



Melissa Lemos

Workflow para Bioinformática

Tese de Doutorado

Tese apresentada ao Programa de Pós-graduação em Informática da PUC-Rio como requisito parcial para obtenção do título de Doutor em Informática.

Orientador: Prof. Marco Antonio Casanova

Co-Orientador: Prof. Antônio Basílio de Miranda

Rio de Janeiro, Setembro de 2004



Melissa Lemos

Workflow para Bioinformática

Tese apresentada como requisito parcial para obtenção do título de Doutor pelo Programa de Pós-Graduação em Informática da PUC-Rio. Aprovada pela Comissão Examinadora abaixo assinada.

Prof. Marco Antonio Casanova

(Orientador)

Departamento de Informática – PUC-Rio

Prof. Antônio Basílio de Miranda

(Co-Orientador)

Departamento de Bioquímica – Fiocruz

Profa. Marta L. Queirós Mattoso

Departamento de Informática – UFRJ

Profa. Ana Maria de Carvalho Moura

Departamento de Informática – IME-RJ

Prof. Paulo Mascarello Bisch

Departamento de Bioquímica – UFRJ

Prof. Luiz Fernando Bessa Seibel

Departamento de Informática – PUC-Rio

Prof. Marcus V. S. Poggi de Aragão

Departamento de Informática – PUC-Rio

Prof. Rubens Nascimento Melo

Departamento de Informática – PUC-Rio

Prof. José Eugênio Leal

Coordenador Setorial do Centro Técnico Científico – PUC-Rio

Rio de Janeiro, 13 de Setembro de 2004

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, da autora e do orientador.

Melissa Lemos

Graduou-se em Engenharia de Computação na PUC-Rio (Pontifícia Universidade Católica do Rio de Janeiro) em 1998. Obteve o título de Mestre em Informática na PUC-Rio em 2000, tendo a Dissertação de Mestrado ênfase em Banco de Dados e Bioinformática. Apresentou trabalhos sobre Bioinformática em congressos. É uma das pesquisadoras responsáveis pelo BioNotes, um sistema de anotação de sequências que está sendo utilizado no Estado do Rio de Janeiro para a pesquisa de genomas de organismos.

Ficha Catalográfica

Lemos, Melissa

Workflow para bioinformática /
Melissa Lemos ; orientador: Marco Antonio
Casanova ; co-orientador: Antônio Basílio de
Miranda. - Rio de Janeiro : PUC, Departamento
de Informática, 2004.

239 f. : il. ; 30 cm

Tese (doutorado) – Pontifícia
Universidade Católica do Rio de Janeiro,
Departamento de Informática.

Inclui referências bibliográficas.

CDD: 004

Ao meu marido, Marcio

Agradecimentos

A Deus, por iluminar-me e abençoar-me sempre.

Ao meu querido marido, Marcio Barroso Cavalière, a quem tanto amo, pelo seu amor, carinho, torcida, companheirismo durante este longo caminho, paciência e compreensão nas inúmeras ocasiões em que eu não pude dar-lhe a merecida atenção e por tornar meus dias mais felizes.

Ao meu pai, Mauro César de Lemos, por todo amor e carinho a mim dedicado e por me ensinar a importância do estudo e do conhecimento.

A minha mãe, Alaide das Graças Lemos, pelo amor, carinho, presença em todos os momentos da minha vida e apoio incondicional em todos os meus sonhos.

Ao meu irmão, Frederico Lemos, amigo que sempre está ao meu lado, torce por mim, acompanha e incentiva meus estudos.

Ao meu orientador, Marco Antonio Casanova, que tem minha eterna gratidão e admiração, pela excelente orientação, por sua dedicação e disponibilidade incansáveis, por suas idéias brilhantes, pela credibilidade depositada em mim, pela serenidade de nossas conversas que não me deixavam desanimar, por seus ensinamentos e pela parceria imprescindível para a realização desta tese.

Ao meu co-orientador, Antônio Basílio de Miranda, pela disponibilidade, atenção e paciência infinita em esclarecer tantas dúvidas em Biologia, contribuindo para a solidificação de uma base necessária para a pesquisa realizada nesta tese.

Ao Luiz Fernando Bessa Seibel, amigo que tem meu reconhecimento e gratidão, pela amizade e otimismo que me incentivaram e fizeram acreditar no valor do meu trabalho, pelas inúmeras sugestões, contribuições e conversas que me ajudaram a ter idéias, a criar e a melhorar a tese, e pelas oportunidades de trabalho, que me permitiram crescer profissionalmente e realizar esta tese.

Ao meu amigo José Antônio Fernandes de Macedo, a quem tenho como exemplo de amizade e sabedoria, agradeço pelo bom humor contagiante, pela disponibilidade, atenção e participação, e por sempre procurar me ajudar e enriquecer este trabalho com suas valiosas críticas e idéias.

Aos amigos Michelle Santos Sá, Flávio Freitas, Rogério Costa e Valéria Bastos, por sempre me passarem alegria, otimismo e pelas sugestões e críticas que tornaram o trabalho melhor.

Ao prof. Warren Gish, que sempre respondeu minhas perguntas por e-mails me esclarecendo dúvidas sobre o BLAST. Sem estes esclarecimentos, não teria sido possível propor várias idéias apresentadas nesta tese.

A todos os profissionais que estiveram envolvidos na especificação e no desenvolvimento do BioNotes, com quem eu tive a felicidade de conviver durante estes anos. Agradeço aos pesquisadores do Consórcio Riogene, em especial ao Paulo Ferreira Cavalcanti, Orlando Martins, Marcelo Bertalan, Shaila Rossle, Paulo Bisch e Ana Coelho. Agradeço aos profissionais em Informática que participaram, sempre com muito entusiasmo e empenho, no desenvolvimento do BioNotes. Meu reconhecimento aos colegas Fernando Mano, Vitor Cruz, Carlos Eduardo Vieira, Roberto Cavalcante, Dan Eisenberg e Andrea Cynthia dos Santos.

Aos professores Paulo Bisch, Marta Mattoso, Ana Maria Moura, Luiz Fernando Seibel, Marcus Poggi de Aragão e Rubens Melo pela participação na banca de avaliação.

A minha família e amigos pelo estímulo, apoio e orações.

Ao Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) pelo auxílio que me foi concedido.

Resumo

Lemos, Melissa; Casanova, Marco Antonio (Orientador) ; Miranda, Antonio Basílio de (Co-orientador). **Workflow para Bioinformática**. Rio de Janeiro, 2004, 239p. Tese de Doutorado – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Os projetos para estudo de genomas partem de uma fase de sequenciamento onde são gerados em laboratório dados brutos, ou seja, sequências de DNA sem significado biológico. As sequências de DNA possuem códigos responsáveis pela produção de proteínas e RNAs, enquanto que as proteínas participam de todos os fenômenos biológicos, como a replicação celular, produção de energia, defesa imunológica, contração muscular, atividade neurológica e reprodução. As sequências de DNA, RNA e proteínas são chamadas nesta tese de biossequências. Porém, o grande desafio destes projetos consiste em analisar essas biossequências, e obter informações biologicamente relevantes. Durante a fase de análise, os pesquisadores usam diversas ferramentas, programas de computador, e um grande volume de informações armazenadas em fontes de dados de Biologia Molecular. O crescente volume e a distribuição das fontes de dados e a implementação de novos processos em Bioinformática facilitaram enormemente a fase de análise, porém criaram uma demanda por ferramentas e sistemas semi-automáticos para lidar com tal volume e complexidade. Neste cenário, esta tese aborda o uso de workflows para compor processos de Bioinformática, facilitando a fase de análise. Inicialmente apresenta uma ontologia modelando processos e dados comumente utilizados em Bioinformática. Esta ontologia foi derivada de um estudo cuidadoso, resumido na tese, das principais tarefas feitas pelos pesquisadores em Bioinformática. Em seguida, a tese propõe um *framework* para um sistema de gerência de análises em biossequências, composto por dois sub-sistemas. O primeiro é um sistema de gerência de workflows de Bioinformática, que auxilia os pesquisadores na definição, validação, otimização e execução de workflows necessários para se realizar as análises. O segundo é um sistema de gerência de dados em Bioinformática, que trata do armazenamento e da manipulação dos

dados envolvidos nestas análises. O *framework* inclui um gerente de ontologias, armazenando ontologias para Bioinformática, nos moldes da apresentada anteriormente. Por fim, a tese descreve instâncias do *framework* para três tipos de ambiente de trabalho comumente encontrados e sugestivamente chamados de ambiente pessoal, ambiente de laboratório e ambiente de comunidade. Para cada um destes ambientes, a tese discute em detalhe os aspectos particulares da execução e otimização de workflows.

Palavras-chave

Bioinformática; Banco de Dados; Workflow; Ontologia; Framework de Software.

Abstract

Lemos, Melissa; Casanova, Marco Antonio (Advisor); Miranda, Antonio Basílio de (Co-Advisor). **Workflow for Bioinformatics**. Rio de Janeiro, 2004, 239p. PhD. Thesis – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Genome projects usually start with a sequencing phase, where experimental data, usually DNA sequences, is generated, without any biological interpretation. DNA sequences have codes which are responsible for the production of protein and RNA sequences, while protein sequences participate in all biological phenomena, such as cell replication, energy production, immunological defense, muscular contraction, neurological activity and reproduction. DNA, RNA and protein sequences are called biosequences in this thesis. The fundamental challenge researchers face lies exactly in analyzing these sequences to derive information that is biologically relevant. During the analysis phase, researchers use a variety of analysis programs and access large data sources holding Molecular Biology data. The growing number of Bioinformatics data sources and analysis programs indeed enormously facilitated the analysis phase. However, it creates a demand for systems that facilitate using such computational resources. Given this scenario, this thesis addresses the use of workflows to compose Bioinformatics analysis programs that access data sources, thereby facilitating the analysis phase. An ontology modeling the analysis program and data sources commonly used in Bioinformatics is first described. This ontology is derived from a careful study, also summarized in the thesis, of the computational resources researchers in Bioinformatics presently use. A framework for biosequence analysis management systems is next described. The system is divided into two major components. The first component is a Bioinformatics workflow management system that helps researchers define, validate, optimize and run workflows combining Bioinformatics analysis programs. The second component is a Bioinformatics data management system that helps researchers manage large volumes of Bioinformatics data. The framework includes an ontology manager that stores Bioinformatics ontologies, such as that previously described. Lastly,

instantiations for the Bioinformatics workflow management system framework are described. The instantiations cover three types of working environments commonly found and suggestively called personal environment, laboratory environment and community environment. For each of these instantiations, aspects related to workflow optimization and execution are carefully discussed.

Keywords

Bioinformatics; Databases; Software Frameworks; Workflow; Ontology.

Sumário

1	Introdução	19
1.1	. Objetivos e Contribuição da Tese	19
1.2	. Organização da Tese	21
2	. Preliminares	23
2.1	Introdução	23
2.2	Contexto Biológico	23
2.3	Programas de Análise	27
2.4	Bancos de Dados	29
2.5	Sistemas de Anotação	30
2.6	Integração	32
2.7	Ontologia	33
2.8	Workflow	34
2.9	Sistema de Gerência de Workflows de Bioinformática	39
2.10	Sistema de Gerência de Dados em Bioinformática	41
2.11	Sistemas de Gerência de Análises em Biossequências	42
2.12	Trabalhos Relacionados	42
2.12.1	myGrid e Proteus	42
2.12.2	Biopipe	43
2.12.3	LabFlow	44
2.12.4	Imagene	45
2.13	Comentários Finais	45
3	Uma Ontologia para Sistemas de Gerência de Análises em Biossequências	47
3.1	Introdução	47
3.2	Workflow	49
3.3	Processos	50
3.3.1	Classes de Processos	50
3.3.2	Propriedades de Processos	54

3.3.2.1 Parâmetros	54
3.3.2.2 Propriedades de Qualidade	55
3.4 Conexões	56
3.5 Contêineres	57
3.5.1 Classes de Contêineres	57
3.5.2 Propriedades de Contêineres	59
3.5.2.1 Parâmetros	59
3.5.2.2 Propriedades de Qualidade	61
3.6 Projetos	61
3.7 Comentários Finais	62
4 Uma Linguagem para Workflows de Sistemas de Gerência de Análises em Biossequências	63
4.1 Introdução	63
4.2 Especificação da Linguagem de Workflow em XML Schema	64
4.3 Comentários Finais	70
5 Um <i>Framework</i> para Sistemas de Gerência de Análises em Biossequências	71
5.1 Introdução	71
5.2 <i>Framework</i> para Sistemas de Gerência de Análises em Biossequências	71
5.3 Controlador	76
5.4 Assistente	77
5.5 Gerente de Ontologia	83
5.6 Gerente de Otimização	84
5.6.1 Validação	85
5.6.2 Otimização	86
5.6.3 Documento de Especificação do Workflow	86
5.7 Gerente de Execução	87
5.8 Gerente de Repositório	88
5.9 Comentários Finais	95

6	Instanciações do <i>Framework</i> para Sistemas de Gerência de Análises de Biossequências	96
6.1	Introdução	96
6.2	Arquiteturas de Hardware dos Ambientes de Trabalho dos Pesquisadores	97
6.2.1	Ambiente Pessoal	97
6.2.2	Ambiente de Laboratório	97
6.2.3	Ambiente de Comunidade	99
6.3	SGWBio no Ambiente Pessoal	99
6.3.1	Implementação de Contêineres	100
6.3.2	Exemplos de estimativas	109
6.3.2.1	Exemplos de estimativas para o tamanho dos contêineres	109
6.3.2.2	Exemplo de estimativa para taxa de consumo	112
6.3.3	Otimização Dinâmica por Pipelining	117
6.3.4	Exemplo de Otimização Dinâmica por Pipelining	135
6.4	SGWBio no Ambiente de Laboratório	140
6.4.1	Implementação de Contêineres	141
6.4.2	Otimização Dinâmica por Pipelining	142
6.4.3	Otimização Dinâmica por Paralelização	154
6.5	SGWBio no Ambiente de Comunidade	158
6.5.1	Implementação de Contêineres	160
6.5.2	Otimização Dinâmica por Pipelining	161
6.5.3	Otimização Dinâmica por Paralelização	165
6.6	Comentários Finais	166
7	Prototipação do SGWBio	171
7.1	Introdução	171
7.2	Assistente	171
7.3	Gerente de Otimização	177
7.4	Gerente de Execução	179
7.5	Gerente de Ontologia	180
7.6	Comentários Finais	181
8	Conclusão	182

8.1 Contribuição	182
8.2 Trabalhos Futuros	186
Referências	188
Anexo 1 – Definição da Ontologia	198
Ontologia em OWL	198
Ontologia em Prolog	221

Lista de Figuras

Figura 1. Esquema do workflow de projeto genoma completo.	36
Figura 2. Esquema de workflow de projeto genoma de ESTs.	38
Figura 3. Diagrama de classes da ontologia.	48
Figura 4. Representação de workflow por grafo bipartido.	50
Figura 5. Processo de controle externo de verificação.	54
Figura 6. Representação da linguagem do workflow de Bioinformática.	64
Figura 7. Exemplo de grafo bipartido de workflow.	69
Figura 8. Diagrama esquemático do SGABio.	72
Figura 9. Diagrama de Classes do SGABio.	72
Figura 10. Diagrama de classes: módulo compartilhamento de objetos. ...	76
Figura 11. Diagrama de classes: módulo controlador.	77
Figura 12. Diagrama de classes: módulo assistente.	78
Figura 13. Diagrama de classes: módulo gerente de ontologia.	83
Figura 14. Diagrama de classes: módulo gerente de otimização.	85
Figura 15. Diagrama de classes: módulo gerente de execução.	88
Figura 16. Modelo lógico simplificado do <i>data warehouse</i>	89
Figura 17. Arquitetura sem compartilhamento.	98
Figura 18. SGWBio em um ambiente de trabalho pessoal.	99
Figura 19. Exemplo 1 - BLASTs compartilhando contêineres.	115
Figura 20. Exemplo 2 - BLASTs compartilhando contêineres.	116
Figura 21. Exemplo: grafo bipartido.	135
Figura 22. Exemplo: grafo bipartido após inicialização.	136
Figura 23. Exemplo: <i>estágio_ideal</i> 1.	137
Figura 24. Exemplo: <i>estágio_real</i> 1.	138
Figura 25. Exemplo: <i>estágio_real</i> 2.	138
Figura 26. Exemplo: <i>estágio_real</i> 3.	139
Figura 27. Exemplo: grafo bipartido do workflow no cenário 1.	139
Figura 28. Exemplo: <i>estágio_ideal</i> 2.	140
Figura 29. SGWBio em um ambiente de trabalho laboratório.	141
Figura 30. SGWBio em um ambiente de trabalho de comunidade.	159

Figura 31. Interface do módulo assistente no protótipo.	172
Figura 32. Protótipo: Configuração de dados.	176
Figura 33. Protótipo: Configuração de parâmetros.	177

Lista de Tabelas

Tabela 1 – Conexão de leitura gradativa e não-gradativa.....	57
Tabela 2 – Conexão de escrita gradativa e não-gradativa.....	57
Tabela 3 – Tipos de contêineres.....	59
Tabela 4 – Propriedade gradativa para conexão de leitura.....	100
Tabela 5 – Propriedade gradativa para conexão de escrita.....	100
Tabela 6 – Tipos de contêineres.....	101
Tabela 7 – Especificação de BufferLimitado.....	102
Tabela 8 – Especificação de BufferIlimitado.....	103
Tabela 9 – Especificação de Arquivo.....	103
Tabela 10 – Especificação de ArquivoBufferLimitado.....	104
Tabela 11 – Especificação de ArquivoBufferIlimitado.....	105
Tabela 12 – Estimativas dos tamanhos dos contêineres.....	106
Tabela 13 – Escolha da forma de implementação de um contêiner.....	109
Tabela 14 – Exemplos para estimativas de tamanho de contêineres.	110
Tabela 15 – Tipos de BLAST.....	116

Lista de Quadros

Quadro 1. Sub-classes dos processos construtivos.....	51
Quadro 2. Sub-classes dos processos de alinhamento de sequências....	52
Quadro 3. Sub-classes de contêineres.....	58
Quadro 4. Instâncias de recursos.....	58
Quadro 5. Sequência de nucleotídeos no formato do Genbank.....	60
Quadro 6. Sequência de nucleotídeos no formato FASTA.....	60
Quadro 7. XML Schema do Workflow.....	66
Quadro 8. Workflow em XML-Schema: Elemento <i>Containers</i>	66
Quadro 9. Workflow em XML-Schema: Elemento <i>Processes</i>	67
Quadro 10. Workflow em XML-Schema: Elemento <i>Parameters</i>	68
Quadro 11. Workflow em XML-Schema: Elemento <i>Connections</i>	69
Quadro 12. Exemplo de documento de especificação de workflow.....	70