

Uma trajetória de pesquisa em engenharia de dados para aplicações em larga escala

Autor: Vítor Silva PESC • COPPE • UFRJ

Orientadores: Marta Mattoso PESC • COPPE • UFRJ

Daniel de Oliveira IC • UFF

Patrick Valduriez INRIA • LIRMM



Vítor Silva

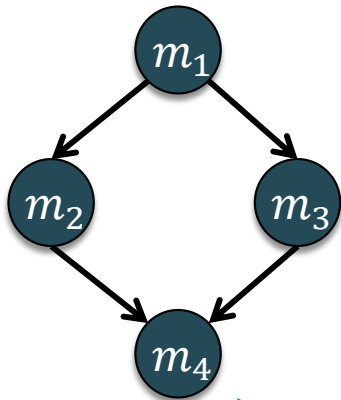


- **Engenheiro de Computação e Informação pela UFRJ**
 - ▣ Iniciação Científica com a Profa. Marta desde o primeiro período
- **Mestre e Doutor em Engenharia de Sistemas e Computação pela COPPE/UFRJ**
- ***Senior Engineering Technologist* na Dell EMC**
- ***Research Engineer* na Snap Inc.**

Cenário de Ciência Computacional e Engenharia (CSE) em Larga Escala

3

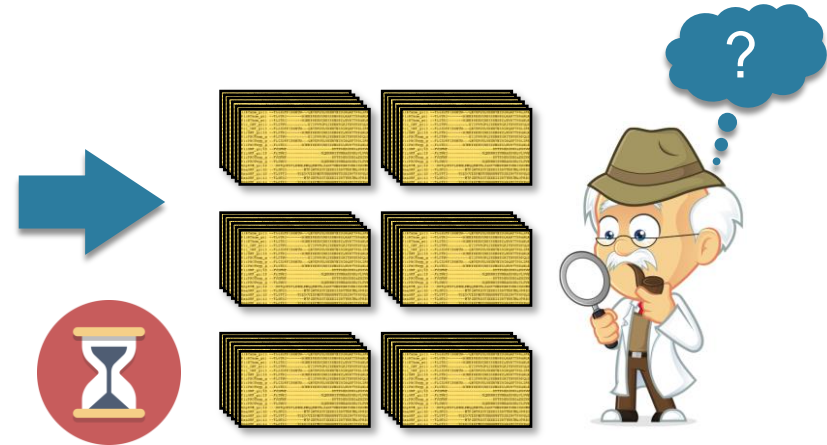
Simulações computacionais



Ambientes de PAD



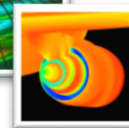
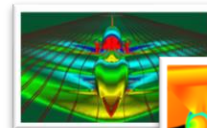
Análise de dados ad-hoc em tempo de execução



Hipótese científica



Parâmetros de entrada



Dados intermediários ou resultados finais

Cenário de Ciência Computacional e Engenharia (CSE)

4

Análise não se resume ao resultado final

- ✓ Analisar parâmetros da simulação e dados intermediários
- ✓ Relacionar dados de diferentes etapas da simulação
- ✓ Realizar ajustes ou mesmo interromper a simulação

Dados intermediários
ou resultados finais

Análise de dados científicos

5

- **Processo ad-hoc e manual**
 - ▣ Desenvolvimento de script para cada análise de interesse

- **Diferentes tipos de consultas:**
 - ▣ Conteúdo de **um arquivo específico** da simulação

 - ▣ Relações de dependência entre **múltiplos arquivos relacionados** pelos programas de simulação
 - Fluxo de arquivos

 - ▣ **Elementos de dados relacionados** a partir de múltiplos arquivos
 - Fluxo de elementos de dados

Alternativas existentes

6

- **Análise de dados científicos** (e.g., FastBit, AQUAdex, PostgresRaw)
 - ▣ Análise *post-mortem* de arquivos de dados “isolados”
(sem proveniência)

- **SGWfC** (e.g., Pegasus)
 - ▣ Análise via captura de dados de proveniência
 - ▣ Conflito com o paralelismo em aplicações de CSE

- **Alternativas baseadas em proveniência** (e.g., noWorkflow, RDataTracker)
 - ▣ Análise via captura de dados de proveniência em granularidade fina
 - ▣ Sobrecarga no desempenho computacional

Alternativas existentes

7

Nenhuma das alternativas permite...

- **a análise de dados de arquivos durante a execução da simulação; e**
- **a extração/indexação de dados científicos**

Problema geral da tese

8

Como acompanhar a execução de simulações computacionais em larga escala?

□ **Desafios**

- Como transformar os dados de simulação em dados passíveis de serem consultados em tempo de execução?
- Como relacionar dados gerados em etapas distintas da simulação?
- Como não interferir no desempenho da simulação?

Hipótese

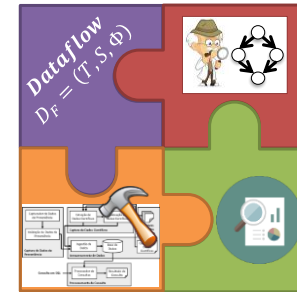
9

Ao prover **uma visão global dos dados** com seus relacionamentos gerados na simulação computacional,
a análise de dados científicos sobre múltiplos arquivos pode ser realizada **durante a execução das simulações**

Solução para o problema geral da tese

10

Como acompanhar a execução de simulações computacionais em larga escala?



□ Desafios

- Como transformar os dados de simulação em dados passíveis de serem consultados em tempo de execução?
- Como relacionar dados gerados em etapas distintas da simulação?
- Como não interferir no desempenho da simulação?



Solução proposta

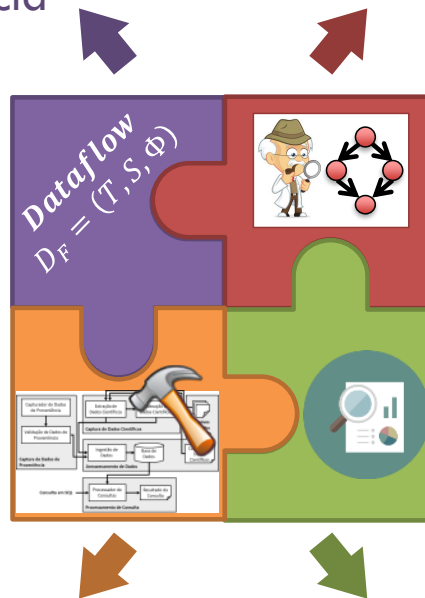
11

Abstração de fluxo de dados

- Nível físico e lógico

Modelo de dados de proveniência
compatível com W3C PROV

**Metodologia para identificar
o fluxo de dados em simulações**



Arquitetura baseada em componentes

- A-Chiron (SGWfC)
- DfAnalyzer (sobrecarga desprezível)

**Extração/indexação de dados e
consulta ao fluxo de dados em
tempo de execução**

Solução proposta aos problemas

12

Abstração de fluxo de dados

- Nível físico e lógico
- Modelo de dados de proveniência compatível com W3C PROV

Como transformar os dados de simulação em dados passíveis de serem consultados em tempo de execução?

Como não interferir no desempenho da simulação?

Arquitetura baseada em componentes

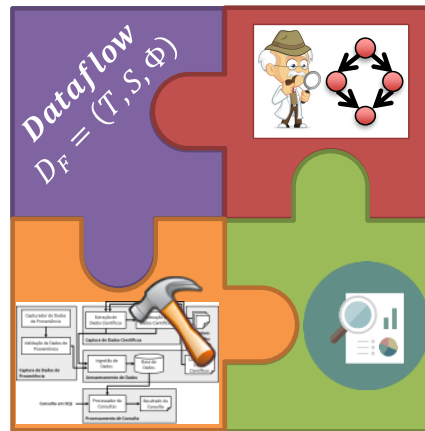
- A-Chiron (SGWfC)
- DfAnalyzer (sobrecarga desprezível)

Metodologia para identificar o fluxo de dados em simulações

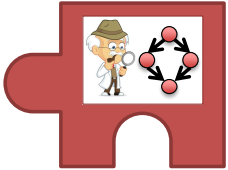
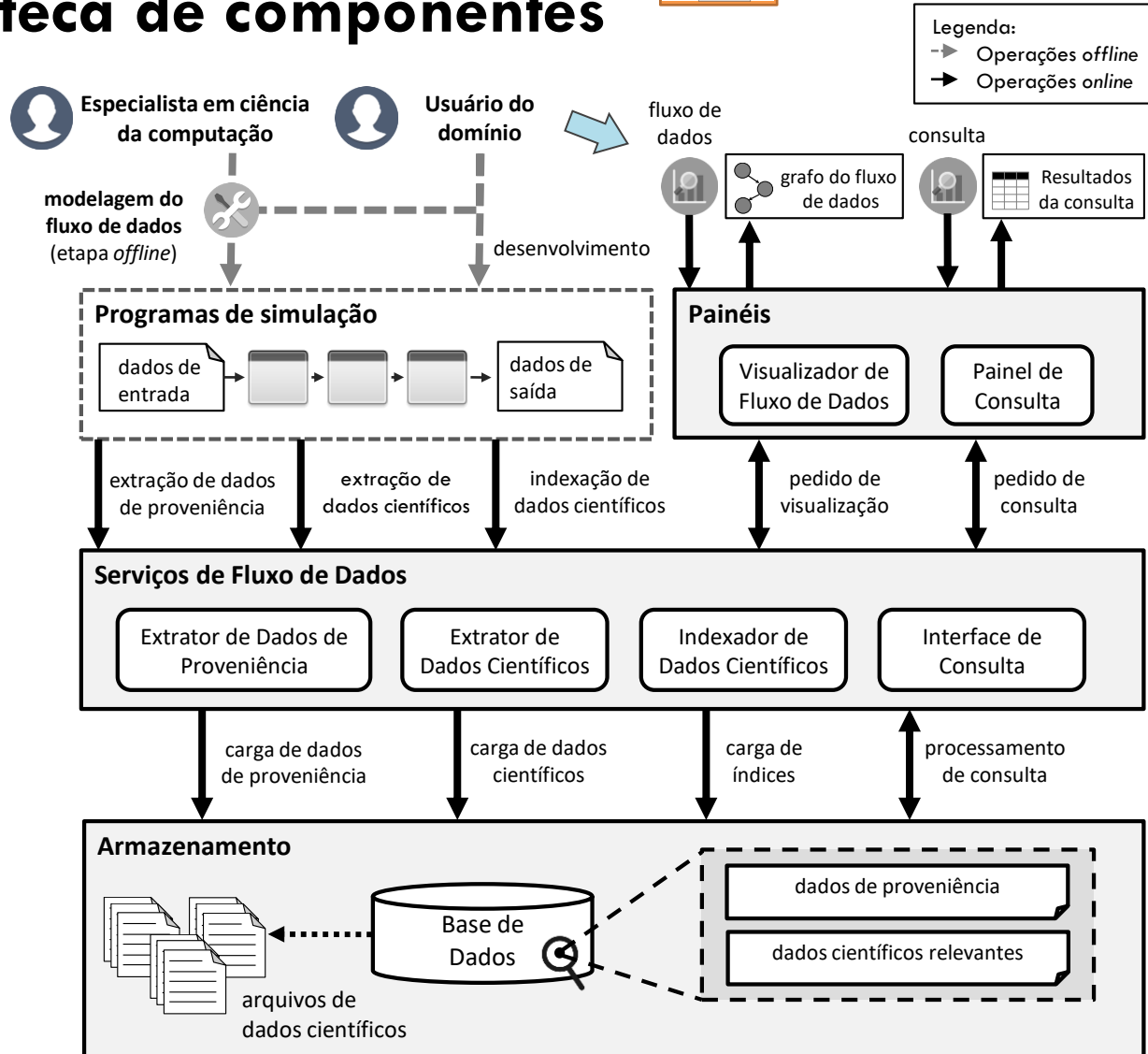
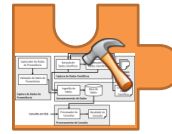
Como transformar os dados de simulação em dados passíveis de serem consultados em tempo de execução?

Como relacionar dados gerados em etapas distintas da simulação

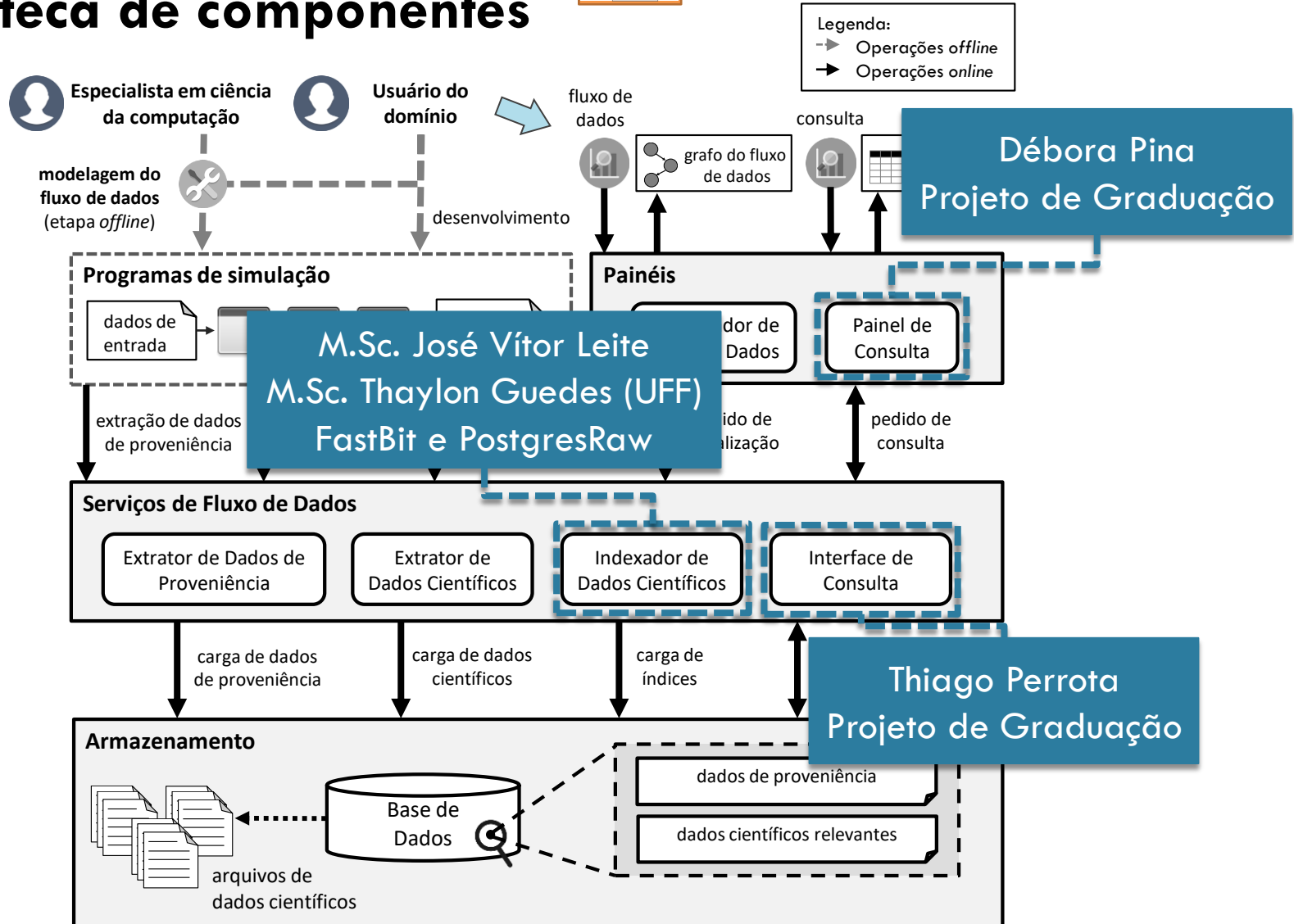
Extração/indexação de dados e consulta ao fluxo de dados em tempo de execução



DfAnalyzer: biblioteca de componentes



DfAnalyzer: biblioteca de componentes

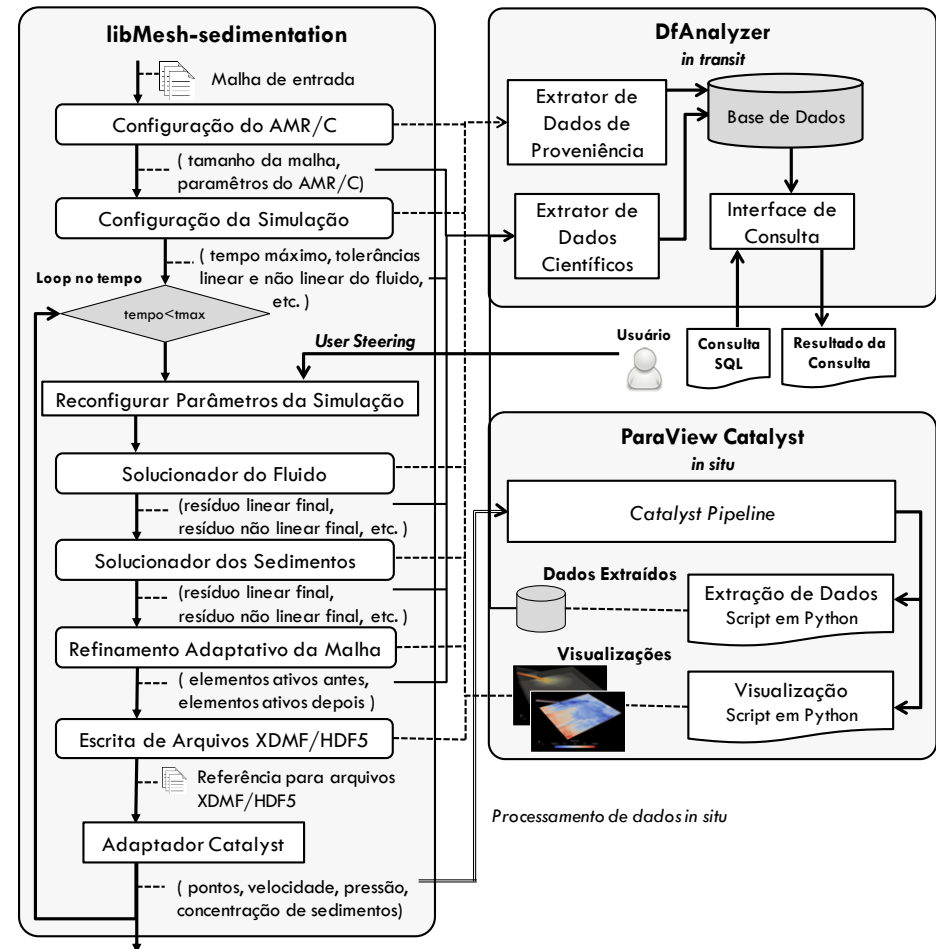
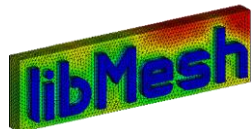
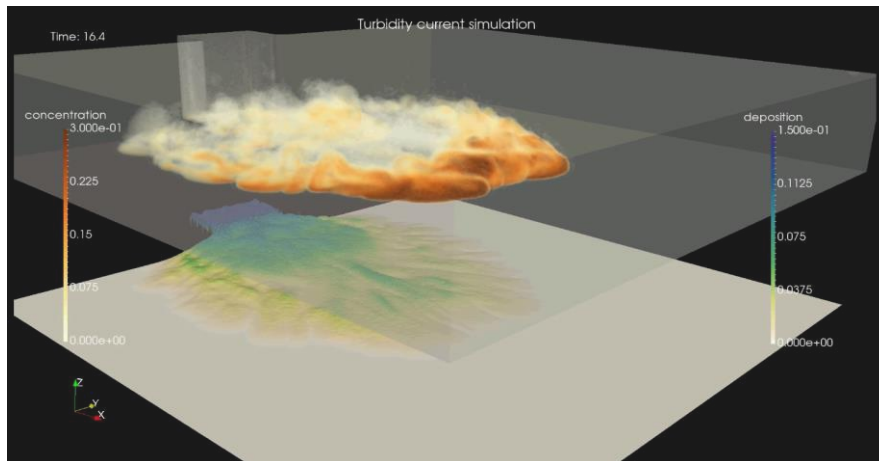


Estudo de caso em correntes turbidíticas usando a DfAnalyzer

15

□ libMesh-sedimentation

- ▣ Tanque real de batimetria
 - Milhões de pontos nas malhas
- ▣ Uso do LoboC (480 cores)



libMesh-sedimentation com DfAnalyzer

16

Contribuição em termos de tempo	Tempo de execução (em segundos)	Tempo de execução (%)
Solucionador do Fluido	75.523,49	50,71%
Solucionador do Sedimento	28.000,50	19,58%
Escrita de dados em arquivos XDMF/HDF5	421,23	0,29%
Extração e visualização de dados <i>in situ</i>	2.175,16	1,52%
Proveniência (DfAnalyzer)	451,70	0,32%
Total	143.029,00	



Tempo de
execução

Sobrecarga desprezível
da DfAnalyzer para a
captura de dados

Tempo de execução
Simulação: 39,73 horas
DfA: 7,53 minutos

Armazenamento
de dados



Visão apenas dos
dados relevantes
para o desenvolvedor
da aplicação

Tipo de dados	Espaço de armazenamento (em GB)	Dados científicos (%)
Arquivos de visualização	0,28	1,21%
Proveniência (DfAnalyzer)	0,38	1,60%
Dados científicos armazenados em arquivos XDMF/HDF5	23,44	-

Conclusões

17

- **Validação da hipótese da tese**
 - ▣ Com a visão global dos dados da simulação, a análise de dados científicos sobre múltiplos arquivos pode ser realizada em tempo de execução

- **Resultados experimentais com a solução desenvolvida**
 - ▣ Interferência desprezível no desempenho da simulação ($< 2\%$)
 - ▣ Análises de dados em múltiplos arquivos em tempo de execução
 - ▣ Adaptações em parâmetros da simulação durante a execução
 - ▣ Confiabilidade dos dados a partir das consultas

Trajetória: Publicações selecionadas

18

Silva, V.; Oliveira, D.; Mattoso, M.
**SciCumulus 2.0: Um Sistema de Gerência
de Workflows Científicos para Nuvens
Orientado a Fluxo de Dados**
SBBD, 2014, Curitiba



Trajetória: Publicações selecionadas

19

Silva, V.; Oliveira, D.; Mattoso, M.
**Exploratory Analysis of Raw Data Files
through Dataflows**
SBAC/PAD Workshop, 2014, Paris. 2014. p. 114



Trajetória: Publicações selecionadas

20

Oliveira, D.; Silva, V.; Mattoso, M.
**How much domain data should be in
provenance databases?**
TaPP, 2015, Edinburgh



Trajetória: Publicações selecionadas

21

Silva, V.; Oliveira, D.; Valduriez, P.; Mattoso, M.

**Analyzing related raw data files
through dataflow**

CCPE, v. 28, p. 2528-2545, 2016



Trajetória: Publicações selecionadas

22

Silva, V.; Leite, J.; Camata, J.; Oliveira, D.;
Coutinho, A.; Valduriez, P.; Mattoso, M.

**Raw data queries during data-intensive
parallel workflow execution**

FGCS, 2017, v. 75, p. 402-422



Trajetória: Publicações selecionadas

23

Camata, J.; Silva, V.; Valduriez, P.; Mattoso, M.; Coutinho, A.

**In situ visualization and data analysis
for turbidity currents simulation**

Computers & Geosciences, 2018,

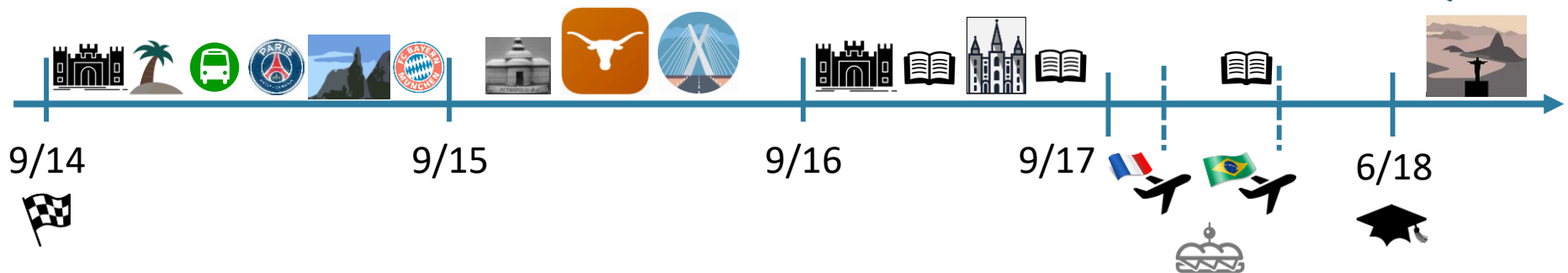
v. 110, p. 23-31



Trajetória: Publicações selecionadas

24

Silva, V.; Oliveira, D.; Valduriez, P.; Mattoso, M.
**DfAnalyzer: Runtime Dataflow Analysis of
Scientific Applications using Provenance**
VLDB Endowment, 2018



Trajectoria: Premiação

25





Agradecimentos



Obrigado!

Uma trajetória de pesquisa em engenharia de dados para aplicações em larga escala

Autor: Vítor Silva PESC • COPPE • UFRJ

Orientadores: Marta Mattoso PESC • COPPE • UFRJ

Daniel de Oliveira IC • UFF

Patrick Valduriez INRIA • LIRMM

