

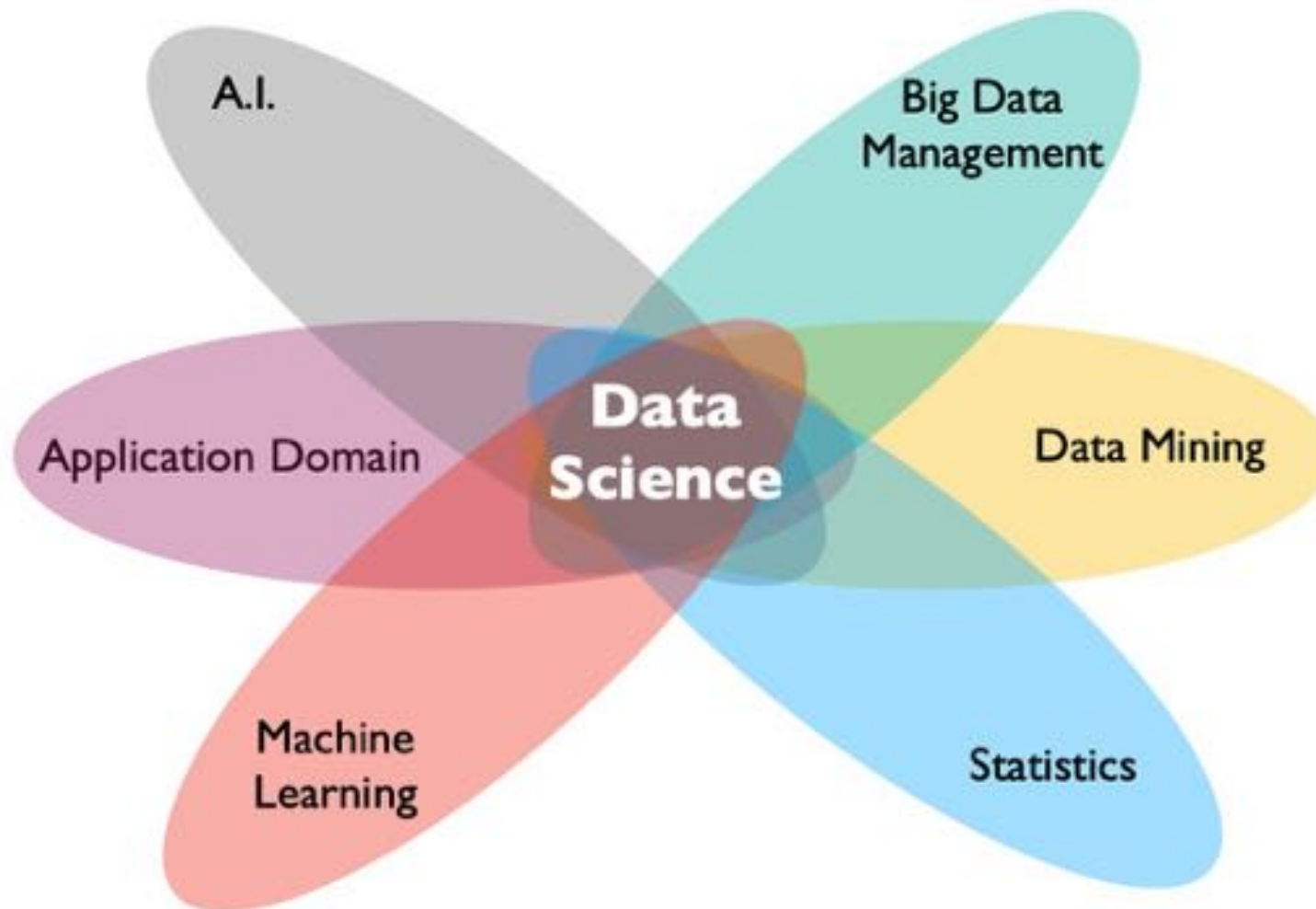






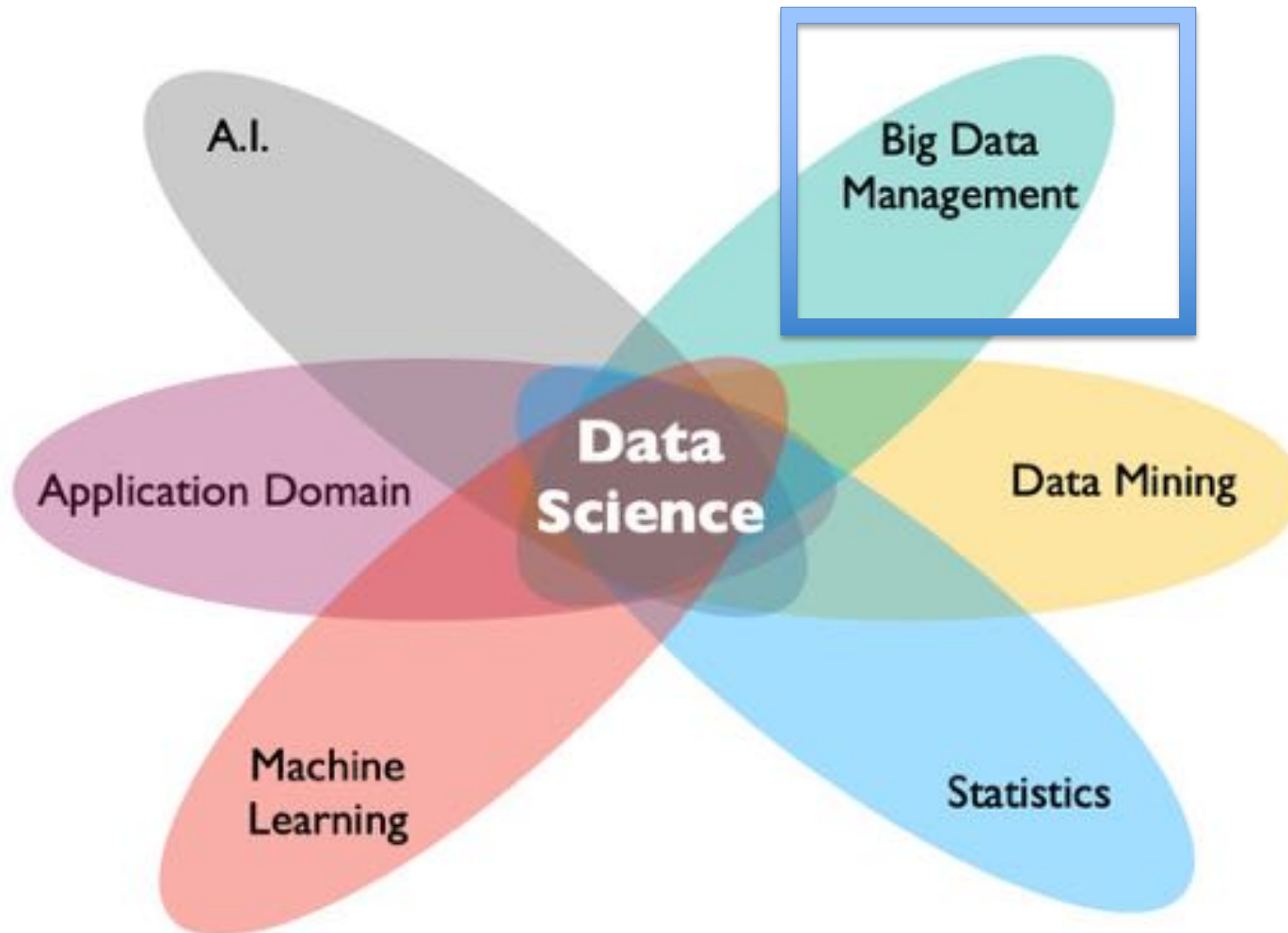


# Data Science $\neq$ Machine Learning



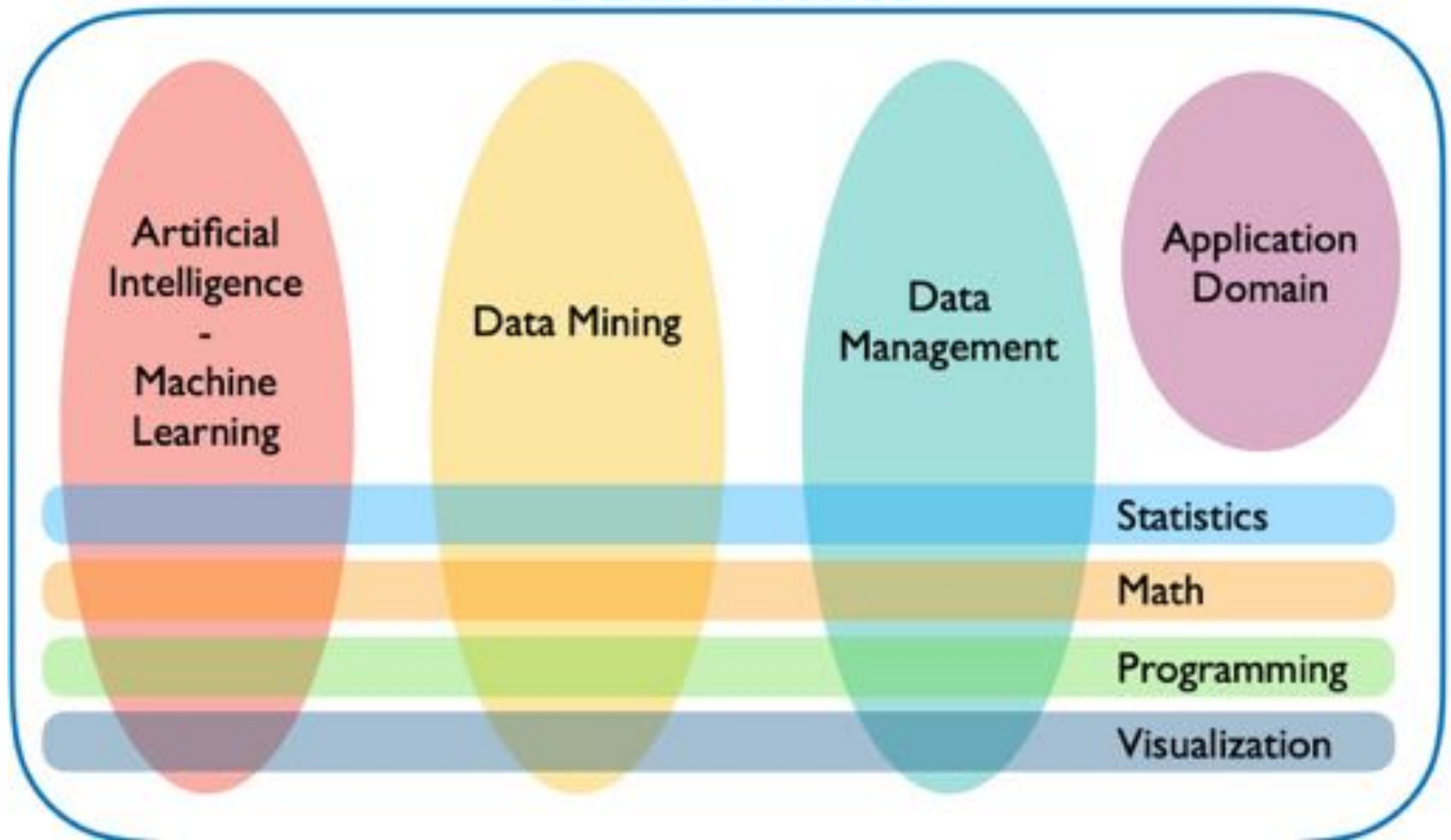
Jens Dittrich "Data Science  $\neq$  Machine Learning: Some Thoughts on the Role of Data Management in the new AI-Tsunami" -- Keynote DEEM@SIGMOD 2018, June 2018

# Data Science $\neq$ Machine Learning



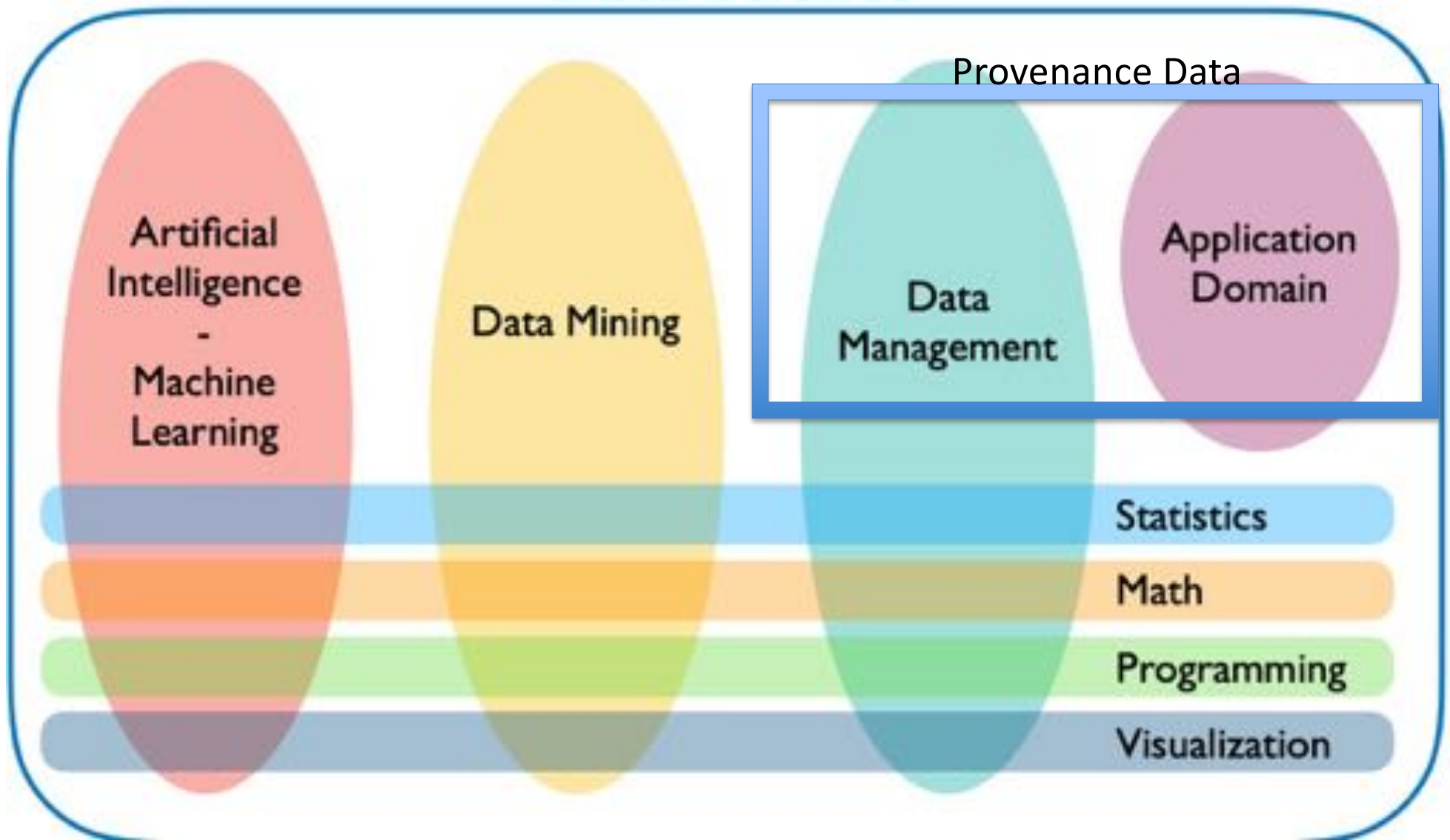
Jens Dittrich "Data Science  $\neq$  Machine Learning: Some Thoughts on the Role of Data Management in the new AI-Tsunami" -- Keynote DEEM@SIGMOD 2018, June 2018

# Data Science



Jens Dittrich, VLDB 2017 (invited talk: “Deep Learning (m)eats Databases”)

# Data Science



Jens Dittrich, VLDB 2017 (invited talk: “Deep Learning (m)eats Databases”)

# Projetos de Pesquisa

- D-Interpret - Gerência de dados para auxiliar a explicação de resultados em aplicações de ciência de dados- CNPq Universal/Faixa C
- MonDataSim - Análise de dados de simulações computacionais por meio de monitoramento da execução, dados de proveniência e intervenções dinâmicas- Faperj Cientista RJ
- SciDISC - Scientific data analysis using Data-Intensive Scalable Computing- França-Brasil, INRIA Associate Teams - Patrick Valduriez
- Bolsa Produtividade de Pesquisa CNPq- 1B



# Trabalho em Equipe



# ECI – aluno de IC 2008

Fernando Chirigati

Research

Awards and Honors

Publications

Talks

Professional Activities

Bio

CV



Fernando Chirigati

Postdoctoral Research Associate

NYU Tandon School of Engineering  
VIDA Research Center  
370 Jay Street, 11th Floor  
Brooklyn, NY 11201

fchirigati [at] nyu [dot] edu

Twitter

LinkedIn

GitHub

## About Me

Currently, I'm a Postdoctoral Research Associate at [NYU Tandon School of Engineering](#), under the supervision of [Prof. Juliana Freire](#). I came to the beautiful - and very busy - city of New York in January 2012 to pursue a Ph.D. degree. Before, I worked as a Research Assistant at Federal University of Rio de Janeiro (UFRJ), under the supervision of [Prof. Marta Mattoso](#). I have a Ph.D. in Computer Science from NYU, and a B.E. in Computer and Information Engineering from UFRJ. To check out my full CV, click [here](#).

I come from the gorgeous city of [Petrópolis](#), in Brazil, where almost all my family and friends still reside. I had the opportunity to study in [Rio de Janeiro](#), the "Marvelous City," where I not only made a lot of good friends, but also started working with research in the database area.

## Research

My research interests are mainly in the area of scientific data management, including **provenance management and analytics, large-scale data analytics, data science, data mining, reproducibility, and data visualization.**

ECI – aluno de MSc (2011) e DSc (2013)



Melhor Tese de  
Doutorado no  
Concurso de Teses  
do XXX Simpósio  
Brasileiro de  
Bancos de Dados,  
SBC.

Jonas Dias · 1st

Data Science Consultant at Dell EMC, Distinguished  
Member of the Technical Staff



ECI – aluno de IC (2012); MSc (2014) e DSc (2018)



Atualmente na Snap, Los Angeles



DCC/IM – aluno de MSc (2015) e DSc previsão defesa 2019



Melhor Dissertação  
de Mestrado no  
Concurso de Teses  
do XXXII Simpósio  
Brasileiro de Bancos  
de Dados, SBC.

**Renan Souza** - 1st

Research Software Engineer at IBM and Computer Science

PhD Candidate at COPPE/UFRJ

# Scicumulus: A lightweight cloud middleware to explore many task computing paradigm in scientific workflows

Authors Daniel de Oliveira, Eduardo Ogasawara, Fernanda Bailão, Marta Mattoso

Publication date 2010/7/5

Conference 2010 IEEE 3rd International Conference on Cloud Computing

Pages 378-385

Publisher IEEE

Description Most of the large-scale scientific experiments modeled as scientific workflows produce a large amount of data and require workflow parallelism to reduce workflow execution time. Some of the existing Scientific Workflow Management Systems (SWMS) explore parallelism techniques - such as parameter sweep and data fragmentation. In those systems, several computing resources are used to accomplish many computational tasks in homogeneous environments, such as multiprocessor machines or cluster systems. Cloud computing has become a popular high performance computing model in which (virtualized) resources are provided as services over the Web. Some scientists are starting to adopt the cloud model in scientific domains and are moving their scientific workflows (programs and data) from local environments to the cloud. Nevertheless, it is still difficult for the scientist to express a parallel computing ...



Total citations Cited by 189



Fonte:  
Google Scholar

# Data Challenges



Computing in  
Continuum



Human-In-the-Loop  
(HIL)



Data Integration

# Especialização: extremos



“a set of special-purpose appliances”



<http://www.ianfoster.org/wordpress/presentations/> (Code in Continuum)



# Exascale “sob medida”



- Google projetou sua própria unidade de processamento de tensores (TPU), projetada inicialmente para alta vazão de operações aritméticas de baixa precisão. Novas gerações de TPU chegam a petaflops além da Edge TPU
  - TPU aumenta o desempenho das redes neurais por trás de Google Search, Street View, Google Photos and Google Translate.

# Code in continuum

## Globus



Ian Foster  
(IPDPS'2019)



<https://www.slideshare.net/ianfoster/coding-the-continuum>

# 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

---

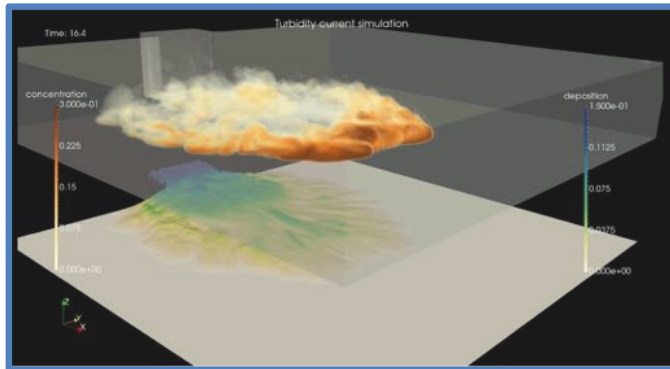
## Big Data and Its Technical Challenges

**Systems were not built to have humans in the loop**

# **PROVENANCE IN DATA ANALYTICS**



# Analyzing sedimentation solver data with provenance data



## Real Case Study:

### ➤ DfAnalyzer tool

- Provenance Gatherer
- Simulation Data Extractor
- Extract data from libMesh
- Query and visualize at runtime

### ➤ TACC computer using 1,040 cores

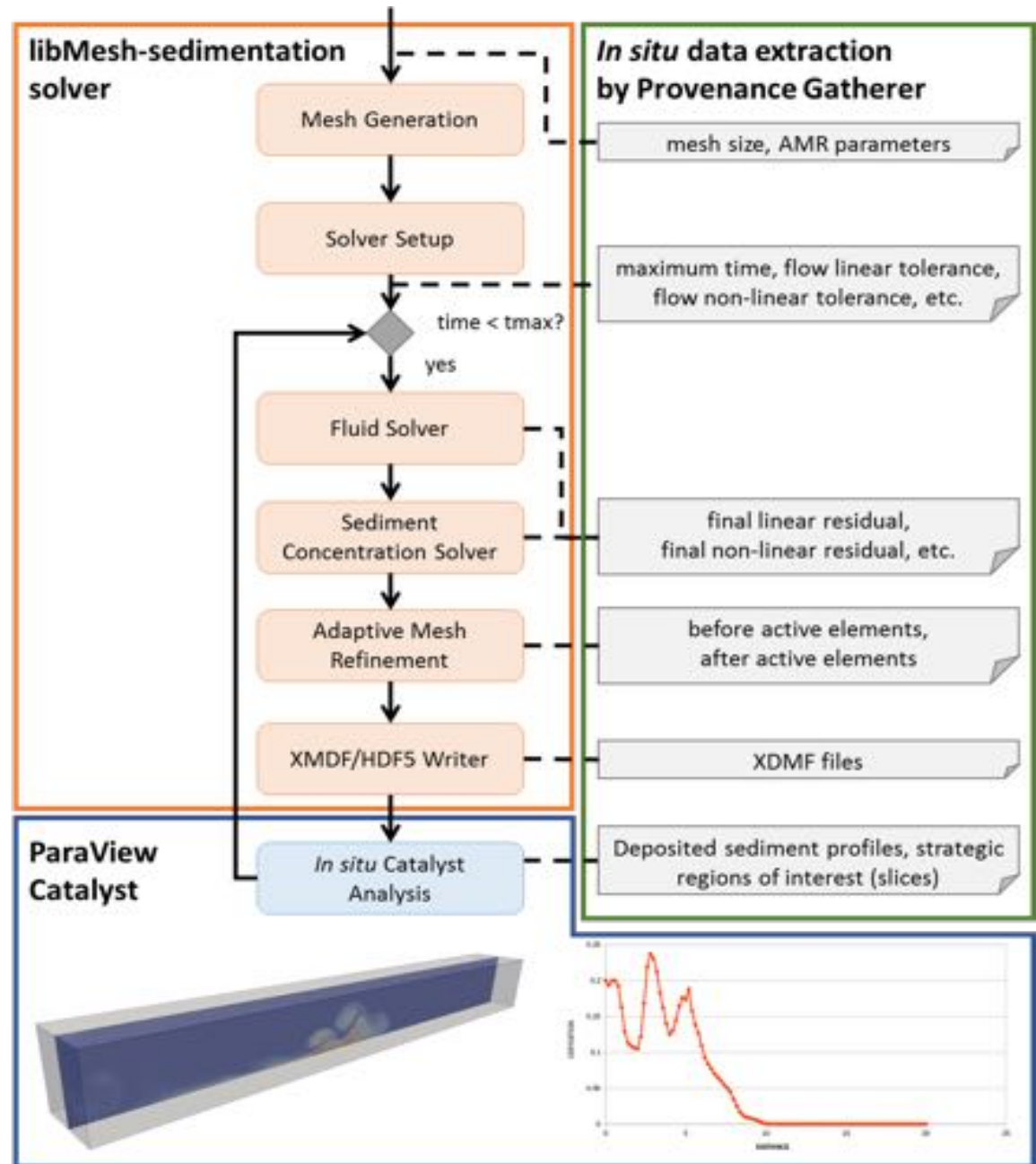
### ➤ Simulation elapsed time: 137.75 min

**Solver:** 136.80 min → 99.31% of total time

**DfAnalyzer w. Catalyst:** 0.95 min → 0.69%

**Initial input mesh:** 480 x 80 x 80 hexahedra

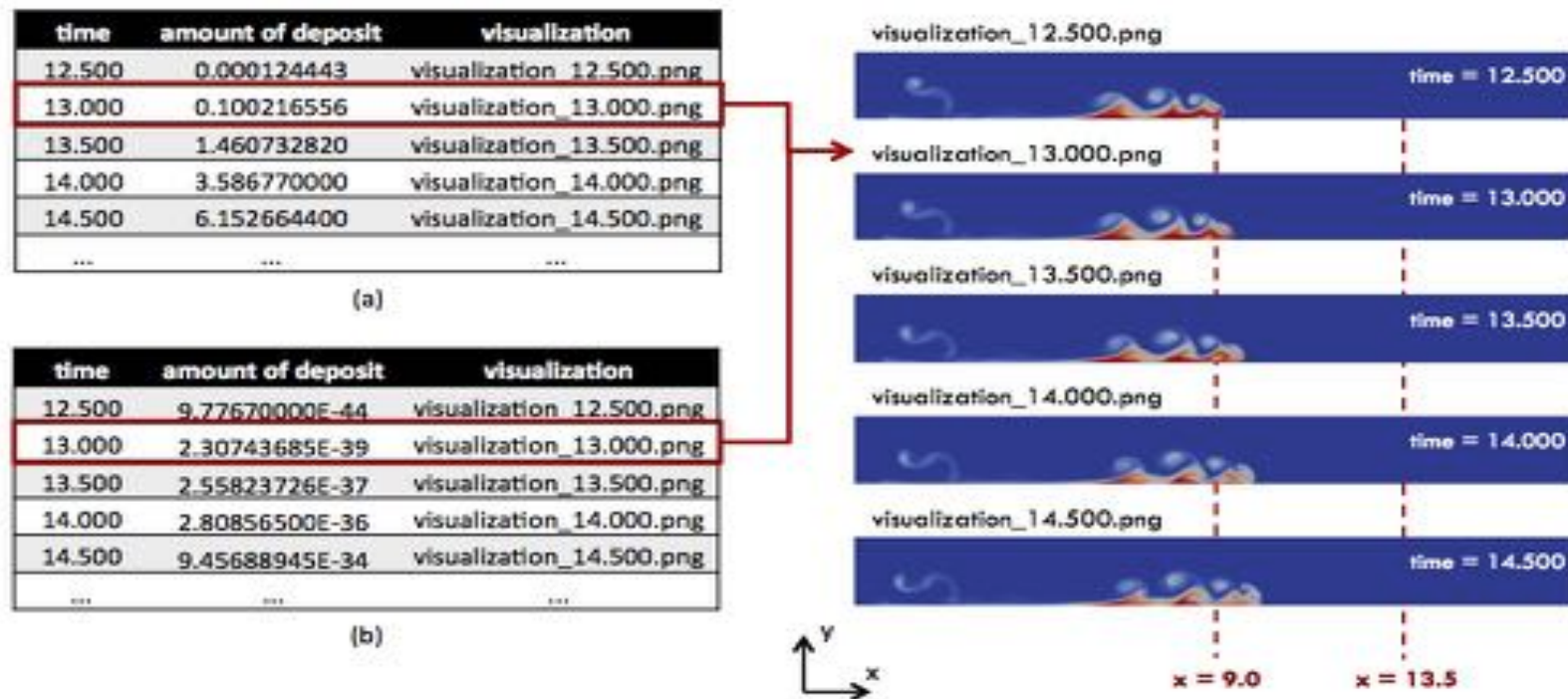
**Total time steps:** 200



Silva, V., de Oliveira, D., Valduriez, P., & Mattoso, M.  
DfAnalyzer: Runtime Dataflow Analysis of Scientific  
Applications using Provenance  
PVLDB 2018.

# Sediment 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

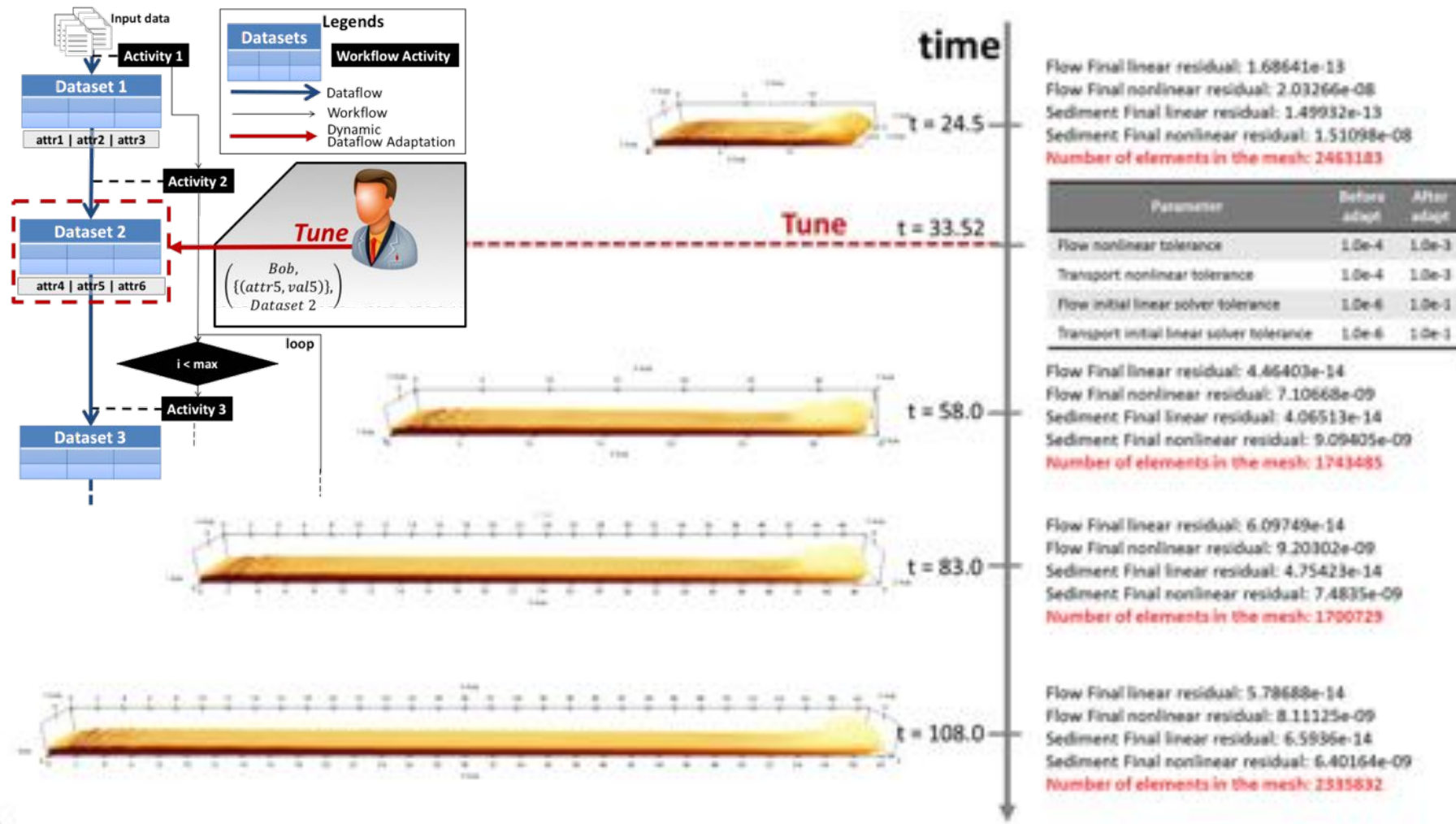
# Query processing

- 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
...	...	...	...

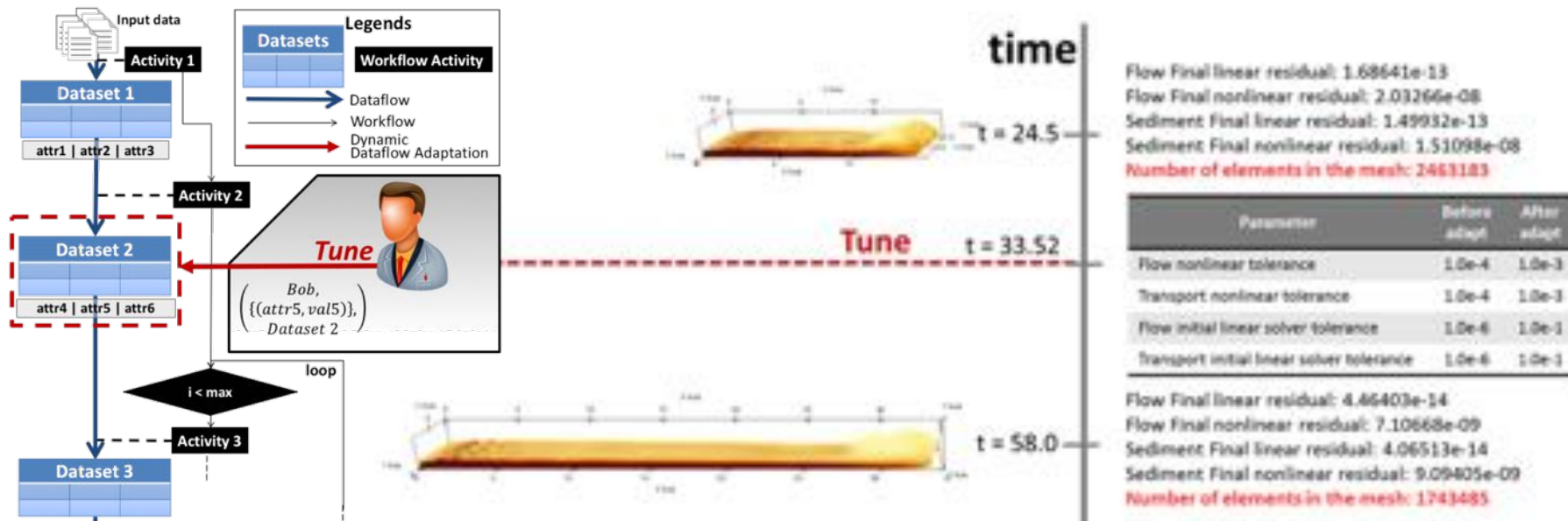
# Turbidity currents simulation data analysis with provenance



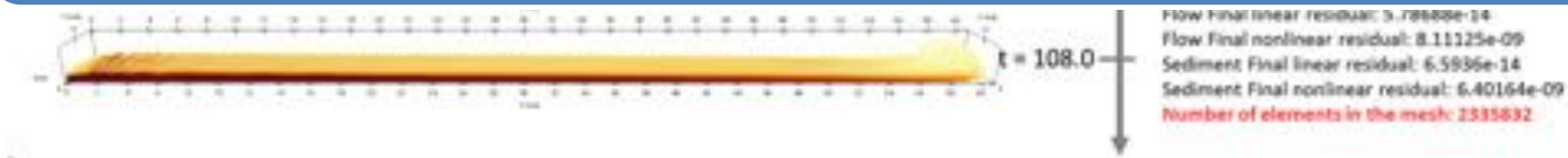
Renan Souza et al. Tracking of Online Parameter Fine-tunings in Scientific Workflows, *WORKS workshop in ACM/IEEE Supercomputing*, 2017. Extended to FGCS Special Issue 2019



# Turbidity currents simulation data analysis with provenance

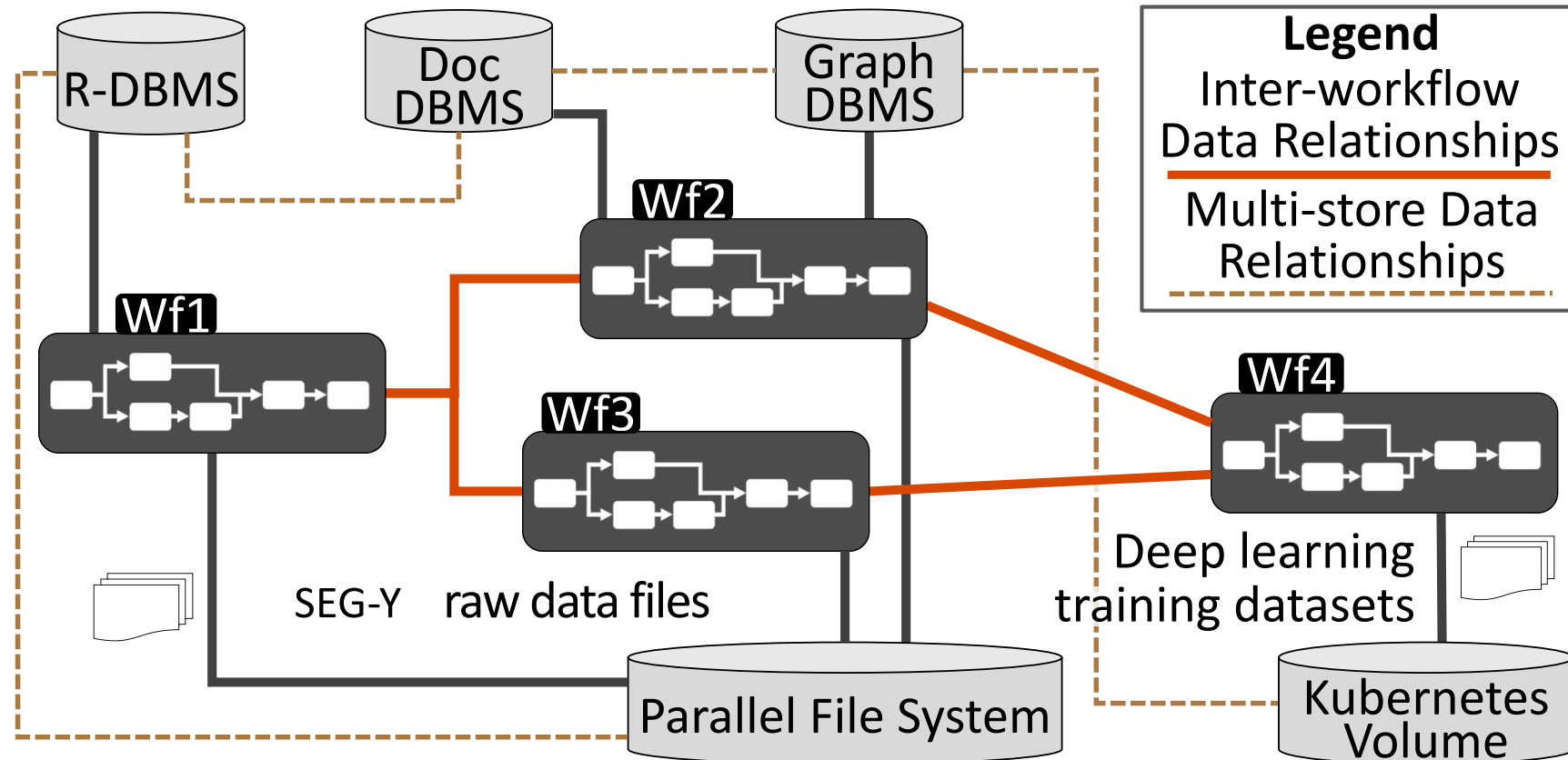


Parameter-tuning reduced the execution time by 10 days (37%), overhead < 1% and allowed job to finish successfully



# Provenance relating data for DL

(IEEE eScience 2019)



*How the geographic coordinates were extracted from the SEG-Y file that is being used to produce training and validation files?*

*What is the spatial resolution between slices in the seismic data?*


# Big data need big theory too

Peter V. Coveney<sup>1</sup>, Edward R. Dougherty<sup>2</sup> and Roger R. Highfield<sup>3</sup>

<sup>1</sup>Centre for Computational Science, University College London, Gordon Street, London WC1H 0AJ, UK

<sup>2</sup>Center for Bioinformatics and Genomic Systems Engineering, Texas A&M University, College Station, TX 77843-31283, USA

<sup>3</sup>Science Museum, Exhibition Road, London SW7 2DD, UK

 PVC, 0000-0002-8787-7256

The current interest in big data, machine learning and data analytics has generated the widespread impression that such methods are capable of solving most problems without the need for conventional scientific methods of inquiry. Interest in these methods is intensifying, accelerated by the ease with which digitized data can be acquired in virtually all fields of endeavour, from science, healthcare and cybersecurity to economics, social sciences and the humanities.

How BD might assist with the struggle of the human mind to overcome three notorious barriers:

- nonlinearity,
- non-locality and
- hyperdimensional spaces.

## Big data: the end of the scientific method?

Sauro Succi<sup>1,2</sup> and Peter V. Coveney<sup>3,4</sup>

<sup>1</sup>Center for Life Nano Sciences at La Sapienza, Istituto Italiano di Tecnologia, viale R. Margherita, 265, 00161, Roma, Italy

<sup>2</sup>Institute for Applied Computational Science, J. Paulson School of Engineering and Applied Sciences, Harvard University, 29 Oxford Street, Cambridge, USA

<sup>3</sup>Centre for Computational Science, Department of Chemistry, University College London, London, UK

<sup>4</sup>Yale University, New Haven, USA

 PVC, 0000-0002-8787-7256

*For it is not the abundance of knowledge, but the interior feeling and taste of things, which is accustomed to satisfy the desire of the soul*

2019, Phil.Trans.R.Soc.A -Special issue 'Multiscale modelling, simulation and computing: from the desktop to the exascale'.




# Big data need big theory too

Peter V. Coveney<sup>1</sup>, Edward R. Dougherty<sup>2</sup> and Roger R. Highfield<sup>3</sup>

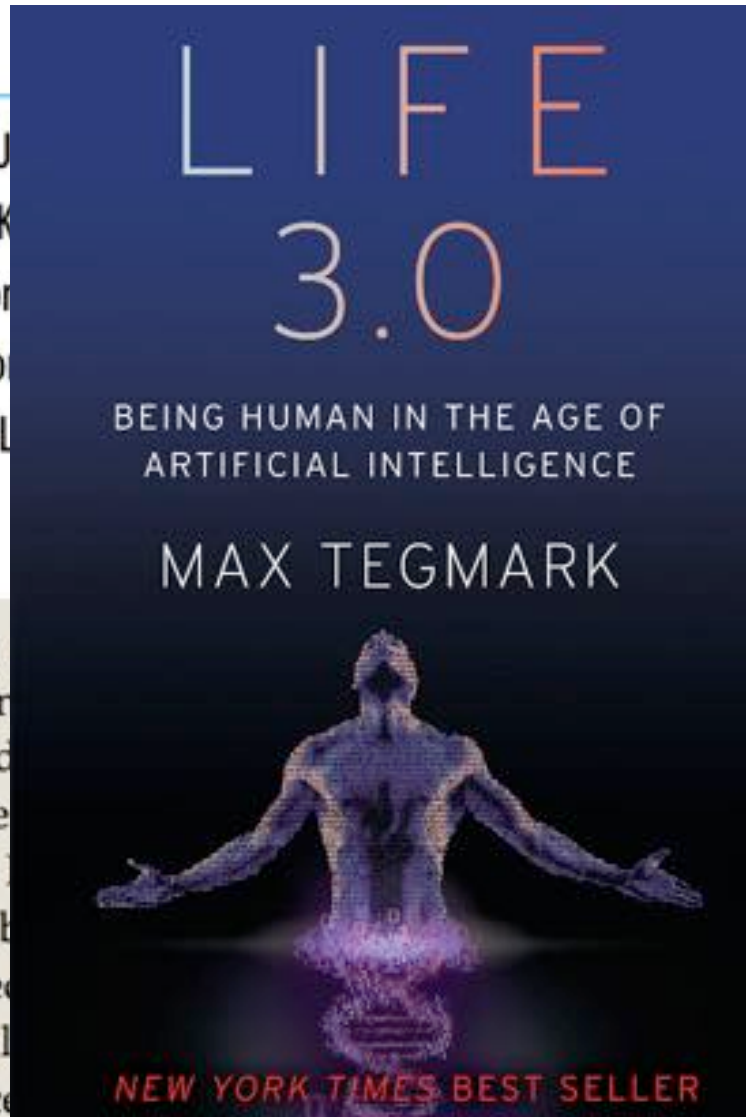
<sup>1</sup>Centre for Computational Science, UCL, Gower Street, London WC1E 6BT, UK

<sup>2</sup>Center for Bioinformatics and Genomics, Texas A&M University, College Station, TX 77843, USA

<sup>3</sup>Science Museum, Exhibition Road, London SW7 2BX, UK

 PVC, 0000-0002-8787-7256

The current interest in big data and data analytics has given the impression that such methods can solve most problems without the need for scientific methods of inquiry. This is intensifying, accelerated by the fact that digitized data can be acquired in an endeavour, from science, health care to economics, social science



How BD might assist with the struggle of the human mind to overcome three notorious barriers:

- nonlinearity,
- non-locality and high-dimensional spaces.

**Big data: the end of the scientific method?**

by Edward R. Dougherty<sup>1,2</sup> and Peter V. Coveney<sup>3,4</sup>

<sup>1</sup>Department of Nano Sciences at La Sapienza, Istituto Italiano di Tecnologia, viale R. Margherita, 265, 00161, Roma, Italy

<sup>2</sup>Department of Applied Computational Science, J. Paulson School of Applied and Applied Sciences, Harvard University, 29 Oxford Street, Cambridge, MA 02138, USA

<sup>3</sup>Centre for Computational Science, Department of Chemistry, UCL, Gower Street, London, UK

<sup>4</sup>Department of Chemistry, Yale University, New Haven, USA

ORCID: 0000-0002-8787-7256

*It is not the abundance of knowledge, but the interior feeling and taste of things, which is accustomed to satisfy the desire of the soul.* (St Ignatius of Loyola).

Due to the bold claims of big data (BD) and its promise to revolutionize science, it is

# Data Science $\neq$ ML

## The Data Science Cake



### Ingredients:

50g statistics  
120g linear algebra  
200g programming  
1kg visualisation  
300g software engineering

### Additional skills:

creativity  
out of the box thinking  
grit  
team spirit

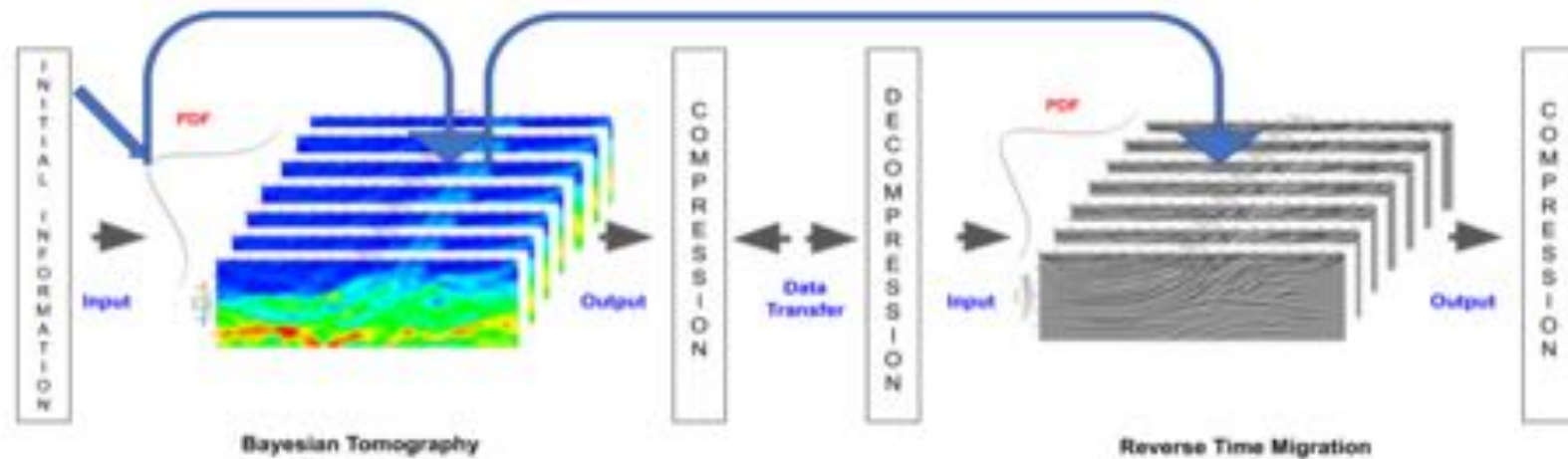
© istock.com sasilsolutions

[Jens Dittrich 2018] - <http://www.youtube.com/user/jensdit>



# Acknowledgements





Obrigada!

Marta Mattoso

