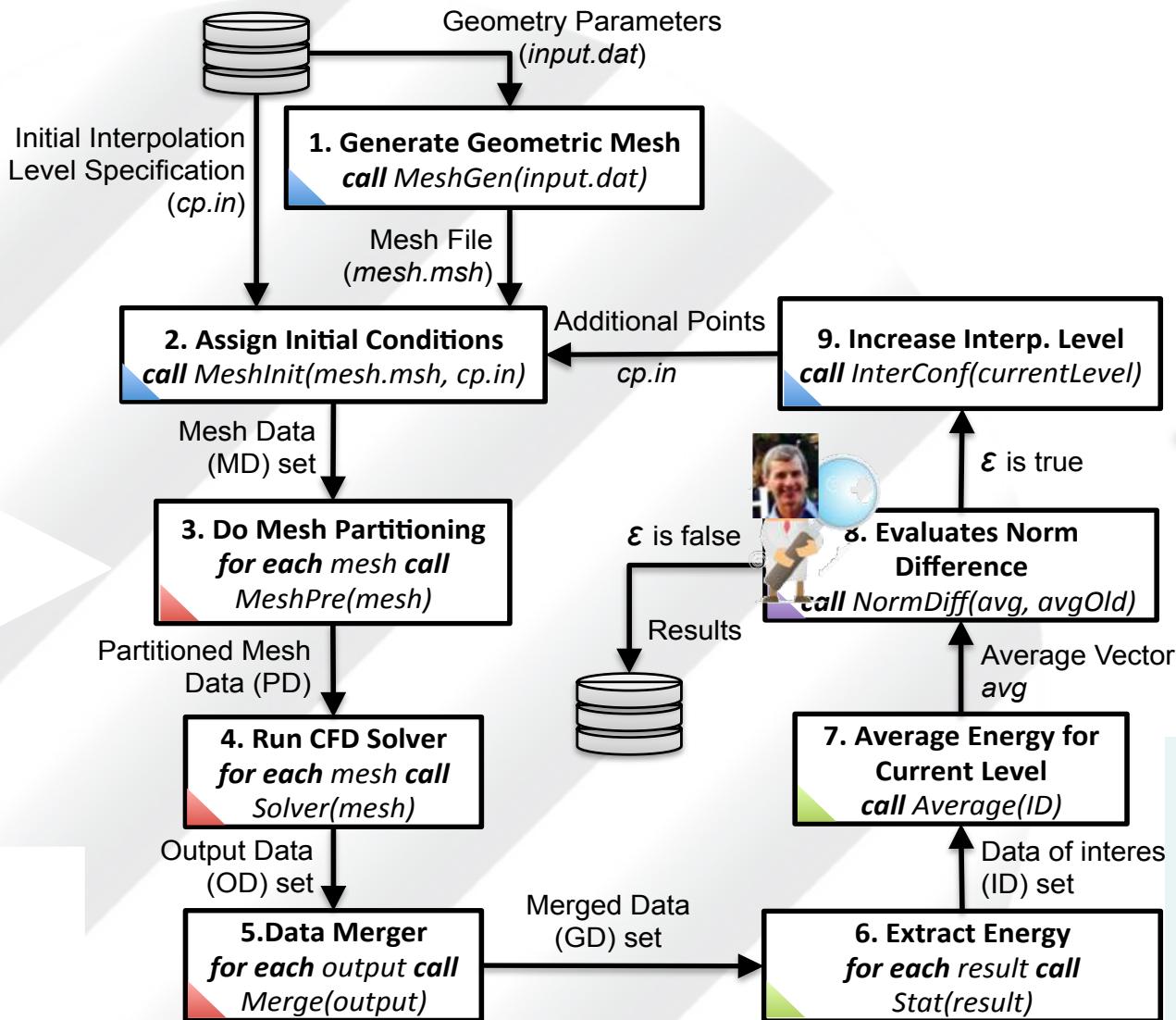


Analyzing Scientific Data with Big Data Techniques

- **Preserving autonomy of scientific data**
- **Extracting relevant data for analysis**
 - Extracting data in-situ is efficient
 - Relationships can be captured and registered at runtime
 - Can be loaded asynchronously in a DBMS (small data)
 - Add provenance data, metadata and performance data
 - Complements viz
- **Using standards to represent data relationships**
 - W3C PROV (provenance for data derivation)



UQ with User Steering & Adaptation



**Change
NormDiff
0.001 to 0.01**

- Compute the velocity at controlling points
- Estimate if the variance at those points are greater than some tolerance previously set
- Avoids aborting to adapt

PREDICTIVE DATA SCIENCE

What is Predictive Data Science?

Predictive data science is a convergence of the fields of Data Science and Computational Science & Engineering.

Predictive data science is needed for high-consequence applications across science, engineering and medicine, where machine learning approaches based on data alone are insufficient.

Data Science

Computational
Science &
Engineering

Predictive Data Science

From: K. Wilcox, Predictive Data Science, www.kiwi.oden.utexas.edu

Challenges

- **High-consequence applications are characterized by complex multiscale multiphysics dynamics**
- **High (and even infinite) dimensional parameters data are relatively sparse and expensive to acquire**
- **Uncertainty quantification in model inference and certified predictions in regimes beyond training data**

Ingredients

- **A physics-based model: conservation laws, that is, conservation of mass, momentum, energy , species, etc.**
- **A projection to define a low-dimensional model**
- **Variable transformations that expose structures in the model**
- **Non-intrusive learning of the reduced model**

See Wilcox at https://youtu.be/fnT__yycb34

The Fundamental Idea: Reduced Order Model

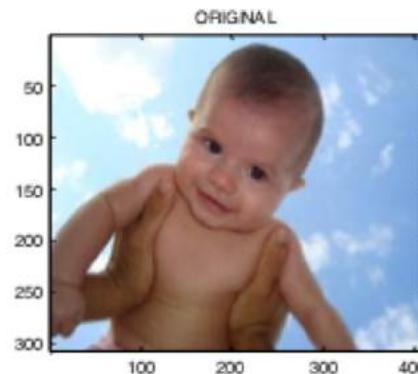
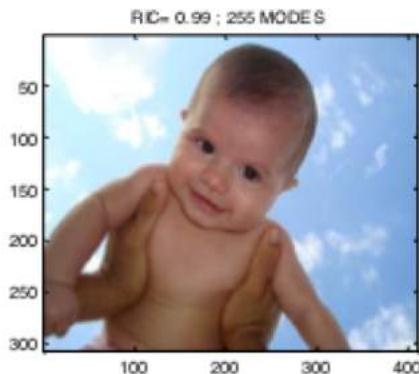
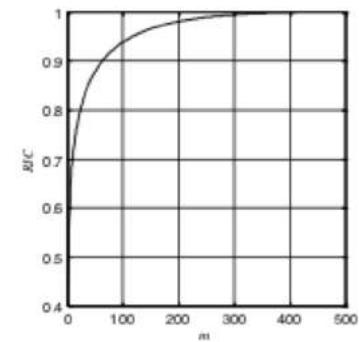
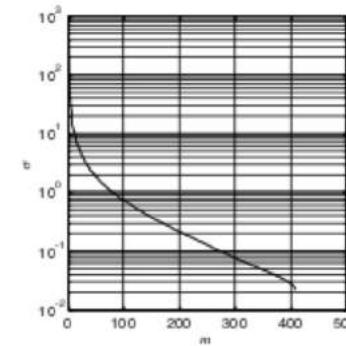
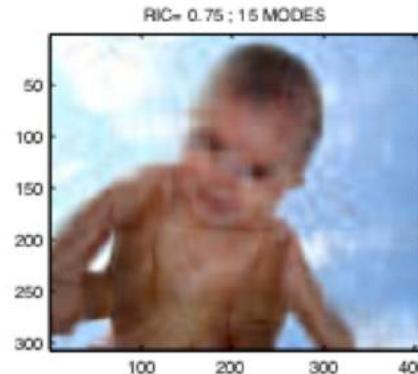
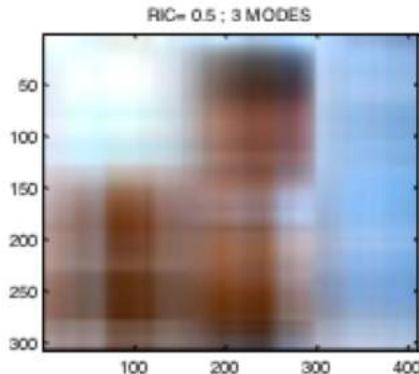


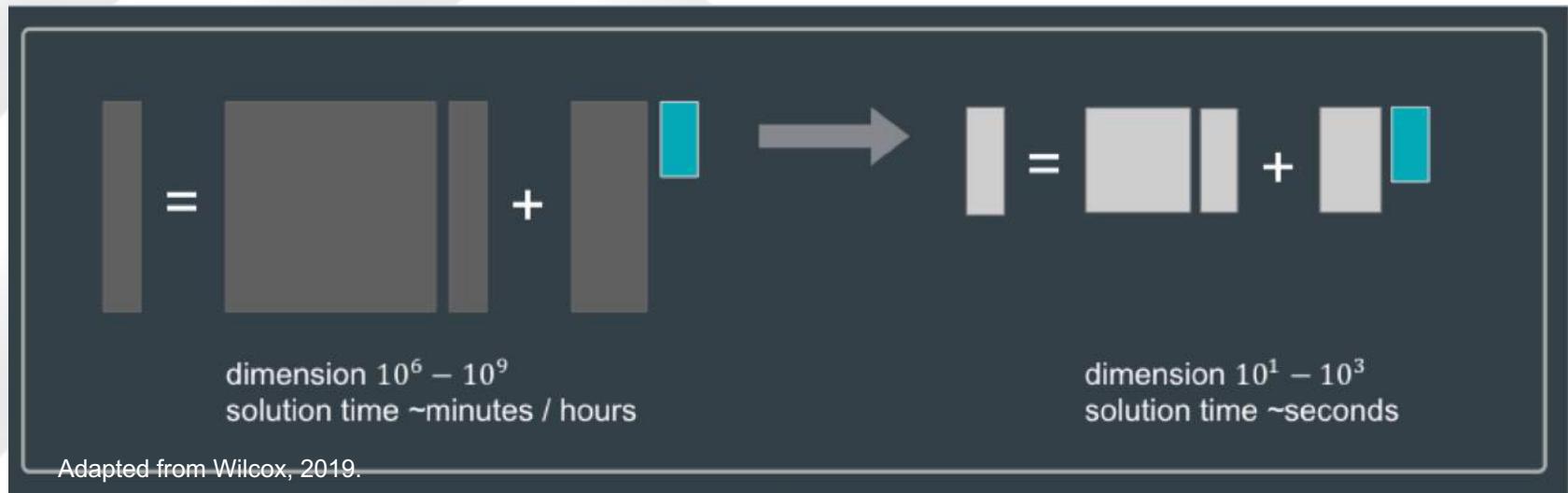
Figure 2.2 – Representation of a little beautiful girl image with reduced set of data.

Figure 2.1 – Singular values (σ) and Relative Information Content (RIC) for an image.

The use of SVD for image processing is an example commonly explored in basic courses of reduced order modeling. Considering an image stored in the RGB format, each pixel has values of Red, Green and Blue scales, which are going to be the data to be reduced.

$$\begin{aligned}
 M &= U \Sigma V^* \\
 m \times n &\quad m \times m \quad m \times n \quad n \times n \\
 U &\quad U^* = I_m \\
 V &\quad V^* = I_n
 \end{aligned}$$

Projection-based Reduced Order Modeling



1. Solve PDE to generate training data: snapshots
2. Get the structure: generate a low dimensional basis from snapshots
3. Training: project PDE into low-dimensional space

ROM for a CFD Problem: The backward facing step

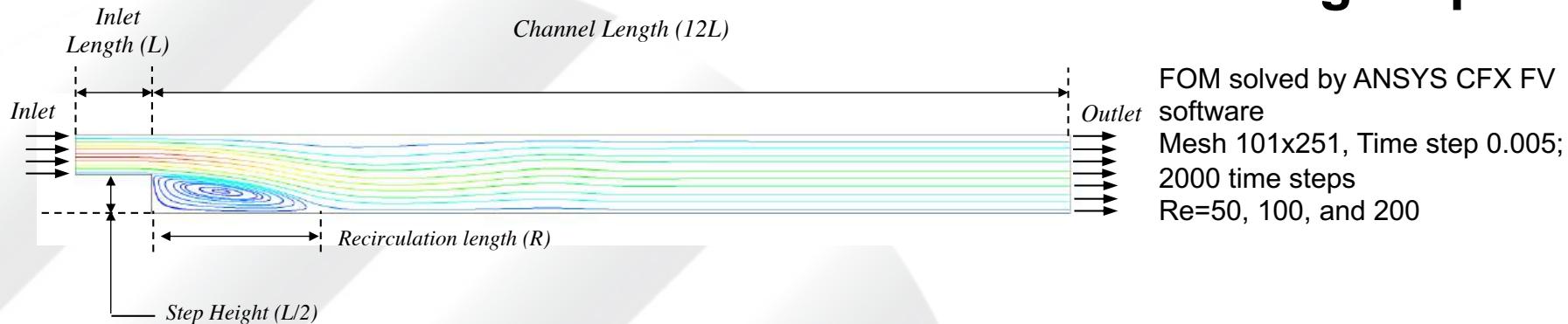


Figure 3.10 – Backward facing step domain.

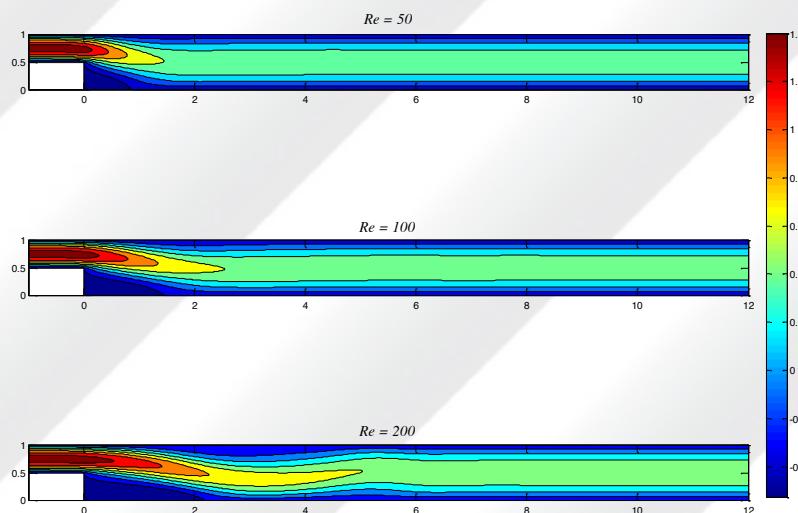


Figure 3.11 –CFX simulations for different Reynolds numbers.

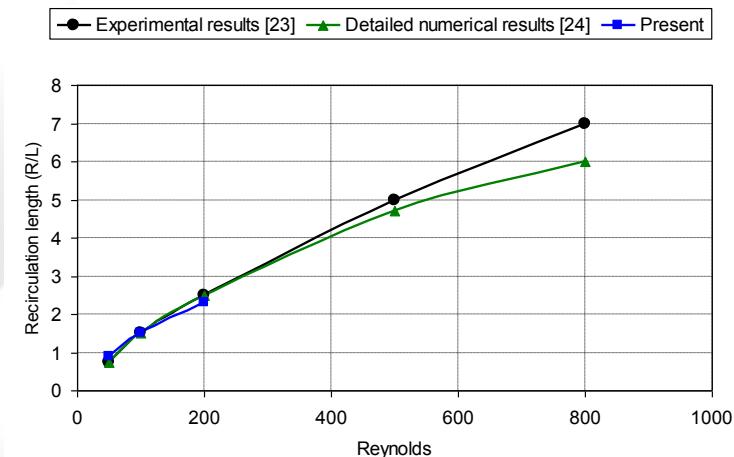


Figure 3.12 – Results from the present CFD model in comparison with results from the literature.
Recirculation length versus Reynolds number.

ROM for a CFD Problem: The backward facing step

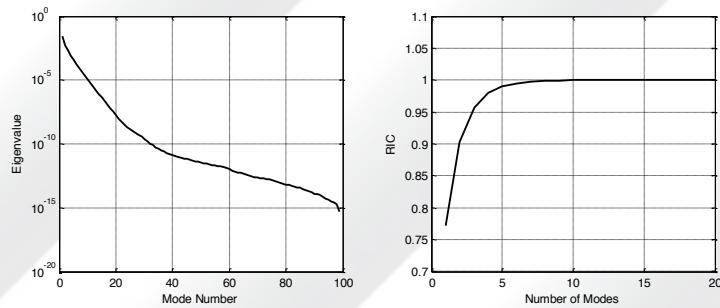
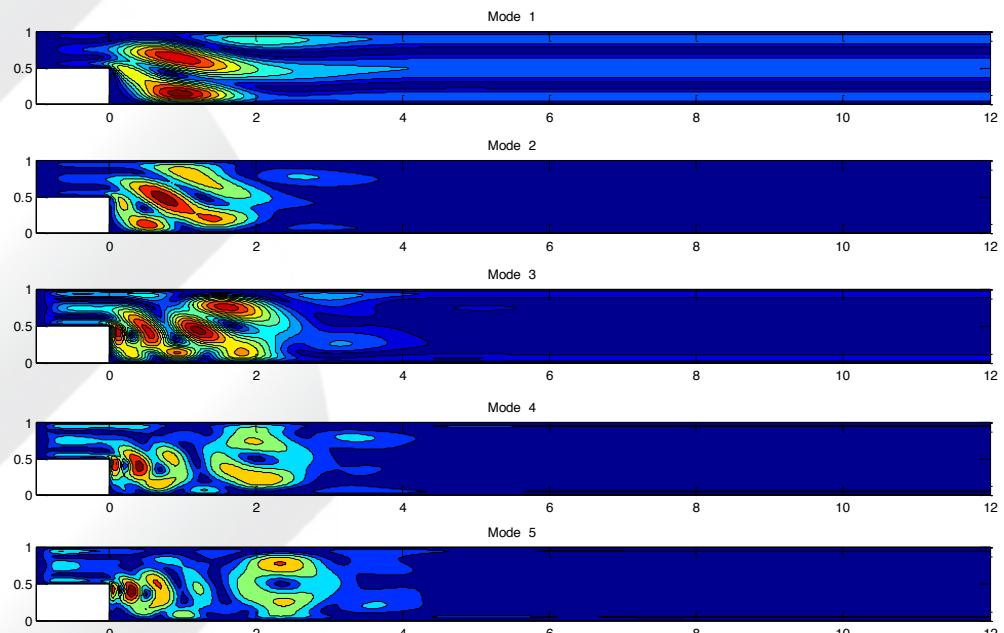


Figure 3.14 – POD modes eigenvalues and RIC for $Re = 100$.



ROM for a CFD Problem: The backward facing step

Table 3.3 – Summary of the errors obtained with ROMs. RMS errors of velocity components and execution time.

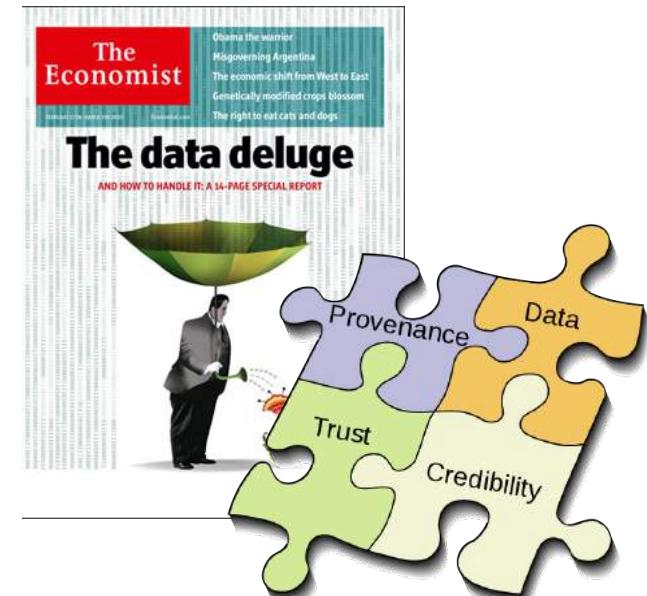
Reynolds ROM	Reynolds Snapshots	U RMS Error	V RMS Error	FOM/ROM Execution Time
50	50	0.0040	0.0011	303
	100	0.0320	0.0114	293
	200	0.0623	0.0226	297
100	50	0.0348	0.0123	488
	100	0.0048	0.0023	468
	200	0.0375	0.0155	475
	50,200	0.0176	0.0070	594
200	50	0.0837	0.0291	446
	100	0.0535	0.0172	435
	200	0.0251	0.0217	432

PROVENANCE IN DATA ANALYTICS

Why Use Provenance ?

Motivation in general

- Determines trust on results
- Ensures reliability, quality of data
- Repeatability/verifiability
- Avoids effort duplication



Why Use W3C PROV ?

Definition

“Provenance is information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form assessments about its quality, reliability or trustworthiness.”
(W3C PROV)

Domain independent

Relates data

Registers data relationships

DFAnalyzer:

https://gitlab.com/ssvitor/dataflow_analyzer/#prototype-of-multiphysics-application

PROVENANCE IN DATA ANALYTICS FE SOLVER WITH FENICS





```

dataflow_tag = "fenics-df"
# Create mesh
mesh = UnitSquareMesh(96, 96)
# Define function spaces
V = FiniteElement("Lagrange", mesh.ufl_cell(), 1)
ME = FunctionSpace(mesh, V*V)
# parts of code were omitted
# ...
# Define Newton solver
solver = NewtonSolver()
solver.parameters["linear_solver"]      = "gmres"
solver.parameters["convergence_criterion"] =
"incremental"
solver.parameters["relative_tolerance"]   = 1e-6

# Output file
file = File("output.pvd", "compressed")
# Step in time
t = 0.0; T = 50*dt; i = 0
prev = t3
while (t < T):
    t += dt; i += 1
# Solver execution
u0.vector()[:] = u.vector()
iter_count, converged_flag = solver.solve(problem,
u.vector())
# Visualization
file << (u.split())[0], t

```

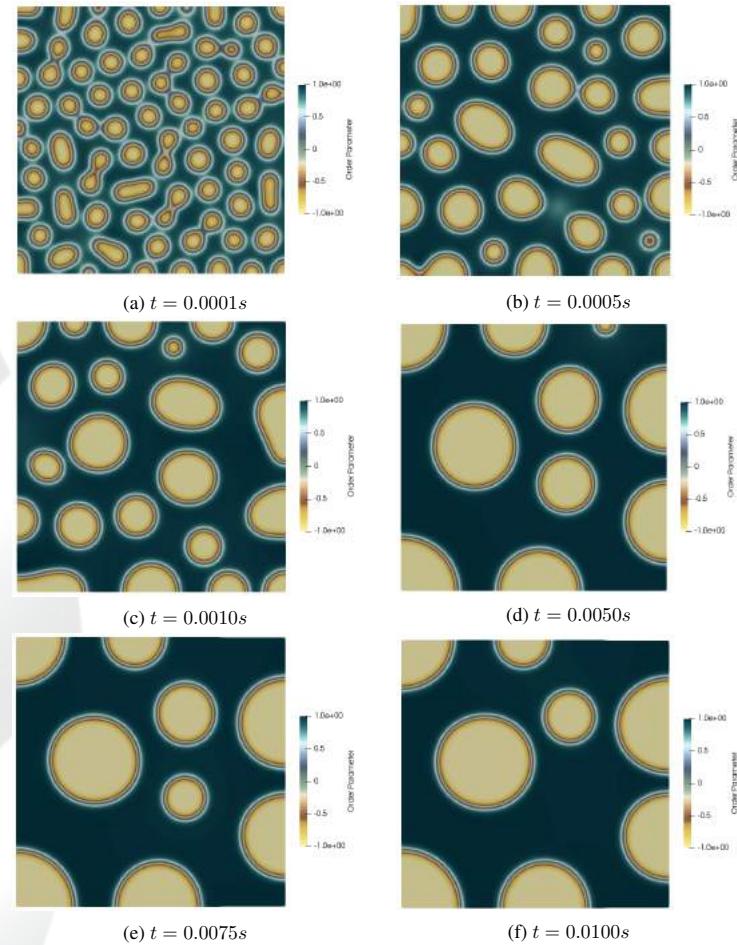


Figure 3.5: Spinodal decomposition using fixed time step and fixed mesh at different times.

Solving a 4h-order PDE governing phase separation: the Cahn-Hilliard equation

```

dataflow_tag = "fenics-df"
t1 = Task(1, dataflow_tag, "MeshCreation")
t1.add_dataset(DataSet("iMeshCreation", [Element([96 , 96])]))
# Create mesh
mesh = UnitSquareMesh(96, 96)
t1.add_dataset(DataSet("oMeshCreation",
[Element([mesh.num_vertices(), mesh.num_cells()])]))
t1.end()
t2 = Task(2, dataflow_tag, "FunctionSpace", dependency=t1)
t2.add_dataset(DataSet("iFunctionSpace",
[Element(["Lagrange", 1])]))
# Define function spaces
V = FiniteElement("Lagrange", mesh.ufl_cell(), 1)
ME = FunctionSpace(mesh, V*V)
t2.add_dataset(DataSet("oFunctionSpace",
[Element(ME.dim())]))
t2.end()
# parts of code were omitted
# ...
t3 = Task(3, dataflow_tag, "NewtonSolver", dependency=t2)
t3.add_dataset(DataSet("iNewtonSolver", [Element(["lu",
"incremental", 1e-6])]))
# Define Newton solver
solver = NewtonSolver()
solver.parameters["linear_solver"]      = "gmres"
solver.parameters["convergence_criterion"] = "incremental"
solver.parameters["relative_tolerance"]   = 1e-6
t3.add_dataset(DataSet("oNewtonSolver", [Element(["gmres",
"incremental", 1e-6])]))
t3.end()
...
# Output file
file = File("output.pvd", "compressed")
# Step in time
t = 0.0; T = 50*dt; i = 0
prev = t3
while (t < T):
    t += dt; i += 1
    current = Task(int(t3._id)+i ,dataflow_tag,"TimeStep",
                   dependency=prev)
    current.add_dataset(DataSet("iTimestep",
                               [Element([t,dt])]))
# Solver execution
u0.vector()[:] = u.vector()
iter_count, converged_flag = solver.solve(problem,
u.vector())
current.add_dataset(DataSet("oTimeStep", [Element
([converged_flag,iter_count,solver.residual()])]))
current.end()
twrite = Task(int(current._id)+1, dataflow_tag,
              "Visualization"+iter_count,
              dependency=current)
twrite.add_dataset(DataSet("iVisualization",
                           [Element(["output.pvd"])])
# Visualization
file << (u.split()[0], t)
# Raw data extraction
extracted_data =
Extractor(ExtractorCartridge.PROGRAM, "output.pvd")
twrite.add_dataset(DataSet("oVisualization",
                           [Element(extracted_data[i-1])]))
twrite.end()

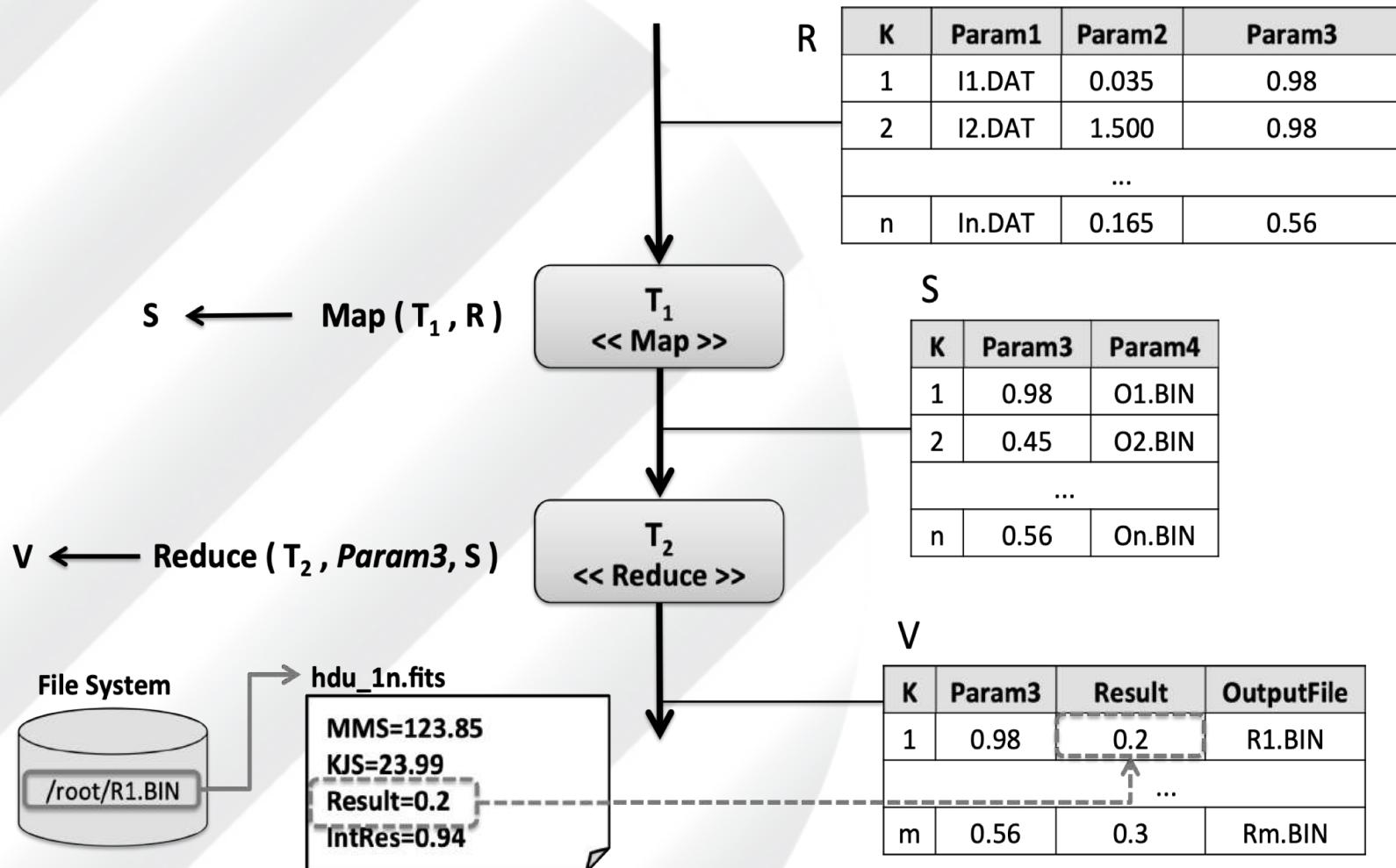
```

Finite element solution of Cahn-Hilliard equation with FEniCS – user steering

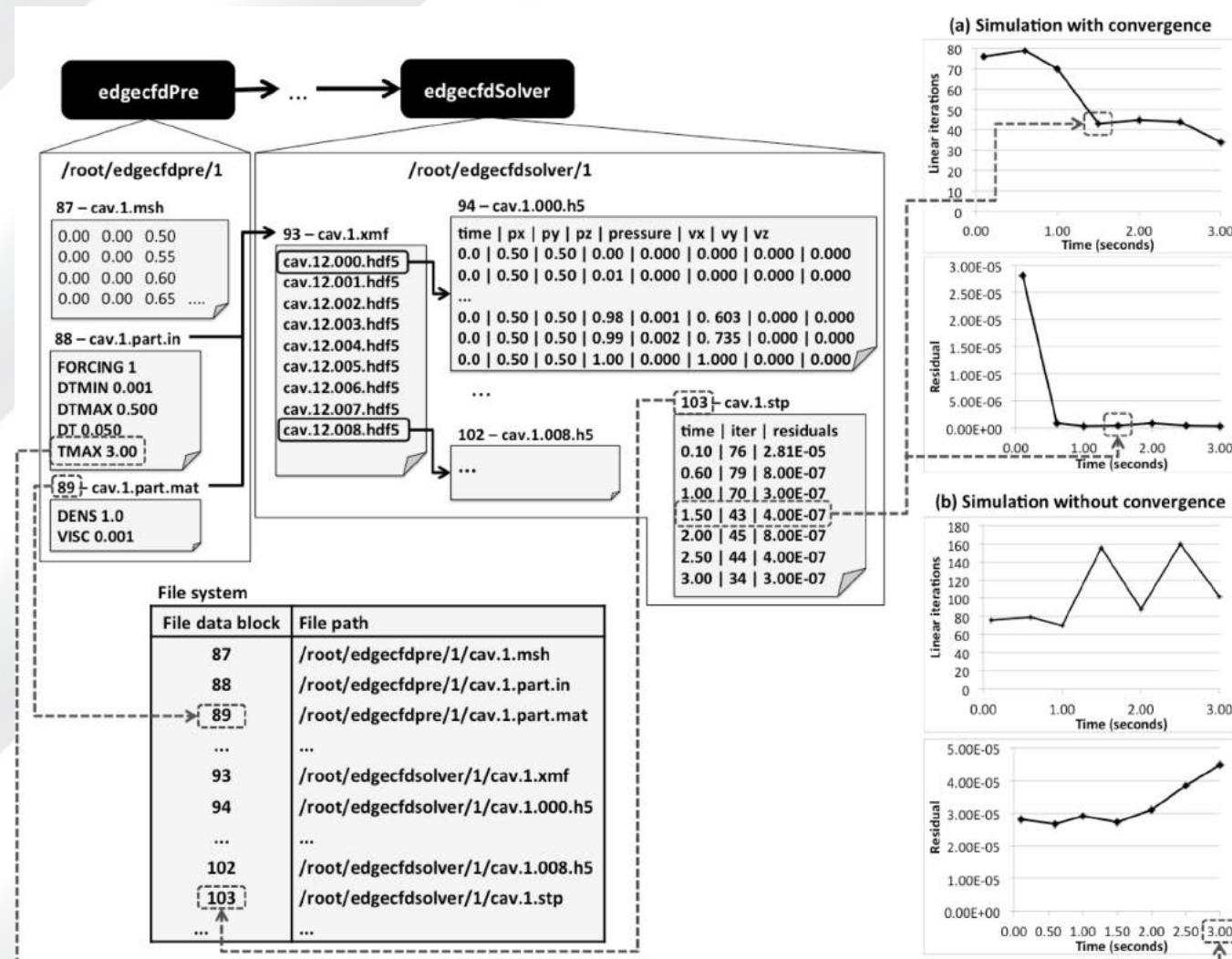
time step	Flow		Transport		Visualization
	Linear residual	Nonlinear residual	Linear residual	Nonlinear residual	
1000	0.0168372	0.0053668	0.000226965	0.0000000	/viz/img_1000.png
2000	0.0643741	1.75431e-06	0.000109237	0.00186711	/viz/img_2000.png
3000	0.1270304	0.0027324	0.000230409	3.265e-06	/viz/img_3000.png
	...				

Query results for numerical analysis with FEniCS to detect possible misbehavior of nonlinear and linear solvers

Data Extraction and relationships in general



Tools for Simulation Data Analysis: Metadata Relates Multiple Files

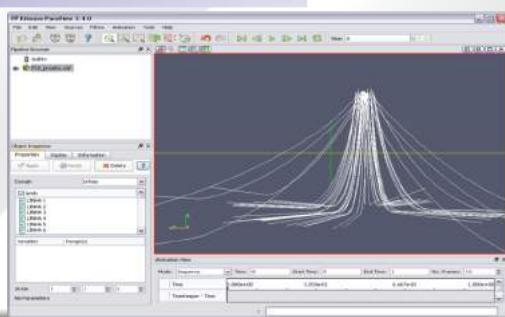


PROVENANCE IN DATA ANALYTICS RISERS FATIGUE ANALYSIS

Risers' Fatigue Analysis : Data Exploration



**Stop criteria,
tolerance**



Parameter Sweep



Estimate risers lifetime

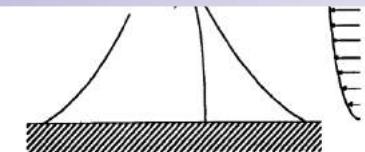
3. Results are analyzed POSFAL to evaluate critical regions

Input Data to simulate Environment conditions:
Waves, wind, currents, bathymetry, etc.

1. Coupled movement Analysis (TPN or Prosim)

**Different
Methodology**

Linha de Acesso
Modelo



Generates large amount of data ...
(finite element meshes,

**Solver
options**



**Iterative
Methods**

2. ... to do Structural Analysis of Risers (ANFLEX)



NACAD

DFAnalyzer:

https://gitlab.com/ssvitor/dataflow_analyzer

PROVENANCE IN DATA ANALYTICS LIBMESH LIBRARY

Turbidity currents

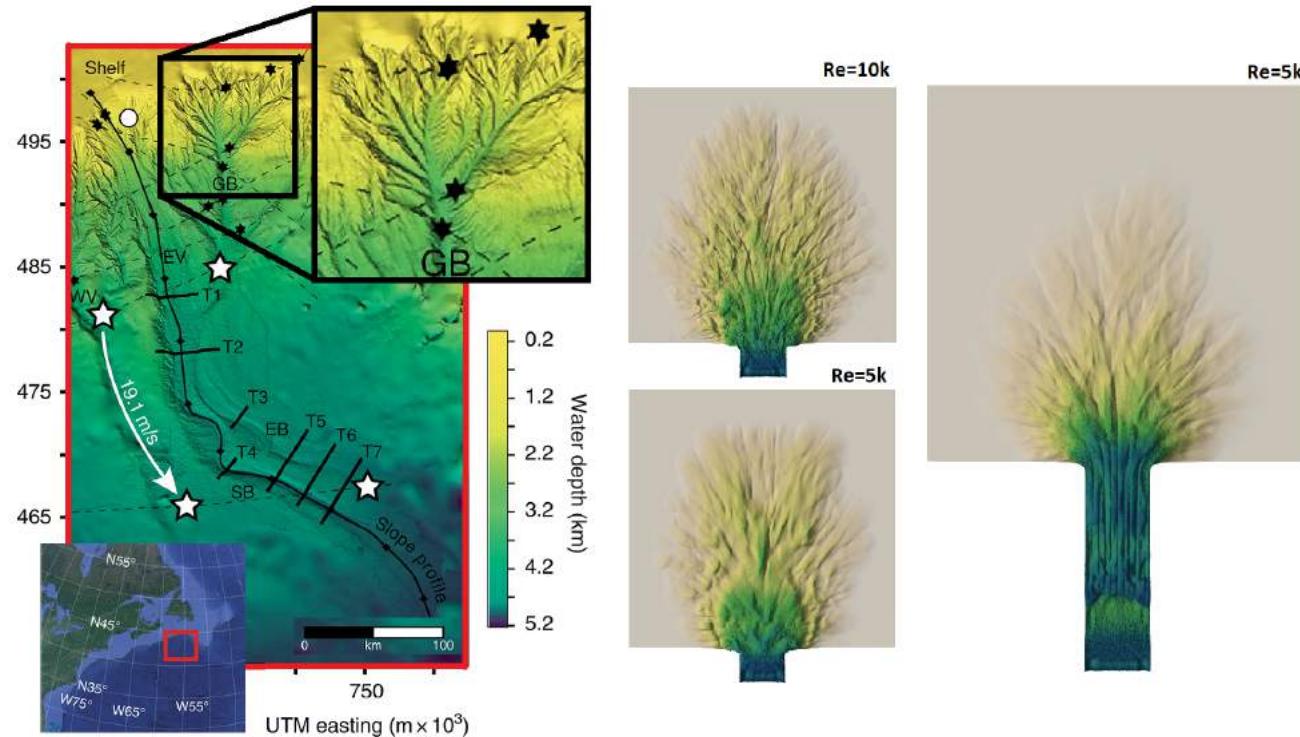
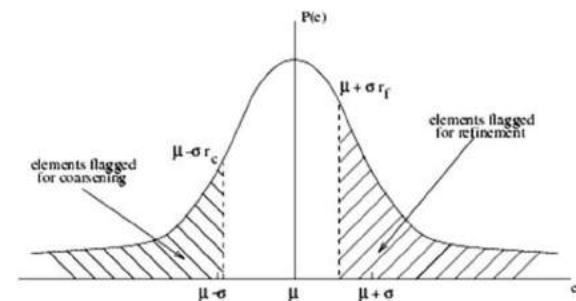
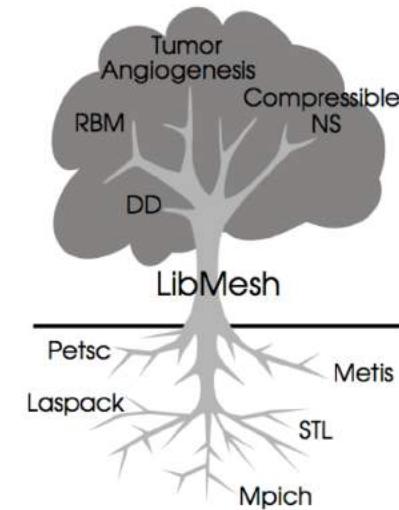


Figure: **Left:** Reconstruction of the 1929 Giant turbidity current deposits; **Top insect:** sediment deposit detail¹; **Right:** Sediment deposits from a turbidity current simulation at different Re (COPPE, 2017)

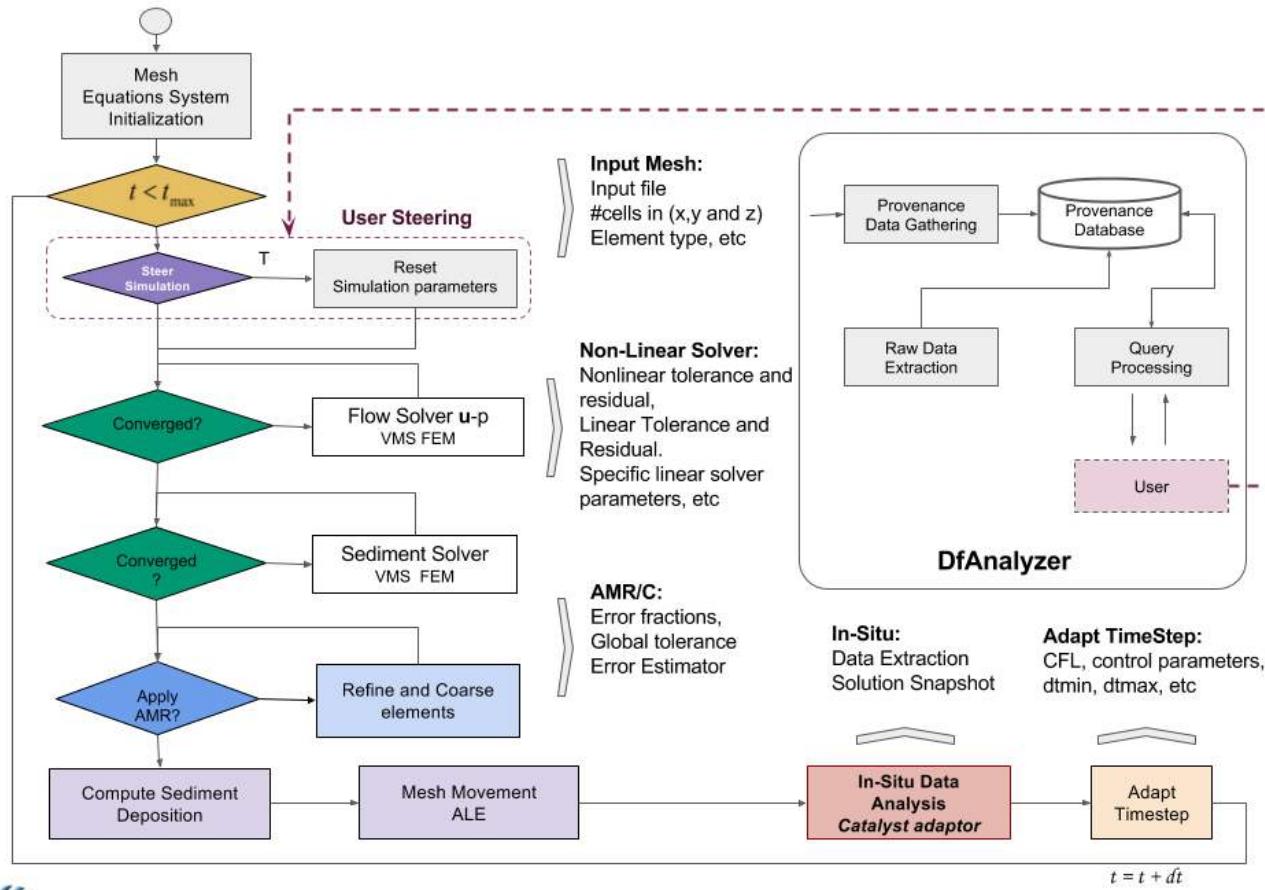
libMesh: A Framework for Finite Element Analysis

- ▶ Open-source library with parallel adaptive mesh refinement and coarsening (AMR/C) support
 - ▶ AMR/C is an optimal strategy for large-scale simulations
 - ▶ libMesh supports h , p and $h-p$ adaptive strategies
- ▶ AMR/C based on a statistical scheme:
 - ▶ the error e (flux jump at the interfaces) have a normal probability density function $P(e)$ with mean μ and standard deviation σ
 - ▶ Development initiated at CFDLab, UT Austin, see Kirk et al, Eng. Comp, 2006



<http://libmesh.github.io/index.html>

Exploring sedimentation solver data with provenance data



Computational Setup

- **Experiments were carried out on *Lobo Carneiro* machine at NACAD/COPPE/UFRJ¹**
 - SGI ICE-X 252 computer nodes
 - Each node with two Intel Xeon E5-2670v3 (Haswell)
 - 64GB of memory per node
 - Network: Infiniband FDR - 56Gb/s (Hypercube)
- **libMesh-Sedimentation solver:**
 - Intel Compilers (version 16.0) with -O3 optimization flag
 - MPI Intel
 - Linked with Paraview Catalyst (v. 5.3) / Offscreen Rendering (Mesa 13.0 version)
 - Linked with Dataflow Analyzer (DfAnalyzer)
- **Run configuration:**
 - Standard nodes for libMesh and Paraview Catalyst
 - A dedicated node for simulation data management
 - *Includes database component from DfAnalyzer*

¹<http://www.nacad.ufrj.br/en/recursos/sgiicex>

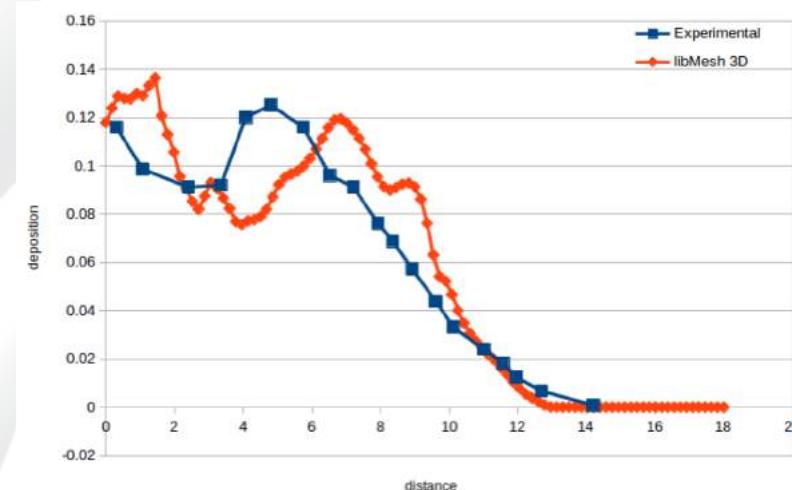
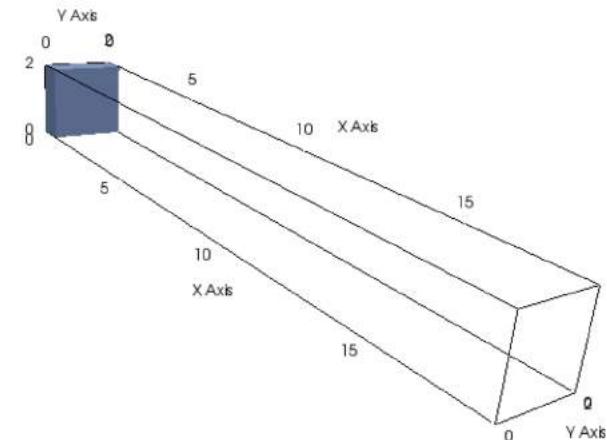
Test Case: de Rooij and Dalziel Sedimentation Tank

- **Simulation Setup**

- Domain size: $20.0 \times 2.0 \times 2.0$
- Initial uniform grid in all directions, grid space 0.075
- One initial uniform refinement, solution with max refine level=3
 - *1.5M HEX8 elements*
- Kelley's error estimator for \mathbf{u} and c
- The lock, in which the fluid initially is at rest, has dimensions $0.75 \times 2.0 \times 2.0$
- Reynolds number $Re=5000$
- Run with 480 cores at Lobo Carneiro

- **Simulation Management**

- Raw data files are written each 50 time steps
- In-situ data extraction are called each 50 time steps
- In-situ visualization generates pngs files each 50 timesteps
- Catalyst data extraction: plot over line filter
- Catalyst viz: slice filter



In-Situ Catalyst data extraction
plot over line filter

Sedimentation Tank:

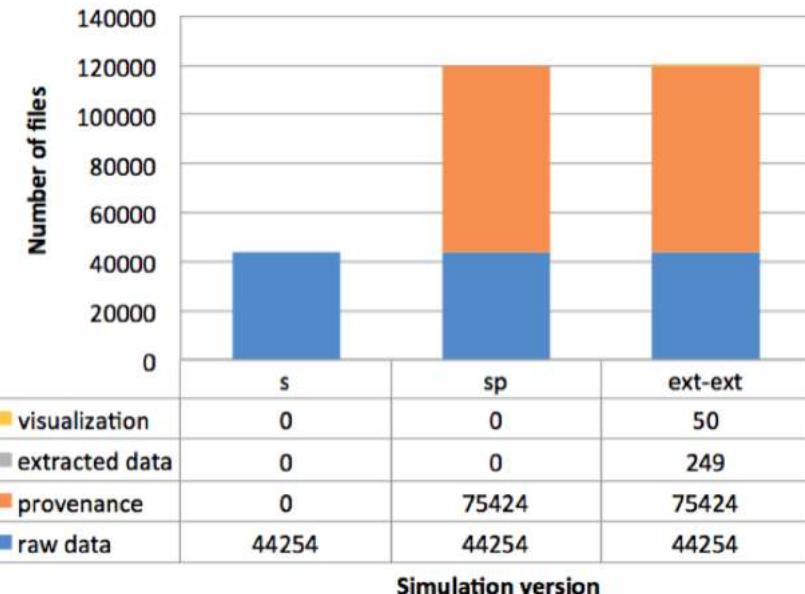
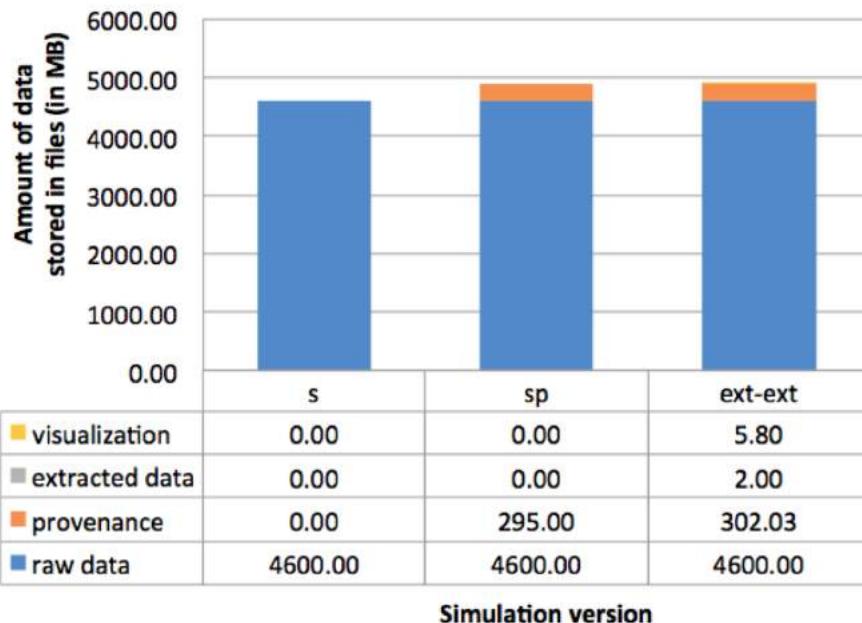
Computation Cost - CPU Time

TABLE 1: Elapsed time for different stages on libMesh-Sedimentation - Sedimentation Tank

Time Contribution	CPU Time (in s)	cost/call	%cost
Flow Solver	23533.21	1.26	37.85%
AMR/C	13122.20	93.73	21.11%
Sediment Solver	4941.66	0.27	7.95%
InSitu Catalyst Viz+Extraction	2065.08	22.45	3.32%
XDMF/HDF5 Raw Data	1329.21	14.45	2.14%
Provenance (DfAnalyzer)	83.38	0.01	0.13%
Others (libMesh)	17096.26		27.5%
Total	62171.00		

- **Remarks:**
 - Provenance adds low overhead to overall simulation costs.
 - In-Situ Viz + extraction cost in relation to Raw data Writer could be offset by disk bandwidth constraints as well as limited disk capacity.
 - CPU time spent in In-Situ Catalyst depends how many filters are applied.

Sedimentation Tank Computation Cost: Raw Data vs. In-Situ and Provenance Data Storage


Legend:

s:: Solver

sp: Solver with Provenance

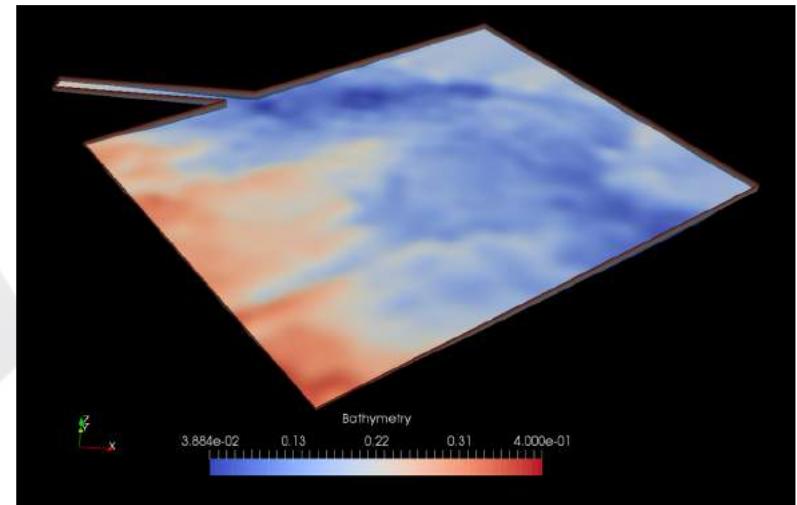
ext-ext: Solver + Provenance + In-situ

Provenance Data: 6%
In-Situ Visualization : 0.12%
Data Extraction: 0.04%
 % in relation of raw data stored

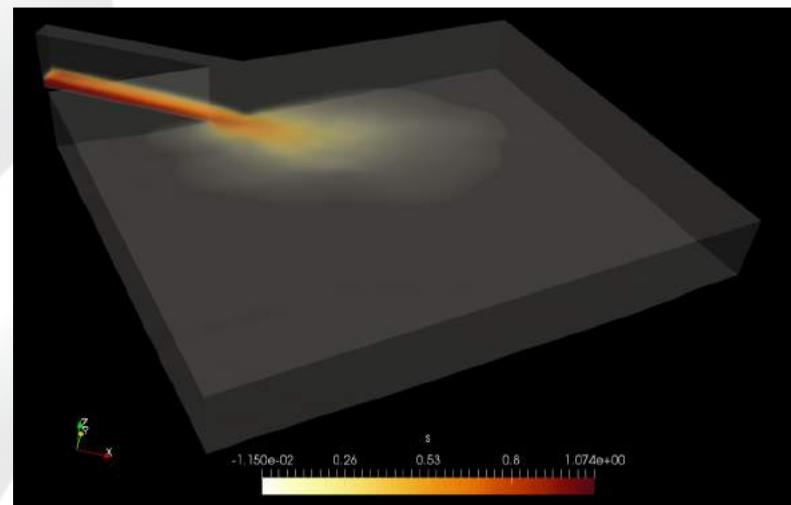
Real Case:

Sedimentation and Deposition with Real Bed Bathymetry

- Run on 480 cores at Lobo Carneiro
- Domain Size:
 - Tank: 14 x 12 x 2
- Fixed Unstructured mesh:
 - 7.6M linear tetrahedral elements
 - 1.4M nodes
- Dimensionless Parameters:
 - Grashof: 10^6 (Reynolds approx. 2000)
 - Monodisperse:
 - Settling velocity: $5.6651E-03$
- Boundaries conditions:
 - Flow:
 - no slip is applied at bottom and channel walls
 - Prescribed velocity at front wall channel
 - Free slip at top
 - Sediment:
 - Concentration prescribed at front channel wall
- Data Analysis
 - In-Situ visualization of sediment deposition
 - Data Extraction: sediments profile over lines



Bed bathymetry



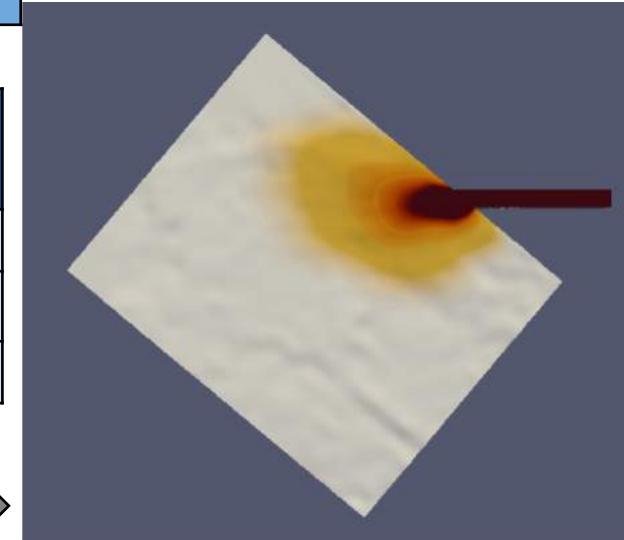
Real Bed Bathymetry Case:

Computation Cost

TABLE 2: Elapsed time for different stages - Real Bed Bathymetry Case - Simulation final time t=100

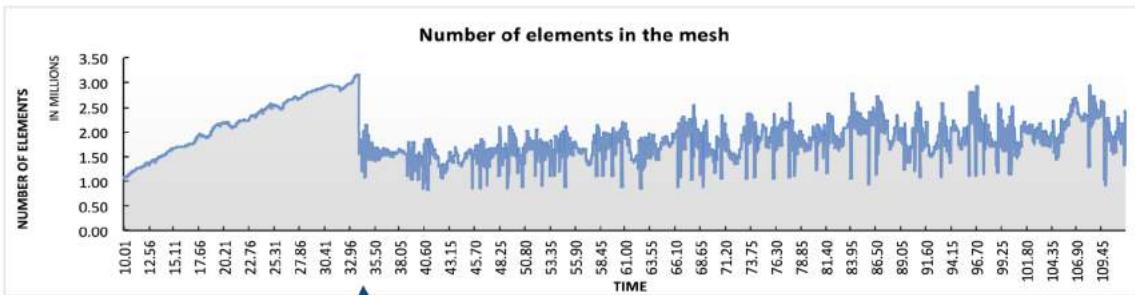
Time Contribution	Elapsed Time	%Cost
Flow Solver	72523.49	50.71%
Concentration Solver	28000.50	19.58%
In-Situ Viz + Extraction	2175.16	1.52%
XDMF/HDF5 Raw data	421.23	0.29%
Provenance	451.70	0.32%
Total	143029.00	

Storage Requirements	Size (in GB)	% Raw Data
XDMF/HDF5 Raw data	23.44	
Provenance + Data Extraction	0.38	1.60%
In-Situ Viz + Extraction	0.28	1.21%

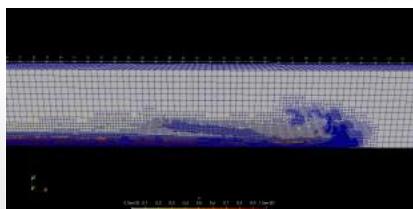


Sedimentation map generated by In-Situ Catalyst →

Parameter Online Tracking and Fine-tuning



Parameters fine-tuning at $t = 33.52$
time



$t = 24.5$

Tune

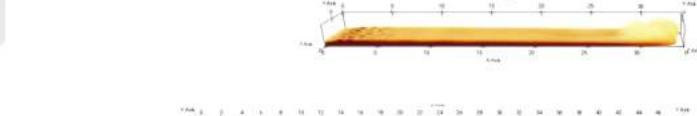
$t = 33.52$

Flow Final linear residual: 1.68641e-13
Flow Final nonlinear residual: 2.03266e-08
Sediment Final linear residual: 1.49932e-13
Sediment Final nonlinear residual: 1.51098e-08
Number of elements in the mesh: 2463183

Parameter	Before adapt	After adapt
Flow nonlinear tolerance	1.0e-4	1.0e-3
Transport nonlinear tolerance	1.0e-4	1.0e-3
Flow initial linear solver tolerance	1.0e-6	1.0e-1
Transport initial linear solver tolerance	1.0e-6	1.0e-1

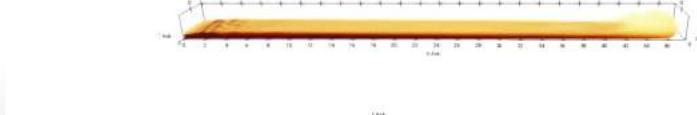
$t = 58.0$

Flow Final linear residual: 4.46403e-14
Flow Final nonlinear residual: 7.10668e-09
Sediment Final linear residual: 4.06513e-14
Sediment Final nonlinear residual: 9.09405e-09
Number of elements in the mesh: 1743485



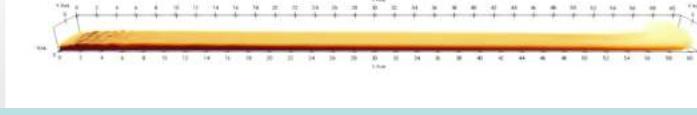
$t = 83.0$

Flow Final linear residual: 6.09749e-14
Flow Final nonlinear residual: 9.20302e-09
Sediment Final linear residual: 4.75423e-14
Sediment Final nonlinear residual: 7.4835e-09
Number of elements in the mesh: 1700729

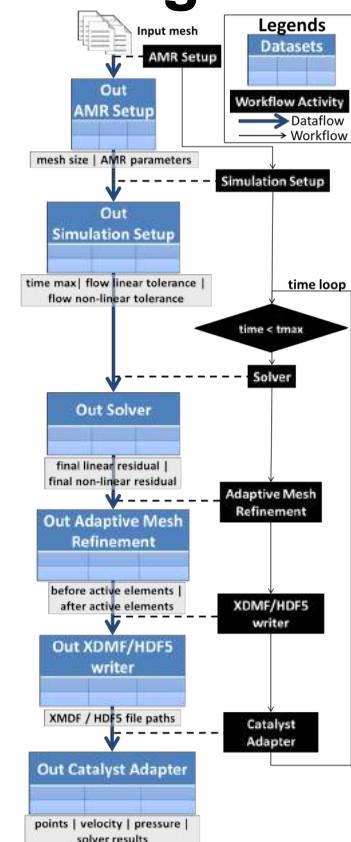


$t = 108.0$

Flow Final linear residual: 5.78688e-14
Flow Final nonlinear residual: 8.11125e-09
Sediment Final linear residual: 6.5936e-14
Sediment Final nonlinear residual: 6.40164e-09
Number of elements in the mesh: 2335823



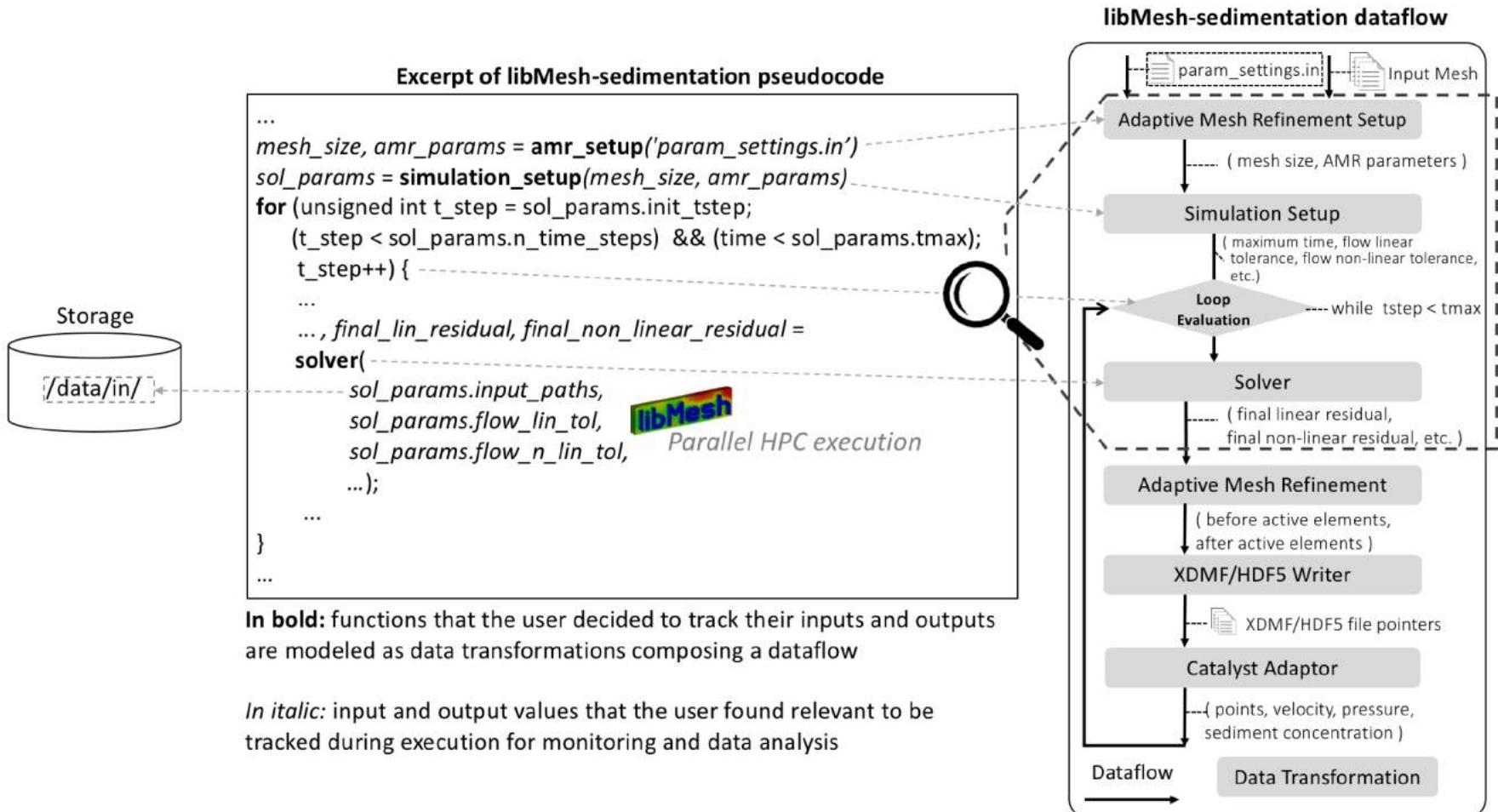
In situ visualization with ParaView Catalyst



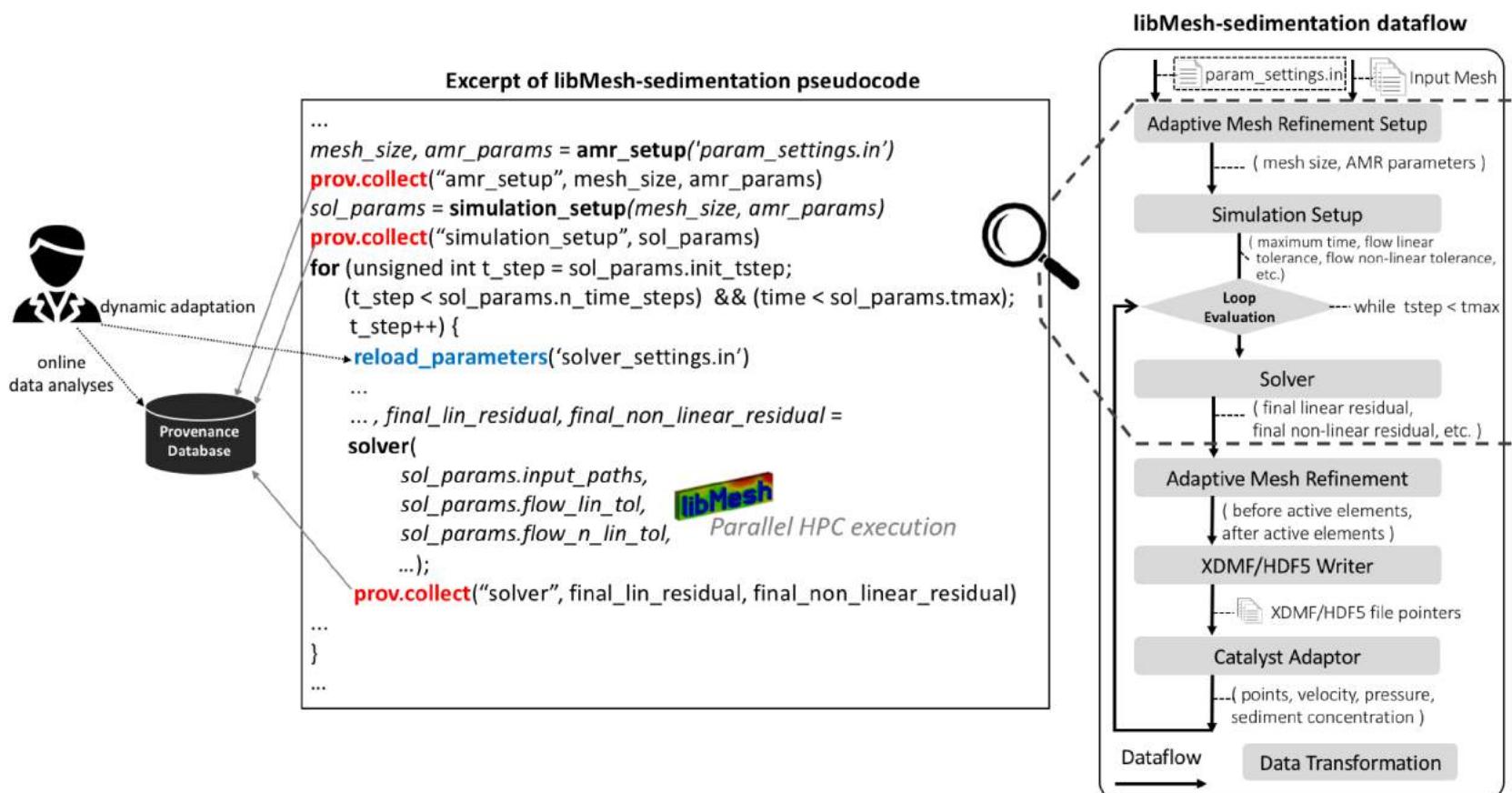
libMesh-sedimentation workflow

	Before	After	Reduction
Avg. Solver Time by iteration	3.82 min	2.21 min	42.14%
Avg. Number of elements	2.4×10^6	1.7×10^6	29.24%
Simulation time (expected)	~ 27 days	~ 17 days	37%

Exploring and adapting sedimentation solver data with provenance data



Exploring and adapting sedimentation solver data with provenance data



Sedimentation provenance data analysis complementing viz tools

- ▶ DfAnalyzer registers deposition along time at predefined locations and pointers to viz files.
- ▶ We can query online with a negligible time (< 500 ms).

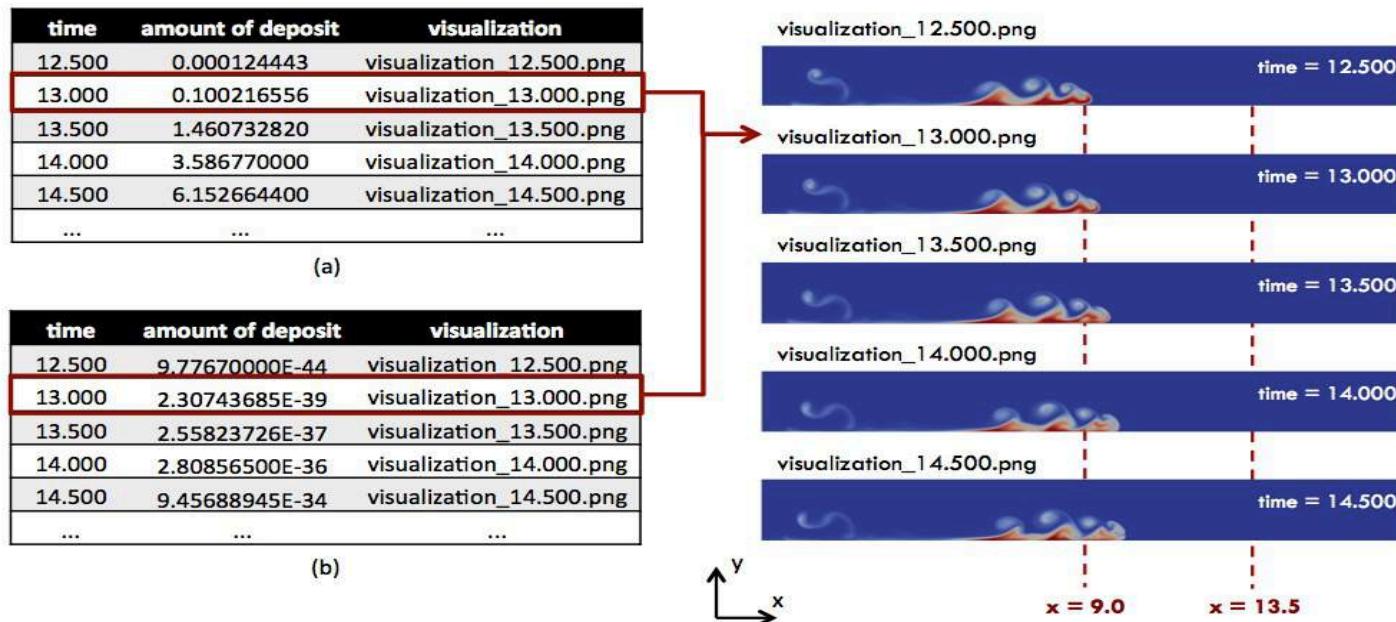


Figure: Sediment deposition monitoring at five time instants at $x = 9.0$ (a) and $x = 13.5$ (b) combining data with in-situ visual information

Query processing

- Analytical queries for...

- **Monitoring:**

- *The appearance of sediments in the domain bottom layer for a specific time step*

- **Debugging and user steering:**

- *Analysis of the algorithm output parameters after the convergence of the solver in the fluid and sediments loops in a specific execution of the sedimentation solver*

time step	x	y	z	d
1	0.000	1.000	0.000	2.00E-04
1	0.180	1.000	0.000	2.00E-04
1	0.360	1.000	0.000	2.00E-04
1	0.540	1.000	0.000	2.00E-04
1	0.720	1.000	0.000	1.99E-04
1	0.900	1.000	0.000	1.19E-04
1	1.080	1.000	0.000	3.04E-08
...

Fluid		Sediments	
linear residual norm	nonlinear residual norm	linear residual norm	nonlinear residual norm
3.98E-06	13.54823207	8.66E-06	0.004445721
4.30E-06	0.390224835	2.00E-09	0.002435432
4.30E-06	0.390224835	1.31E-05	0.016144017
7.00E-09	0.002712742	2.00E-09	0.002435432
7.00E-09	0.002712742	1.31E-05	0.016144017
...

Queries users can now ask

Inspecting parameter tunings (“*who*”, “*when*”, “*what*”, “*which*”)

- How many tunings?
- Which parameters did I change?
- What were the old and new values?
- When did each adaptation happen?

Understanding consequences of a tuning (“*how*”)

- In parameter tuning 3, how was the average of Solver output values 10 iterations before and after?

Data reduction (“*how*”, “*which*”)

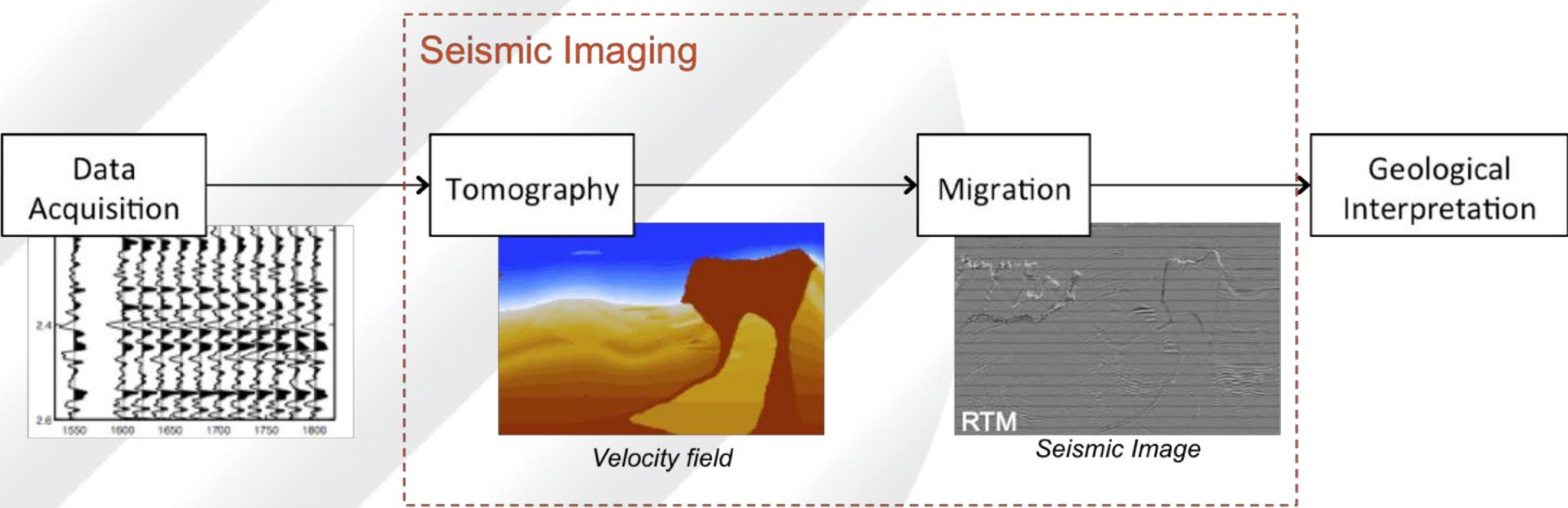
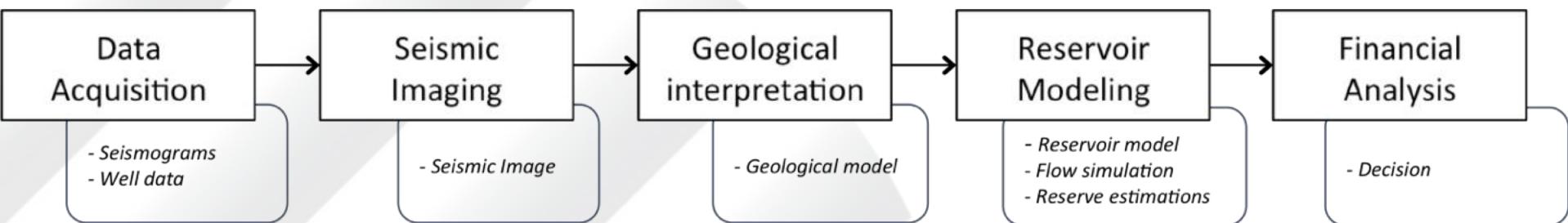
- On average, how long iterations were lasting before and after I reduced input files from the input data? Which files were affected?

DFAnalyzer:

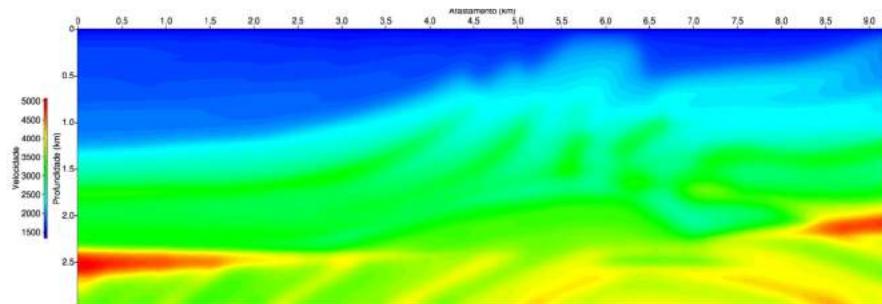
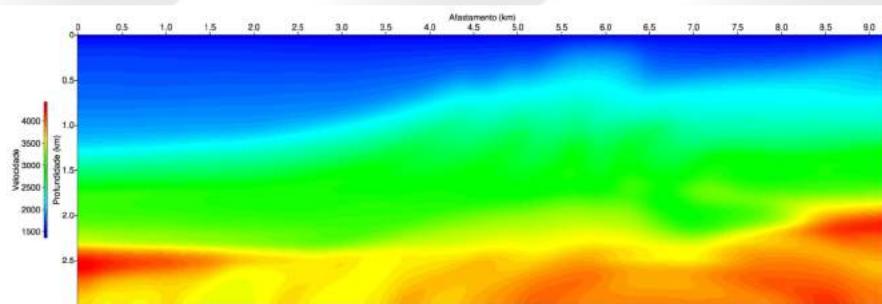
https://gitlab.com/ssvitor/dataflow_analyzer

PROVENANCE IN DATA ANALYTICS UQ IN SEISMIC IMAGING

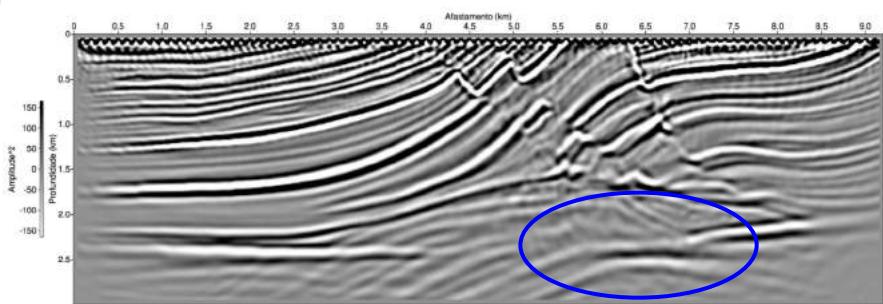
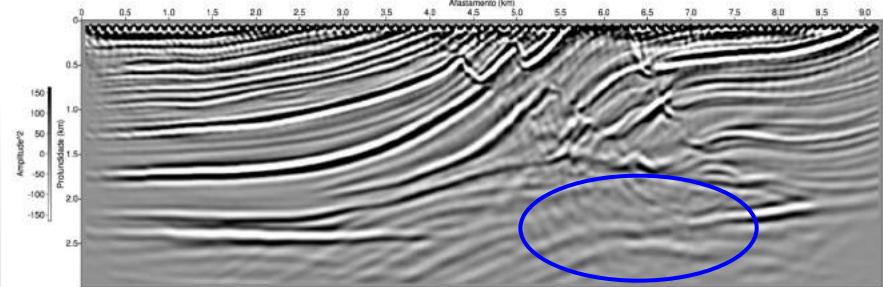
Seismic Imaging



Uncertainty in Seismic Imaging



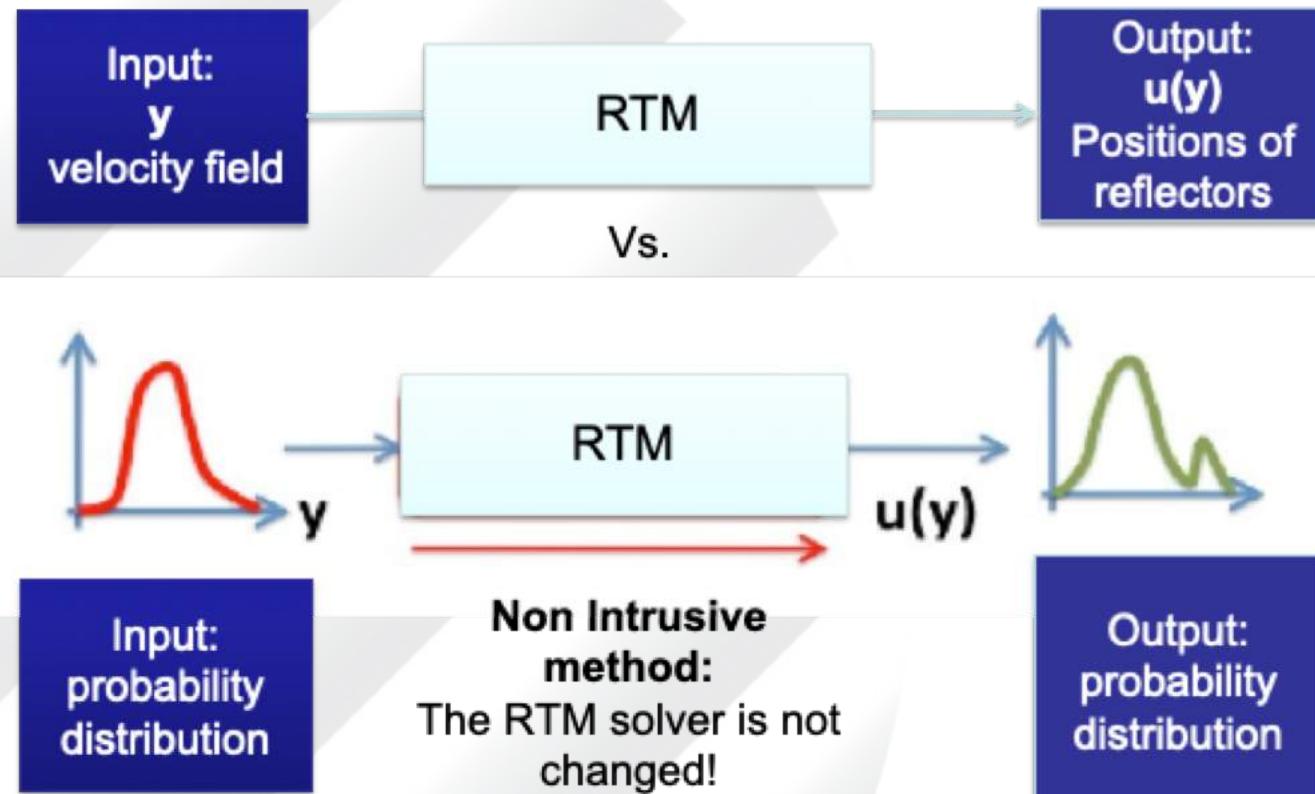
(a) Possible velocity fields



(b) Corresponding Seismic images

- + Fomel, S., Landa, E. Structural Uncertainty Of Time-migrated Seismic Images. *Journal of Applied Geophysics*, vol. 101, 27-30, 2014.
- + Pawelec, I., Uncertainty Quantification in Seismic Imaging, Dissertation. Colorado School of Mines, Colorado, US, 2018.

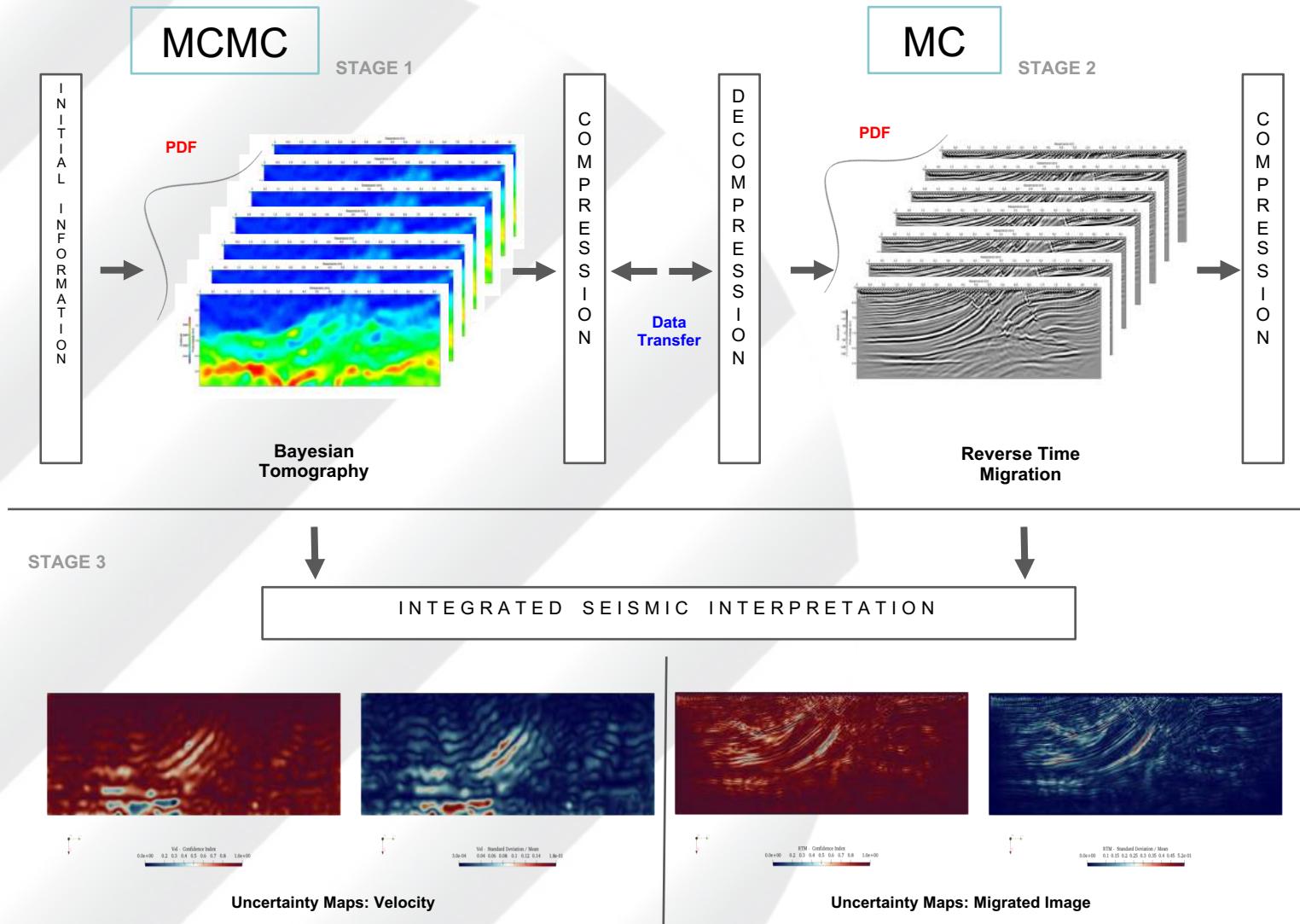
Uncertainties in Seismic Images



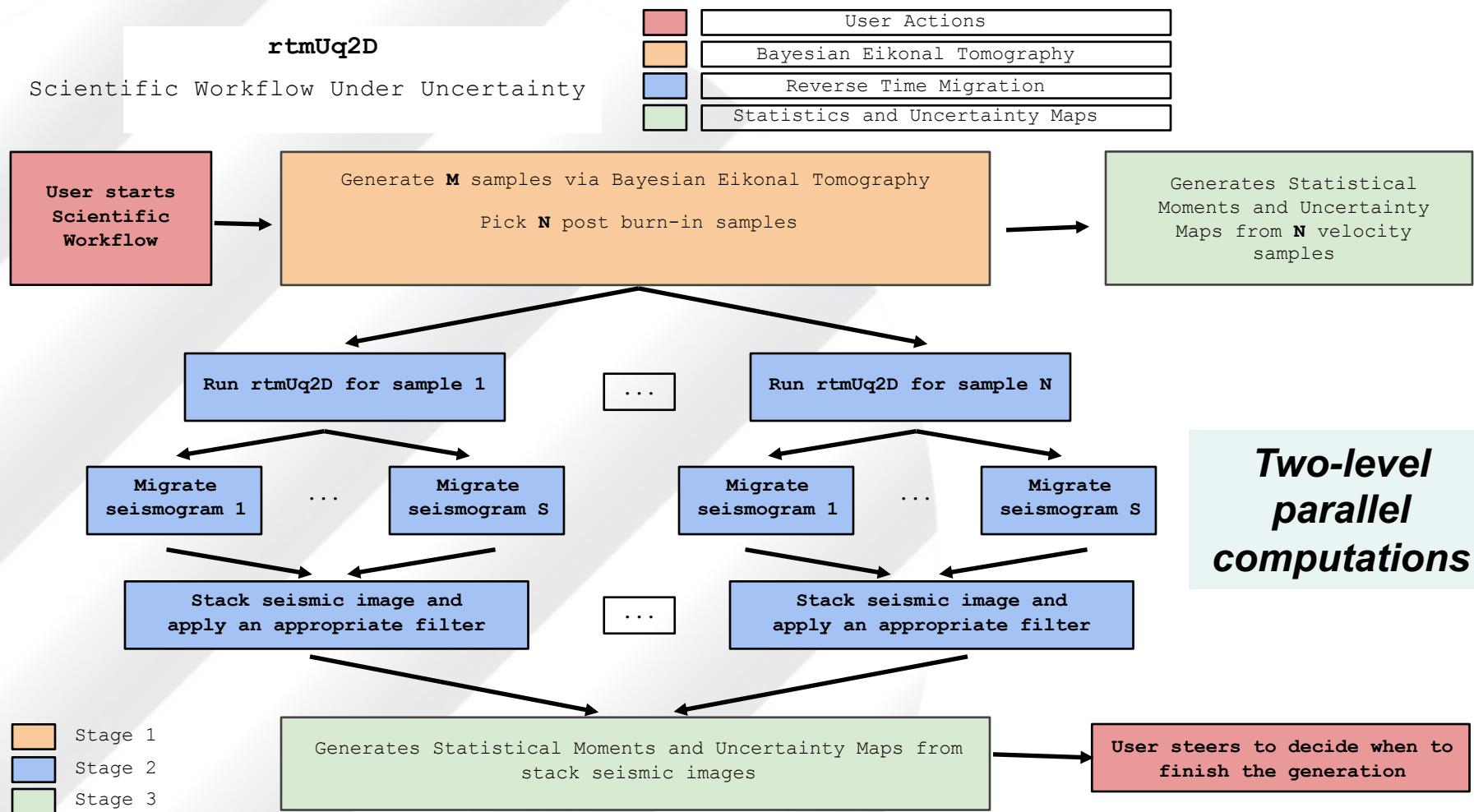
Remark: RTM = 2 migrations -> UQ on arrival time for each model

- + Lavril, T., J. P., *A Probabilistic Framework for Uncertainty Quantification in Large-Scale Simulations: Application in Seismic Imaging*. Dissertação, Universidade Federal do Rio de Janeiro, COPPE/UFRJ, Rio de Janeiro, 2015.

ENSEMBLE BASED UNCERTAINTY ANALYSIS



RTM under UQ Workflow



Provenance Data in the RTM-UQ Workflow

- **File correlation as seismic images are generated**
 - What is the image corresponding to the “m-velocity field?
 - There are too many files for a manual search. Tiresome procedure and prone to errors.
- **Image statistics**
 - What are the 5 images closest to the mean?
 - Which (and how many) images are within 2 standard deviation
 - Which images are within 95% confidence interval?
- **RTM processing time and image sizes**
 - What are the mean and standard deviation of RTM CPU times?
 - What are the most compressed images?

Provenance Data in the Workflow

```

int task_id = 1;
Task task_runrtm= Task(dataflow.get_tag(), runrtm.get_tag(), task_id);

vector<string> irunrtm_values = {velocityFile, waveletFile, std::to_string(Nx), std::to_string(Nz)};
Dataset& ds_irunrtm = task_runrtm.add_dataset(irunrtm.get_tag());
ds_irunrtm.add_element_with_values(irunrtm_values);

// RTM per shot...
for (shotID = 1; shotID <= n; shotID++) {
    printf("SHOT %d\n", shotID);
    // arithmetic progression
    shotLocation[shotID] = firstLoc + (shotID - 1) * step;
    task_runrtm.begin();
    isotropicAcousticModeling(shotLocation, modifiedVelocityModel, wavelet);
    sprintf(seismogramFile_02, seismogramFile, shotID);
    vector<string> iseisogram_values = {seismogramFile_02};
    Dataset& ds_iseisogram = task_runrtm.add_dataset(iseisogram.get_tag());
    ds_iseisogram.add_element_with_values(iseisogram_values);
    task_runrtm.begin();
    reading_file(seismogramFile_02 , seismogram, numberofReceivers, numberoftimestep);
    adjointModeling(shotID, shotLocation, modifiedVelocityModel, wavelet, seismogram);
    if (shotID== numberOfShots)
        task_runrtm.end();
    else task_runrtm.save();
}
if (shotID== numberOfShots)
    task_runrtm.end();
else task_runrtm.save();

```

+-----+ id shot shotlocation seismogram_path cross_correlation_path
+-----+ id shot shotlocation seismogram_path cross_correlation_path
1 1 4 ./INPUTS/seismogram_00001.bin ./OUTPUTS/crossCorrelation_00001.bin
2 2 12 ./INPUTS/seismogram_00002.bin ./OUTPUTS/crossCorrelation_00002.bin
3 3 20 ./INPUTS/seismogram_00003.bin ./OUTPUTS/crossCorrelation_00003.bin
4 4 28 ./INPUTS/seismogram_00004.bin ./OUTPUTS/crossCorrelation_00004.bin
5 5 36 ./INPUTS/seismogram_00005.bin ./OUTPUTS/crossCorrelation_00005.bin
6 6 44 ./INPUTS/seismogram_00006.bin ./OUTPUTS/crossCorrelation_00006.bin
7 7 52 ./INPUTS/seismogram_00007.bin ./OUTPUTS/crossCorrelation_00007.bin
8 8 60 ./INPUTS/seismogram_00008.bin ./OUTPUTS/crossCorrelation_00008.bin
9 9 68 ./INPUTS/seismogram_00009.bin ./OUTPUTS/crossCorrelation_00009.bin
10 10 76 ./INPUTS/seismogram_00010.bin ./OUTPUTS/crossCorrelation_00010.bin

PROV challenges

- CAPTURE (GRANULARITY)
- REPRESENTATION
- STORAGE
- QUERIES

PROV is domain-agnostic
can be specialized with domain data

Provenance & Data Extraction for analysis at runtime



Prepare data to be queried, before being produced



Capture provenance data without interference on execution control



Negligible-low overhead on the execution time



Data analysis support during and after execution



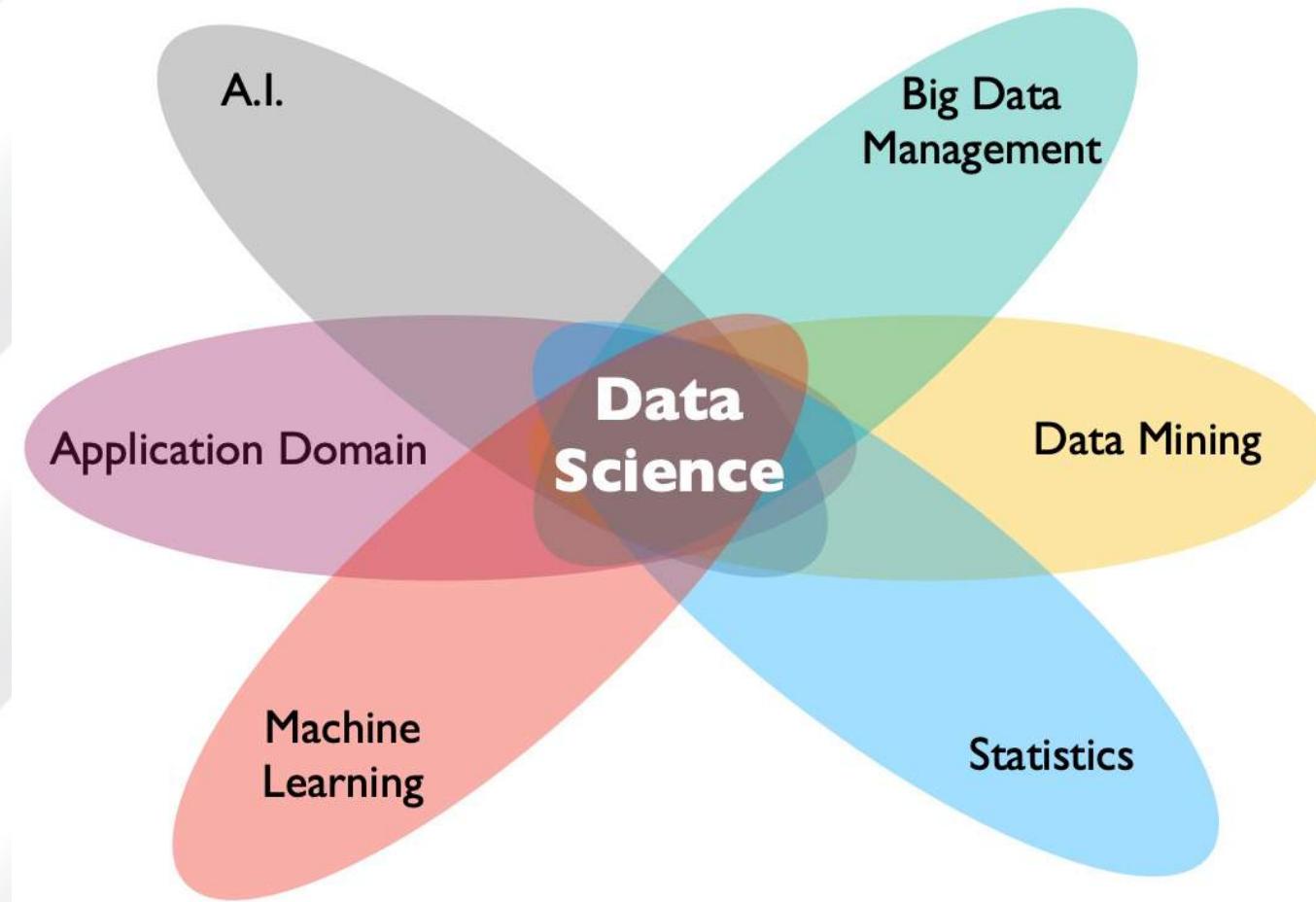
Register user steering actions

Take home lesson:

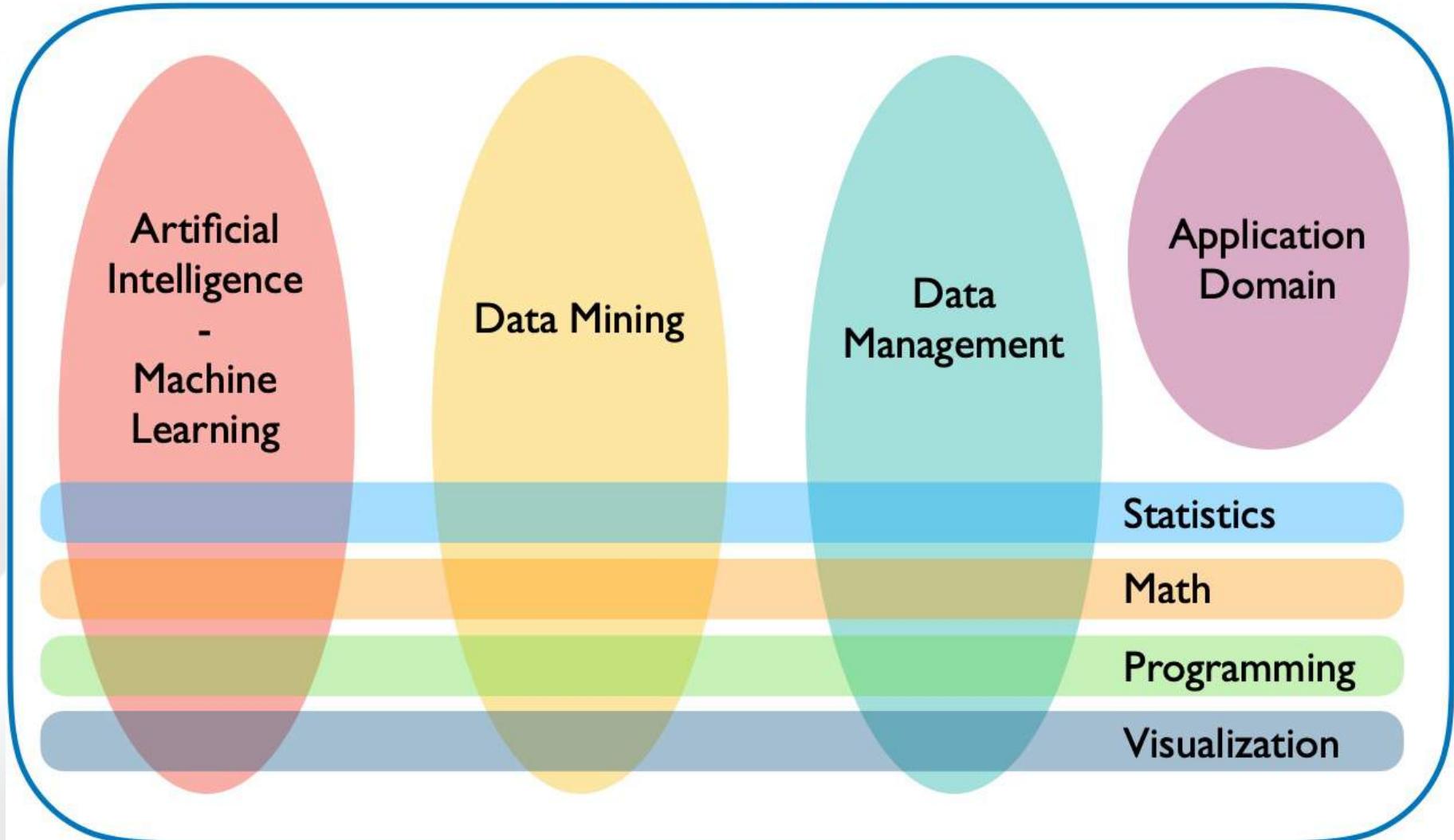
**MANAGING DATA IS NEVER SIMPLE
WITHOUT PROVENANCE EVEN HARDER**

DATA SCIENCE AND MACHINE LEARNING

Data Science ≠ Machine Learning



Data Science



Data Science and Machine Learning Point of View

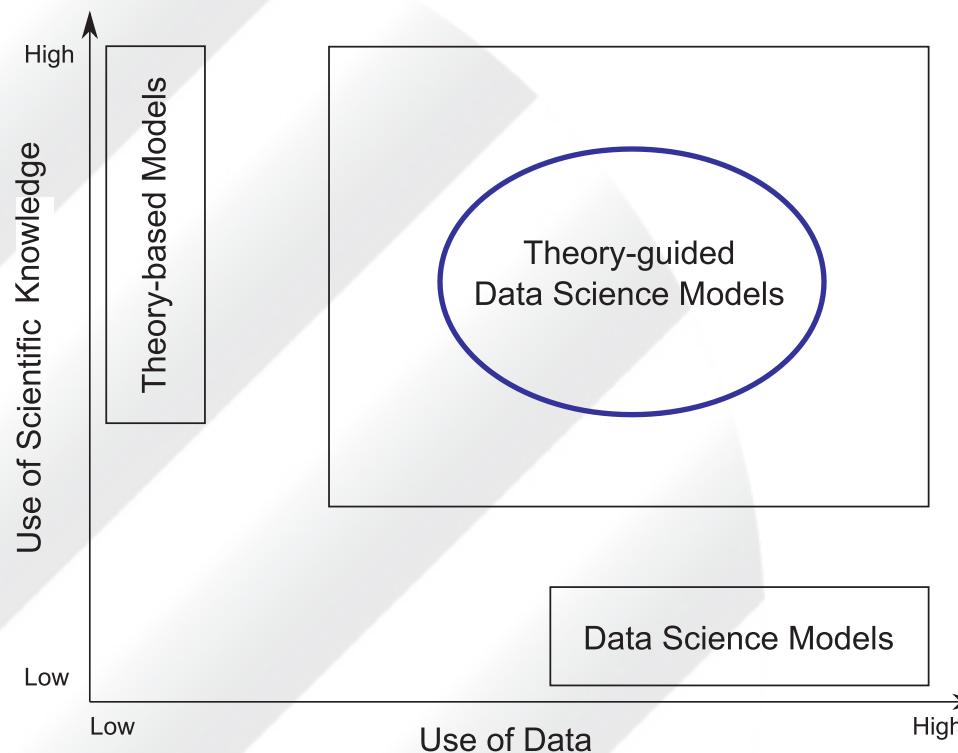
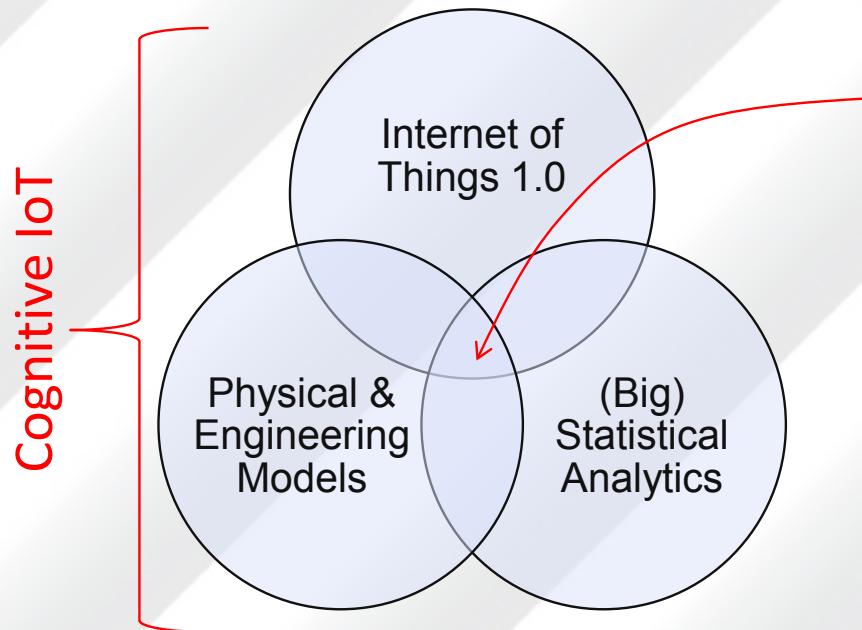


Fig. 1. A representation of knowledge discovery methods in scientific applications. The x -axis measures the use of data while the y -axis measures the use of scientific knowledge. Theory-guided data science explores the space of knowledge discovery that makes ample use of the available data while being observant of the underlying scientific knowledge.

Industry Point of View

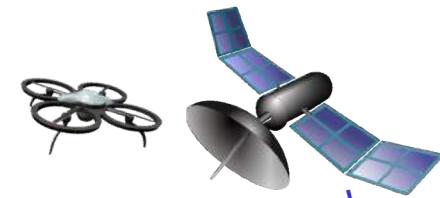
IBM Physical Analytics



Highly interdisciplinary with many interesting research topics

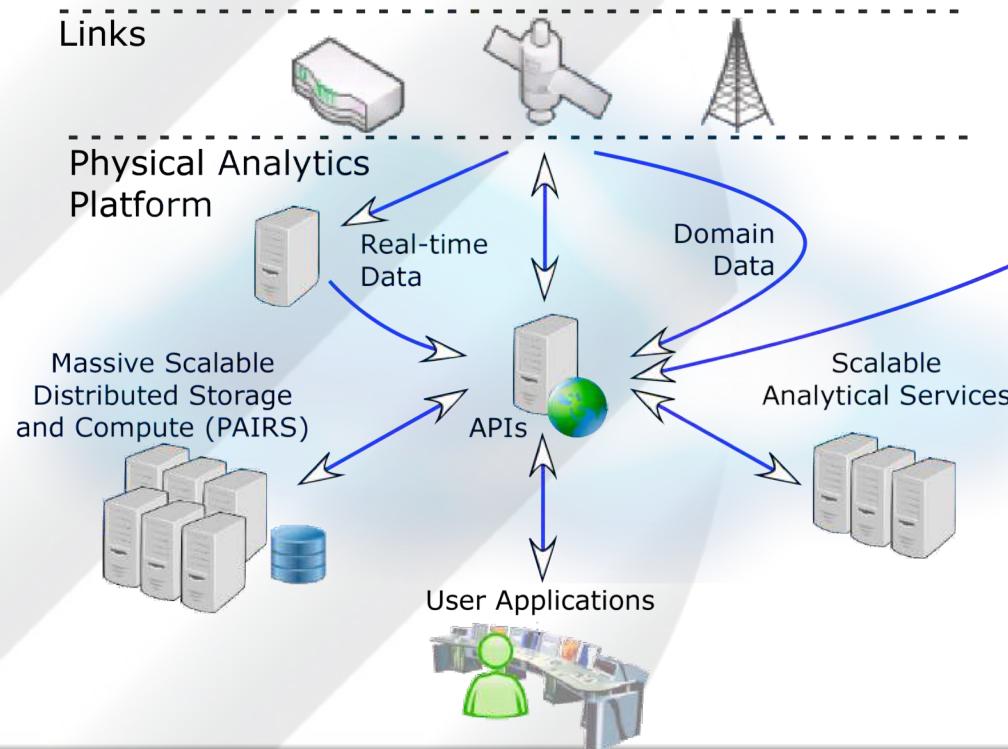
- Statistical learning under physical constraints
- Model complexity reduction
- Situation-dependent, machine-learnt multi-model blending
- Graph theory and statistical physics
- Parallelizing physical models for data-intensive computation
-
- Feed learning back to improve understanding of the underlying physics

IBM Physical Analytics



Physical systems in different Industries

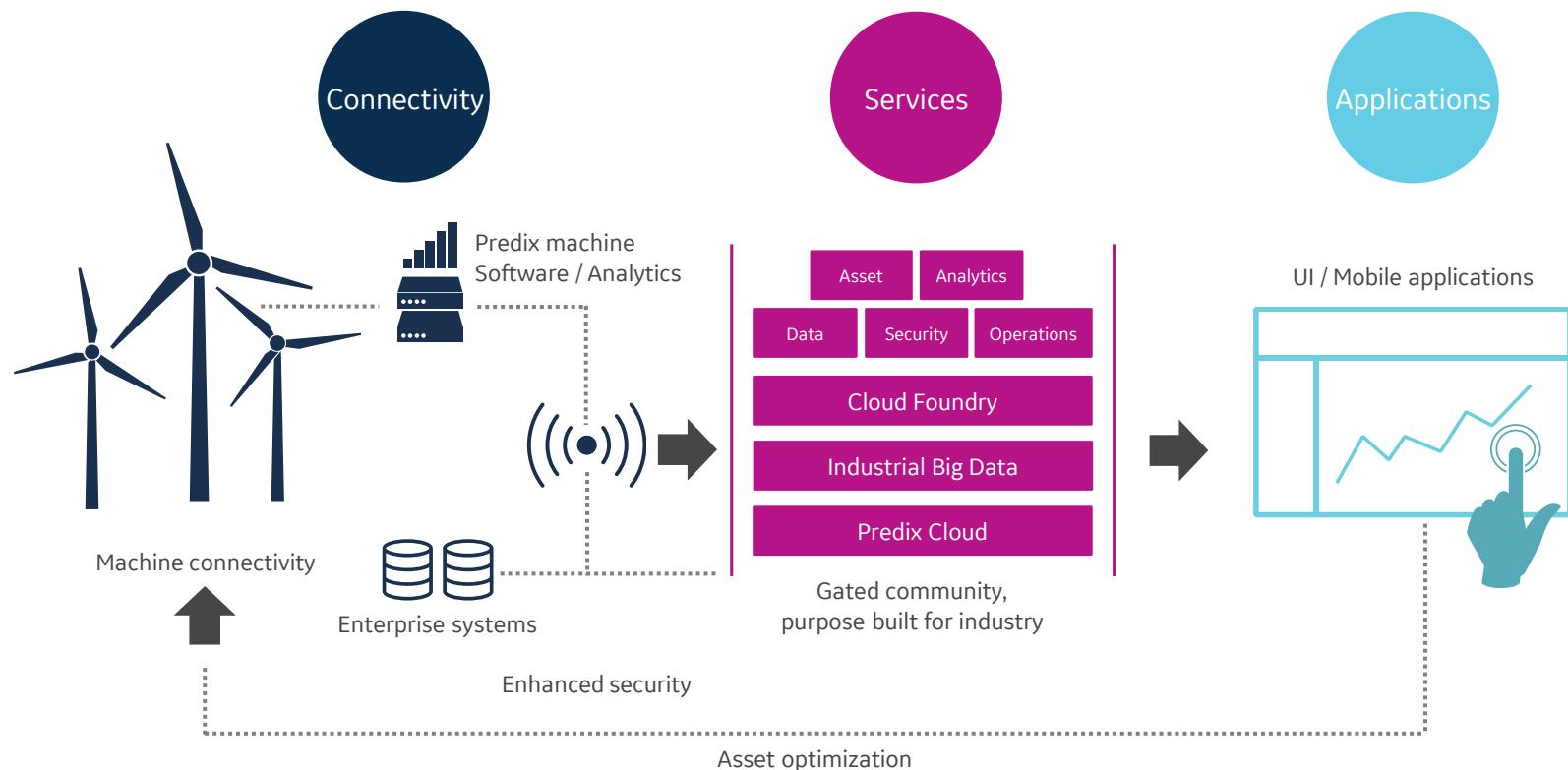
Cross Industry	Telecom	Oil and Gas	Health Care	Transportation & Travel	Utility	Agriculture	Public
Data Centers 	High value Buildings 	Network Offices 	Pipelines & Fracking Operations 	Hospitals 	Bridges / Infrastructure 	Solar farms 	Vineyards, Wineries, Greenhouses



Industry Point of View

GE 'Digital Foundries', where Physics + Analytics Intersect

Predix – A cloud-based platform for industry



A Word of Caution in Machine Learning

- **Machine Learning (ML) are *Universal Approximators***
- **Huge success in text mining, image recognition, etc.**
- **ML finds best solution minimizing some norm between *input* and *output* – standard solutions may be trapped within a local minimum, thus stochastic algorithms are preferred**
- **ML does not know about the problem being solved**
 - In case of physical systems ML may reach unphysical solutions, that is, a solution for incompressible flows where velocity does not satisfy $\nabla \cdot \mathbf{u} = 0$
- **ML learning should know about the physics: adding constraints to the formulation**

Physics-Constrained Data Models

Karthik Duraisamy, Data-enabled, Physics-constrained Predictive Modeling of Complex Systems, SIAM NEWS, July-Aug, 2017

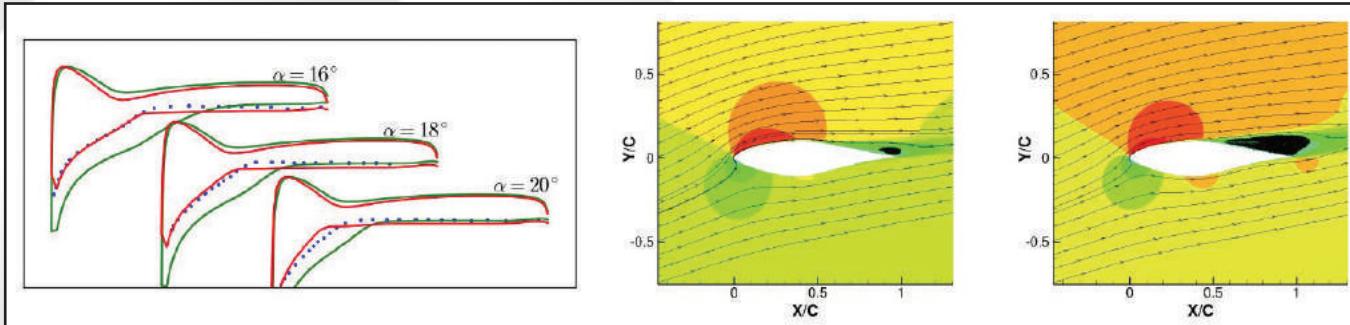
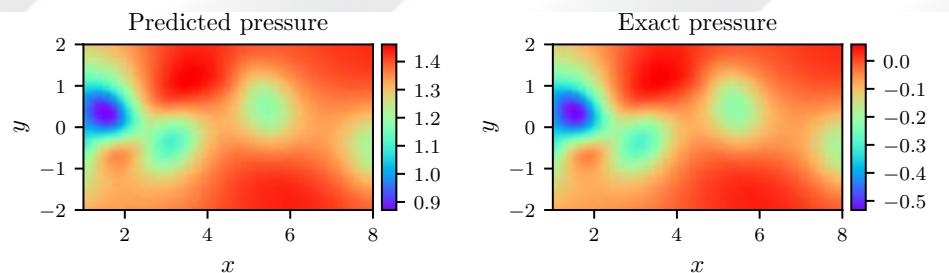


Figure 1. Example of data-augmented, physics-based modeling applied to turbulent flow prediction. Predictive improvement is achieved based on inferring force data over another airfoil and constructing machine-learned model augmentations. **Left.** Pressure over airfoil surface. Green: baseline physics model. Red: machine learning-augmented physics model. Blue: experimental measurements. **Middle.** Baseline flow prediction (pressure contours and streamlines). **Right.** Flow prediction using machine learning-augmented physics model. Image adapted from [8].



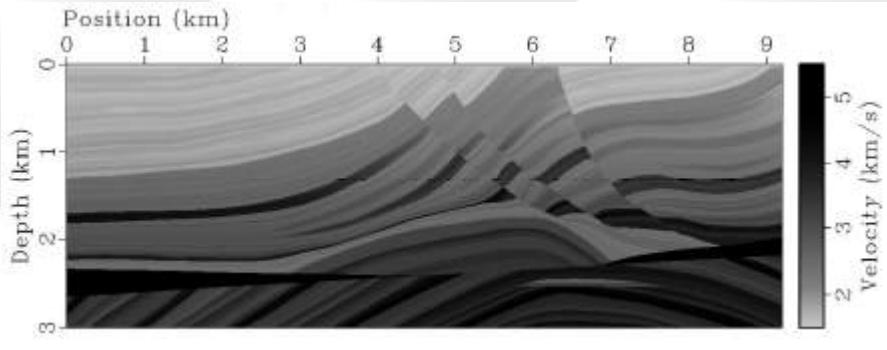
Correct PDE	$u_t + (uu_x + vu_y) = -p_x + 0.01(u_{xx} + u_{yy})$ $v_t + (uv_x + vv_y) = -p_y + 0.01(v_{xx} + v_{yy})$
Identified PDE (clean data)	$u_t + 0.999(uu_x + vu_y) = -p_x + 0.01047(u_{xx} + u_{yy})$ $v_t + 0.999(uv_x + vv_y) = -p_y + 0.01047(v_{xx} + v_{yy})$
Identified PDE (1% noise)	$u_t + 0.998(uu_x + vu_y) = -p_x + 0.01057(u_{xx} + u_{yy})$ $v_t + 0.998(uv_x + vv_y) = -p_y + 0.01057(v_{xx} + v_{yy})$

Navier-Stokes equation: Top: Predicted versus exact instantaneous pressure field $p(t,x,y)$ at a representative time instant.. This remarkable qualitative agreement highlights the ability of **physics-informed neural networks** to identify the entire pressure field, despite the fact that no data on the pressure are used during model training. Bottom: Correct partial differential equation along with the identified one obtained by learning λ_1 , λ_2 and $p(t, x, y)$.

M Raissi, P Perdikaris, GE Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations
Journal of Computational Physics 378, 686-707, 2019

Understanding Uncertainty in a Model Problem: 1D Acoustic Wave Equation

Marmousi velocity model



1D wave equation

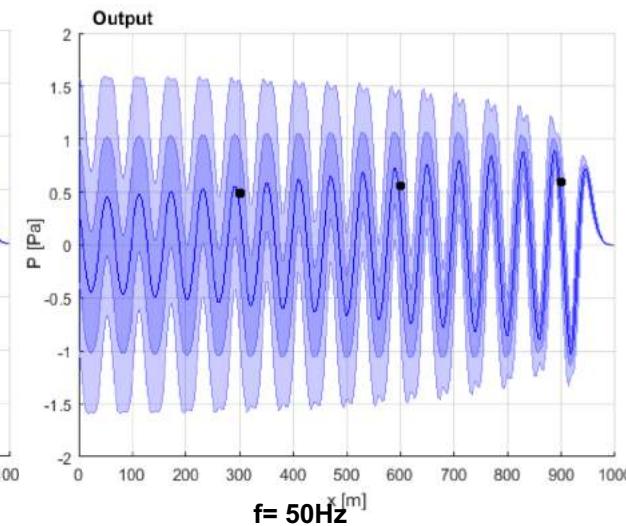
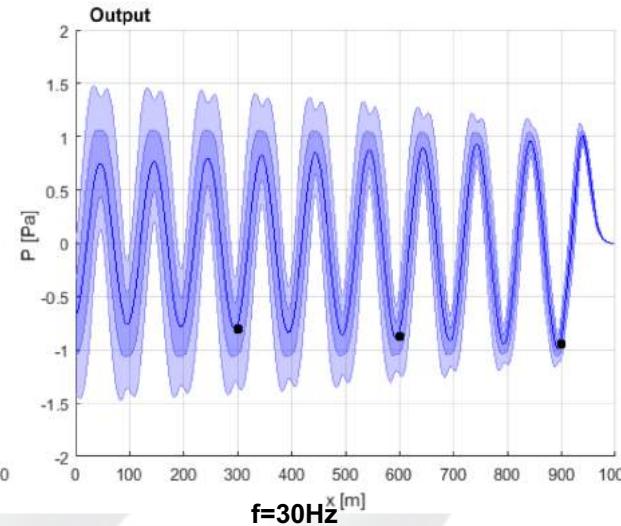
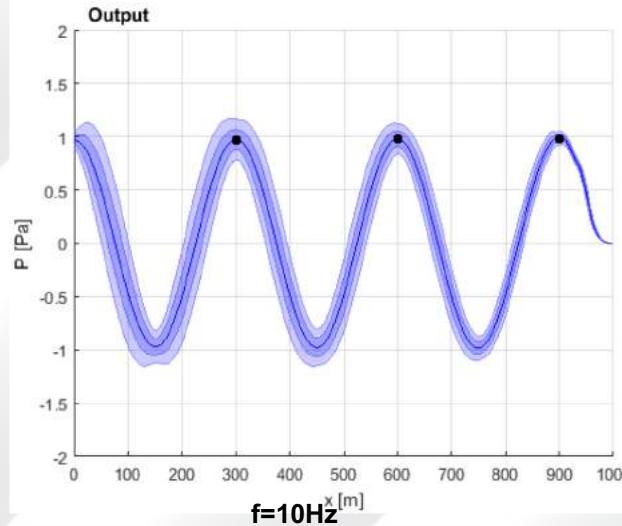
$$\frac{\partial^2 P}{\partial t^2} = c^2 \frac{\partial^2 P}{\partial x^2} \quad \text{with} \quad \frac{\partial P(x,0)}{\partial t} = 0, \quad P(x, 0) = 0$$

Initial condition

$$P(0,t) = \sin(2\pi ft) \quad \text{and} \quad P(L,t) = 0$$

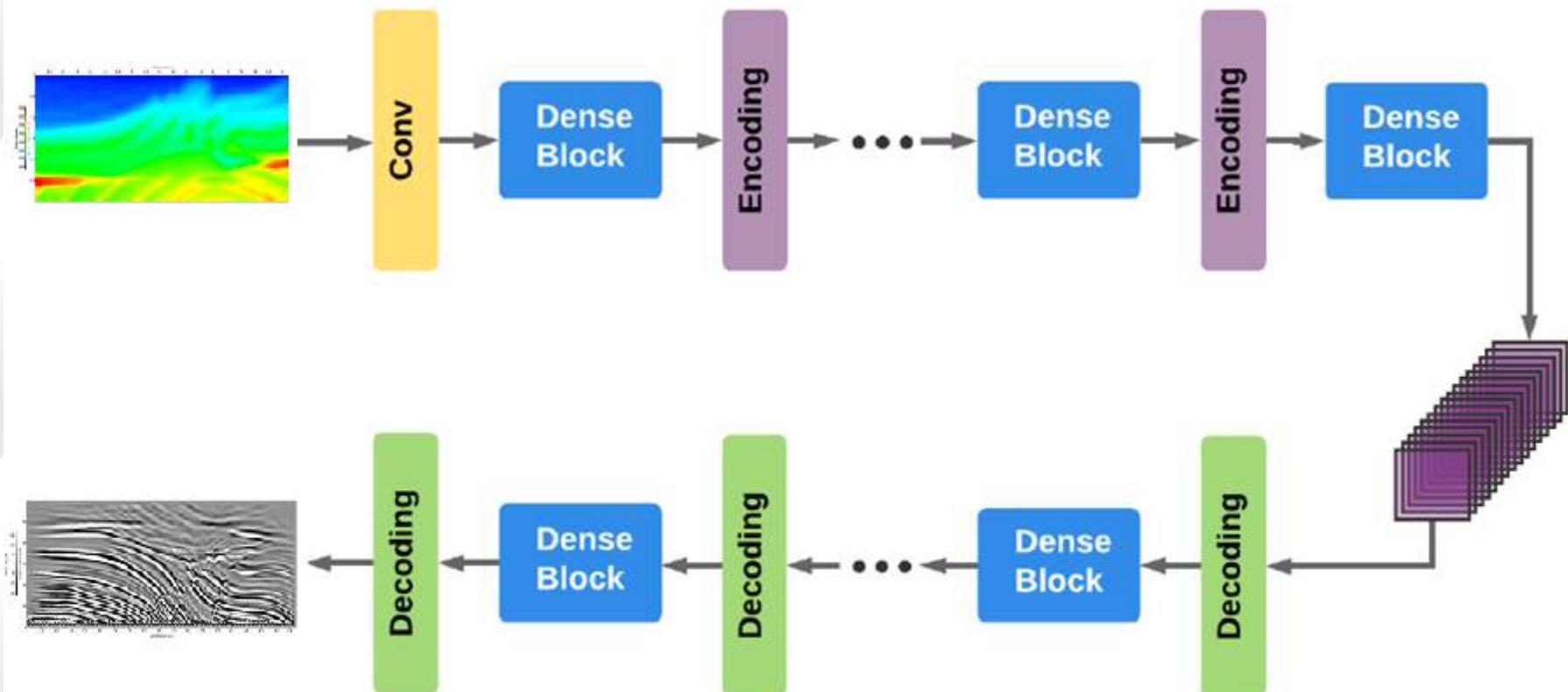
Velocity field representation via Karhunen-Loeve expansion

$$c(x, \omega) \approx \xi_0 E(c(x)) + \sigma(c(x)) \sum_{i=1}^M \xi_i(\omega) \sqrt{\lambda_i} f_i(x)$$



A (Deep) Machine Learning Surrogate

- Aim: Decrease the computational effort of the Monte Carlo Method



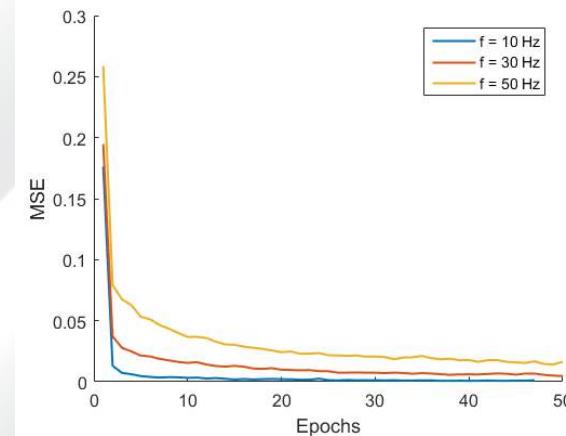
- Adapted from: Zhu and Zabaras, *Bayesian deep convolutional encoder-decoder networks for surrogate modelling and uncertainty quantification*, JCP, 366 (2018) 415–447

DNN Surrogate: Encoder-Decoder

Neural network architecture and training details: Output denotes the number of output feature maps and Dimension is the spatial dimension of the feature map.

Layers	Output	Dimension
Input	1	301
Convolution	48	148
Dense Block	144	148
Encoding	72	74
Dense Block	168	74
Encoding	84	37
Dense Block	180	37
Decoding	1	3

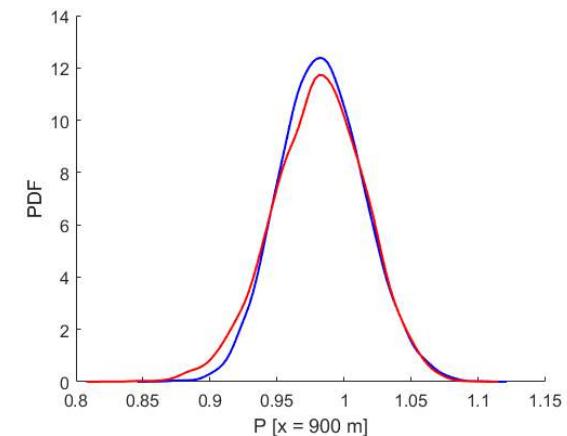
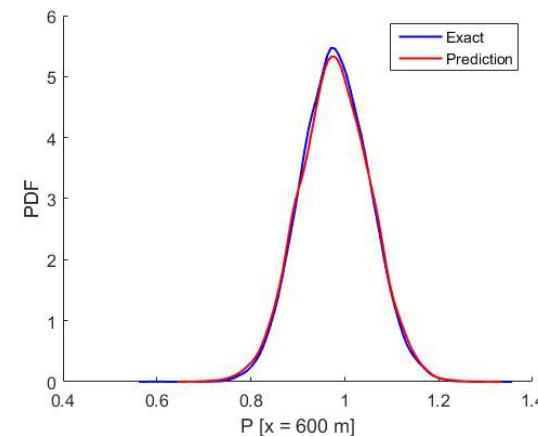
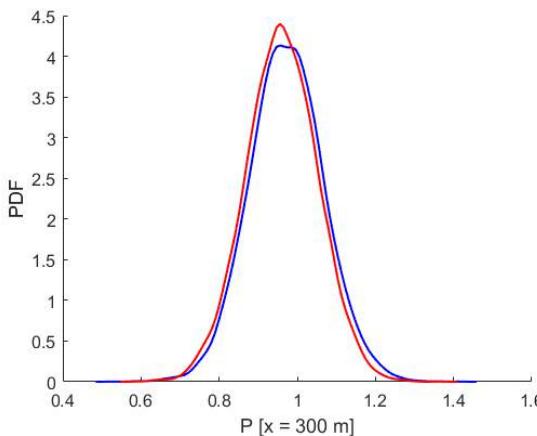
- Activation ReLU
- Optimizer ADAM
- Learning Rate 1e-3
- Epochs 50
- Training Data 10000
- Mini-batch size 64
- Loss function Mean-squared-error



- The neural network is trained using the **TensorFlow** software library (version r1.13) on a 2.30 GHz Xeon E5-2630 CPU, where it took around **17 minutes** to train 50 Epochs.

Predicting Uncertainties with the Surrogate

- 10K samples, different from those used in training, are extrapolated to evaluate the surrogate model performance in predicting the pressure in 3 different points in the spatial domain.
- The computational cost of the surrogate model is around **10 seconds** (**Original model has a cost of 4hrs!**)
- For **low frequencies** of the source term, it was found that in regions where the variability of the solution is lower, the substitutive model presented a better prediction, with a **mean relative error of 1%**.
- For **high frequencies** of the source term, in regions where the variability of the solution is higher, the mean relative error increases significantly, **around 10%**.
- *frequency of source term – 10 Hz*



Final Remarks and Discussion

- ❑ We see a synergy between HPC, CSE, and ML. We have learned that new machines open new possibilities → **what will happen when we have exascale machines?**
- ❑ **Predictive Computational Science and Predictive Data Science are the new paradigms**
- ❑ CSE and Data Science: **we need them both**
- ❑ Machines and Apps are becoming increasingly complex, how to manage all this? **HPC, Storage, Networking and Visualization are becoming more integrated**
- ❑ What is ML role? How it will be integrated with PDE-based solvers? How to deal with the huge amount of data?
 - Los Alamos 2ND PHYSICS INFORMED MACHINE LEARNING, Jan 21-25, 2018
- ❑ International collaboration, like BR-EU H2020 will strengthen the field
- ❑ **Digital Twins** in Energy Industry
 - GE is promoting; Shell joined a JIP for digital twins in Oil&Gas in July, 2017
- ❑ The Brazilian government has sustained a successful R&D policy funding in the Energy area directly promoting partnerships involving the private sector, universities and research labs, supervised by the regulatory agencies ANP (oil & gas) ANEEL (electricity).

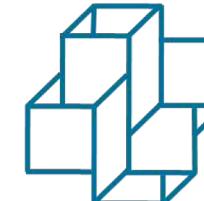
Our Team

- Faculty at COPPE
 - Civil Engineering: J. Alves, R. Elias, L. Landau
 - Mechanical Engineering: F. Rochinha
 - Computer Science: M. Mattoso, A. Cortes, G. Travassos
 - Visitors: R. Cottreau (CNRS-Centrale), P. Valduriez (INRIA)
- Pos-Docs, Research Staff and Students:
 - C. Alves, A. Aveleda, C. Barbosa, G. Barros, O. Caldas, J. Camata, D. Costa, H. Costa, M. Costa, A. Cruz, J. Dias, R. Freitas, I. Ghisi, M. Goncalves Jr, M. Grave, G. Guerra, L. Gesenhues, L. Kunstmann, T. Lavril, E. Lins, T. Miras, E. Ogasawara, D. Oliveira, A. Rossa, F. Seabra, B. Silva, C. Silva, D. F.C. Silva, R. Silva, V. Silva, D. Vasconcelos, S. Zio
- Funding: Petrobras, ANP, MCTIC, MEC
- Computer Resources: NACAD/COPPE/UFRJ, SDumont/LNCC, TACC/UT Austin, Occigen/CINES, France, MN4-BSC, Spain

Our Main Academic Partners



FOR COMPUTATIONAL ENGINEERING & SCIENCES



Laboratório
Nacional de
Computação
Científica



NACAD

Our Industrial Partners



ORACLE®



EMC²



Itautec

Hewlett Packard
Enterprise

HALLIBURTON



Schlumberger



NACAD