



Gerência de Experimentos Científicos em Nuvens de Computadores: Oportunidades de Pesquisa em Banco de Dados

Daniel de Oliveira - danielcmo@ic.uff.br Marta Mattoso - marta@cos.ufrj.br





Quem nós somos

- Professora associada 4 no Programa de Engenharia de Sistemas e Computação da COPPE/UFRJ
- Bolsista de produtividade em pesquisa do CNPq nível 1C.
- Publicou mais de 200 artigos completos em revistas e congressos.
- Foi diretora de publicações da SBC no período de 2005 a 2007.
- Coordena diversos projetos de pesquisa com financiamento do CNPq, Capes, INRIA e FAPERJ
- Atua principalmente nas áreas de distribuição e paralelismo em bancos de dados, workflows científicos em ambientes de paralelismo e gerência de dados de proveniência.

Profa. Marta Mattoso, D.Sc.





Quem nós somos

- Professor adjunto do Instituto de Computação da Universidade Federal Fluminense (UFF) desde 2013
- Recebeu o grau de Doutor em Ciências pela UFRJ em 2012.
- Publicou mais de 50 artigos em periódicos indexados e em congressos nacionais e internacionais.
- Seus interesses de pesquisa incluem bancos de dados, computação em nuvem, gerência de workflows científicos, paralelismo de dados, bioinformática e mineração de dados.
- É membro da ACM, IEEE e SBC.

Prof. Daniel de Oliveira, D.Sc.





Grupo de Pesquisas

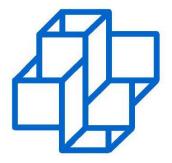
- Marta Mattoso, D.Sc. COPPE/UFRJ
- Daniel de Oliveira, D.Sc. IC/UFF
- Eduardo Ogasawara, D.Sc. CEFET/RJ
- Pós-doutorandos
 - Kary Ocaña, D.Sc. COPPE/UFRJ
- Alunos de D.Sc.
 - Jonas Dias
 - Flavio Costa

Colaborações Atuais

- INRIA
 - Patrick Valduriez



- LNCC
 - Fabio Porto
 - Luiz Gadelha Jr.



Laboratório Nacional de Computação Científica

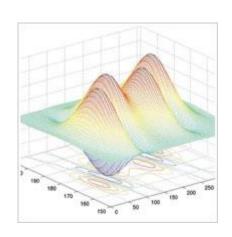
Roteiro do Tutorial

- Motivação
- Workflows Científicos
- · Nuvens de Computadores
- Proveniência de Dados
- Máquinas de Workflow para Nuvem
- · Aplicação de Proveniência em e-Science
- Demo

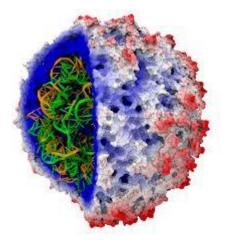


e-Science

 "A combinação de pesquisa em computação e modelagem matemática que proporciona a aceleração da pesquisa em outras áreas do conhecimento"

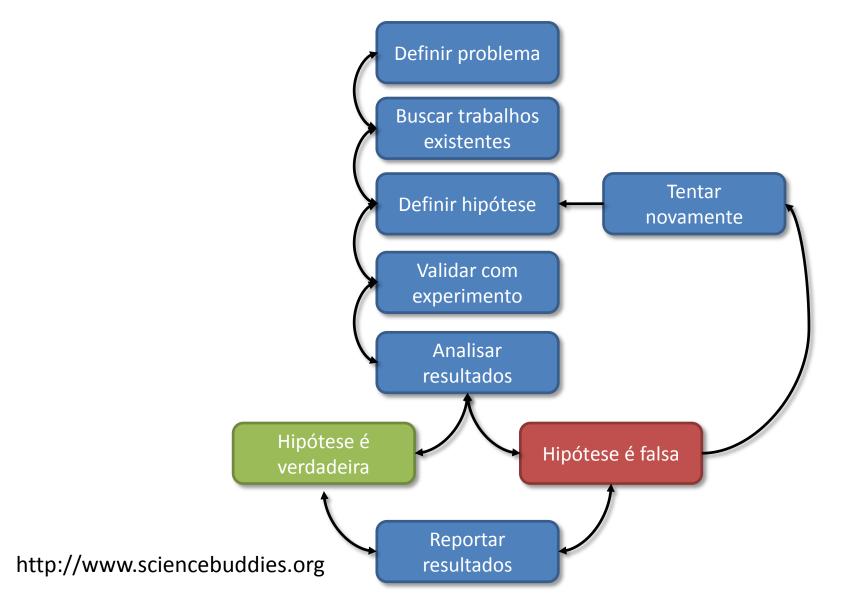






Fonte: www2.bioetanol.org.br/escibioenergy

Método Científico

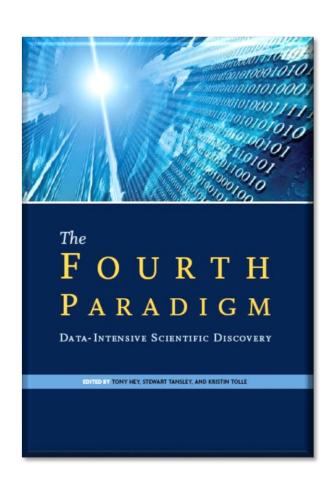


Método Científico in silico

- "Ten years ago, the problem was how to obtain data."
- Nos dias de hoje o gargalo se encontra em determinar quais estratégias computacionais serão utilizadas para que o cientista possa analisar esse grande volume de dados produzido e coletado
- Dados geralmente são heterogêneos e distribuídos

Fonte: www2.bioetanol.org.br/escibioenergy

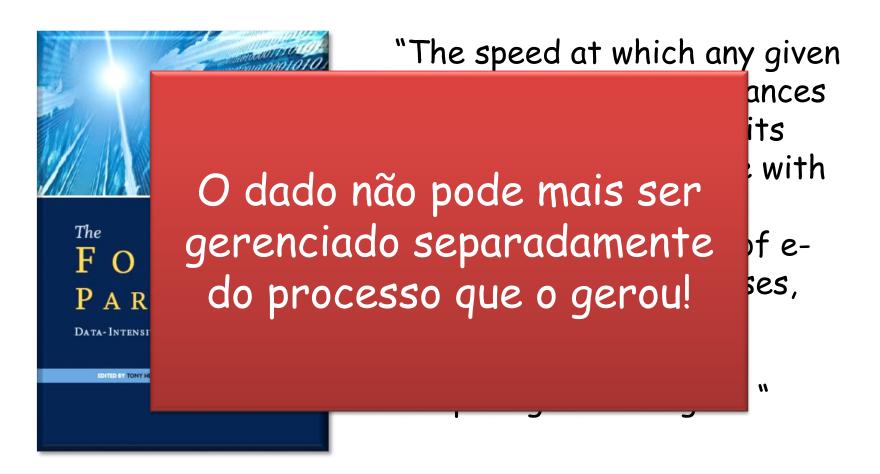
O quarto paradigma



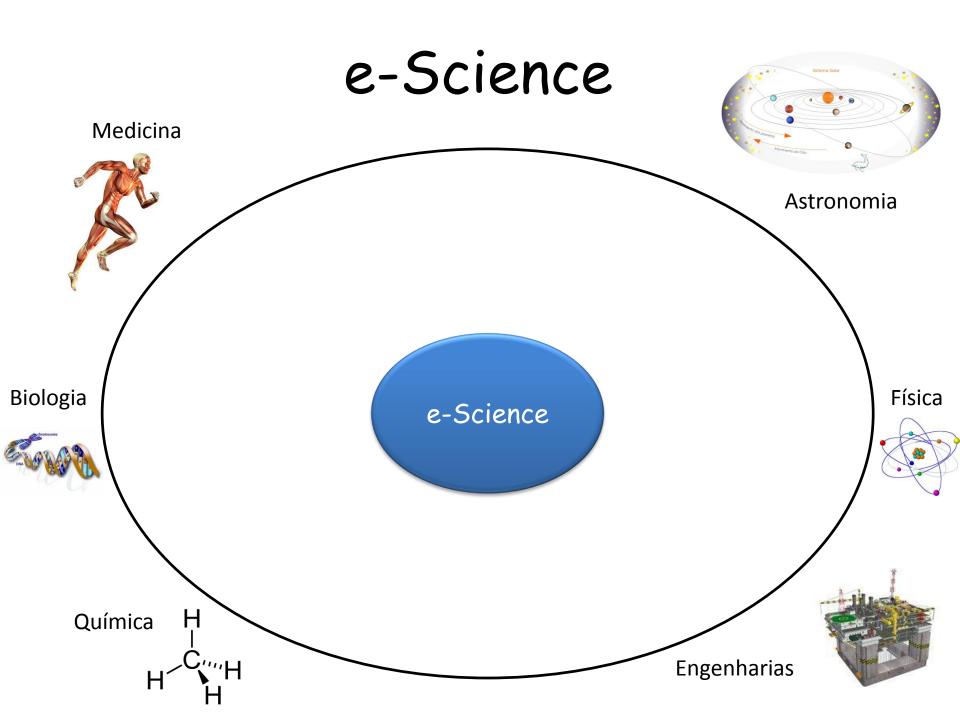
"The speed at which any given scientific discipline advances will depend on how well its researchers collaborate with one another, and with technologists, in areas of e-Science such as databases. workflow management, visualization, and Cloud computing technologies."

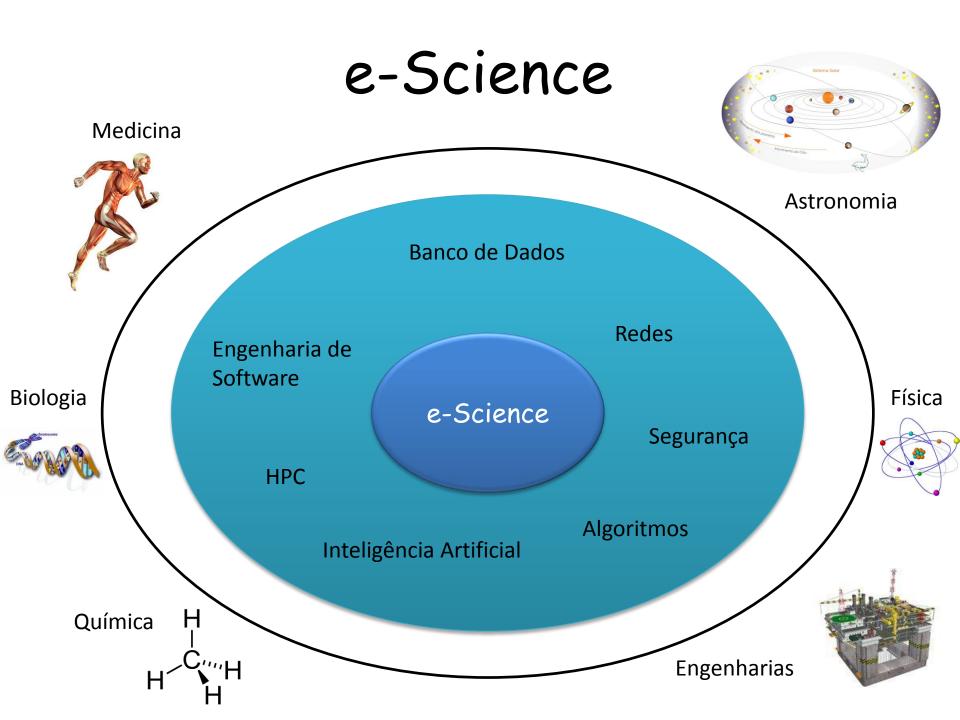
http://research.microsoft.com/en-us/collaboration/fourthparadigm/

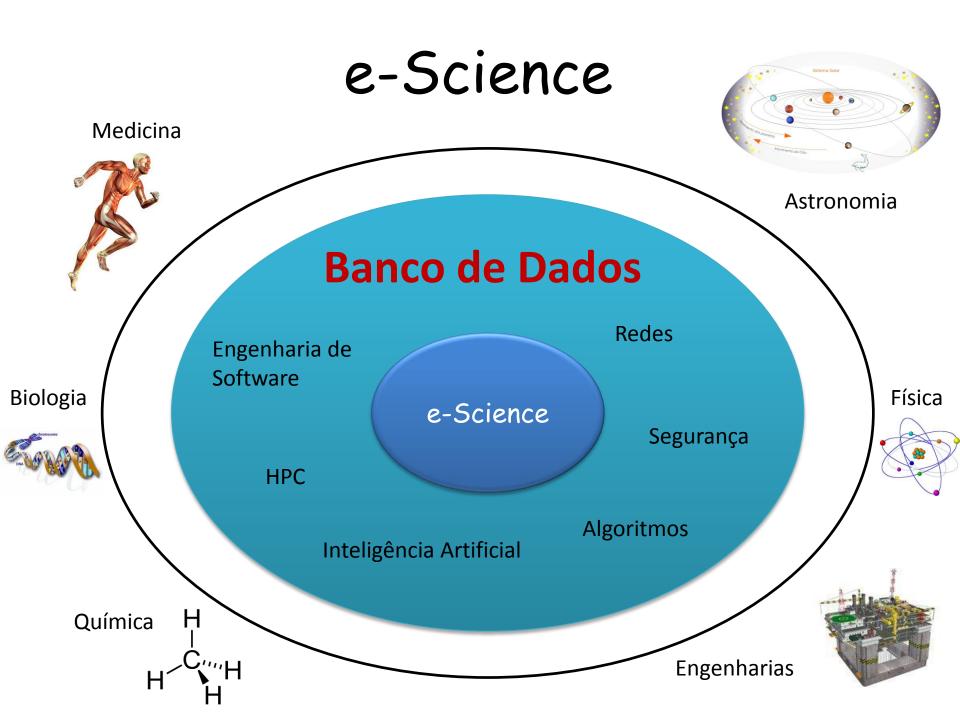
O quarto paradigma



http://research.microsoft.com/en-us/collaboration/fourthparadigm/









e-Science no Reino Unido

- E-Science é vital para a exploração do poderoso aparato científico da próxima geração
 - Telescópios
 - Sensores
 - Satélites







- Cada um desses aparatos podem gerar vários terabytes de dados diariamente
- Os usuários podem estar dispersos geograficamente ao redor do globo
- A e-Science deve construir a infraestrutura necessária para este exploração



e-Science no Brasil

 São Paulo School of Advanced Science on e-Science for Bioenergy Research



- Grandes Desafios da Computação no Brasil
 - http://www.gd2.ufam.edu.br/
- BreSci Brazilian e-Science workshop
 - http://www.ic.ufal.br/csbc2013/noticias/e-science



e-Science no Brasil



Página inicial » A Instituição » Eventos

Latin American eScience Workshop 2013

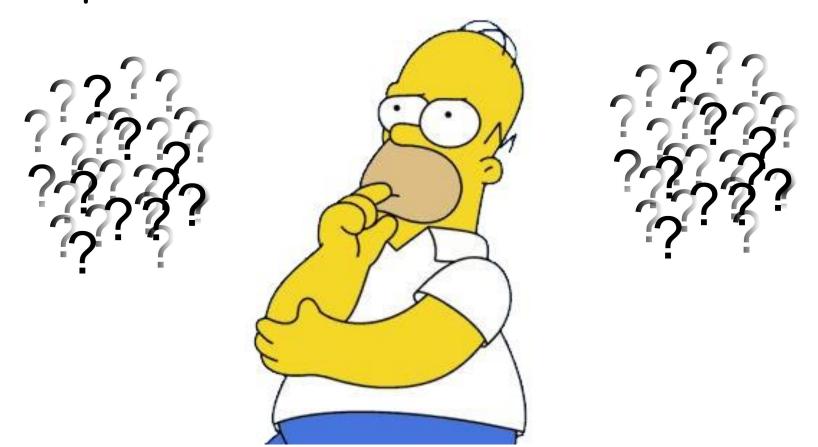


Sponsored jointly by Microsoft Research and FAPESP

May 13-15, 2013

e-Science

 Ok, mas por que e-Science é importante?

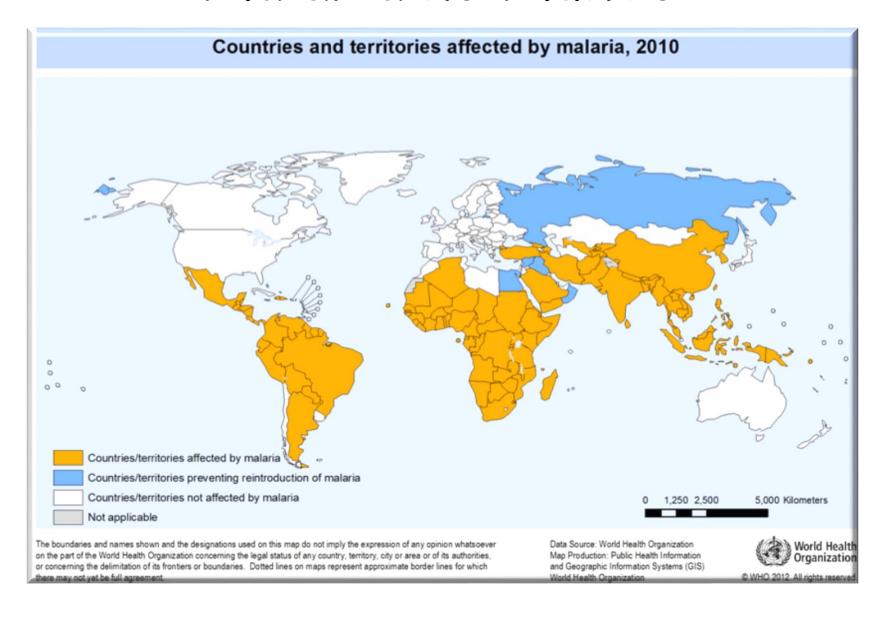


Experimento de Motivação

- DTN compreendem infecções parasitárias, virais e bacterianas que atacam especialmente as populações de baixa renda nas regiões em desenvolvimento da África, Ásia e América Latina
- 17 DTN em 149 países afetam 1 bilhão de pessoas (1/7 da população mundial)
 - > Doença de Chagas
 - > Leishmaniose
 - > Hanseníase
 - > Esquistossomose
 - > Infecções por trematódeos transmitidas por alimentos
 - > Helmintíases transmitida pelo solo

- Cisticercose
- > Dengue
- > Dracunculíase
- > Equinococose
- Oncocercose
- > Raiva
- > Tracoma
- > Úlcera de Buruli
- Tripanossomíase humana africana
- A OMS reporta progressos contra a malária como uma prioridade

Malária no Mundo



Malária no Brasil

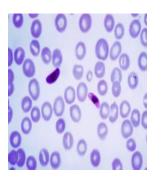


Fonte: Sistema de Informação de Vigilância Epidemiológica-Malária/Secretaria de Vigilância em Saúde/Ministério da Saúde

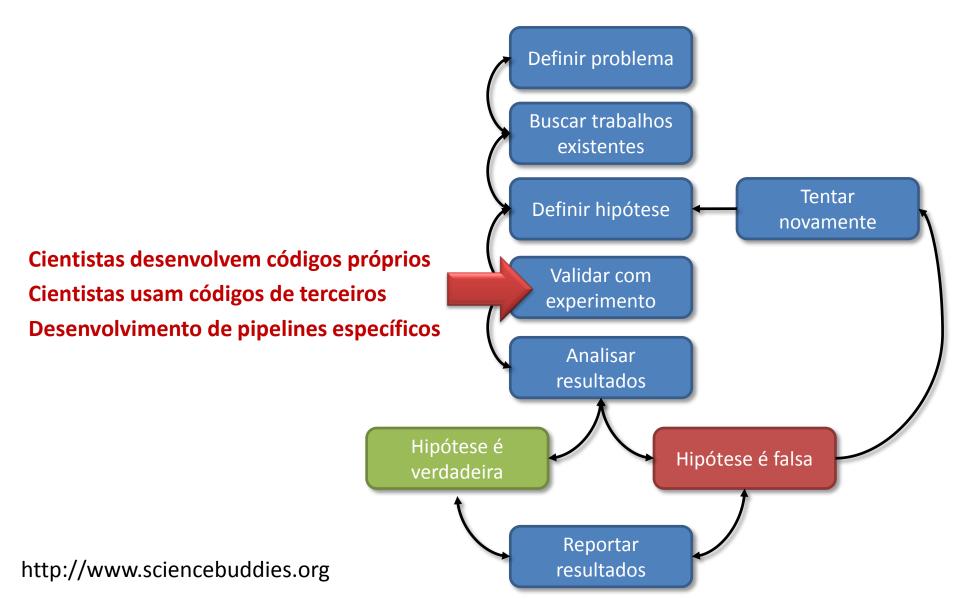
Mapa de risco da malária por Município de Infecção. Brasil 2011.

Cisteíno Proteases em Plasmodium

- · Plasmodium falciparum gera os casos mais graves de malária
- As drogas antimaláricas apresentam uma crescente resistência
- Evidências indicam que cisteíno proteases (CP) têm um papel essencial no tratamento médico da malária
- Projetos de descoberta de drogas
 - →alvos fármaco-terapêuticos promissores
 - →inibidores antimaláricos
 - >papel importante no ciclo de vida do parasita
 - →papaína e ubiquitina

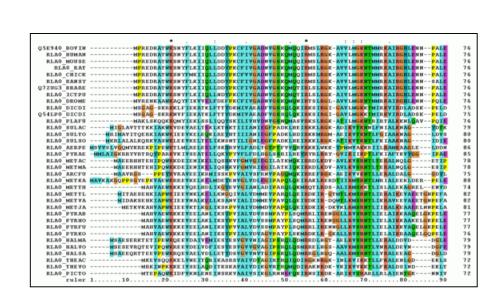


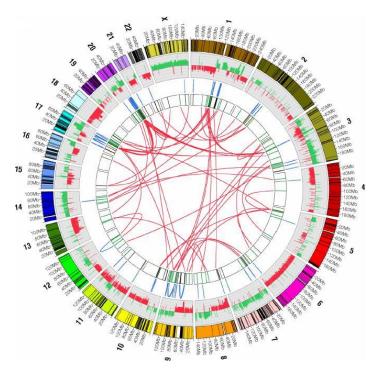
Testes para Análise de Cisteíno Proteases



1º Passo - Genômica Comparativa

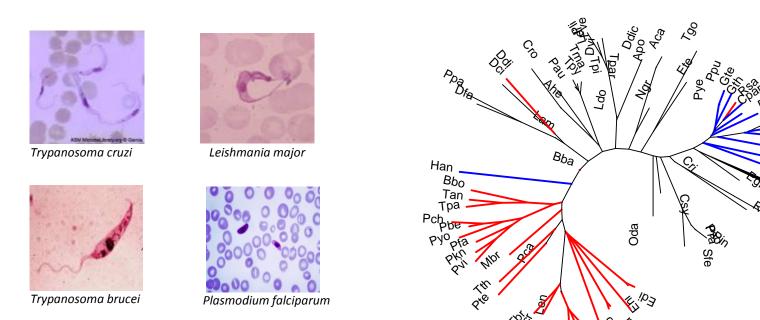
 Um alinhamento múltiplo de sequência (AMS) identifica regiões de similaridade que podem ser consequência de relações funcionais, estruturais ou evolutivas entre as sequências





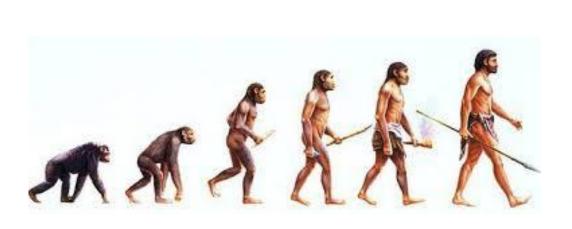
2º Passo - Filogênia

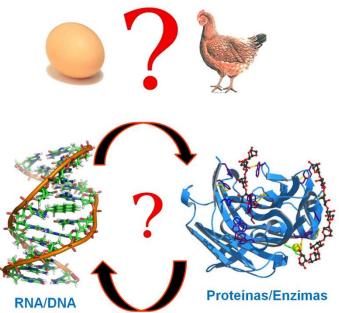
 Faz uso extensivo de AMS para a construção de árvores filogenéticas usados para inferir relações evolutivas entre genes homólogos



3º Passo - Evolução Molecular

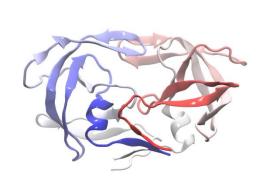
 Fornece base para inferências evolutivas e biológicas, e evoluem devido ao crescimento explosivo de dados genômicos e novas abordagens computacionais e estatísticas

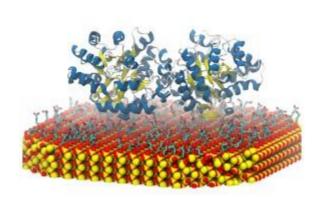


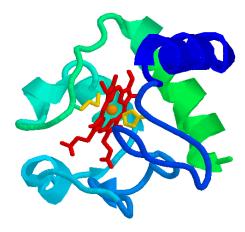


4º Passo - Modelagem Molecular

- Construção do modelo de resolução atômico (modelo 3D) da sequência da proteína (alvo) e uma estrutura 3D experimental de proteínas homólogas (molde)
- Uso de métodos de Genômica Estrutural e.g., Computer-Aided Drug Design (CADD)







Pipelines Farmacofilogenômica

(A) Genômica Comparativa

Construção do AMS Sequências do MAFFT **GOLD** Conversão de Formato ReadSea Busca por Ortologia OrthoSelect do AMS Não Rejeitar Comparação de Perfis **HHSearch** Seleção do Modelo ModelGenerator de Substituição **MEROPS** Busca por Ortologia Rejeitar **PSIBLAST** Construção da Árvore RAxML, Phylip. Refseq Filogenética **MrBayes** Rejeitar Verificação da Anotação Verificação Visualização da Árvore **MEGA** Filogenética 2 Marcadores de Genes

Construção d

Não

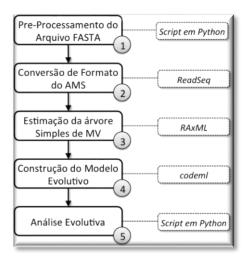
e Refinamen

Avaliação da Predição

Modelo

(B) Filogenia

(C) Evolução Molecular



Atribuição das Dobras
e Seleção do Molde

Alinhamento das
Sequências Alvo-Mold

(D)

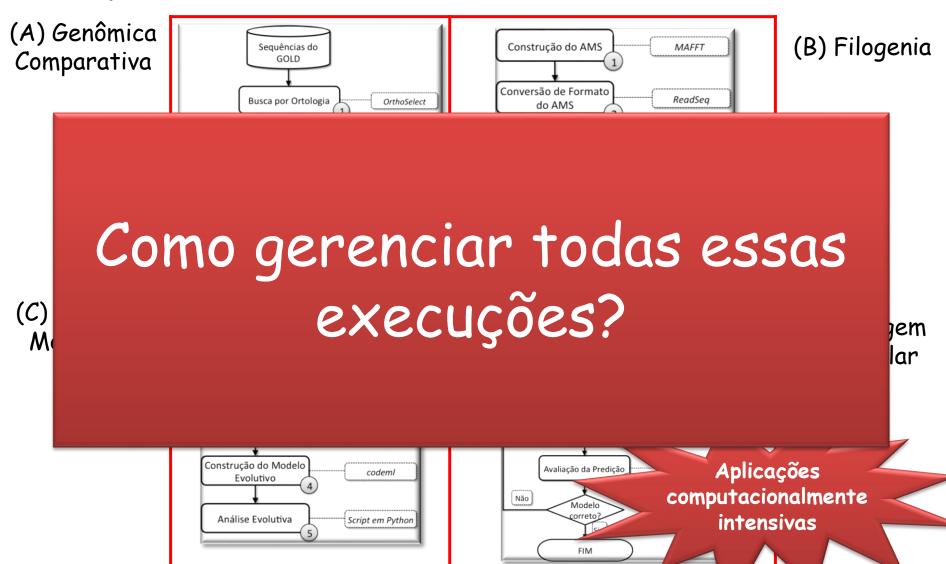
Modelagem

Molecular

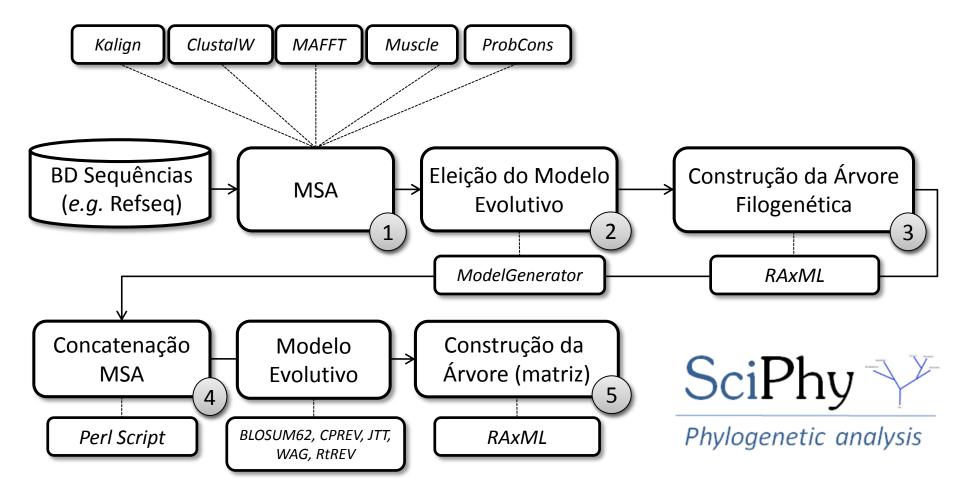
MODELLER

Aplicações computacionalmente intensivas

Pipelines Farmacofilogenômica



Análise de Cisteíno Proteases Modelada como Workflow Científico



Fonte: Ocaña K.A.C.S., Oliveira D., Ogasawara E., Dávila A.M.R., Lima A.A.B., and Mattoso M. SciPhy: A Cloud-Based Workflow for Phylogenetic Analysis of Drug Targets in Protozoan Genomes. Springer, pp. 66-70, 2011.

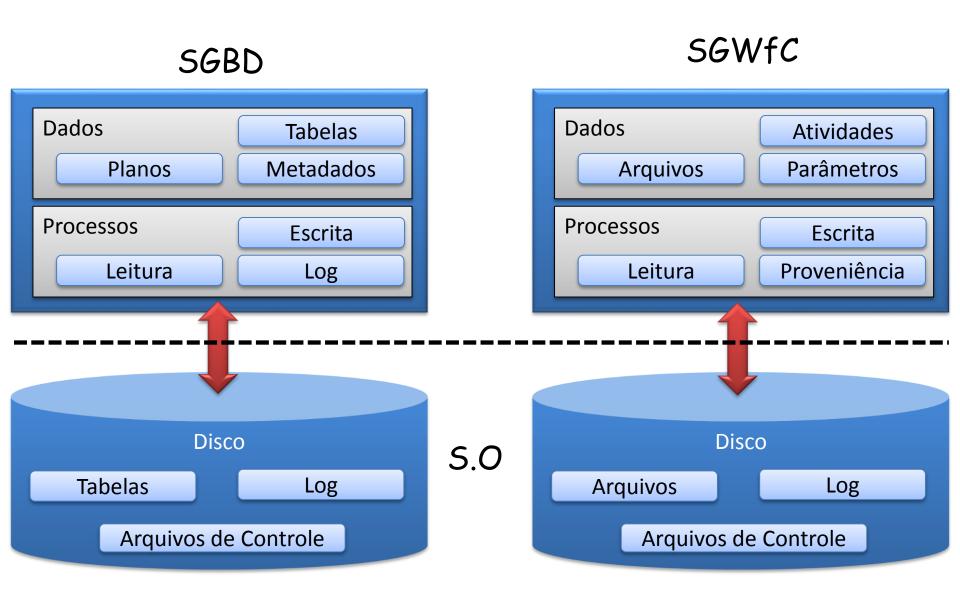
Workflows Científicos

- Workflows científicos são abstrações que representam a cadeia de atividades dentro de um experimento científico
- Eles são gerenciados por Sistemas de Gerência de Workflows Científicos (SGWfC)
- Diversas variações do workflows dentro do mesmo experimento
- Essas variações incluem a alteração de dados de entrada e parâmetros
- · Proveniência é uma questão fundamental

Sistemas de Gerência de Workflows Científicos

- Definem e executam os workflows
- · Proveem execução eficiente
- Controlam falhas garantindo a integridade
- Acessam/Armazenam/Consultam dados usando SGBD
- · Rastreiam a proveniência dos dados

Analogia com Banco de Dados



Abordagens Existentes para Execução de Workflows Científicos

 Modelo de execução com paralelismo nativo Ex.: Swift, Pegasus, Triana, SciCumulus









Motor sequencial
 Ex.: Kepler, VisTrails, Taverna







 Motor sequencial e camada de paralelismo Ex.: VisTrails+Hadoop, Kepler+Hadoop



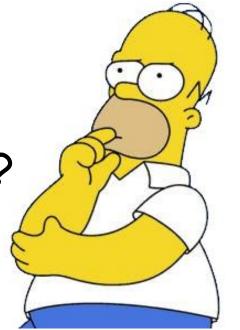
Execução de uma Análise de Cisteíno Protease

- Execução demanda grande tempo de processamento
- · Cientistas necessitam variar:
 - Parâmetros
 - Métodos
 - Imagens
 - Arquivos de Dados
 - Configurações
- Como executar essas análises?

Execução de uma Análise de Cisteíno Protease

- Execução demanda grande tempo de processamento
- · Cientistas necessitam variar:
 - Parâmetros
 - Métodos
 - Imagens
 - Arquivos de Dados
 - Configurações
- Como executar essas análises?

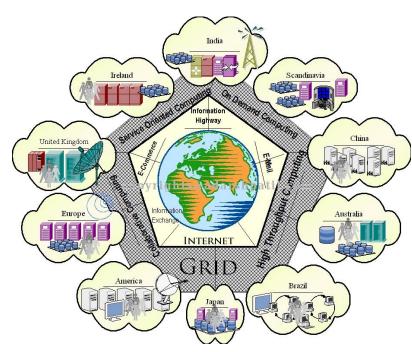
Paralelismo!



Em que ambiente executar?



Clusters



Grades



Desktop

Em que ambiente executar?

- Clusters podem ser financeiramente custosos
- Grades demandam um esforço grande de configuração
- Desktops podem demandar um tempo exageradamente grande para a execução
- Tempo de execução versus Tempo de Análise dos Dados

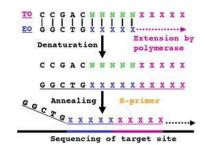
Alternativa: Análise de Cisteíno Protease em Nuvens



Dados de entrada para Realizar análises filogenéticas: Sequências de DNA e RNA de diversos organismos...



1. Alinhamento de sequências (MAFFT, ProbCons, ClustalW, Muscle e Kalign) e conversão



Produz uma grande quantidade de dados ... (sequências alinhadas)...





3. Geração das árvores filogenéticas (RA×ML)

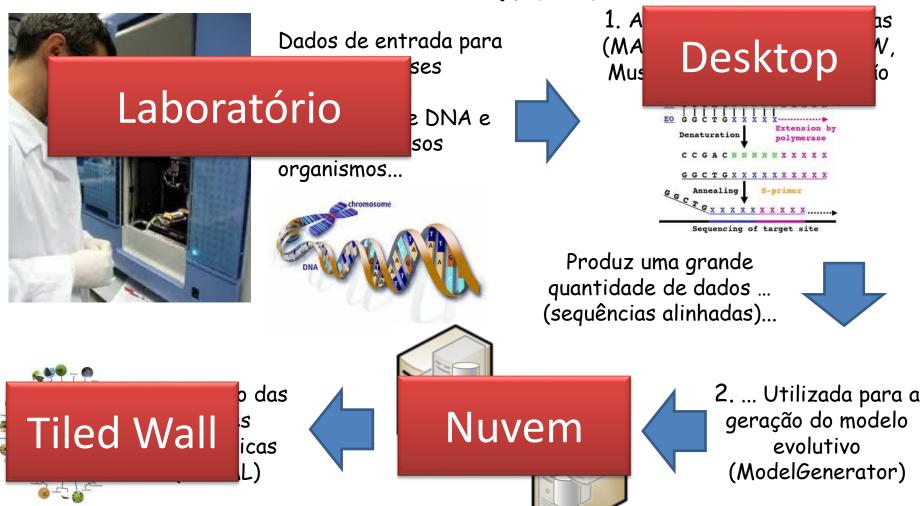






2. ... Utilizada para a geração do modelo evolutivo (ModelGenerator)

Alternativa: Análise de Cisteíno Protease em Nuvens



Execução de uma Análise de Cisteíno Protease em Nuvens

- Utilização da nuvem faz com que cientistas não tenham que comprar equipamentos
- Execução paralela reduz o tempo de processamento
- Mas.....





Nem tudo que reluz é ouro.



Problemas na Execução de Workflows em Nuvens



Nova Visão: Análise de Cisteíno Protease em Nuvens



Dados de entrada para Realizar análises 1. Alinhamento de sequências (MAFFT, ProbCons, ClustalW, Muscle e Kalian) e conversão

Qual árvore corresponde a uma determinada sequência de entrada?

O programa MAFFT é melhor que o ProbCons? Qual e-value foi escolhido?

Como recuperar uma execução que falhou? Como monitorar uma execução a distância?

árvores filogenéticas (RAxML)



geração do modelo evolutivo (ModelGenerator)

'a a

Execução de uma Análise de de Cisteíno Protease

- Execução paralela reduz o tempo de processamento
- · Utilização da nuvem faz com que cientistas não

Como relacionar todos esses recursos?

- Arquivos de Dados
- Configurações

Importância da Proveniência na Ciência

- Interpretar e reproduzir dados
- Verificar se um determinado experimento foi executado seguindo os procedimentos
- Identificar as entradas e saídas dos experimentos
- Garantir qualidade dos dados
- Rastrear quem executou um experimento e quem é responsável pelos resultados

"Provenance is as (or more) important as the results" (Juliana Freire e Susan Davidson, SIGMOD 2008)

Desafios na Execução de *Workflows* em Nuvens

- Ambiente "caixa-preta" e em constante mudança
 - Variação de desempenho das máquinas virtuais
 - Mudanças constantes no ambiente
 - Falhas constantes nas máquinas virtuais
- Diversos tipos de máquinas virtuais passíveis de utilização
 - Qual tipo utilizar?
- Transferência de dados
 - Como evitar sobrecarga na transferência
- Como monitorar os workflows?

Questão

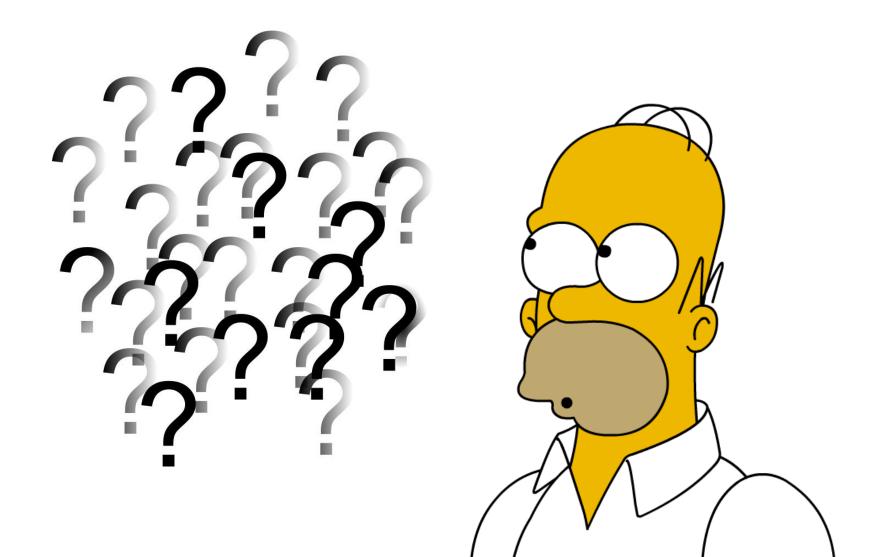
• "...considerando os desafios na execução de workflows em nuvem é possível prover capacidade de alto desempenho e a gerência distribuída do workflow científico por meio da adoção de soluções que considerem as características únicas do ambiente, distribuam as atividades e monitorem a execução das mesmas em várias máquinas virtuais?"

Hipótese

• "Se for adotada uma infraestrutura adaptativa para o escalonamento, despacho, monitoramento e captura de proveniência distribuída das atividades paralelas de um workflow científico então podemos prover a capacidade de alto desempenho, monitoramento e controle que são necessárias aos cientistas."



Em que podemos nos inspirar na área de Banco de Dados?



Em que podemos nos inspirar na área de Banco de Dados?

- Representação da consulta por meio de álgebra relacional ≅ representação do workflows por meio de álgebra de workflows
- Otimização do plano de execução de consultas ≅ otimização de workflows científicos

Em que podemos nos inspirar na área de Banco de Dados?

- Escalonamento das atividades no ambiente de nuvem de acordo com uma álgebra de workflows
- · Distribuição dos dados de proveniência
- Monitoramento



Roteiro do Tutorial

- Motivação
- · Workflows Científicos
- · Nuvens de Computadores
- Proveniência de Dados
- Máquinas de Workflow para Nuvem
- · Aplicação de Proveniência em e-Science
- Demo



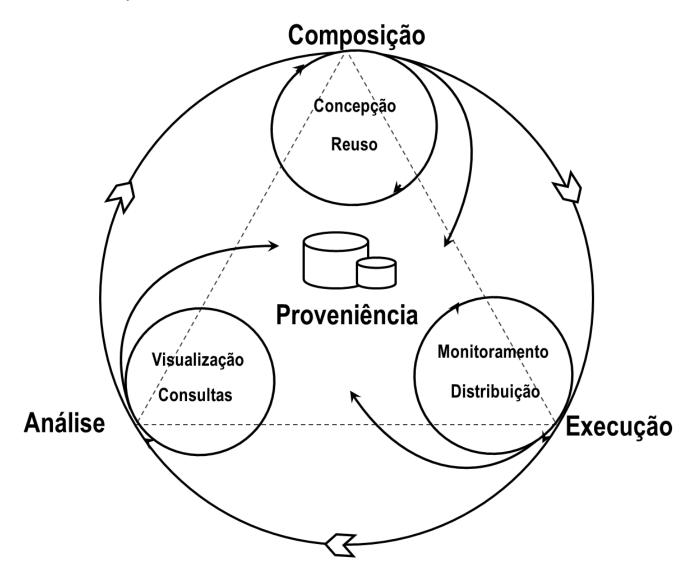
Experimento Científico

- Uma das formas utilizadas pelos cientistas para apoiar a formulação de novas teorias
- Possui um ciclo de vida com três etapas:
 - Composição
 - Execução
 - Análise





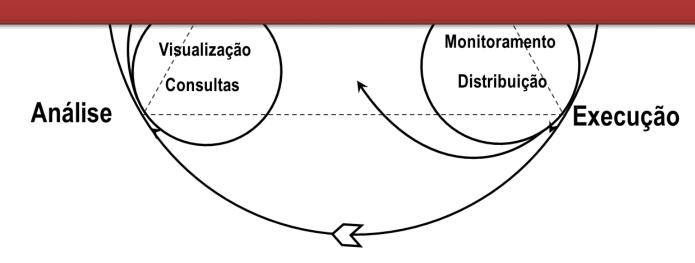
Nossa visão do ciclo de vida do experimento científico



Nossa visão do ciclo de vida do experimento científico



O experimento pode ser composto pela execução de múltiplos workflows científicos



Workflow

A automação de um processo de negócio, completo ou apenas parte dele, através do qual, documentos, informações ou tarefas são transmitidas de um participante a outro por ações, de acordo com regras procedimentais. (WFMC, 1995)

 Provê a abstração necessária para descrever uma série de processos estruturados e suas atividades, oferecem um contexto robusto de resolução de problemas e promovem o uso efetivo e otimizado dos recursos computacionais

Workflow

- Uma coleção de atividades organizadas para acompanhar algum experimento (processo de negócio).
- As atividades ou tarefas são os componentes de software independentes que implementam alguma funcionalidade e são executadas por um ou mais sistemas de softwares.
- Exemplos de atividades incluem executar um programa, transformar um arquivo ou atualizar um banco de dados.

Workflow

- Um workflow define a ordem de execução dessas atividades ou as condições em que essas atividades serão executadas e a sua eventual sincronização.
- Os dados de entrada e saída das atividades (variáveis) são definidos como o fluxo de dados do workflow.

SGWf e Máquina de Workflow

- Sistema de gerência de workflows SGWf (WfMS - Workflow Management Systems)
 - Software que fornece toda a infraestrutura para definir, executar e monitorar workflows.
- Máquina de Workflow
- Na utilização deste tipo de software é importante para separar a definição do workflow, da máquina (i.e. engine) encarregada de executar o mesmo

Especificação do Workflow

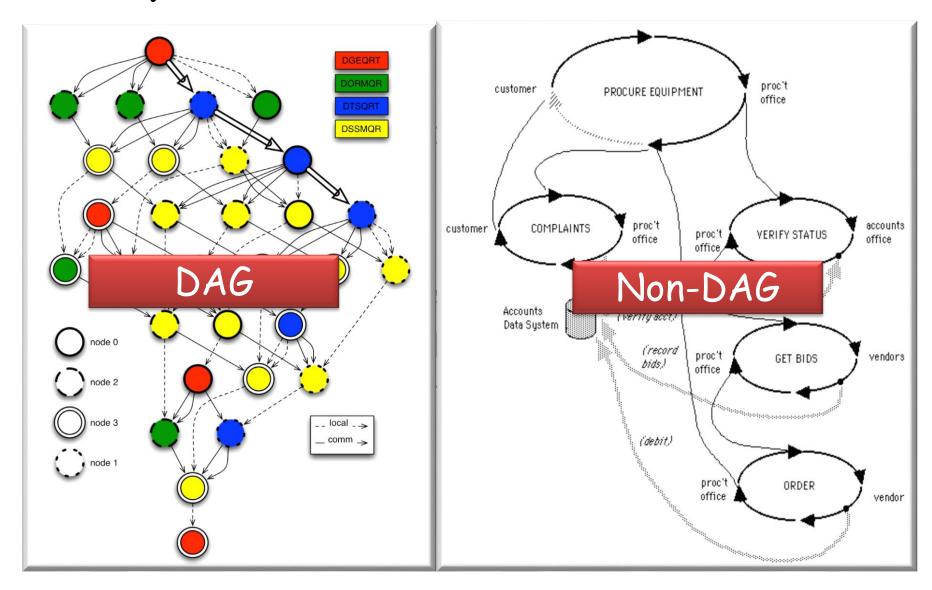
· DAG

 Contém estruturas do tipo sequenciais, paralelas ou livres.

Non-DAG

- Incluem estruturas de iteração
- Permitem que os workflow executem tarefas ou sub-workflows de forma repetida.

Especificação do Workflow



Diversidade de Sistemas de Gerência de Workflows

- Askalon
- Chiron
- Kepler
- Pegasus
- SciCumulus
- Swift
- Taverna
- VisTrails















Diversidade de Sistemas de Gerência de Workflows

- Askalon
- Chiron
- Kepler
- Pegasus
- SciCumulus
- Swift
- Taverna
- VisTrails















SGWf...um caldeirão de tecnologias

- Muitos Modelos...
 - Computação control flow, data flow, pipelines baseados em scripts
 - Interação- interativos, batch
 - Representação de Wf DAG, non-DAG
 - Adaptabilidade dinâmicos, estáticos
- · Muitos Ambientes...
 - Componentes grid, cluster, **nuvem**, etc
 - Natureza- open-source, comerciais
 - Escala-long running, data streaming, etc...
- Muitas tecnologias....
 - centralizado, distribuído, local, Web...



Roteiro do Tutorial

- Motivação
- Workflows Científicos
- · Nuvens de Computadores
- Proveniência de Dados
- Máquinas de Workflow para Nuvem
- · Aplicação de Proveniência em e-Science
- Demo



Computação em Nuvem

- Uma das opções para prover um ambiente de alto desempenho para cientistas é executar seus experimentos em nuvens de computadores
- Se baseia na ideia de prover recursos (software e hardware) utilizando o ambiente Web para uma gama (virtualmente infinita) de usuários
 - Elasticidade de Recursos
 - Virtualização de Recursos
 - Personalização
 - Disponibilidade



Quem já usa Computação em Nuvem?



52% das grandes empresas já utilizam cloud computing, diz Nicholas Carr



75% das grandes empresas no Brasil já usam cloud computing, aponta estudo



Uso de computação na nuvem cresce 40% em 12 meses no Brasil



FAPEMIG lança edital para *cloud computing* 30 de Dezembro de 2011

Quem já usa Computação em Nuvem?

♦ Home > Tecnologia > Cloud veio para ficar, diz IDC

Cloud veio para ficar, diz IDC



19/10/2010 9:15, Marcelo Bernstein, com agências - de Nova York

TI tem investido em cloud e virtualização, não em PCs

Embora o mercado de computadores tenha crescido 2,6%, o número é menor que a previsão inicial de 2,9%, aponta a IDC.

Shane O'Neill, CIO/EUA

Publicada em 15 de julho de 2011 às 12h01

Brasil ainda está na primeira fase de adoção da nuvem

Convergência Digital - Hotsite Cloud Computing

:: Por Fábio Barros* :: 28/07/2011

TI | 05/10/2011 14:06

>> Compartilhar:









Computação em nuvem poderá movimentar US\$ 1,45 bi em 2015

Segundo a empresa de pesquisa de mercado Gartner, a receita das empresas de serviços de armazenamento aumentou 56% em 2011

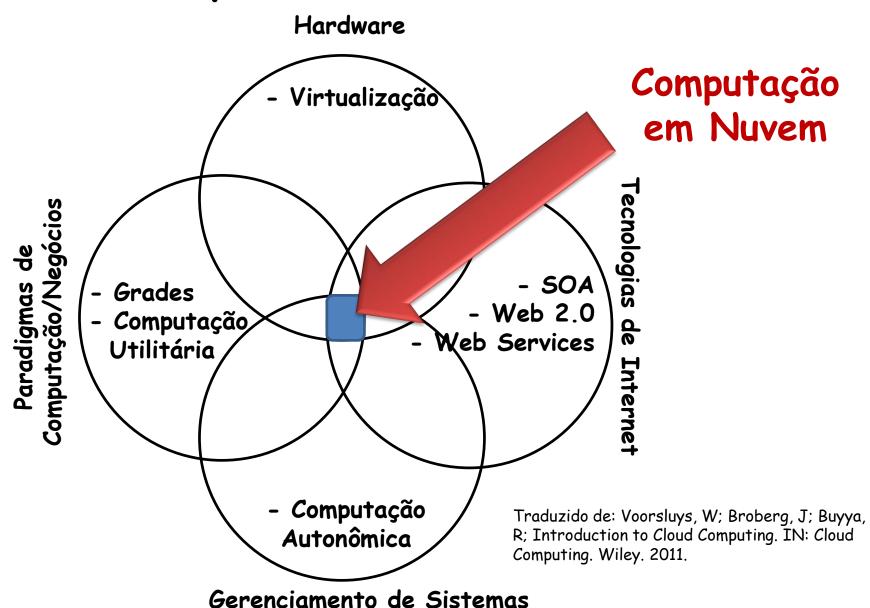
Cloud computing reduzirá consumo de energia de data centers em 31% até 2020

Convergência Digital - Hotsite Cloud Computing :: Da redação :: 21/09/2011

Computação em Nuvem: Definição

- Diversas definições: não há consenso!
- Mais de 20 definições em Vaquero (2009)
- Segundo Ian Foster:
 - "...um paradigma de computação em larga escala que possui foco em proporcionar economia de escala, em que um conjunto abstrato, virtualizado, dinamicamente escalável de poder de processamento, armazenamento, plataformas e serviços são disponibilizados sob demanda para clientes externos através da Internet."

Computação em Nuvem



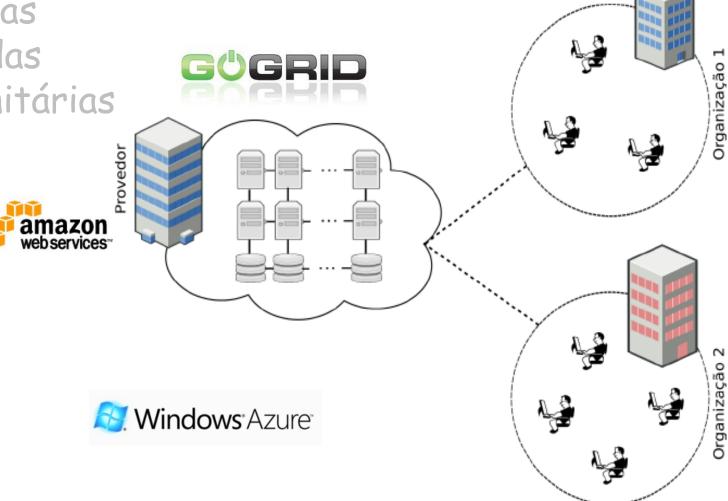
Tipos de Nuvem

Públicas

Privadas

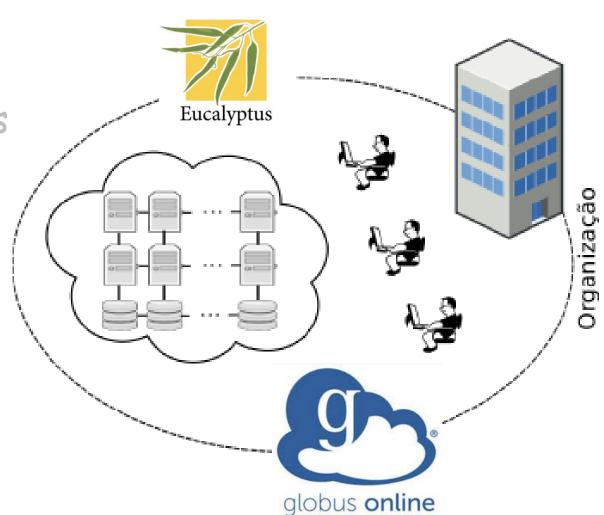
Híbridas

Comunitárias



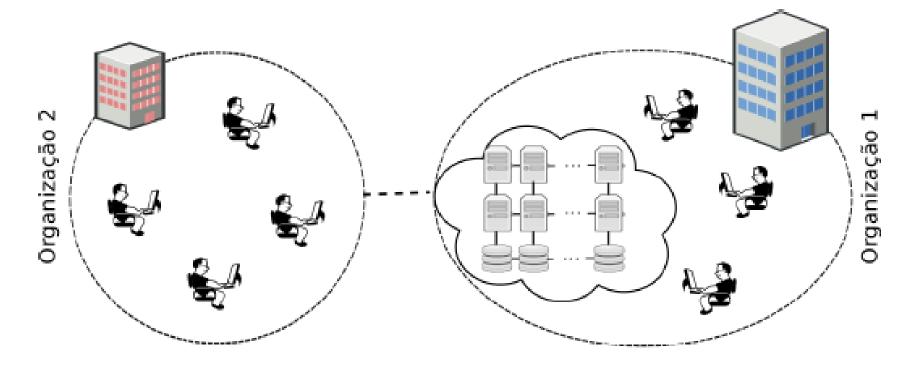
Tipos de Nuvem

- Públicas
- Privadas
- Híbridas
- Comunitárias



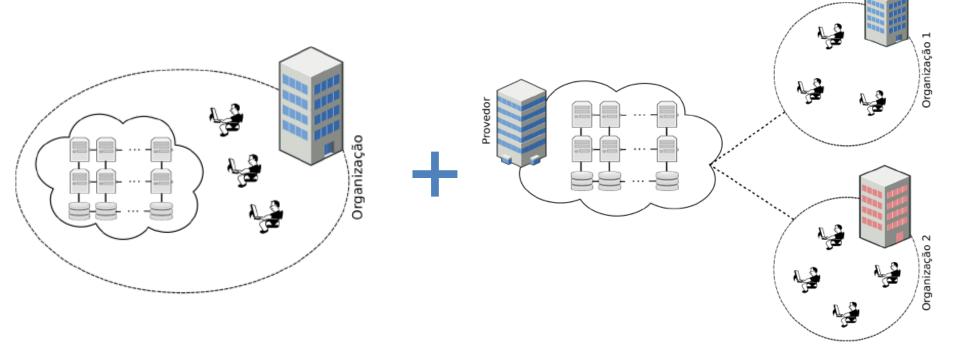
Tipos de Nuvem

- Públicas
- Privadas
- Comunitárias
- Híbridas

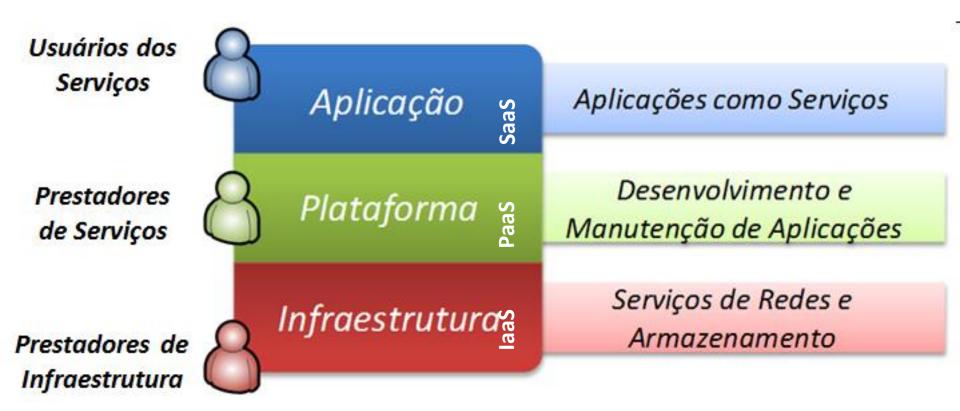


Tipos de Nuvem

- Públicas
- Privadas
- Comunitárias
- Híbridas



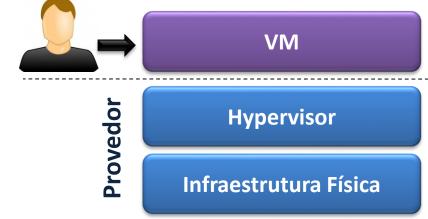
Taxonomia da Nuvem



Fonte: Lamia Youseff, Maria Butrico and Dilma Da Silva Toward a Unified Ontology of Cloud Computing.", Grid Computing Environments Workshop, 2008. GCE '08

IaaS - Infrastructure as a Service

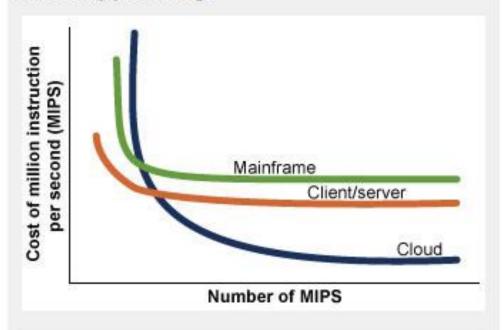
- Oferece infraestrutura de hardware
 - Normalmente por meio de virtualização
- Funciona como um aluguel de recursos:
 - Equipamentos de Rede
 - Memória
 - CPU
 - Armazenamento
- Infraestrutura deve ser escalável
 - Aumentar ou diminuir recursos de acordo com a necessidade do cliente



IaaS - Custo Financeiro

Moving to the cloud can potentially save companies money

Consolidating computing power has the potential to lead to more efficient computing, and to drive down the cost of processing.



Source: Microsoft.

IaaS: Exemplos







Sign in to the AWS Management Console



Products & Services ▼

Amazon Elastic Compute Cloud (Amazon EC2)

Instâncias on demand padrão	
Pequena (padrão)	\$0.060 por hora
Médio	\$0.120 por hora
Grande	\$0.240 por hora
Extragrande	\$0.480 por hora
Instâncias on demand padrão de segunda geração	
Extragrande	\$0.500 por hora
Dupla extragrande	\$1.000 por hora
Microinstâncias on demand	
Micro	\$0.020 por hora
Instâncias on demand com mais memória	
Extragrande	\$0.410 por hora
Dupla extragrande	\$0.820 por hora
Quádrupla extragrande	\$1.640 por hora
Instâncias on demand com CPU de alta performance	
Médio	\$0.145 por hora
Extragrande	\$0.580 por hora
Instâncias de computação em cluster	
Quádrupla extragrande	\$1.300 por hora
Óctupla extragrande	\$2.400 por hora
Instâncias on demand de cluster com mais memória	
Óctupla extragrande	\$3.500 por hora

Instâncias on demand padrão	
Pequena (padrão)	\$0.060 por hora
Médio	\$0.120 por hora
Grande	\$0.240 por hora
Extragrande	\$0.480 por hora
Pequena (padrão) Médio Grande Extragrande Instâncias on demand padrão de so stâncias de uso geral Extragrande	
Extragrande	\$0.500 por hora
Dupla extragrande	\$1.000 por hora
Microinstâncias on demand	
Micro	\$0.020 por hora
Instâncias on demand com mais memória	
Extragrande	\$0.410 por hora
Dupla extragrande	\$0.820 por hora
Quádrupla extragrande	\$1.640 por hora
Instâncias on demand com CPU de alta performance	
Médio	\$0.145 por hora
Extragrande	\$0.580 por hora
Instâncias de computação em cluster	
Quádrupla extragrande	\$1.300 por hora
Óctupla extragrande	\$2.400 por hora
Instâncias on demand de cluster com mais memória	
Óctupla extragrande	\$3.500 por hora

Instâncias on demand padrão	
Pequena (padrão)	\$0.060 por hora
Médio	\$0.120 por hora
Grande	\$0.240 por hora
Extragrande	\$0.480 por hora
Instâncias on demand padrão de segunda geração	
Extragrande	\$0.500 por hora
Dupla extragrande	\$1.000 por hora
Microinstâncias on demand	test
Micro : as de	\$0.020 por hora
Dupla extragrande Microinstâncias on demand Micro Instâncias on demand com mais memória Extragrande	
Extragrande	\$0.410 por hora
Dupla extragrande	\$0.820 por hora
Quádrupla extragrande	\$1.640 por hora
Instâncias on demand com CPU de alta performance	
Médio	\$0.145 por hora
Extragrande	\$0.580 por hora
Instâncias de computação em cluster	
Quádrupla extragrande	\$1.300 por hora
Óctupla extragrande	\$2.400 por hora
Instâncias on demand de cluster com mais memória	
Óctupla extragrande	\$3.500 por hora

Instâncias on demand padrão	
Pequena (padrão)	\$0.060 por hora
Médio	\$0.120 por hora
Grande	\$0.240 por hora
Extragrande	\$0.480 por hora
Instâncias on demand padrão de segunda geração	
Extragrande	\$0.500 por hora
Dupla extragrande	\$1.000 por hora
Microinstâncias on demand	
Micro	\$0.020 por hora
Instâncias on demand com mais memória	
Extragrande	\$0.410 p
Dupla extragrande	\$n nen!
Quádrupla extragrande	hora
Instâncias on demand com CPU de alta performance	140 des
Médio	\$0.145 por hora
Extragrande 600	\$0.580 por hora
Instâncias de computação em cluster	
Quádrupla extragrande	\$1.300 por hora
Óctupla extragrande	\$2.400 por hora
Extragrande Dupla extragrande Quádrupla extragrande Instâncias on demand com CPU de alta performance Médio Extragrande Instâncias de computação em cluster Quádrupla extragrande Óctupla extragrande Instâncias on demand de cuatra de computação em cluster Instâncias on demand de cuatra de computação em cluster Octupla extragrande Octupla extragrande	
Óctupla extragrande (NSTA	\$3.500 por hora

Free Tier*

As part of AWS's Free Usage Tier, new AWS customers can get started with Amazon EC2 for free. Upon sign-up, new AWS customers receive the following EC2 services each month for one year:

750 hours of EC2 running Linux/Unix Micro instance usage

750 hours of Elastic Load Balancing plus 15 GB data processing

10 GB of Amazon Elastic Block Storage (EBS) plus 1 million IOs, 1 GB snapshot storage, 10,000

snapshot Get Requests and 1,000 snapshot Put Requests

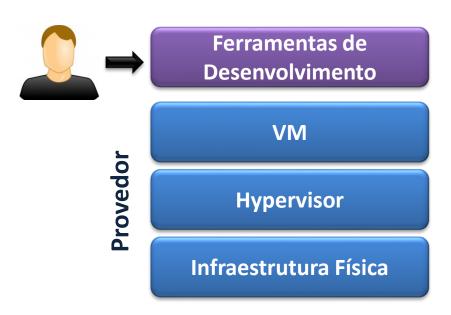
15 GB of bandwidth out aggregated across all AWS services

1 GB of Regional Data Transfer



PaaS - Platform as a Service

- Modelo onde se fornece recursos para a construção de aplicações e serviços para a Internet
- Os recursos incluem:
 - Ferramentas de desenvolvimento
 - Teste
 - Hospedagem
 - Banco de Dados
 - Segurança
 - Controle de versões



PaaS: Exemplos

Google App Engine

Página inicial

<u>Documentos</u>

Perguntas frequent



Execute os seus aplicativos da web na infraestrutura do Google.

Fácil de criar, fácil de manter e fácil de escalar.

Uma visão inicial do suporte à linguagem Java™ Novo!

Recentemente, o Google App Engine está descobrindo a sua segunda linguagem: Java. Essa versão inclui uma visão inicial da nossa execução em Java, a integração com o Google Web Toolkit e um Plug-in do Google para o Eclipse, fornecendo a você uma solução Java completa para aplicativos AJAX da web. Nosso suporte à linguagem Java ainda está em desenvolvimento e estamos ansiosos pela sua ajuda e comentários. A execução em Java está disponível para qualquer pessoa usar. Experimente e nos envie o seu feedback.

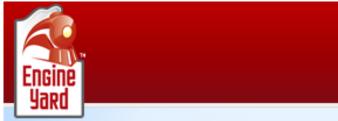








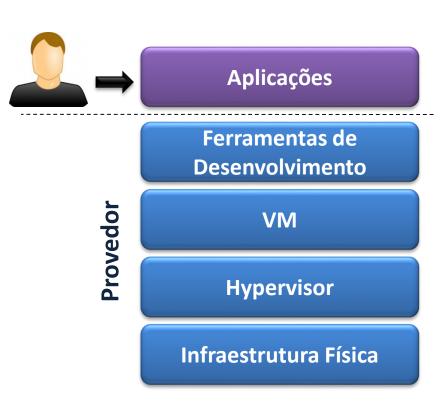




Cloud Features

SaaS - Software as a Service

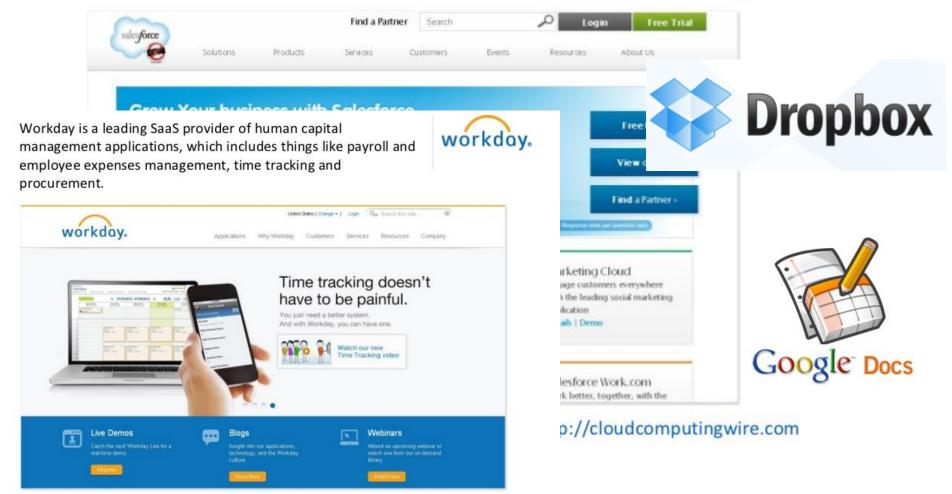
- Modelo no qual uma aplicação é armazenada em um servidor
- Usuários a acessam via Internet
 - Não há a necessidade de dar suporte à aplicação



SaaS: Exemplos

Salesforce is a leader is SaaS computing – and is best known for its on-demand Customer Relationship Management (CRM) solutions.





Casos de Sucesso

- · The New Hork Times
 - Amazon EC2 e 53 (100 instâncias)
 - Conversão de 11 milhões de artigos (4TB)





- Processamento de vídeos do Big Brother Brasil
- Evita sobrecarga nos servidores locais

Casos na área científica

Análise Filogenômica

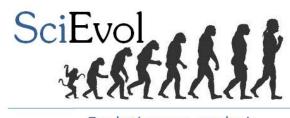
OLIVEIRA, D., Ocaña, K. A. C. S., Ogasawara, E., Dias, J., Goncalves, J., Mattoso, M., (2012), "Cloud-based Phylogenomic Inference of Evolutionary Relationships: A Performance Study". 2nd International Workshop on Cloud Computing and Scientific Applications (CCSA), Ottawa.



Phylogenomic analysis

Estudos Evolutivos

Ocaña, K. A. C. S., OLIVEIRA, D. de, Horta, F., Dias, J., Ogasawara, E. Mattoso, M., (2012), "Exploring Molecular Evolution Reconstruction Using a Parallel Cloud-based Scientific Workflow". Brazilian Symposium or Bioinformatics (BSB 2012), Campo Grande, MT.



Evolutionary analysis

Genômica Comparativa

OCANA, K.; OLIVEIRA, D.; DIAS, J.; OGASAWARA, E.; MATTOSO, M. L. Q. . Optimizing Phylogenetic Analysis Using SciHmm Cloud-based Scientific Workflow. 7th IEEE e Science conference. IEEE Computer Society, 2011.



Roteiro do Tutorial

- Motivação
- Workflows Científicos
- · Nuvens de Computadores
- · Proveniência de Dados
- Máquinas de Workflow para Nuvem
- · Aplicação de Proveniência em e-Science
- Demo



O que é proveniência?

- Dictionary

prov•e•nance | prävənəns |

noun

the place of origin or earliest known history of something: an orange rug of Iranian provenance.

 the beginning of something's existence; something's origin: they try to understand the whole universe, its provenance and fate.

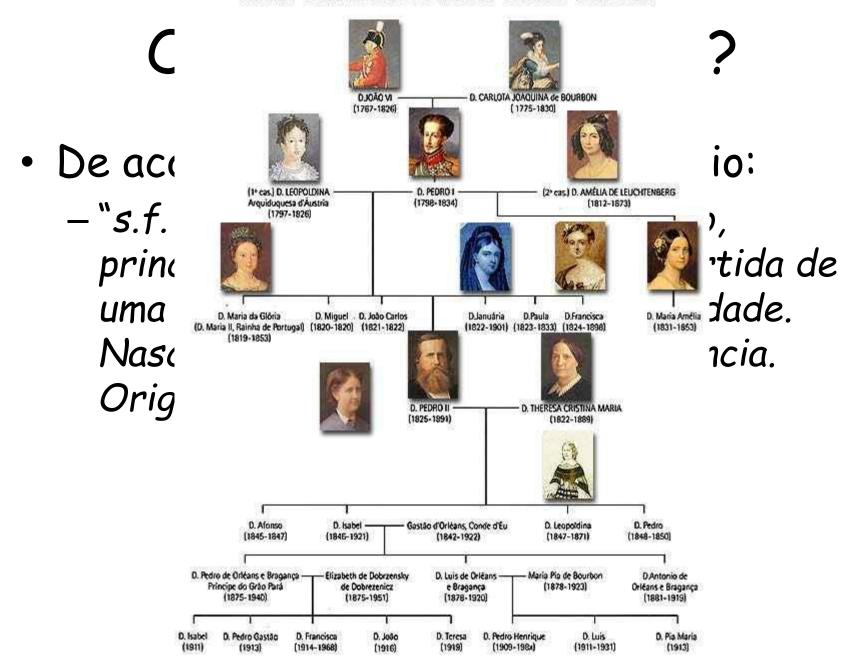
See note at ORIGIN.

 a record of ownership of a work of art or an antique, used as a guide to authenticity or quality: the manuscript has a distinguished provenance.

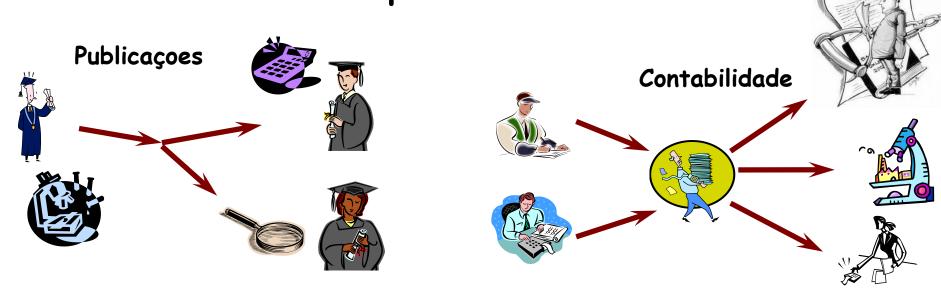
ORIGIN late 18th cent.: from French, from the verb provenir 'come or stem from,' from Latin provenire, from pro- forth' + venire 'come.'

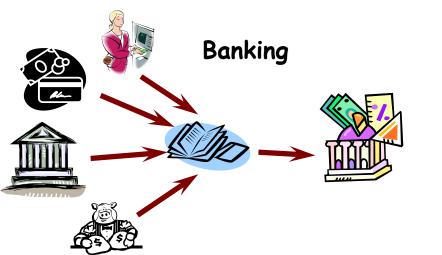
O que é proveniência?

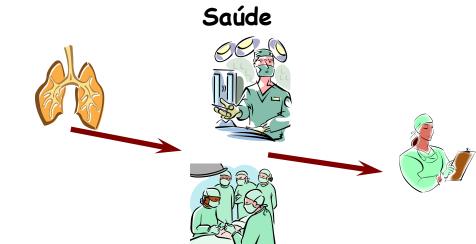
- · De acordo com o dicionário Aurélio:
 - "s.f. Primeira manifestação; começo, princípio. Procedência; ponto de partida de uma nação, de uma família; naturalidade. Nascimento, proveniência; ascendência. Origem."



Onde a proveniência pode ser aplicada?







Proveniência nas artes

ARTE E CULTURA

09 de Setembro de 2013 * 10h22

Após um século escondido, quadro "falso" é confirmado como Van Gogh





Uma paisagem francesa que passou um século guardada em um sótão por ser considerada uma tela falsa de Vincent van Gogh, na verdade, é um trabalho autêntico, segundo o resultado de um novo estudo divulgado nesta segunda-feira.

A tela "Ocaso em Montmajour", que mostra azinheiras retorcidas e uma distante ruína banhada pela luz do entardecer, foi pintada em 1888, quando Van Gogh vivia em Arles, no sul da França.

A pintura, pertencente a um colecionador privado, será exposta a partir deste mês, durante um

ano, no Museu

- Pesquisa a autoria de obras de arte em acervos de museus ou coleções particulares
 Dados de proveniência
- Dados de proveniência podem incluir restaurações realizadas nas obras
- Ajudam especialistas a verificar autenticidade de obras

Adaptação dos slides de Freire e Davidson – SIGMOD 2008

Proveniência nas artes





- Pesquisa a autoria de obras de arte em acervos de museus ou coleções particulares
- Dados de proveniência podem incluir restaurações realizadas nas obras
- Ajudam especialistas a verificar autenticidade de obras

Adaptação dos slides de Freire e Davidson – SIGMOD 2008

Proveniência na indústria



- Garante a procedência do alimento
- É possível saber em que fazenda, granja etc o alimento foi produzido ou colhido

Adaptação dos slides de Freire e Davidson – SIGMOD 2008

Proveniência na medicina

Prontuário Ministrio de Sados Fundado Osendro Cour RISTITUTO FERNANCIES PROJERA NONE DATAE PORM	28k6/6 Sinuso Socal Jac comparati plantic Cantic 8 100 kg Re 185 au 27hillio Peso 8 100 kg Re 185 au
35/6/5/C D. S Bers n. 280 10. 2.66 10. 2.6	Ca the Son what free, a interestive on the Co. of the form of the form of the form of the form of the control o
Expunds to un on inition to gention of the control	Joseph Pessing (1965) 10.50 - 10.50 pm 1 pm . Candyn auch i els pronto. . Co : FO patent. . Uso nominal (sic) Jone / 170 acompanhada pelo points: Em un a Bathum De voto mudasel Engalines i ma "para har landheado (pacco) pala japan' e mania. Ao Hame : Ottim, Enjalones

fonte: A epistemologia narrativa e o exercício clínico do diagnóstico, Maria Helena Cabral de Almeida Cardoso, Kenneth Rochel de Camargo Jr., Juan Clinton Llerena Jr., Ciência & Saúde Coletiva

- Garante o rastreio das mudanças do paciente
- Cirurgias
 realizadas,
 procedimentos
 realizados,
 medicamentos
 prescritos

Adaptação dos slides de Freire e Davidson – SIGMOD 2008

Proveniência na Ciência

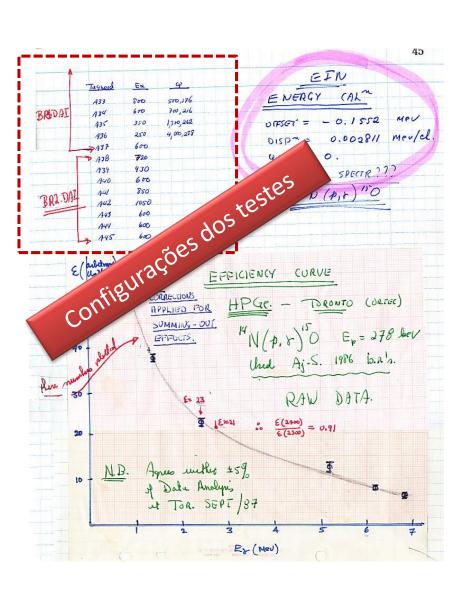
- Reprodução e interpretação de um resultado
- Compreensão do raciocínio empregado no experimento
- Verificação se o experimento segui procedimentos específicos
- · Rastreio de responsabilidades
 - Quem fez o que e quando?

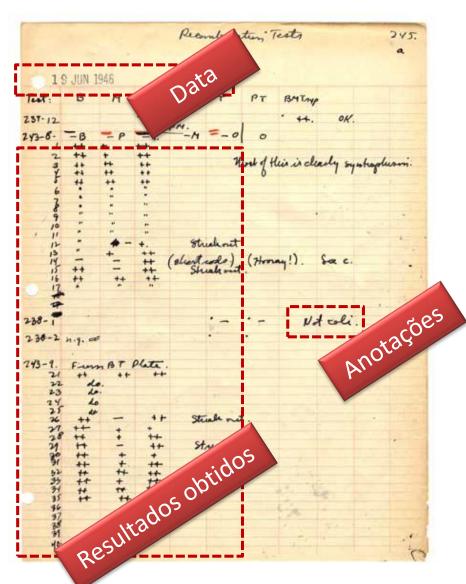
Proveniência na Ciência

- Reprodução e interpretação de resultado
- Compreensão do raci experimento
- "Provenance is as (or more) important as the results"

 (Juliana Freire e Susan Davidson, SIGMOD 2008) Verifica m fez o que e quando?

Proveniência na Ciência





Por que a curadoria de dados científicos é importante em e-Science?

- · Faz parte do processo de pesquisa
- Agrega valores Intrínsecos e Extrínsecos à pesquisa
- Aumenta o potencial de criar "novos conhecimentos" a partir de dados préexistentes
- Preserva o contexto da descoberta científica

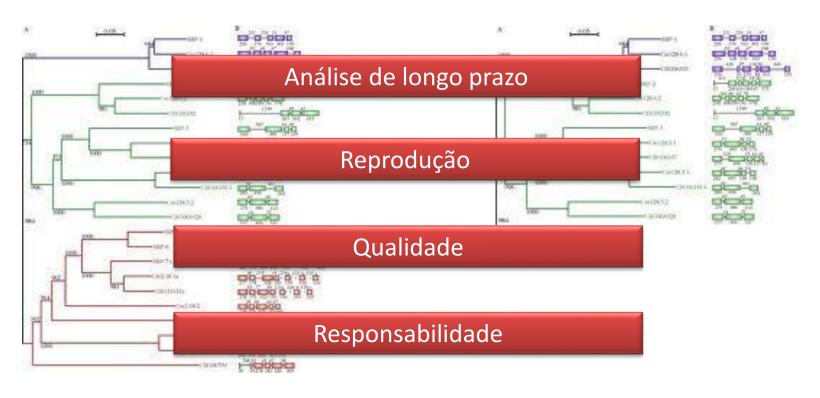
Qualidade dos dados

- A proveniência dos dados permite avaliar a qualidade deles para uma aplicação
- Erros introduzidos por defeitos nos dados tendem a se tornar mais graves quando propagados
- O nível de detalhe da proveniência determina com que grau a qualidade dos dados pode ser estimada

Proveniência em experimentos in silico

iii-201102281EHWP25EX8.tcoffeeStockholm

iii-201102281EHWP25EX9.tcoffeeStockholm



Formas da proveniência

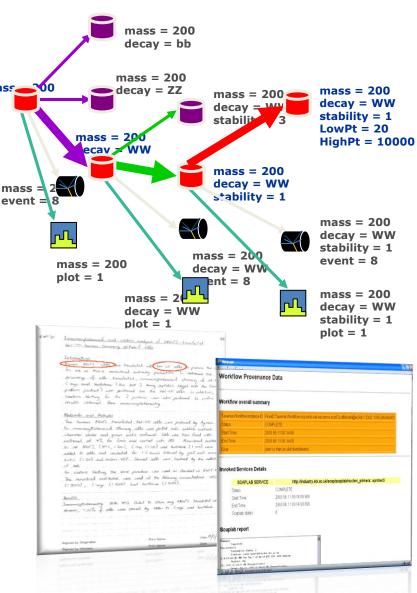
Derivação

- Representação de um caminho, workflow, roteiro, ou consulta.
- Associação entre itens (grafos)
- Geralmente é uma explicação (2W1H = when, who, how)
- Centrada no processo

Anotação

- Anexado aos items ou coleções de dados (estruturado ou semi ou texto livre)
- Geralmente é uma explicação (5W1H = why, when, where, who, what, how)
- Centrada no dado

Fonte: Knowledge and Provenance (Goble, 2003)



Proveniência em Tempo Real

- Uso da proveniência em tempo real possibilita análises mais precisas e uma resposta mais rapida durante a execução do experimento
 - As derivações que ocorrem durante a execução do workflow devem estar disponíveis assim que os dados forem produzidos

Dados específicos de domínio também são necessários

A Survey of Data Provenance in e-Science

Yogesh L. Simmhan Beth Plale Dennis Gannon Computer Science Department Indiana University, Bloomington, IN 47 405 (ysimmhan, piale, gannon)@cs.indiana.edu

ABSTRACT

Data scale

R 900

abun produ

05 SQ E perve

denty

origin

In th

chara

e-scie

аррес cates

RI COD ю pro à. Th

10 Sec.

The

R 900

provi

allow

in co

the n

work rich. make

prov

which

and c

inten

Scien

and

form

data

(DOI

that i

amaly

data

beyou

acces

infor

meta

descr

of their application [2]. Materials engineers choose

Provenance of e-Science Experiments - experience from Bioinformatics.

Mark Greenwood, Carole Goble, Robert Stevens, Jun Zhao, Matthew Addis, Darren Marvin, Luc Moresu, Tom Cinn EPSRC e-Science Pilot Project *** Grid http://www.mygrid.org.uk

COMPUTATIONAL PROVENANCE

Provenance for Computational Tasks: A Survey

Provenance in Scientific Workflow Systems

Susan Davidson. Sarah Cohen-Boulakia, Anat Eyal University of Pennsylvania {susan, sarabeb, anate }@cis.upenn.edu

Bertram Ludäscher, Timothy McPhillips Shawn Bowers, Manish Kumar Anand University of California, Davis {ludaesch, tmcphillips, showers, maanand}@ucdavis.edu

Juliana Freire University of Utah juliana@cs.utah.edu

CONCURRENCY AND COMPUTATION: PRACTICE AND EXPERIENCE Conserving Computat.: Pract. Exper. 2007; V OLUME:1-9869.

Prepared sating operations (Verstern)

Abstract

ching and storage of provenance information promises to be a major advantage

Tracking Provenance in a

1 Introd Like experimen

2006/03/23 12:005

Library

reduced vel

information

perferred:

e-Science e

securitation.

कार्य एक शहर

per sentation

data associated

reduced value i

the movemence

michadata, recor

experiments are

or even by thos

also importanti

Our investigation

project, which

sonibe our under

present how we

Section 4 and or

of results.

Tipos de proveniência?

Computation Institute, University of Chicago, Chicago, IL

ments for e-Seis iments as well ^a Math & Computer Science Division, Argume National. ments by scient

Laboratory, Argonna, IL 600 30, USA USC Information Sciences Institute Between Experiments and Data

COMPUTATIONAL

PROVENANCE

Scientific applications are often structured as workflows compiled from abstract experiment

Briefings in Bioinformatics Advance Access published May 14, 2007 BRIDRINGS IN BIOINFORMATICS, pay 1 of II.

step often hides details about how rovenance helps scientists connect Co. 10.000/16/16/16/05

too experiment SUMMARY

> The virtual data model allow their physical materializatik Language (VDL) and associa query, and retrieval of virtmaterialization of virtual da exercise to illustrate the pow by these tools, which for a computational 'procedure(s) runtime log(s) produced by t annotation(i) that associate

Using provenance to manage knowledge of In Silico experiments

Robert Stevens, Jun Zhao and Carole Goble Schoolseli 64 january 2007, resolve & Shirifarch 2007

This article offers a briefing in one of the knowledge management issues of in silky experimentation in bidinformatics. Recording of the province of an experiment—what was done; where, how and why, etc. is an important aspect of scientific best practice that should be excended to insilito experimentation. We will do this in the context of escience which has been part of the move of biginformatics towards an industrial setting. Despite the computational nature of biblinformatics, these analyses are scientific and thus necessitate their own versions of typical adentific rigour. Just as recording who, what, why, when, where and how of an experiment is central to the adentific describe computations, their parameters, I/O data, and data or control dependencies between them, and software systems manage these workflows by following dependencies and executing reating the computations on the desired data. Workflow technologies not only help automate complex scientific analyses, but they also provide the opportunity to capture transformations performed on the data.

To ilkistrate the complexity of today's analyses, we examine a popular astronomy application, called Montage,4 which produces science-grade mosaics of the sky on demand. We can structure this application as a workflow that takes several

parameter ght of as a e modules

d Chimera

s. In such

logical se-

rs between (b), Nodes

ch consists

material for to reuse any

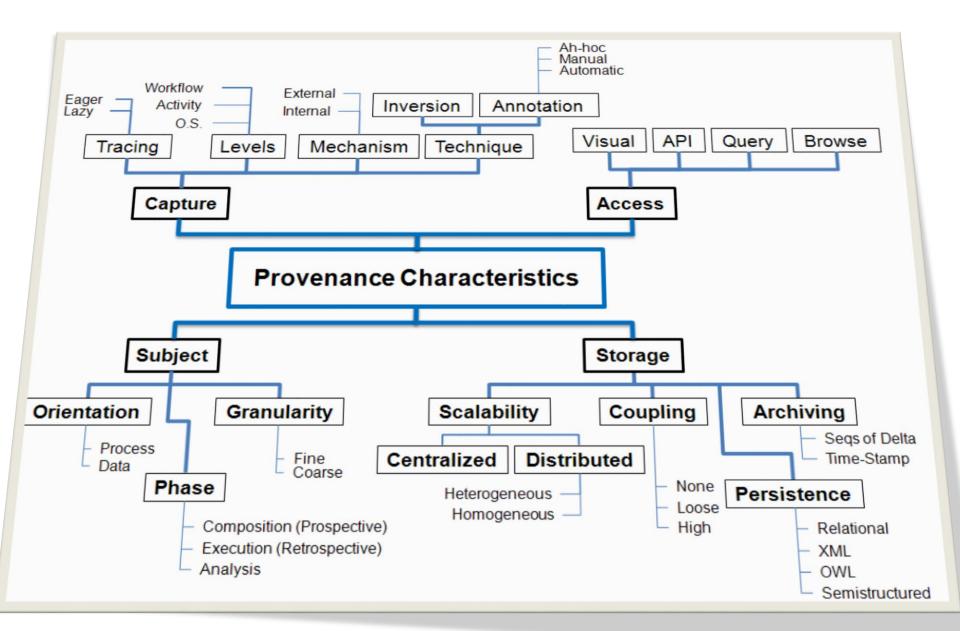
The my

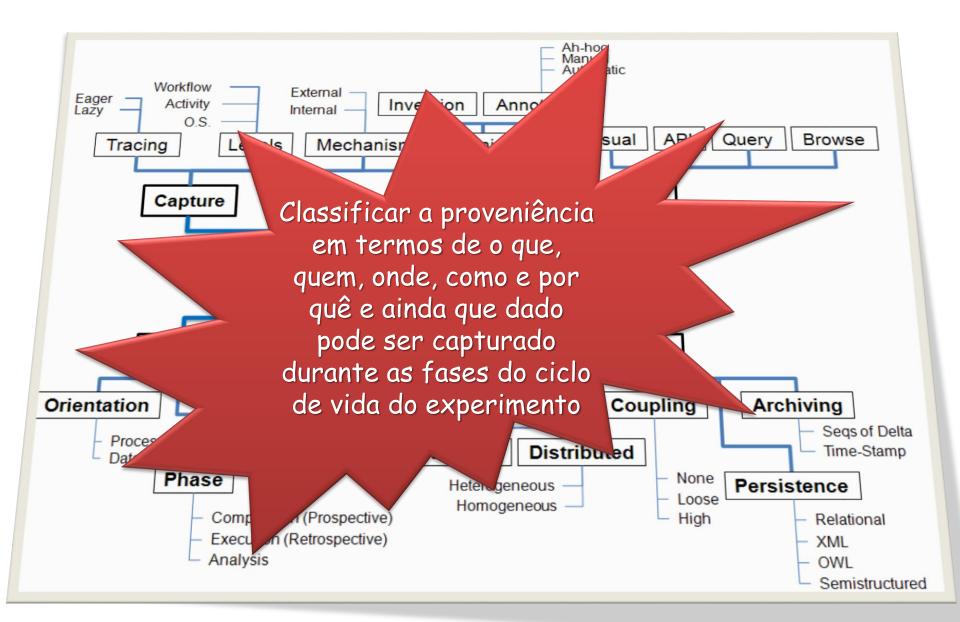
e-Science is the ments, sensors, computers - by lacce distributes tific problems. that uses control computational summary, search

fact. The PYCh

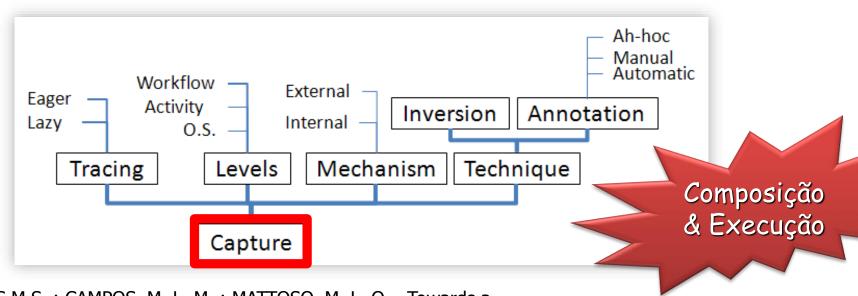
in molecsees. This id supplies ment. The sequences. M4, which determines atput from d to create of modules

Taxonomia de Proveniência

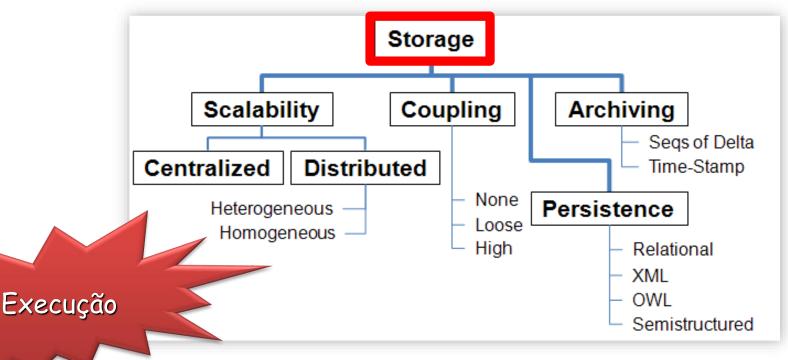




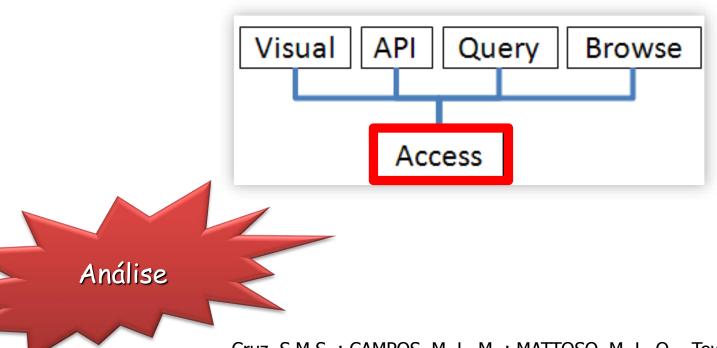
 Classificar como os dados de proveniência podem ser capturados nos sistemas de proveniência.



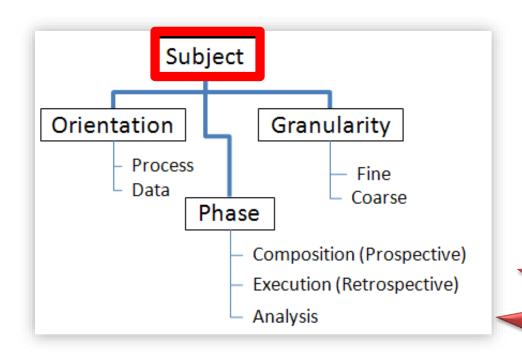
 Formas de registro e armazenamento dos dados de proveniência.



 Descreve as formas de acesso aos dados e repositórios de proveniência.



 Representação da proveniência em termos de assunto e granulosidade.



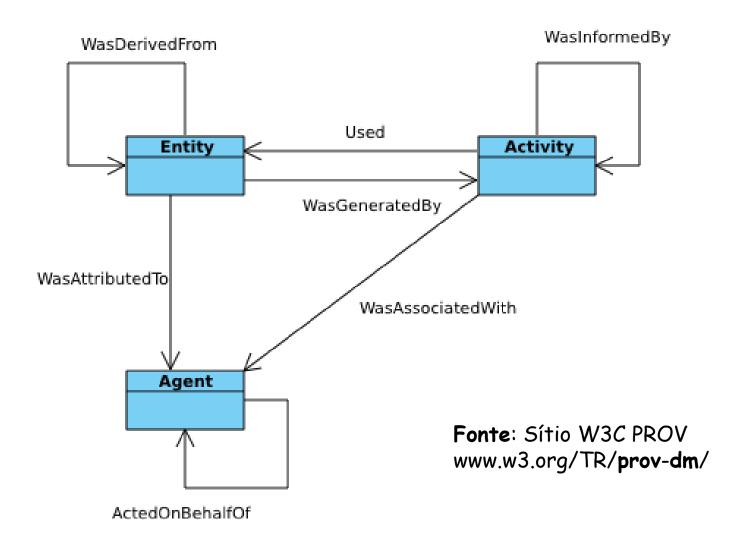
Ortogonal à todas as fases

Representação de proveniência com o PROV

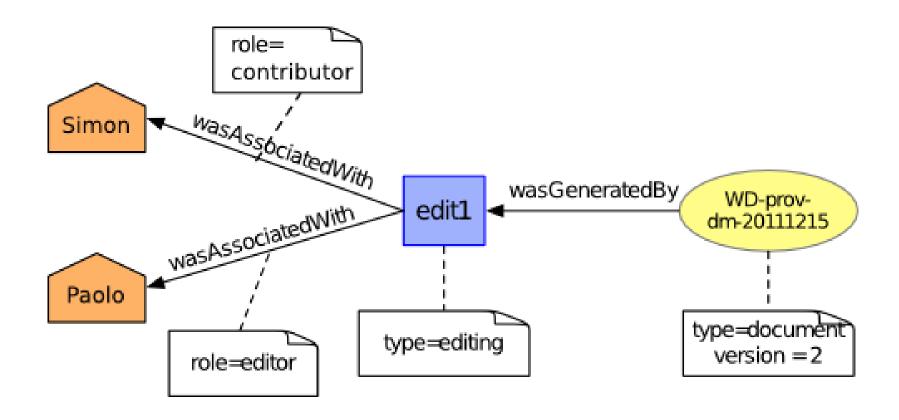
- Padrão W3C para representação de proveniência
- Baseado em 3 elementos principais:
 - Entidade (Entity)
 - Agente (Agent)
 - Atividade (Activity)
- "aims the inter-operable interchange of provenance information in heterogeneous environments such as the Web"



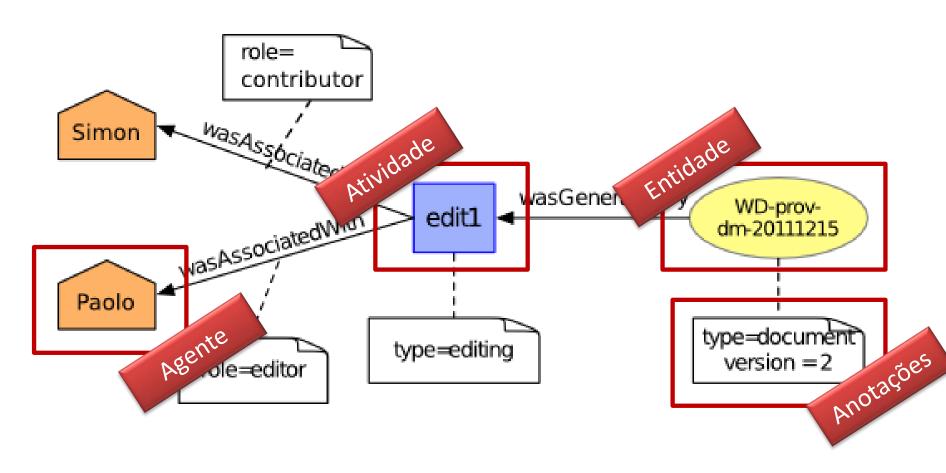
O modelo de dados do PROV



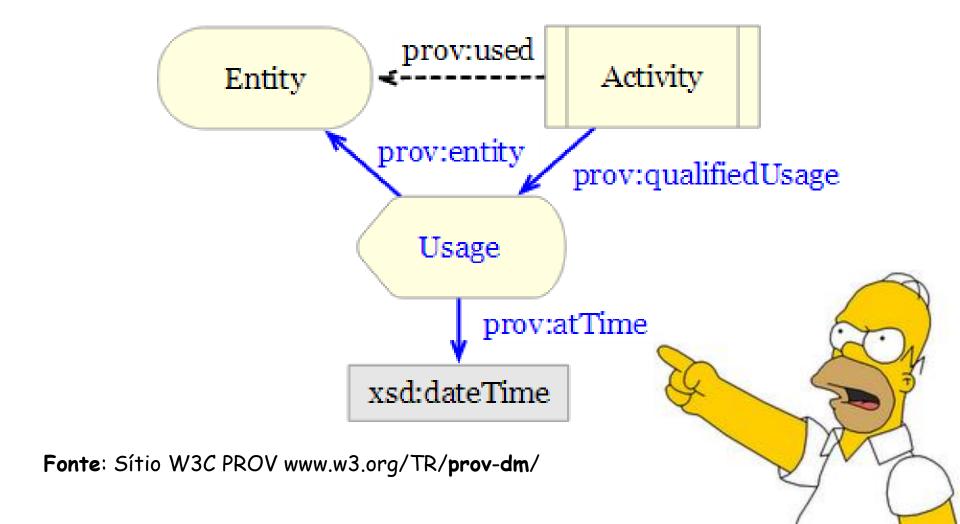
Exemplo de utilização do PROV



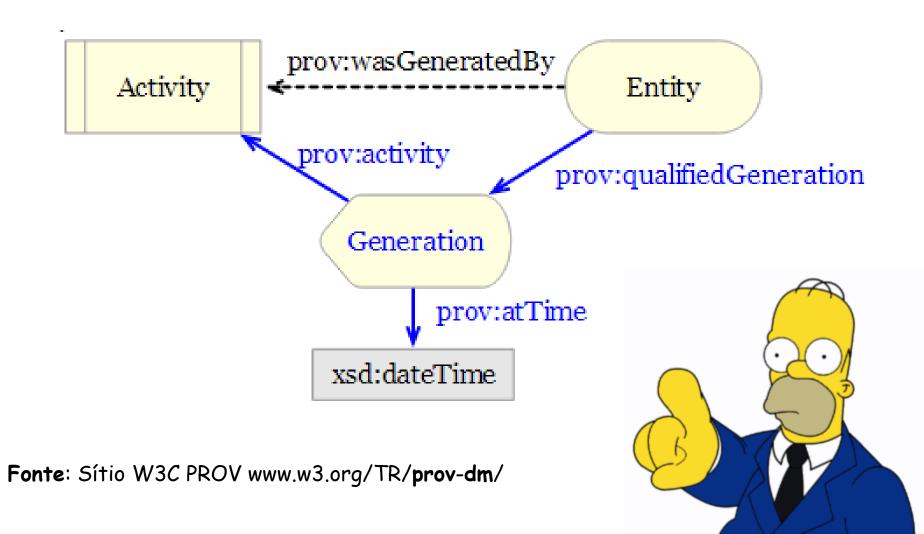
Exemplo de utilização do PROV



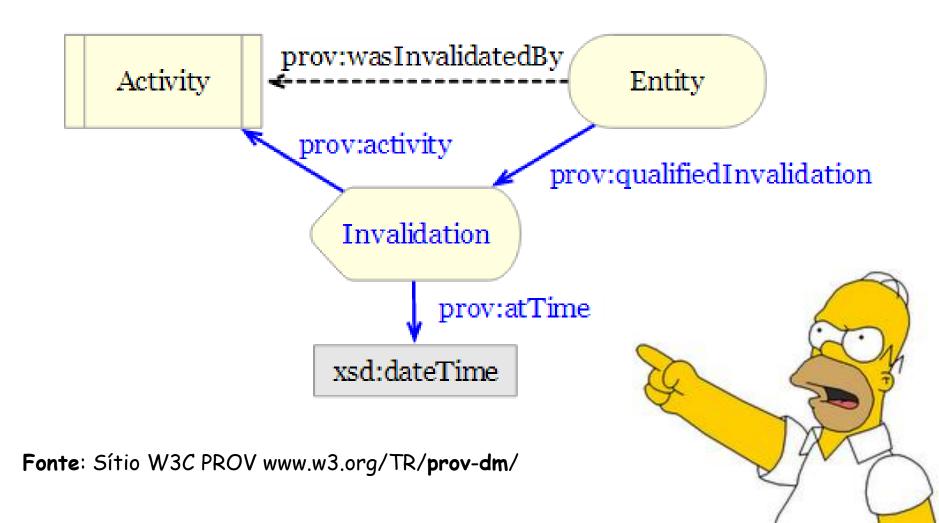
Relacionamento used



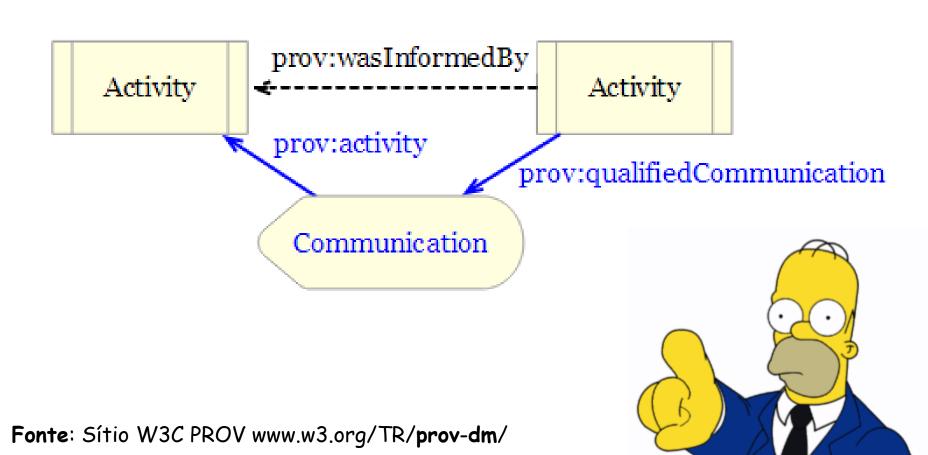
Relacionamento wasGeneratedBy



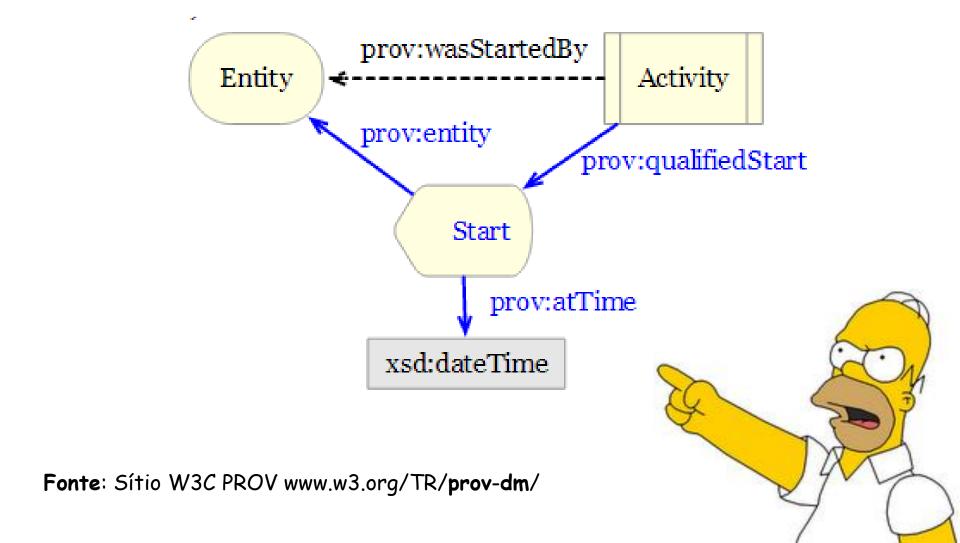
Relacionamento was Invalidated By



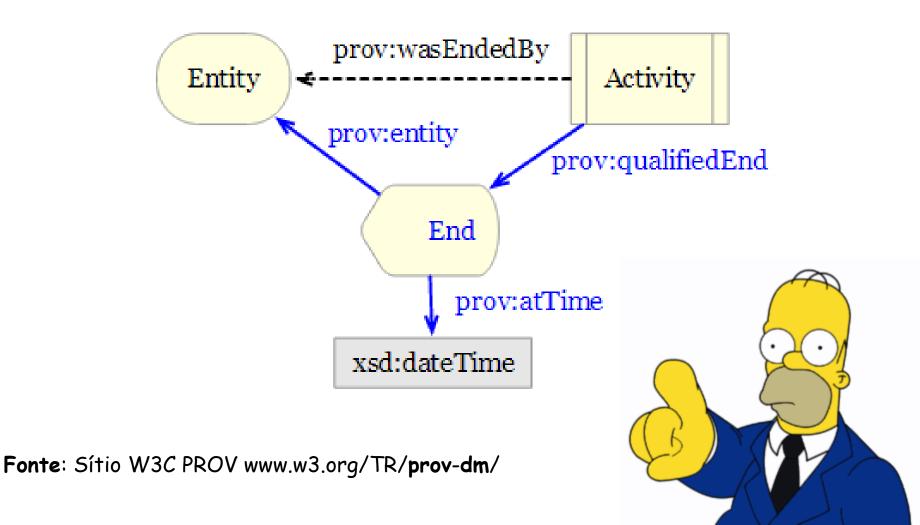
Relacionamento was Informed By



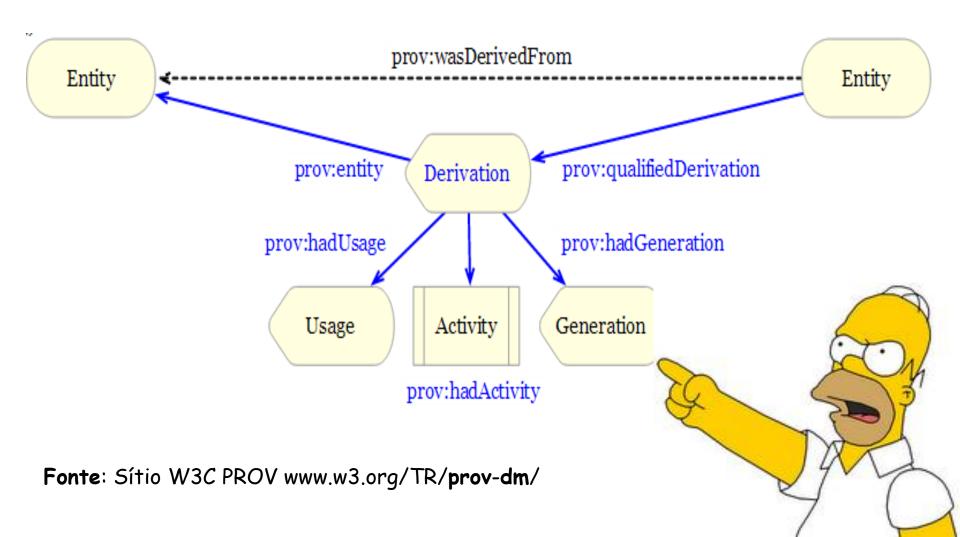
Relacionamento wasStartedBy



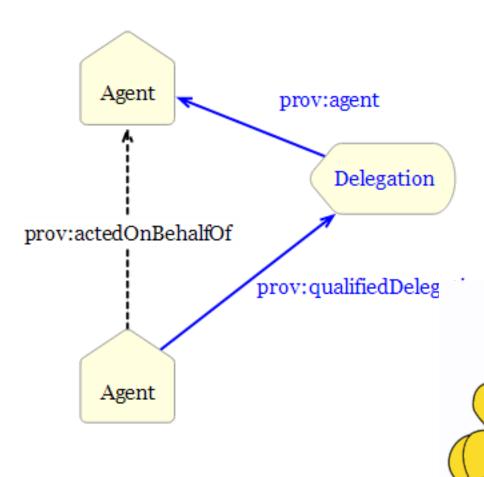
Relacionamento was Ended By



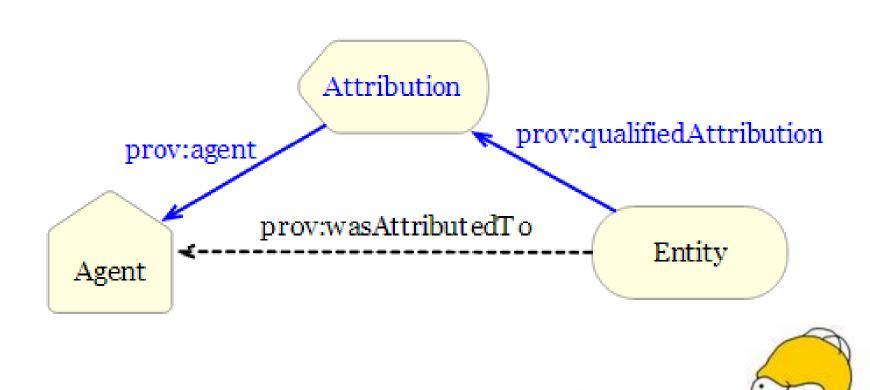
Relacionamento was Derived From



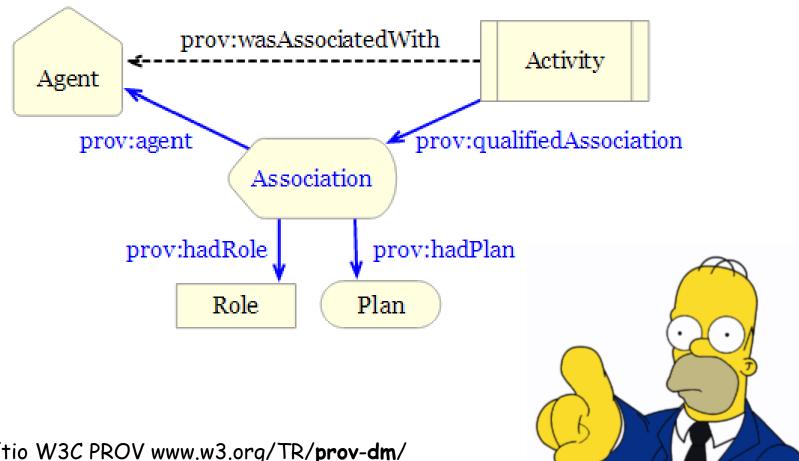
Relacionamento actedOnBehalfOf



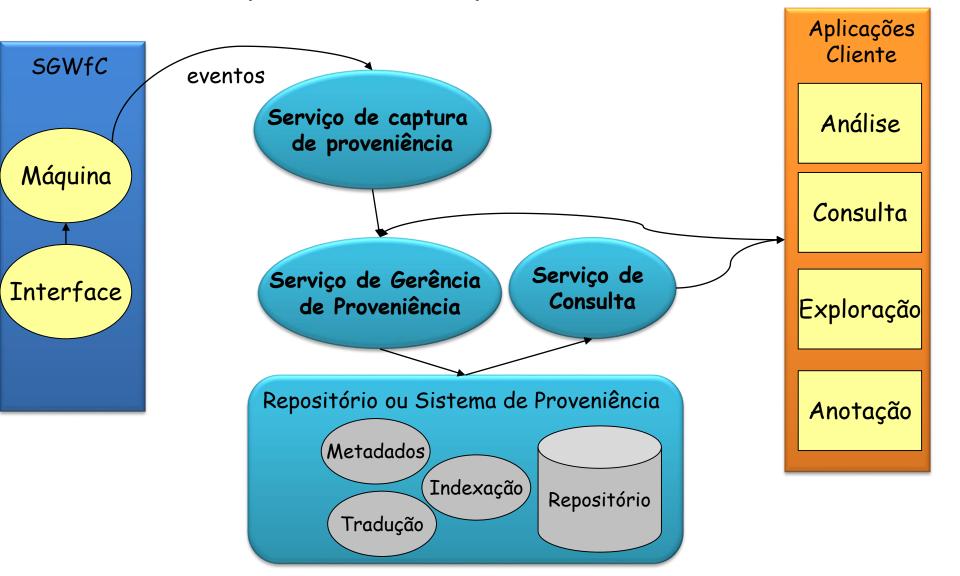
Relacionamento was Attributed To



Relacionamento was Associated With



Arquitetura genérica para captura de proveniência



Mecanismos de Captura

- Atrelados a máquina de execução de workflows
 - Taverna
 - Kepler
 - VisTrails
 - SciCumulus

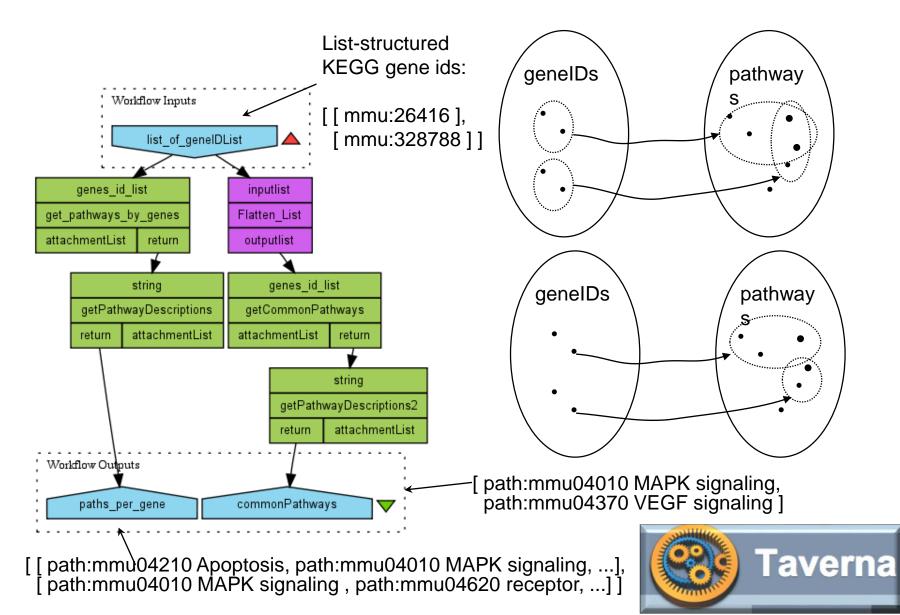
Centralizados

Distribuído

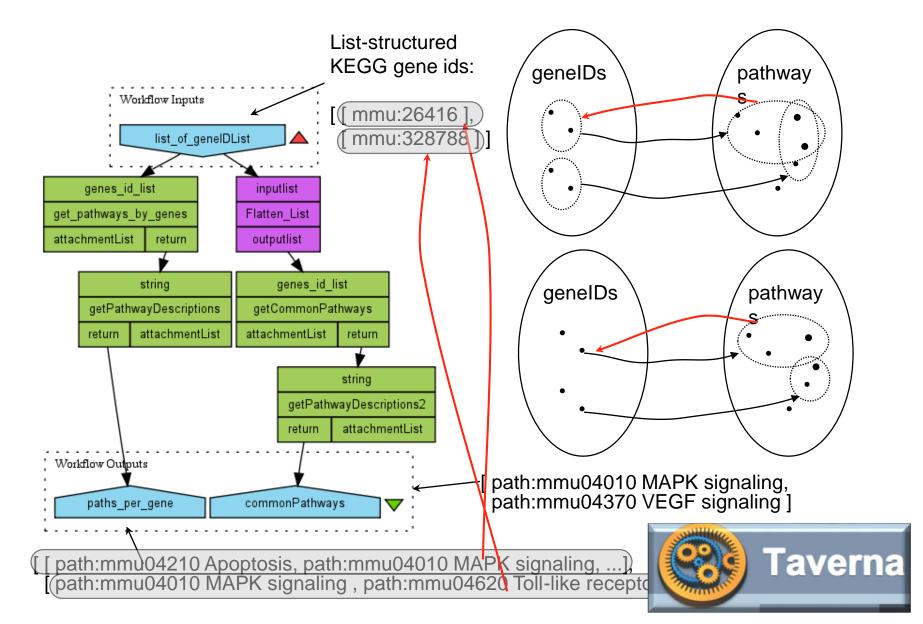
- Desacoplados da máquina de workflows
 - PASS
 - PASOA
 - Karma

Distribuídos

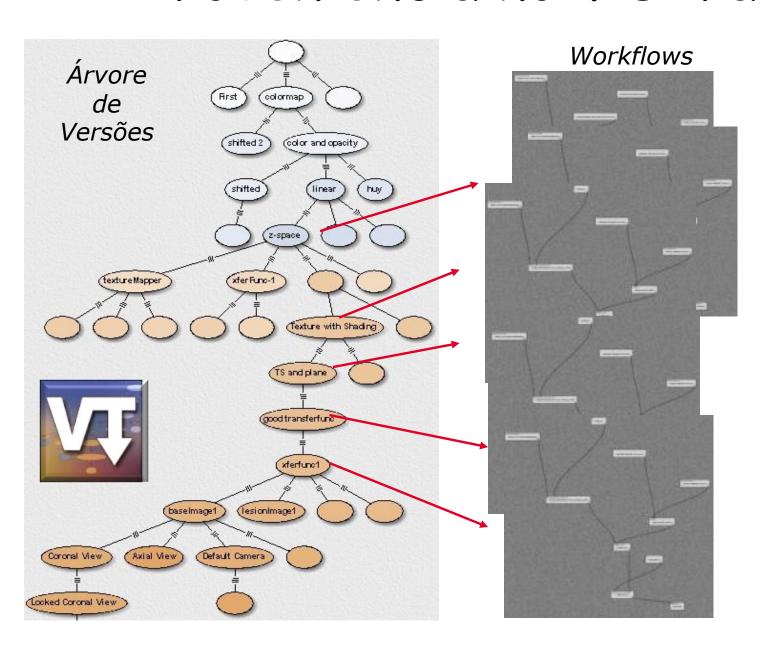
Proveniência no Taverna



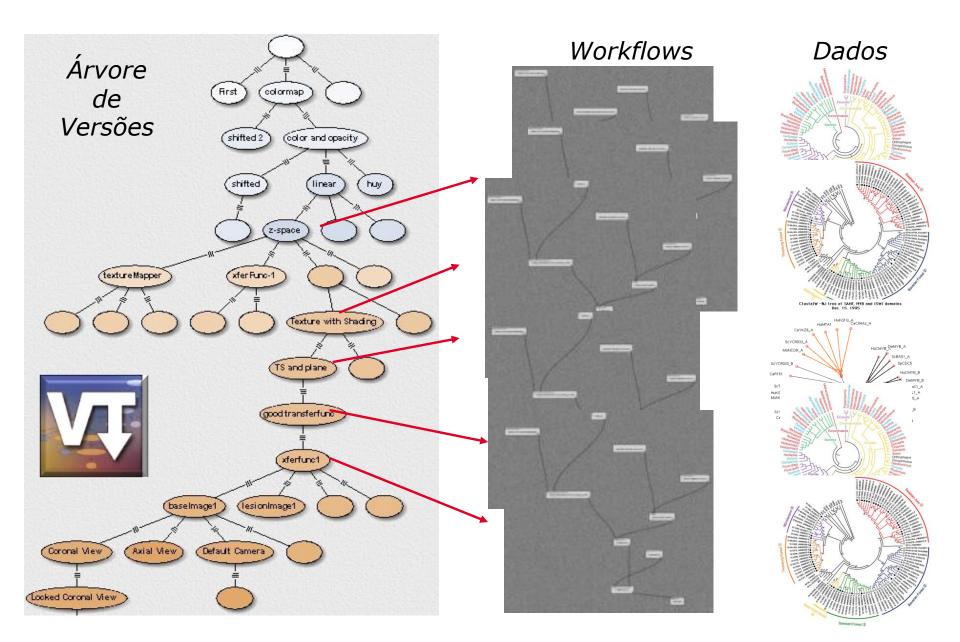
Proveniência no Taverna



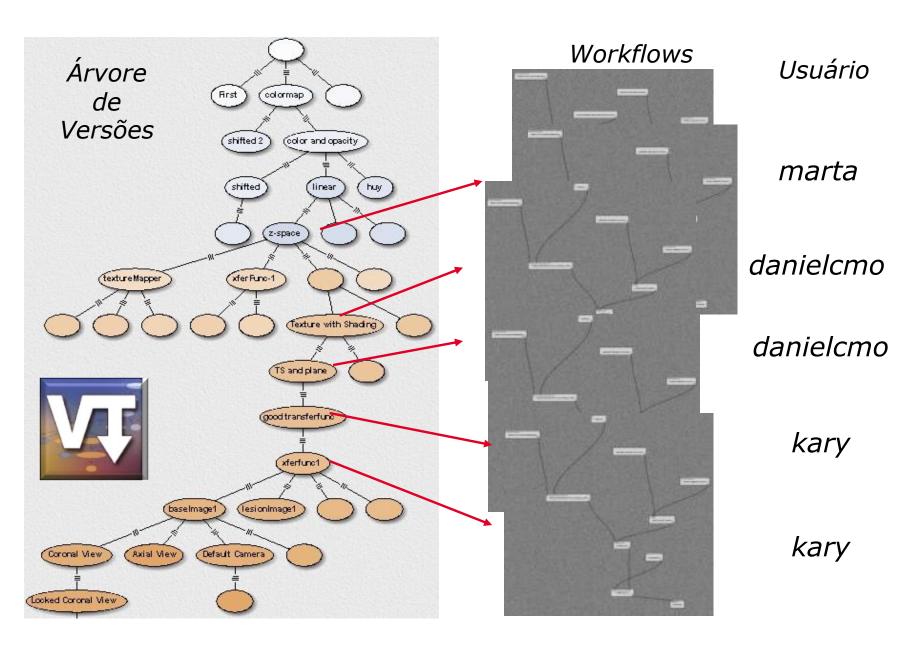
Proveniência no VisTrails



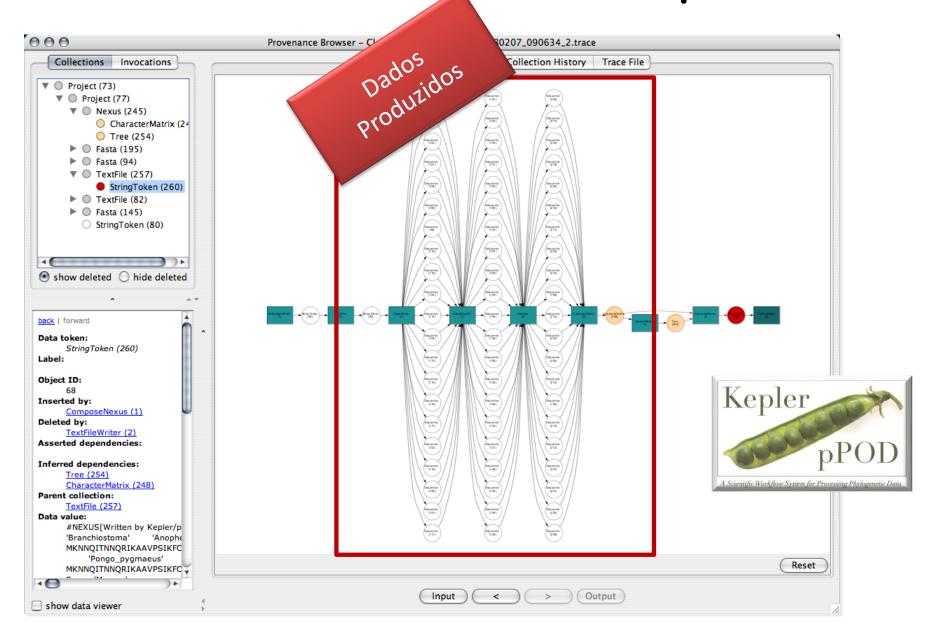
Proveniência no VisTrails



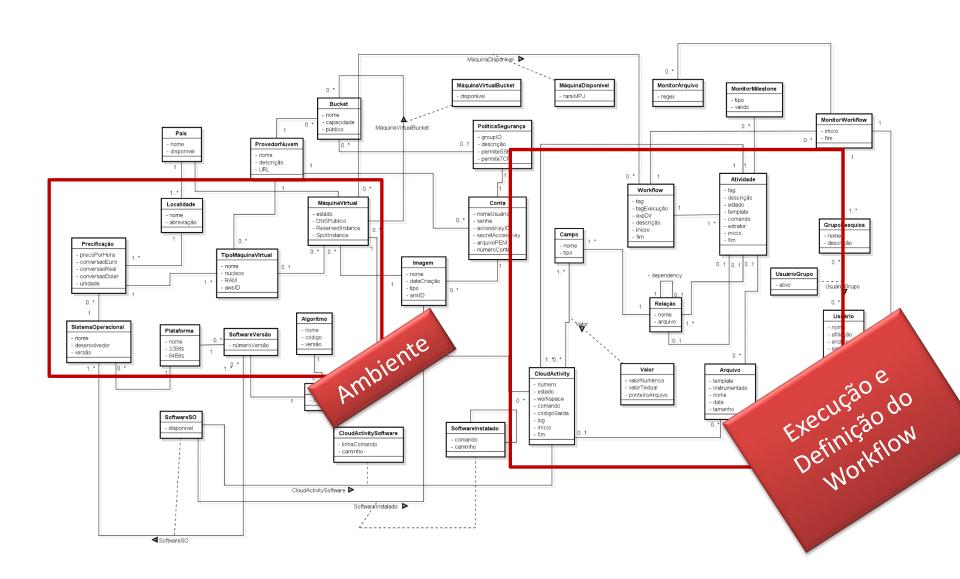
Proveniência no VisTrails



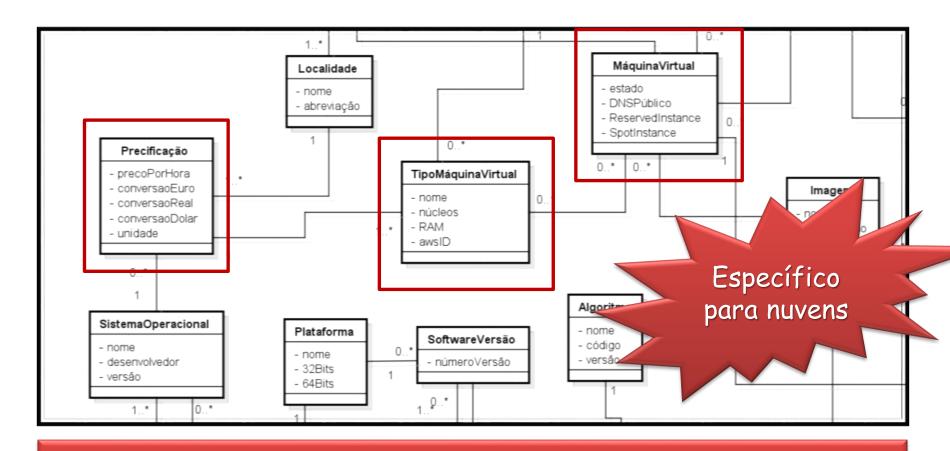
Proveniência no Kepler



Proveniência no Chiron/SciCumulus

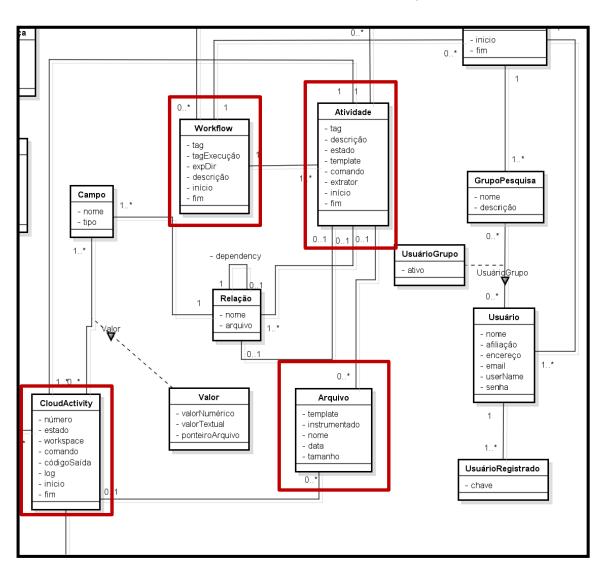


Proveniência no Chiron/SciCumulus



Este tipo de informação pode ser usado para fins além da reprodução e validação

Proveniência no Chiron/SciCumulus



Consultas de Proveniência no SciCumulus

 "Recuperar, por ordem crescente de identificador dos workflows as atividades relacionadas a cada workflow e as cloud activities associadas a cada atividade"

w.description
a.tag,
t.workspace
from hworkflow w, hactivity a, hactivation t
where w.wkfid = a.wkfid
and a.actid = t.actid
order by w.wkfid

SciPhy-60003	Phylogeny using RAXML-kalign	raxml	/root/exp/raxml/4/
SciPhy-60003	Phylogeny using RAXML-kalign	raxml	/root/exp/raxml/3/
SciPhy-60003	Phylogeny using RAXML-kalign	raxml	/root/exp/raxml/2/
SciPhy-60003	Phylogeny using RAXML-kalign	raxml	/root/exp/raxml/1/
SciPhy-60003	Phylogeny using RAXML-kalign	readseq	/root/exp/readseq/5/
SciPhy-60003	Phylogeny using RAXML-kalign	readseq	/root/exp/readseq/4/
SciPhy-60003	Phylogeny using RAXML-kalign	readseq	/root/exp/readseq/3/
SciPhy-60003	Phylogeny using RAXML-kalign	readseq	/root/exp/readseq/2/
SciPhy-60003	Phylogeny using RAXML-kalign	readseq	/root/exp/readseq/1/
SciPhy-60003	Phylogeny using RAXML-kalign	mafft	/root/exp/mafft/2/
SciPhy-60003	Phylogeny using RAXML-kalign	mafft	/root/exp/mafft/3/
SciPhy-60003	Phylogeny using RAXML-kalign	mafft	/root/exp/mafft/4/
SciPhy-60003	Phylogeny using RAXML-kalign	mafft	/root/exp/mafft/1/
SciPhy-60003	Phylogeny using RAXML-kalign	mafft	/root/exp/mafft/5/

Consultas de Proveniência no SciCumulus

 "Recuperar, por ordem crescente de execuções dos workflows, as datas de início e término, tags dos workflows, bem como o nome de todas as atividades associadas e que não contenham nenhuma cloud activity que executou com erro".

```
SELECT w.tag,
       a.tag,
       t.exitstatus,
       t.processor,
       t.workspace,
       t.status,
       t.endtime,
       t.starttime,
       extract ('epoch' from (t.endtime-t.starttime))||',' as duration
from hworkflow w, hactivity a, hactivation t
where w.wkfid = a.wkfid
and a.actid = t.actid
and not exists (select * from hactivation a2
               where a2.actid = a.actid
               and a2.exitstatus <> 0)
order by w.wkfid
```

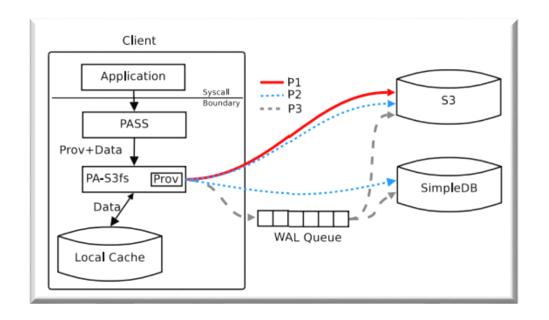
mafft5	0	1	/root/exp/mafft5/47/	FINISHED	- 1
mafft5	0	1	/root/exp/mafft5/62/	FINISHED	
mafft5	0	2	/root/exp/mafft5/77/	FINISHED	
mafft5	0	1	/root/exp/mafft5/105/	FINISHED	
mafft5	0	2	/root/exp/mafft5/120/	FINISHED	
mafft5	0	1	/root/exp/mafft5/136/	FINISHED	
mafft5	0	2	/root/exp/mafft5/149/	FINISHED	
mafft5	0	2	/root/exp/mafft5/164/	FINISHED	
mafft5	0	1	/root/exp/mafft5/186/	FINISHED	
mafft5	0	2	/root/exp/mafft5/200/	FINISHED	
mafft5	0	1	/root/exp/mafft5/1/	FINISHED	
mafft5	0	1	/root/exp/mafft5/3/	FINISHED	
mafft5	0	2	/root/exp/mafft5/2/	FINISHED	
mafft5	0	1	/root/exp/mafft5/4/	FINISHED	
	mafft5	mafft5 0 mafft5 0	mafft5 0 1 mafft5 0 2 mafft5 0 1 mafft5 0 2 mafft5 0 1 mafft5 0 2 mafft5 0 2 mafft5 0 1 mafft5 0 2 mafft5 0 2 mafft5 0 1 mafft5 0 1 mafft5 0 2	mafft5 0 1 /root/exp/mafft5/62/ mafft5 0 2 /root/exp/mafft5/77/ mafft5 0 1 /root/exp/mafft5/105/ mafft5 0 2 /root/exp/mafft5/120/ mafft5 0 1 /root/exp/mafft5/136/ mafft5 0 2 /root/exp/mafft5/149/ mafft5 0 2 /root/exp/mafft5/164/ mafft5 0 1 /root/exp/mafft5/186/ mafft5 0 2 /root/exp/mafft5/200/ mafft5 0 1 /root/exp/mafft5/1/ mafft5 0 1 /root/exp/mafft5/1/ mafft5 0 1 /root/exp/mafft5/3/ mafft5 0 1 /root/exp/mafft5/3/ mafft5 0 2 /root/exp/mafft5/3/	mafft5 0 1 /root/exp/mafft5/62/ FINISHED mafft5 0 2 /root/exp/mafft5/77/ FINISHED mafft5 0 1 /root/exp/mafft5/105/ FINISHED mafft5 0 2 /root/exp/mafft5/120/ FINISHED mafft5 0 1 /root/exp/mafft5/136/ FINISHED mafft5 0 2 /root/exp/mafft5/149/ FINISHED mafft5 0 2 /root/exp/mafft5/164/ FINISHED mafft5 0 1 /root/exp/mafft5/186/ FINISHED mafft5 0 2 /root/exp/mafft5/200/ FINISHED mafft5 0 1 /root/exp/mafft5/1/ FINISHED mafft5 0 1 /root/exp/mafft5/3/ FINISHED mafft5 0 1 /root/exp/mafft5/3/ FINISHED mafft5 0 1 /root/exp/mafft5/3/ FINISHED

Karma

- http://www.extreme.indiana.edu/karma
- Mecanismo de Captura: cada serviço envia dados de proveniência a um serviço central
- · Modelo de Dados:
 - Retrospectiva → grafo representado em XML
 - Prospectiva → BPEL
- Minimiza sobrecarga de captura
 - Desacoplado da máquina de workflow

PASS

- http://www.eecs.harvard.edu/syrah/pass/
- Mecanismo de Captura: cada serviço envia dados de proveniência a um serviço central





Proveniência em tempo real

Por que dados de proveniência em tempo real são importantes?



Proveniência



Especialmente em ambientes de nuvem

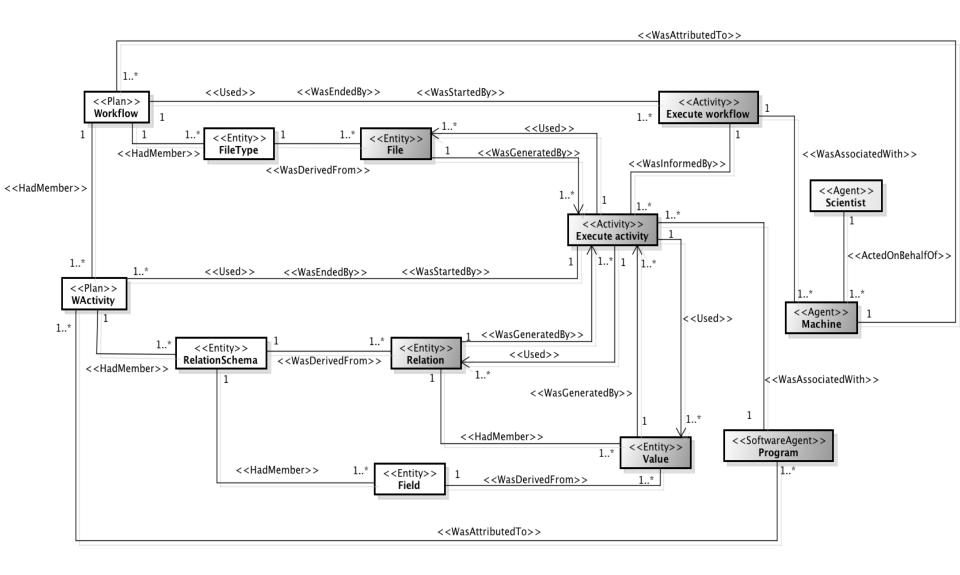


PROV-Wf

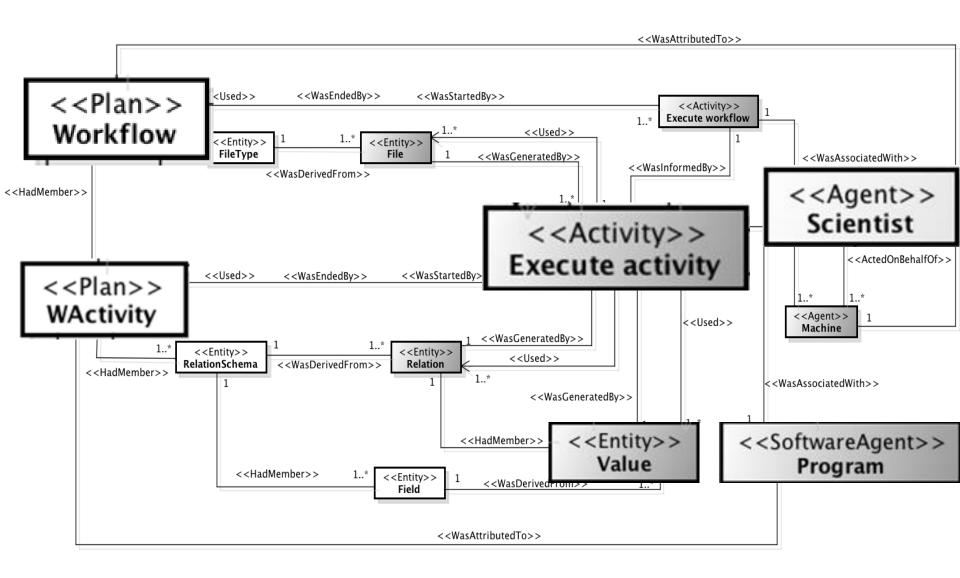
- Extensão do PROV para representação de proveniência em workflows científicos
- Relaciona elementos do modelo estendido com os elementos do PROV por meio de estereótipos da UML
- Proposto no workshop BigProv do EDBT/ICDT 2013



PROV-Wf



PROV-Wf



Mapeamento no PROV-Wf

VisTrails	SciCumulus	PROV-Wf
VisTrail.ID	Workflow.ExecTag	Execute Workflow.ID
VisTrail.Name	Workflow.Tag	Workflow. Name
Module.Name	Activity.Tag	WActivity. Name
Module.ID	Activation.ID	Execute Activity.ID
Parameter. Alias	Field.Name	Field.Name
Parameter.Type	Field.Type	Field.Type
Parameter.Val	Value.Val	Value.Val

Uso do PROV-Wf

SELECT f1.name, v1.value, f2.name, v2.value

FROM workflow w, activity a, executeActivity ea, relationschema rs, relation r,

field f1, field f2, Value v1, Value v2

WHERE a.wkf_id = w.id AND ea.act_id = a.id

AND r.rel_id = rs.id **AND** r.ea_id = ea.id

AND f1.rel_id = rs.id **AND** f2.rel_id = rs.id

AND v1.uf_id = f1.id **AND** v2.uf_id = f2.id

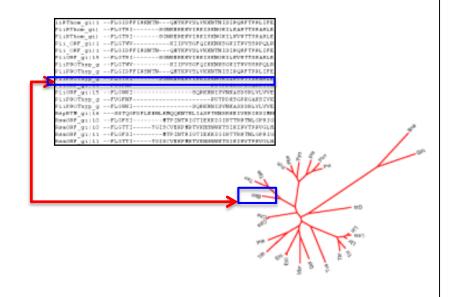
AND rs.name = "Relation MSA"

AND f1.name = "BIO SEQ"

AND f2.name = "TAXON"

AND a.name = %PARAM_ACT_NAME%

AND w.tag = %PARAM_WF_NAME%



Roteiro do Tutorial

- Motivação
- Workflows Científicos
- · Nuvens de Computadores
- · Proveniência de Dados
- · Máquinas de Workflow para Nuvem
- · Aplicação de Proveniência em e-Science
- Demo



Máquinas baseadas em Hadoop

- Framework para execução de aplicações em grandes clusters (virtuais ou não)
- Fornece funcionalidades tanto de confiabilidade e movimentação de dados.
- Implementa o paradigma computacional chamado map / reduce
 - Cada aplicação é "dividida" em muitos pequenos fragmentos, cada um dos quais podem ser executadas ou reexecutado em qualquer nó no cluster.

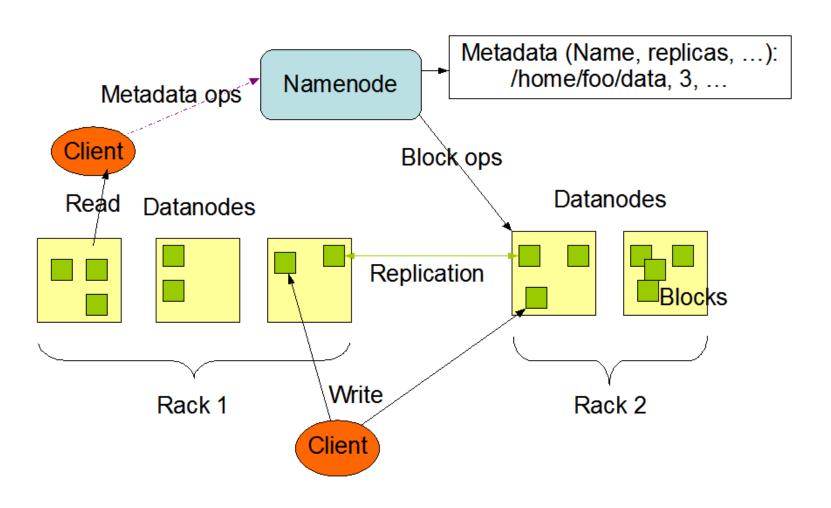
Hadoop



- · Desacoplado do conceito de workflow
- Escalonamento não leva em consideração a estrutura do workflow
- Não captura proveniência de forma nativa
 - Complementos para captura são necessários

Arquitetura do Hadoop

HDFS Architecture

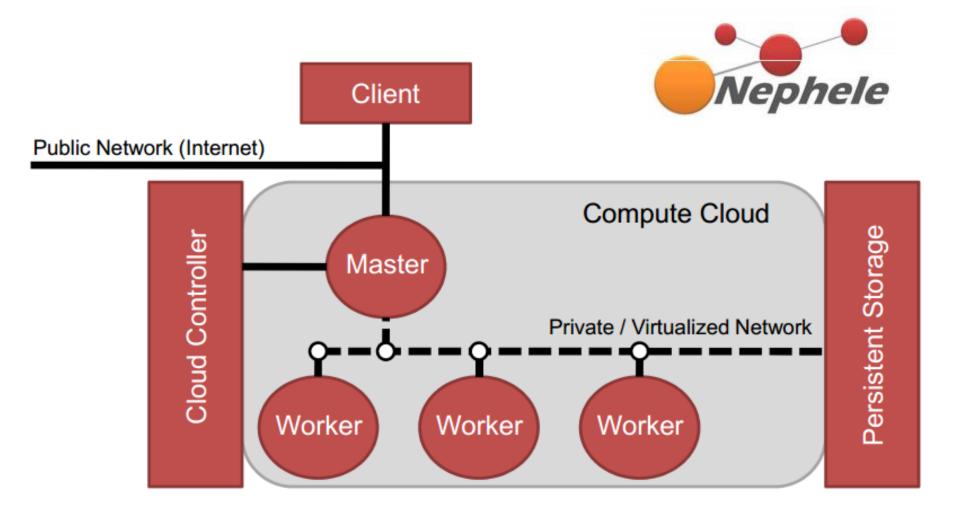


Nephele

- Arcabouço para processamento paralelo em nuvens
- Aloca e desaloca máquinas virtuais dependendo do comportamento do Job
- · Não possui proveniência associada
- Paralelização explícita
 - Tarefas são designadas como "parallelizable"



Nephele



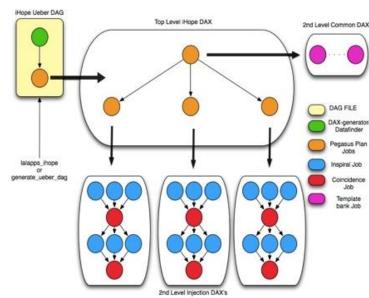
Pegasus

- •Financiado pela NSF/OCI em colaboração com o grupo do escalonador CONDOR de UW Madison
- Distribui atividades em paralelo em ambientes distribuídos
- ·Captura os dados de proveniência
- ·Escalabilidade
 - Big Data (kB to TB)
 - Quantidade de tarefas (1...106 tarefas)

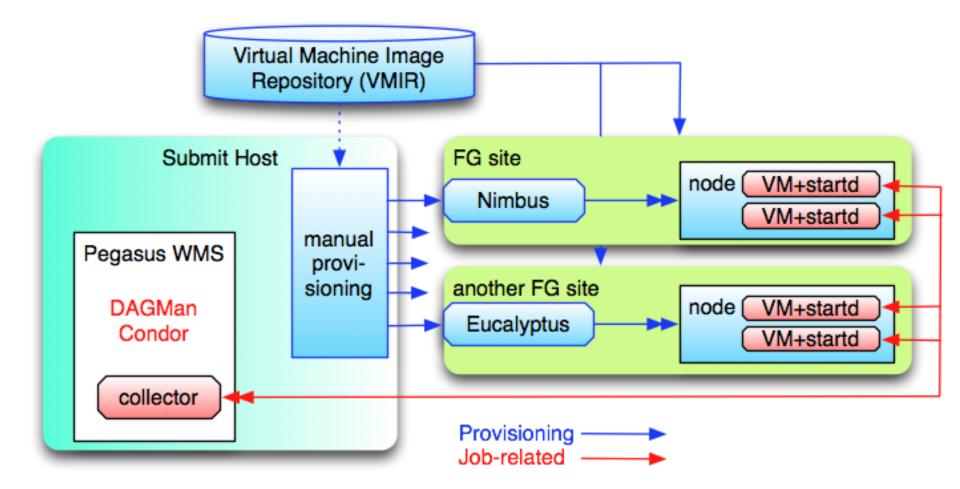


Pegasus

- Re-executa tarefas que falharam automaticamente
- Pode executar nos mais diversos ambientes
 - Laptop
 - campus cluster
 - grid
 - nuvem



Pegasus



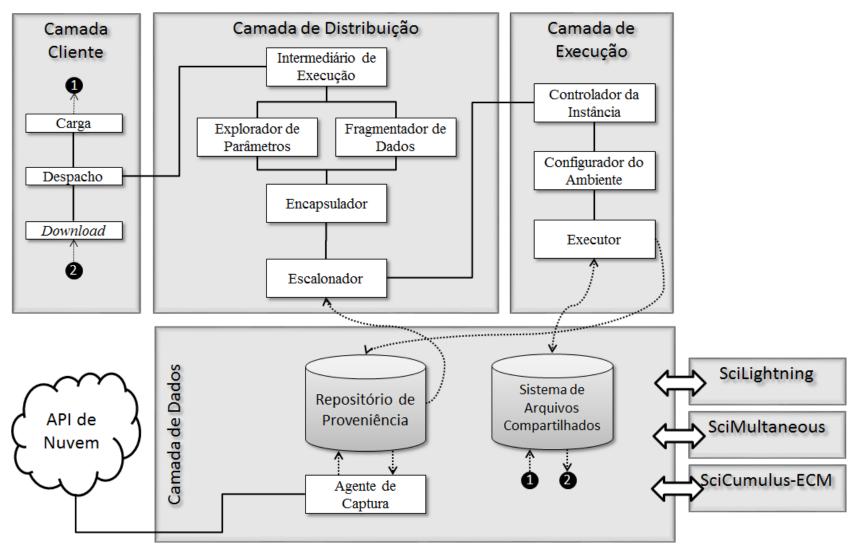
Fonte: pegasus.isi.edu/presentations/

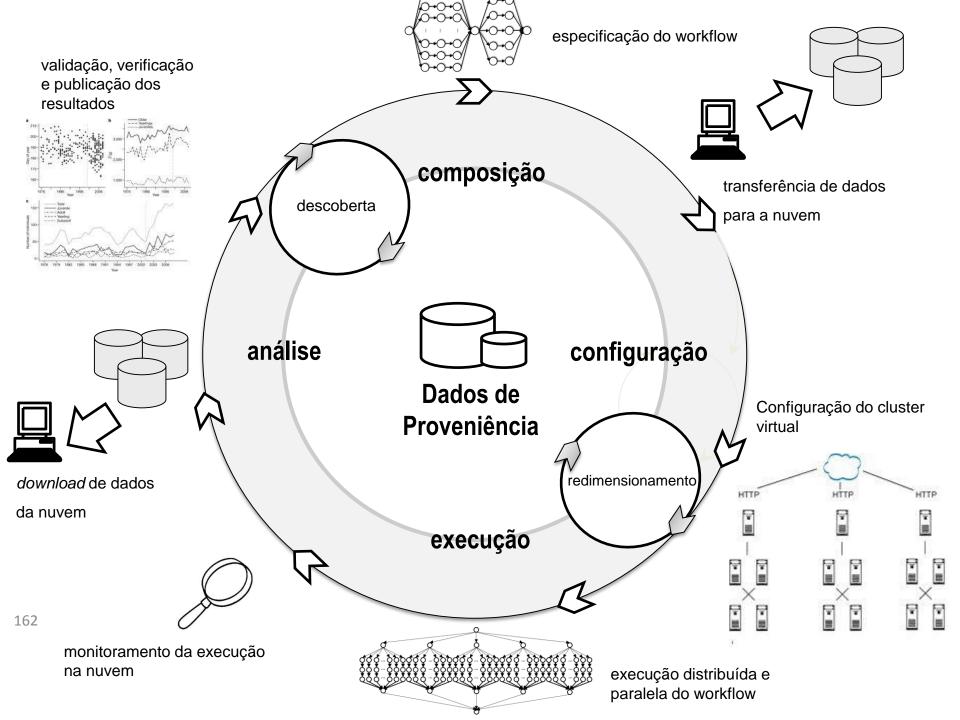
SciCumulus

- Uma máquina para execução de workflows para nuvem que tem como objetivo conectar o SGWfC com o ambiente de nuvem para oferecer apoio à execução de workflows científicos
 - Middleware Hydra
 - Motor de execução Chiron

OLIVEIRA, D.; OGASAWARA, E.; BAIAO, F.; MATTOSO, M. L. Q. . **SciCumulus: A Lightweight Cloud Middleware to Explore Many Task Computing Paradigm in Scientific Workflows**. In: The 3rd IEEE CLOUD 2010, Miami. Proc. of the 3rd IEEE CLOUD, 2010. p. 378-385

Arquitetura





Modelo Algébrico

- O SciCumulus segue o modelo algébrico proposto por Ogasawara et al. (2011)
 - Atividades consomem relações
 - Cada relação é composta de tuplas e atributos (parâmetros das atividades)
 - Execuções de atividades consomem tuplas
- · Os operadores da álgebra Invocam programas:
 - Map (1:1)
 - SplitMap (1:n)
 - Reduce (n:1)
 - Filter (1:0-1)

OGASAWARA, E.; DIAS, J.; OLIVEIRA, D.; PORTO, F.; VALDURIEZ, P.; MATTOSO, M. L. Q. . **An Algebraic Approach for Data-Centric Scientific Workflows**. Proceedings of the **VLDB** Endownment, v. 4, p. 1328-13369, 2011.

Relações como o modelo de dados para consumo e produção

- Relações são definidas como um conjunto de tuplas de tipos primitivos (inteiro, string, etc) ou tipos complexos (e.g. ponteiros para arquivos)
- Exemplo: R(R)

RID	<u>CaseStudy</u>	sdat	ddat
1	U-125	U-125S.DAT	U-125D.DAT
1	U-127	U-127S.DAT	U-127D.DAT
2	U-129	U-129S.DAT	U-129D.DAT

• R = (RID: Integer, CaseStudy: String; SDat: FileRef, DDat: FileRef)

Split Map (SplitMap) $T \leftarrow SplitMap(Y, a, R)$

R	<u>RID</u>	RdZip	
— 1 Project1.zip		Project1.zip	
	2 Project2.zip		

T ← SlipMap(extractRD, 'RdZip', R)

Т	RID	<u>Study</u>	sdat	ddat
\longrightarrow	1	U-125	U-125S.DAT	U-125D.DAT
	1	U-127	U-127S.DAT	U-127D.DAT
	2	U-129	U-129S.DAT	U-129D.DAT

Reduce Activity (Reduce) $T \leftarrow Reduce(Y, g_A, R)$

R	RID	<u>Study</u>	SsSai	DdSai	MEnv
	1	U-125	U-125Ss.SAI	U-125Dd.SAI	U-125.ENV
	1	U-127	U-127Ss.SAI	U-127Dd.SAI	U-127.ENV
	2	U-129	U-129Ss.SAI	U-129Dd.SAI	U-129.ENV

 $T \leftarrow Reduce(CompressRD, \{'RID'\}, R)$

T

RID	RdResultZip	
1	ProjectResult1.zip	
2	ProjectResult2.zip	

SRQuery Activity $T \leftarrow SRQuery(qry, R)$

D
ĸ
••

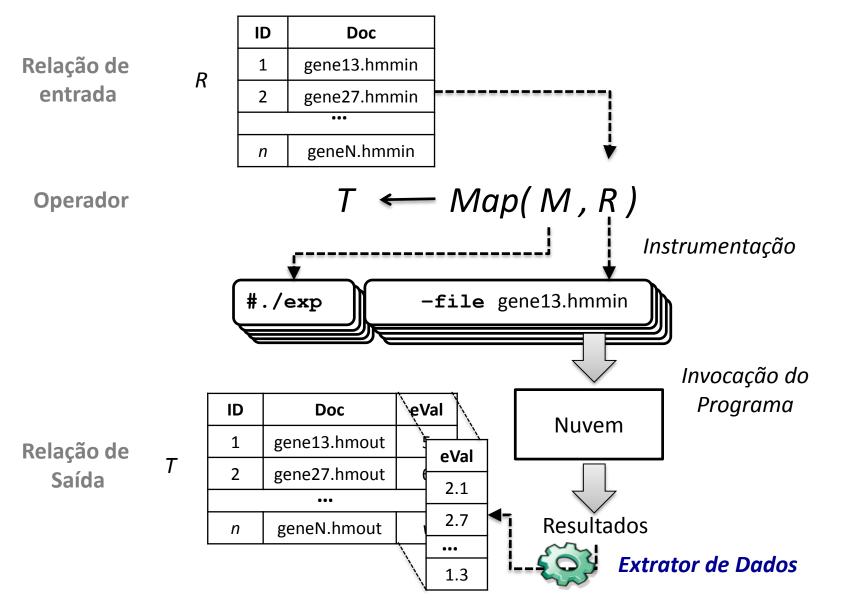
RID	<u>Study</u>	SsSai	Curvature
1	U-125	U-125Ss.SAI	1.5
1	U-126	U-126Ss.SAI	0.9
1	U-127	U-127Ss.SAI	1.2

 $T \leftarrow SRQuery(\pi_{RID, Study, SsSai, Curvature}(\sigma_{Curvature>1}(R)), R)$

T

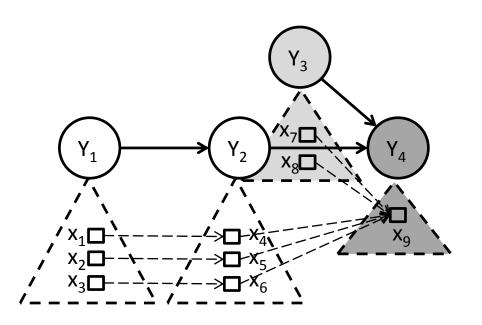
RID	<u>Study</u>	SsSai	Curvature
1	U-125	U-125Ss.SAI	1.5
1	U-127	U-127Ss.SAI	1.2

Execução baseada em relações



Estratégias de Execução

- First Tuple First (FTF)
- First Activity First (FAF)



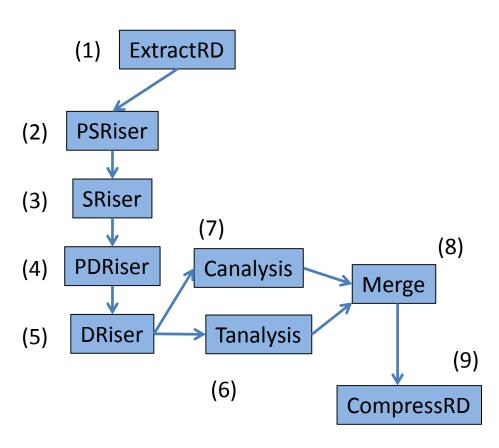
FTF:

$$\{\langle x_1, x_4 \rangle, \langle x_2, x_5 \rangle, \langle x_3, x_6 \rangle\}$$

FAF:

Workflow como uma expressão da álgebra

Workflow



Expressões algébricas

```
T_1 \leftarrow SplitMap(ExtractRD, R_1)
T_2 \leftarrow Map(PSRiser, T_1)
T_3 \leftarrow Map(SRiser, T_2)
T_4 \leftarrow Map(PDRiser, T_3)
T_5 \leftarrow Map(DRiser, T_4)
T_6 \leftarrow Filter(Tanalysis, T_5)
T_7 \leftarrow Filter(Canalysis, T_5)
T_8 \leftarrow SRQuery(T_6 \bowtie T_7, \{T_6,
\mathsf{T}_7\})
T_9 \leftarrow Reduce(CompressRD, T_8)
```

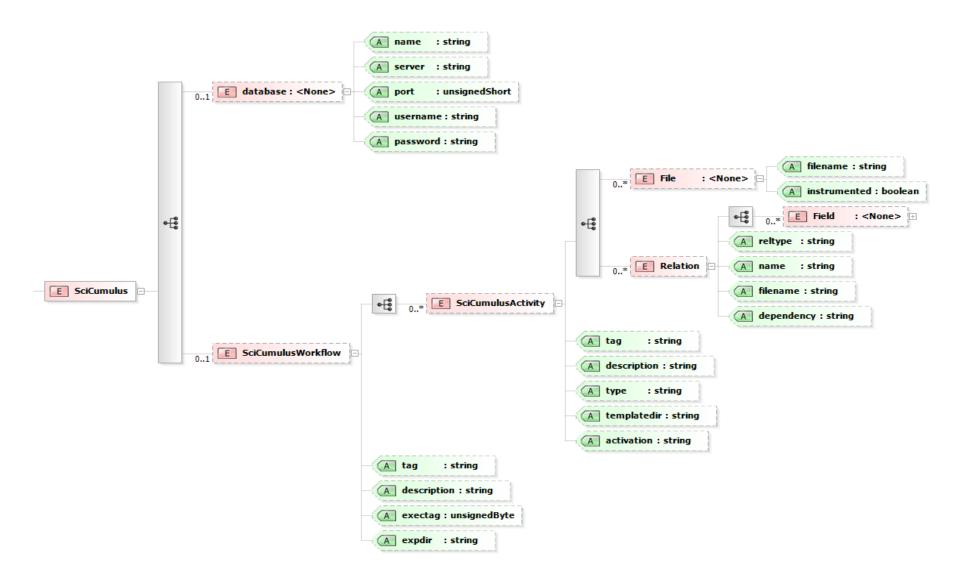
Detalhes de Implementação

- Java Versão 6 Update 15
- HSQLDB
- Drivers JDBC
- MPJ arcabouço de paralelismo em Java
- PostgreSQL versão 8.4.6
- · Código disponibilizado em:
 - https://sourceforge.net/projects/scicumulus/

Representação dos Workflows

```
<?xml version="1.0" standalone="no"?>
<SciCumulus>
  <database name="scicumulus_adaptive" server="mp4-4.dyndns.info"</pre>
    port="5432" username="scicumulus" password="********** />
  <SciCumulus Workflow tag="filogenia" description="This is a test using anflex."
    exectag="Experimento Kary/200 organismos - Adaptive" expdir="/root/exp">
    <SciCumulus Activity tag="mafft" description="mafft" type="MAP"
           templatedir="/root/exp/template_mafft" activation="experiment.cmd">
      <Relation reltype="Input" name="A" filename="parameter.txt" />
      <Field name="NAME" ype="string"/>
      <Field name="FASTA_FILE" type="string" />
      <Relation reltype="Output" name="C" file name="output_mafft.txt" />
      <Field name="NAME" type="str NAME;FASTA_FILE</pre>
                                     GI;ORTHOMCL2033
      <Field name="NUM_SEQ" type=
                                      G2:ORTHOMCL1895
      <Field name="FASTA_FILE" ty
                                      G3;ORTHOMCL2034
      <File filename="experiment.cmd"
                                     G4:ORTHOMCL1896
    </SciCumulusActivity>
                                      G5;ORTHOMCL2035
                                      G6;ORTHOMCL1897
  </SciCumulus Workflow>
                                      G7;ORTHOMCL2036
</SciCumulus>
                                      G8;ORTHOMCL1898
                                      G9;ORTHOMCL2037
                                      G10;ORTHOMCL1899
```

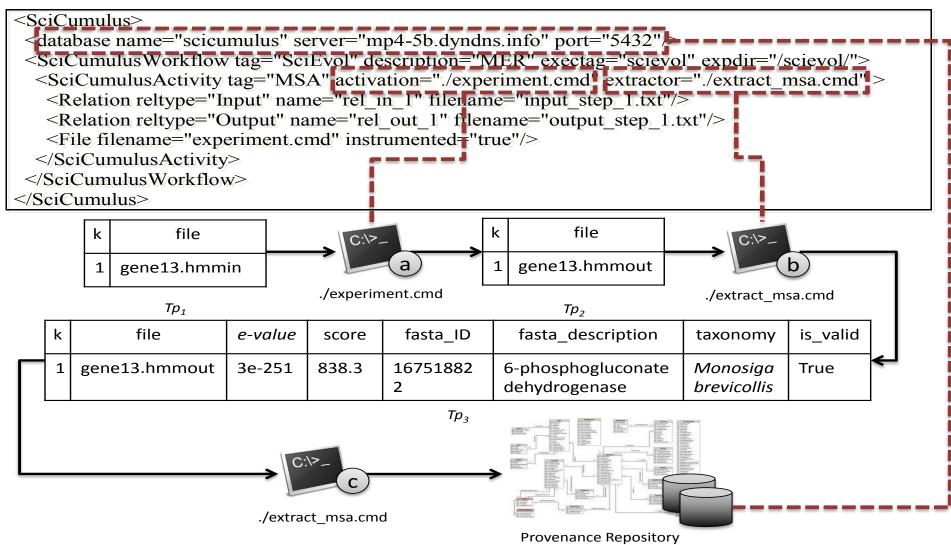
Representação dos Workflows



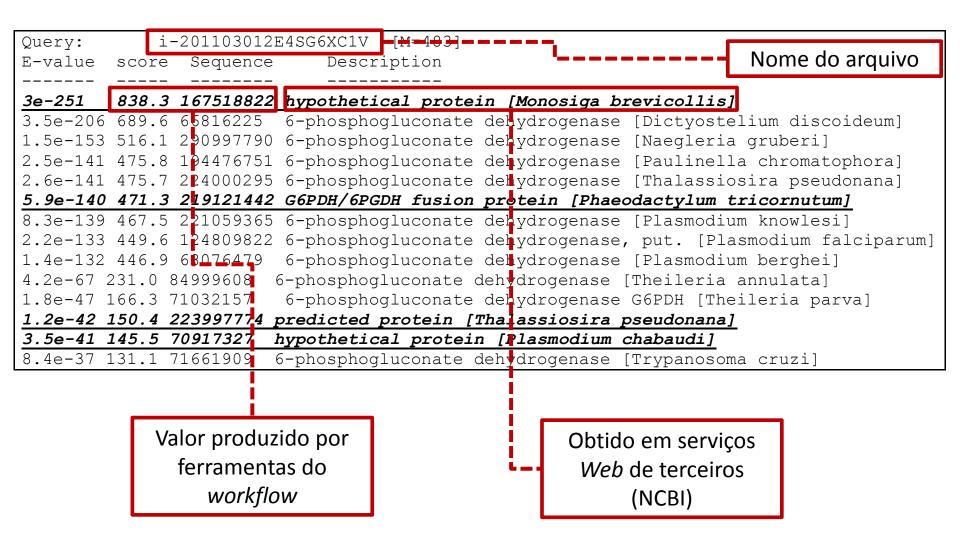
Componente extrator de dados

- Como extrair dados de diferentes formatos (e.g. binário, Fasta, HDF5, etc.)?
- Componente extrator (baseado na álgebra)
 - Invoca um programa externo (definido pelo usuário) que analisa os arquivos produzidos e extrai os metadados de interesse
 - Encapsula regras de extração que são dependents de domínio
 - Relações possuem dados da execução com dados do domínio

Extração de dados de domínio



Extração de Dados de Domínio



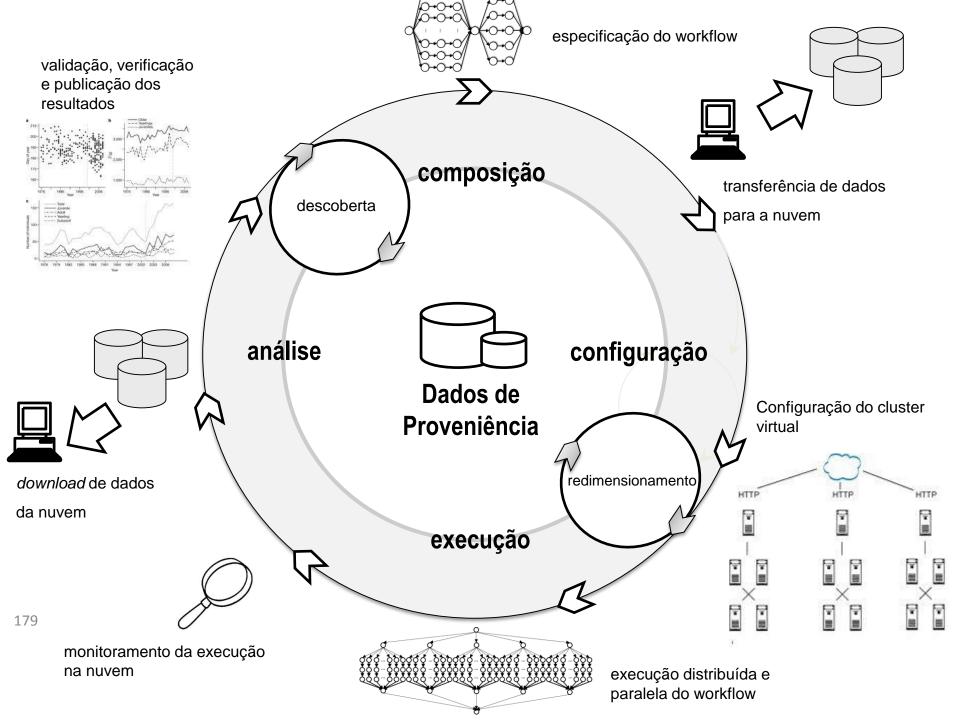
Consulta utilizando dados de domínio

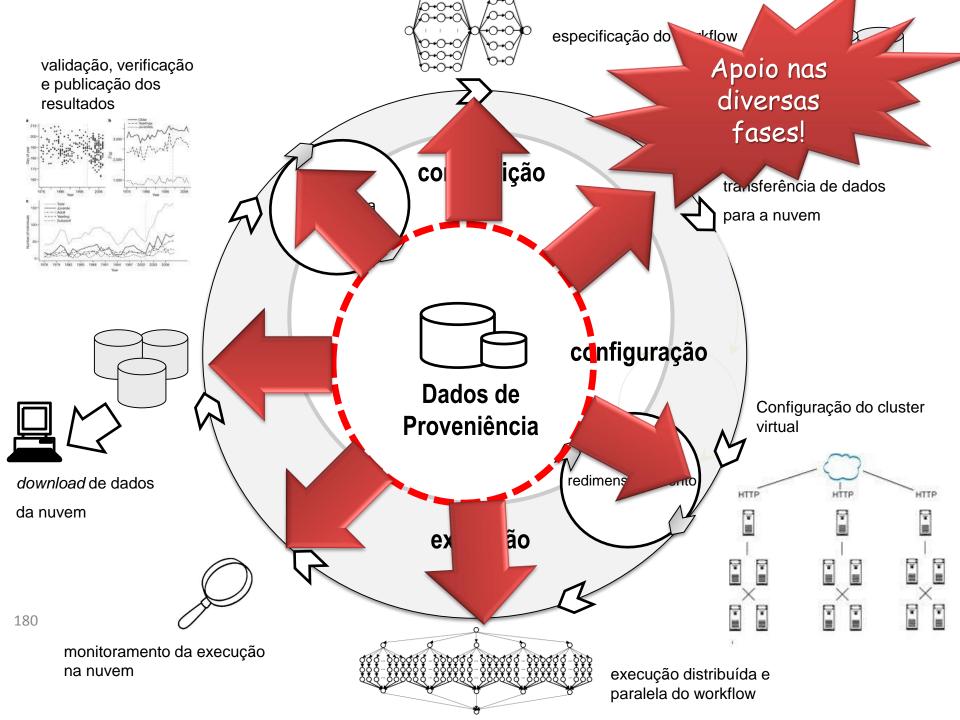
 Quais sequencias dadas como entrada não fazem parte de um gene específico?

```
WHERE T.taskid = H.taskid AND H.msacid = M2.msacid AND M2.msaid = M1.msaid AND M1.sgid = SG.sgid AND SG.sqid = S.sqid AND S.orgid = O.orgid AND O.spid = SP.spid and SP.genid = G.genid AND T.exitStatus = 0 /* No error */
AND G.genus = "PLASMODIUM"
```

Roteiro do Tutorial

- Motivação
- Workflows Científicos
- · Nuvens de Computadores
- Proveniência de Dados
- Máquina de Execução SciCumulus
- · Aplicação de Proveniência em
- e-Science na Nuvem
- Demo





O uso da proveniência



Dimensionamento do Ambiente

- · Recursos são instanciados sob demanda
- Vários tipos de máquinas virtuais passíveis de uso
 - Capacidade de armazenamento
 - CPU
 - Memória
 - Largura de banda
- Variação de provedor para provedor

· Recursos são instanciados sob demanda

Qual tipo utilizar?
Quantas máquinas instanciar?
Força bruta?

- Memória
- Largura de banda
- · Variação de provedor para pro





Dimensioning the Virtual Cluster for Parallel Scientific Workflows in Clouds

Daniel de Oliveira¹, Vitor Viana³ 1IGUEE danielcmo@ic.uff.br

Eduardo Ogasawara² 2CFFFT/B.I. eogasawara@cefet-ri.br

Kary Ocaña³, Marta Mattoso³ {vviana,kary,marta}@cos.ufrj.br

4th Workshop on Scientific Cloud Computing (ScienceCloud) 2013

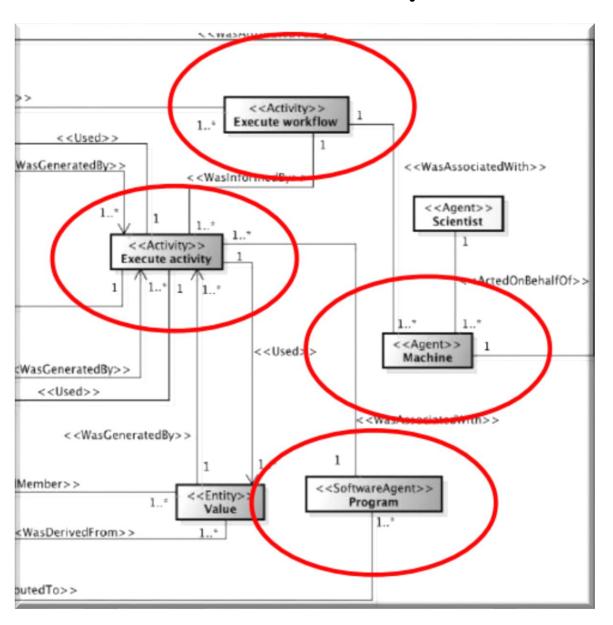
Co-located with ACM HPDC 2013 New York City, NY, USA -- June 17th, 2013 workflow in parallel may be a hard task to accomplish, i.e. it is hard to define and adapt the optimal number of virtual machines to be used. Most systems follow this manual configuration of the scientist for the whole workflow execution, using adaptive techniques only in the presence of failures. Due to the huge number of options (virtual machine types) to configure a cloud environment, the configuration task commonly becomes impractical to be performed

executed, generating idle computing time along the experiment lifecycle. This elastic allocation of resources, on demand, provides a new dimension for HPC applications.

Some Scientific Workflow Management Systems (SWfMS) [11,12] bridge the gap between cloud environments and the management of scientific experiments, such as, SciCumulus engine [4,13], Pegasus [14], Swift [15], and specifically for bioinformatics, the combination of Taverna and Galaxy named Tavaxy [16]. In these SWfMS a virtual cluster is created in the cloud to execute the

- Dimensionamento leva em consideração características do workflow
- · Histórico de execução é utilizado
 - Proveniência retrospectiva
 - Características do ambiente
- Redução do tempo de execução e dos custos
 - Evita o super e o sub dimensionamento!

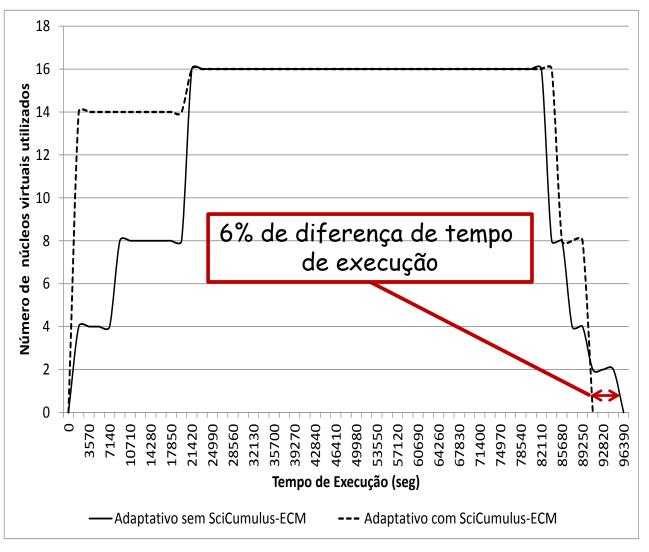
Quais dados são importantes?



SciCumulus-ECM

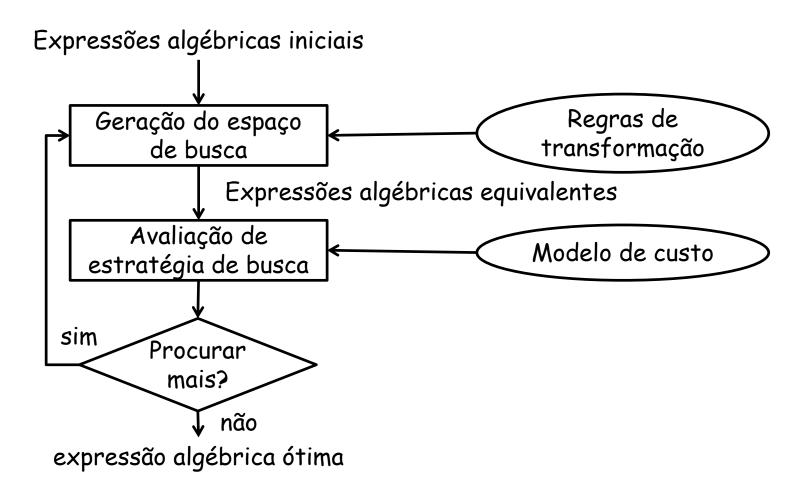
- · Baseado em algoritmos genéticos
- Alimenta o algoritmo com os dados históricos de proveniência
- Com uma base de proveniência de tamanho médio (aprox. 100 execuções de workflows) o algoritmo apresenta uma convergência rápida
- Utilização do arcabouço JGAP

Análise da Execução Adaptativa Otimizando a Quantidade de Máquinas Virtuais



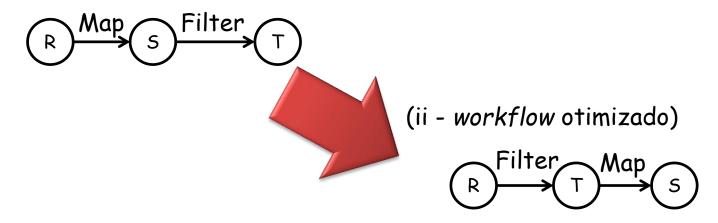
- 400 arquivos multi-fasta
- Cenário com foco no tempo de execução
- Orçamento limite de US\$ 75,00

Processo de otimização de workflows



Transformações algébricas

(i - workflow - perspectiva de relações)



- Garantia do esquema das atividades
- Similar a otimização de consultas na álgebra relacional
- Antecipação do Filtro

 Antecipação da seleção

Escalonamento Adaptativo

- O ambiente de nuvem possui instabilidades
 - Falhas nas VMs
 - Aplicação de atualizações
 - Reinicializações
 - Heterogeneidade das VMs
- O escalonamento deve levar essas características em consideração

Escalonamento Adaptativo

```
-f(a_i, VM_j) =
\alpha_1 \times (\text{Tempo de Execução} + \text{Tempo de Transferência}) +
\alpha_2 \times Custo de Confiabilidade
+
\alpha_3 \times (\text{Vh} \times [\text{Tempo de Execução} + \text{Tempo de Transferência}] + \text{Vt} \times \text{Quantidade de Dados}
+ \text{Transferidos})
```

• A a_i escolhida é aquela que satisfaz $f(a_i, VM_j) = min \ \forall VM_j \in VM \{f(a_i, VM_j)\}$

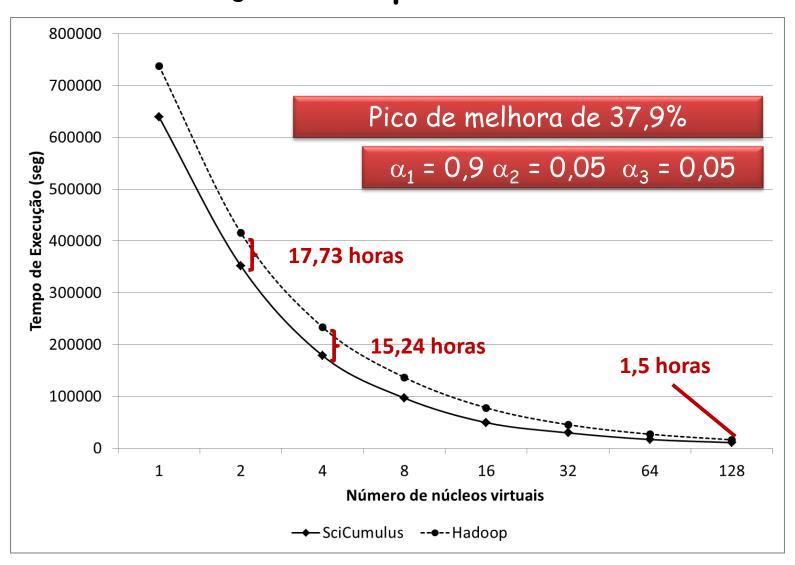
Escalonamento Adaptativo

```
α<sub>2</sub>× Custo de as informações são

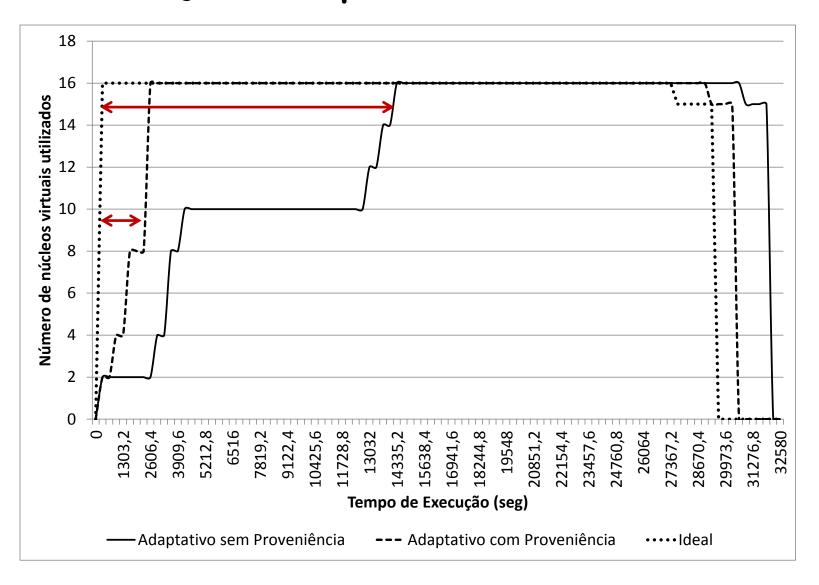
+
Todas as informações são

todas as informações são
               buscadas no repositório de
Todas as interespositoriem em todas as no repositoriem algebrica buscadas no repositoriem algebrica buscadas no repositoria e levam algebrica proveniencia e peração as proveniencia e peração a operação a operação a que satisfaz f(a_i, VM_j) = consider a f(a_i, VM_j)}
```

Execução Adaptativa - sem variação na quantidade de VMs

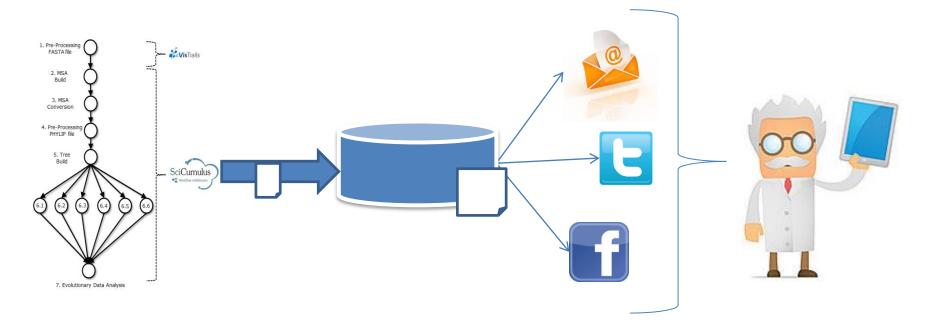


Execução Adaptativa - variação na quantidade de VMs



Monitoramento

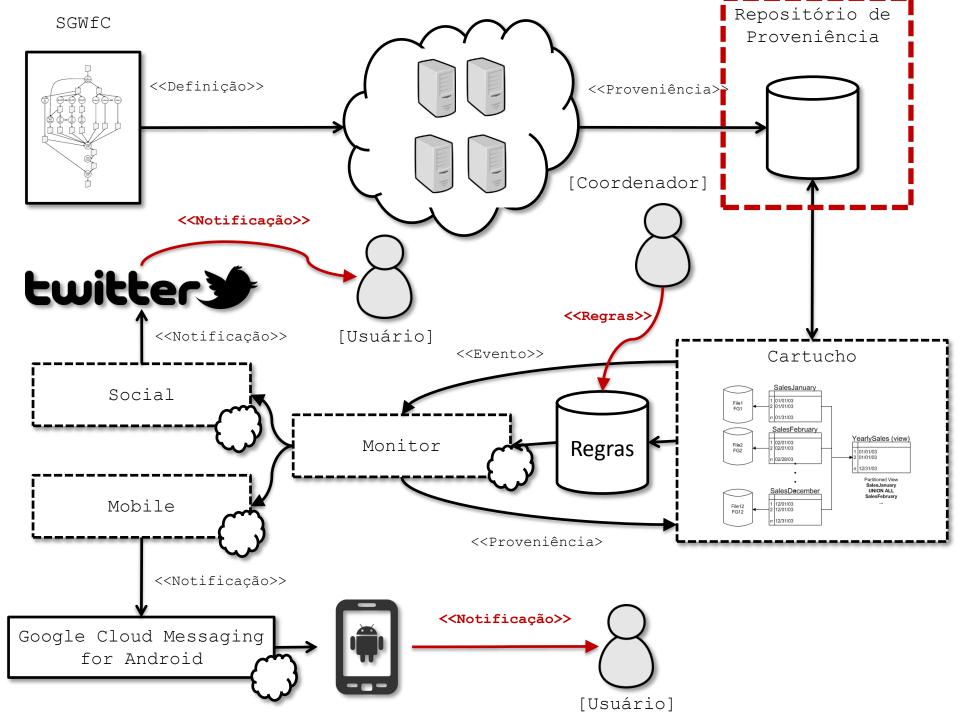
- Execuções podem durar semanas ou meses
- Monitoramento in situ é inviável
- Utilização de redes sociais e computação móvel



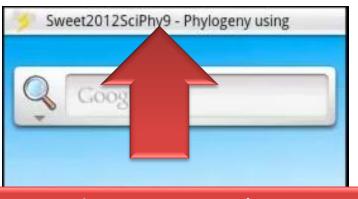
SciLightning

- Serviço de monitoramento de execução de workflows científicos em paralelo
- 3 componentes principais:
 - SciLightning Monitor
 - SciLightning Social
 - SciLightning Mobile
- Todo o monitoramento é baseado nos dados de proveniência gerados pelo SciCumulus

PINTAS, J.; OLIVEIRA, D.; OCANA, K.; DIAS, J.; MATTOSO, M. L. Q. . Monitoramento em Tempo Real de Workflows Científicos Executados em Paralelo em Ambientes Distribuídos. In: VI e-Science Workshop (em conjunto com CSBC 2012), 2012, Curitiba. Anais do VI e-Science Workshop (em conjunto com CSBC 2012). Porto Alegre: Sociedade Brasileira de Computação, 2012.



SciLightning



Recebimento de aviso assíncrono





Recuperação de Falhas

- Falha de VMs é um problema constante
- O provedor garante a reposição de uma VM, mas não de seu contexto
- O cientista não deve re-executar atividades que já foram executadas anteriormente
- Proveniência fornece esse tipo de informação de contexto da execução

SciMultaneous

- Serviço de re-execução de atividades
- Explora características da nuvem como elasticidade por meio de duas heurísticas
 - H1: Redundância de atividades
 - H2: Monitoramento contínuo

COSTA, F.; OLIVEIRA, D.; OCANA, K.; OGASAWARA, E.; MATTOSO, M. L. Q. . **Enabling Re-Executions of Parallel Scientific Workflows Using Runtime Provenance Data**. In: 4th International Provenance and Annotation Workshop, 2012, Santa Barbara, CA. Proceedings of the 4th International Provenance and Annotation Workshop. Heildelberg: Springer Verlag, 2012.

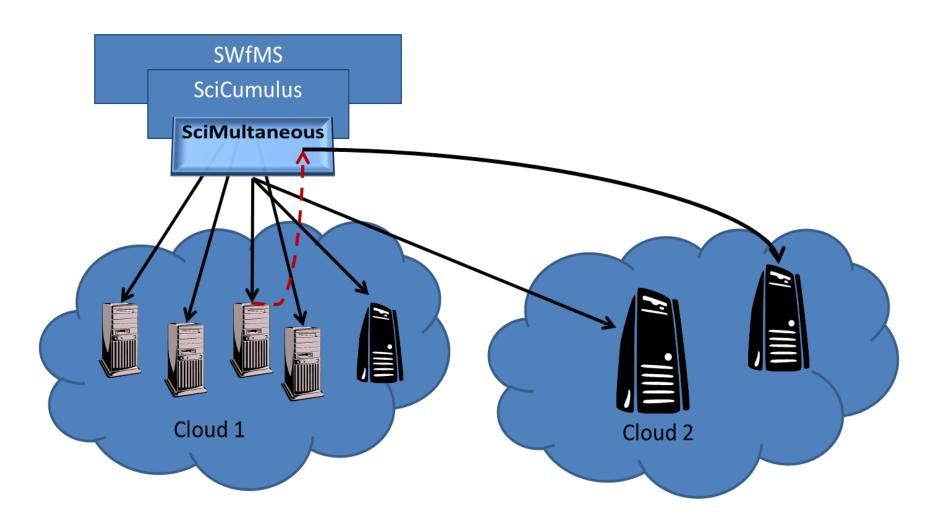
Consulta de Proveniência

 "Recuperar, por ordem crescente de execuções dos workflows, as datas de início e término, tags dos workflows, bem como o nome de todas as atividades associadas e que contenham alguma execução com erro".

SELECT W +20

Select w.tag,	4					
a.tag,	MafftAdaptive	mafft5	0	1	/root/exp/mafft5/47/	FINISHED
t.exitstatus,	MafftAdaptive	mafft5	0	1	/root/exp/mafft5/62/	FINISHED
t.processor,	daptive	mafft5	0	2	/root/exp/mafft5/77/	FINISHED
t.workspace,		mafft5	0	1	/root/exp/mafft5/105/	FINISHED
t.status,	CL.	mafft5	0	2	/root/exp/mafft5/120/	FINISHED
t.endtime,	76	∘fft5	0	1	/root/exp/mafft5/136/	FINISHED
t.starttime,			0	2	/root/exp/mafft5/149/	FINISHED
extract ('epoch' from (t.endtime-t.starttime)) ',' as duration		YC2	0	2	/root/exp/mafft5/164/	FINISHED
extract (epoch from (trematime tistarttime)///// as daration					1/100DEXD/Mattt5/104/	FINISHED
from hworkflow w, hactivity a, hactivation t	Matte.	Joe -	· ·	1	/root/exp/mafft5/186/	FINISHED 2
, , , , , , , , , , , , , , , , , , , ,	Mafft Adaptive	roes		1 2		
from hworkflow w, hactivity a, hactivation t	Matte	roes	COL	1 2	/root/exp/mafft5/186/	FINISHED 2
from hworkflow w, hactivity a, hactivation t where w.wkfid = a.wkfid	MafftAdaptive C	roes	Com	1 2	/root/exp/mafft5/186/ /root/exp/mafft5/200/	FINISHED 2 FINISHED 2
from hworkflow w, hactivity a, hactivation t where w.wkfid = a.wkfid and a.actid = t.actid	MafftAdaptive	ro es	Com	1 2	/root/exp/mafft5/186/ /root/exp/mafft5/200/ /root/exp/mafft5/1/	FINISHED 2 FINISHED 2 FINISHED 2
from hworkflow w, hactivity a, hactivation t where w.wkfid = a.wkfid and a.actid = t.actid and not exists (select * from hactivation a2	MafftAdaptive MafftAdaptive MafftAdaptive	maffts maffts	Com	1 2	/root/exp/mafft5/186/ /root/exp/mafft5/200/ /root/exp/mafft5/1/ /root/exp/mafft5/3/	FINISHED 2 FINISHED 2 FINISHED 2

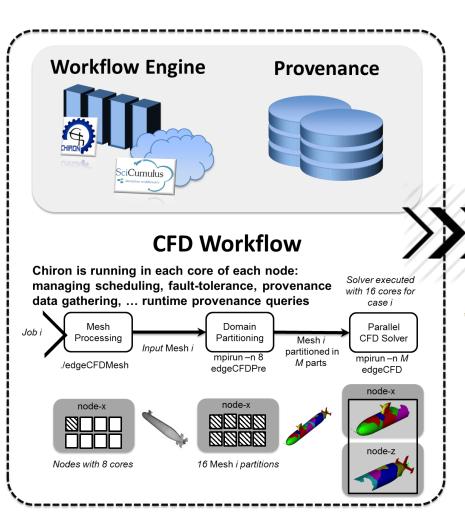
SciMultaneous



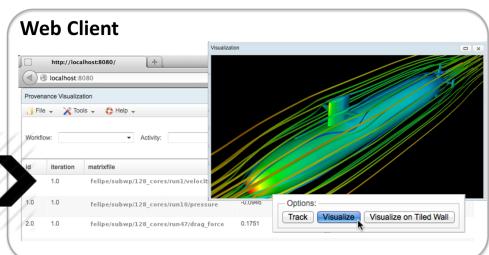
Visualização

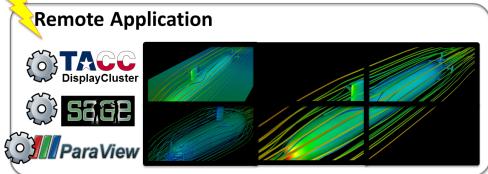
- A visualização dos dados dos experimentos é uma tarefa fundamental de análise
- Grande volume de dados
- Necessidade de associar metadados a imagens e vídeos produzidos
- Necessidade de utilização de ambientes de visualização como tiled walls displays

PROV-Vis



Prov-Vis





PROV-Vis



Remote Application

Roteiro do Tutorial

- Motivação
- Workflows Científicos
- · Nuvens de Computadores
- · Proveniência de Dados
- Máquina de Execução SciCumulus
- · Aplicação de Proveniência em e-Science
- Demo



Demo SciCumulus

- Utilização do ambiente Amazon EC2
 - Necessário que o usuário possua uma conta
- Imagem pública: ami-6e1a8907
- Linux CentOS 5.5
- SciCumulus
 - java jar scicumulus-core.jar 0/root/expSciPhy/machines.conf niodev MPI2 /root/expSciPhy/scicumulus.xml

Download

Home / Browse / Science & Engineering / Bio-Informatics / SciCumulus / Support





Brought to you by: eogasawara, jonasdias, phydoop

Looking for the latest version? Download SciCumulusCore-0.1-SNAPSHOT.jar (6.8 MB)

Home

Name	Modified	Size	Downloads	
SciCumulus_compact.backup	2013-03-20	11.6 MB	0	
SciCumulusCore-0.1-SNAPSHOT.jar	2013-03-20	6.8 MB	2	
SciCumulusDesktop-1.0-SNAPSHOT.jar	2013-03-20	18.3 kB	1	
sciCumulus.xml	2013-03-20	687 Bytes	1	

```
-rw-r--r 1 root root 2.1K Sep 23 15:04 SciCumulus.xml
drwxr-xr-x 2 root root 4.0K Sep 23 13:01 input
-rw-r--r 1 root root 148 Aug 6 17:10 machines.conf
drwxr-xr-x 2 501 games 12K Aug 8 05:00 output
drwxr-xr-x 8 root root 4.0K Sep 23 13:10 result_bioscicumulus
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_modelgenerator
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_raxml
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_readseq
```

```
-rw-r-r-- 1 root root 2.1K Sep 23 15:04 SciCumulus.xml
drwxr-xr-x 2 root root 4.0K Sep 23 13:01 Input
-rw-r-r-- 1 root root 148 Aug 6 1
drwxr-xr-x 2 501 games 12K Aug 8 0
drwxr-xr-x 8 root root 4.0K Sep 23 1
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_modelgenerator
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template_readseq
```

```
-rw-r--r-- 1 root root 2.1K Sep 23 15:04 SciCumulus.xml
drwxr-xr-x 2 root root 4.0K Sep 23 13:01 input
-rw-r--r-- 1 root root 148 Aug 6 17:10 machines.conf
drwxr-xr-x 2 501 games 12K Aug 8 05:00 output
drwxr-xr-x 8 root root 4.0K Sep 23 13:10 result_bioscicumulus
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
drwxr-xr-x 2 501 games 4.0K Feb 16 2013 template mafft
```

```
<Hydra>
    <database name="scicumulus" password="livre2008!" port="5432" server="ec2-50-17-107-164.compute-1.amazonaws.com" username="pos</pre>
    <HydraWorkflow description="Phylogenetic tree construction" exectag="exsciphy-2" expdir="/root/exp_SciPhy" tag="SciPhy">
        <HydraActivity activation="./experiment.cmd" description="alignment" tag="mafft" templatedir="/root/exp_SciPhy/template_ma</p>
                <Relation filename="parameter.txt" name="A" reltype="Input" />
                <Relation filename="output_mafft.txt" name="C" reltype="Output" />
                <File filename="experiment.cmd" instrumented="true" />
                <Field input="A" name="FASTA_FILE" output="C" type="string" />
        </HydraActivity>
        <HydraActivity activation="./experiment.cmd" description="Format Alignment" tag="readseg" templatedir="/root/exp_SciPhy/te</pre>
                <Relation dependency="mafft" filename="output_mafft.txt" name="A" reltype="Input" />
                <Relation filename="output_readseq.txt" name="C" reltype="Output" />
                <File filename="experiment.cmd" instrumented="true" />
                <Field input="A" name="FASTA_FILE" output="C" type="string" />
        </HydraActivity>
        <HydraActivity activation="./experiment.cmd" description="Search for models" tag="modelgenerator" templatedir="/root/exp S</pre>
                <Relation dependency="readseg" filename="output readseg.txt" name="A" reltype="Input" />
                <Relation filename="output modelgenerator.txt" name="C" reltype="Output" />
                <File filename="experiment.cmd" instrumented="true" />
                <Field input="A" name="FASTA_FILE" output="C" type="string" />
        </HydraActivity>
        <HydraActivity activation="./experiment.cmd" description="Phylogeny with RAxML" tag="raxml" templatedir="/root/exp_SciPhy/</p>
                <Relation dependency="modelgenerator" filename="output modelgenerator.txt" name="A" reltype="Input" />
                <Relation filename="output raxml.txt" name="C" reltype="Output" />
                <File filename="experiment.cmd" instrumented="true" />
                <Field input="A" name="FASTA FILE" output="C" type="string" />
        </HydraActivity>
        </HydraWorkflow>
</Hydra>
```

Representação dos Workflows

```
<?xml version="1.0" standalone="no"?>
<SciCumulus>
  <database name="scicumulus_adaptive" server="mp4-4.dyndns.info"</pre>
    port="5432" username="scicumulus" password="********** />
  <SciCumulus Workflow tag="filogenia" description="This is a test using anflex."
    exectag="Experimento Kary/200 organismos - Adaptive" expdir="/root/exp">
    <SciCumulus Activity tag="mafft" description="mafft" type="MAP"
           templatedir="/root/exp/template_mafft" activation="experiment.cmd">
      <Relation reltype="Input" name="A" filename="parameter.txt" />
      <Field name="NAME" ype="string"/>
      <Field name="FASTA_FILE" type="string" />
      <Relation reltype="Output" name="C" file name="output_mafft.txt" />
      <Field name="NAME" type="str NAME;FASTA_FILE</pre>
                                     GI;ORTHOMCL2033
      <Field name="NUM_SEQ" type=
                                      G2:ORTHOMCL1895
      <Field name="FASTA_FILE" ty
                                      G3;ORTHOMCL2034
      <File filename="experiment.cmd"
                                     G4:ORTHOMCL1896
    </SciCumulusActivity>
                                      G5;ORTHOMCL2035
                                      G6;ORTHOMCL1897
  </SciCumulus Workflow>
                                      G7;ORTHOMCL2036
</SciCumulus>
                                      G8;ORTHOMCL1898
                                      G9;ORTHOMCL2037
                                      G10;ORTHOMCL1899
```

```
# Number of Processes
1
# Protocol switch limit
131072
# Entry in the form of machinename@port@rank
ec2-54-242-8-34.compute-1.amazonaws.com@22000@0
```

```
drwxr-xr-x 2 root root 4.0K Sep 23 13:10 1
drwxr-xr-x 2 root root 4.0K Sep 23 13:10 2
drwxr-xr-x 2 root root 4.0K Sep 23 13:10 3
drwxr-xr-x 2 root root 4.0K Sep 23 13:10 4
drwxr-xr-x 2 root root 4.0K Sep 23 13:10 5
```

```
-rw-r--r-- 1 root root 2.6K Sep 23 13:10 ENZ1.fasta
-rw-r--r-- 1 root root 1.4K Sep 23 13:10 ENZ1.fastaNumbered
-rw-r--r-- 1 root root 1.5K Sep 23 13:10 ENZ1.mafft
-rw-r--r-- 1 root root 1.3K Sep 23 13:10 H_Err.txt
-rw-r--r-- 1 root root 24 Sep 23 13:10 H_Relation.txt
-rw-r--r-- 1 root root 0 Sep 23 13:10 H_Result.txt
-rwxr-xr-x 1 root root 383 Sep 23 13:10 experiment.cmd
-rwxr-xr-x 1 root root 1.2K Sep 23 13:10 numberFasta.pl
```

Quer saber mais sobre banco de dados e nuvens?

- Não perca!
- "Análise em Big Data e um Estudo de Caso utilizando Ambientes de Computação em Nuvem"
 - Ticiana L. C. da Silva, Flávio R. C. Sousa,
 José Antônio F. de Macêdo e Javam C.
 Machado
 - Quinta-feira 03 de outubro 09:00 Sala
 Imperial

Agradecimentos







à Pesquisa do Estado do Rio de Janeiro











Obrigado!

Daniel de Oliveira - danielcmo@ic.uff.br Marta Mattoso - marta@cos.ufrj.br



