

## INFERÊNCIA DIFUSA COMO SUPORTE À DESCOBERTA DE POSSÍVEIS SNPs EM SEQUÊNCIAS DE CDNA

Wagner Arbox, Michel Eduardo Beteza Yamagishi, Marcos Vinícius G. Barbosa da Silva

Empresa Brasileira de Pesquisa Agropecuária, Rua Eugênio do Nascimento, 610, 36038-330, Juiz de Fora, MG

E-mail: arbox@cpnpj.embrapa.br, michel@cpnpj.embrapa.br, marcos@cpnpj.embrapa.br

Luiz Alfredo Vidal de Carvalho

Universidade Federal do Rio de Janeiro, Centro de Tecnologia - Bloco H-319, 21945-970, Rio de Janeiro, RJ

E-mail: alfredo@cos.ufrj.br

Diferenças pontuais entre pares de bases de diferentes sequências alinhadas, são o tipo mais comum de variabilidade genética. Tais diferenças, conhecidas como polimorfismos base única (*single nucleotide polymorphisms* - SNPs), são importantes no estudo da variabilidade das espécies, pois podem provocar alterações funcionais ou fenotípicas, que, por sua vez, podem implicar em consequências evolutivas ou bioquímicas nos indivíduos das espécies. A descoberta de SNPs por algoritmos computacionais é uma prática bastante difundida e, nessa área, dois *scripts* se destacam pelo amplo uso: Polyphred [1] e Polybayes [2].

O Polyphred, analisa diretamente os sinais expressos no sequenciamento do material genético e detecta SNPs a partir da variação dos sinais de fluorescência dos cromatogramas, procurando por reduções nas regiões do pico do sinal. Se for encontrada uma redução, onde uma segunda base foi detectada, então esse ponto é identificado como potencial heterozigoto. Após o alinhamento das sequências (*reads*), as bases dessa seção transversal, que inclui *reads* e consenso, são comparadas. O Polybayes analisa as bases geradas a partir da "leitura" dos cromatogramas - feita por *base-calling* [3], que nomeia e atribui um valor de qualidade para cada base (*Phred quality score* - PQS) - e utiliza um algoritmo de inferência Bayesiana, que procura por seções transversais onde os *reads* alinhados apresentam bases diferentes entre si. O Polybayes considera o número de *reads* e, ainda, a taxa a *priori* de pontos polimórficos, como sendo  $(\frac{1-0,003}{4})$ , ou seja, um SNP para cada 333 pares de bases, dividido pelo número de possíveis diferentes bases - A, T, C ou G - em um ponto. Deve ser notado que, esses dois *scripts*, têm seus resultados influenciados pelo PQS, obtido durante a leitura dos cromatogramas.

Os referidos *scripts* trabalham com diferentes metodologias, sobre diferentes atributos, contudo, espera-se que apresentem resultados similares, ao tratar um mesmo conjunto de seqüências, mas, não é incomum fornecerem resultados diferentes, o que produz incerteza na tomada de decisão, quando os resultados são discordantes.

O presente texto apresenta um modelo que se baseia em lógica difusa (*fuzzy logic*) para, a partir dos resultados do Polyphred e do Polybayes, auxiliar na tomada de decisão, no caso em se as informações sejam divergentes e, também, na confirmação de informações coincidentes. Ou seja, utiliza a lógica difusa para dar suporte à decisão, avaliando os resultados gerados por dois diferentes métodos e, ainda, incluindo, explicitamente, o PQS das bases do consenso, como um "valorizador" adicional, que reduz os efeitos específicos de cada um dos *scripts*.

A metodologia aqui apresentada não define nenhum limiar de "corte", no que se refere ao PQS, pois, o modelo de inferência difusa, automaticamente, elimina os pontos de baixa qualidade, não classificando-os como SNPs. Os critérios para a definição das variáveis linguísticas (conjuntos difusos), seus qualificadores e das funções de pertinência (expressões 2, 3 e 4), fundamentaram-se:

1. no índice atribuído pelo Polyphred (*Polyphred score* - PPS), que estabelece seis classes com intervalos *crisp*s, variando de 1, que indica um PPS  $\leq 49$  e uma taxa de verdadeiros positivos de 1%, sendo improvável a existência de SNPs. Até 6, que indica PPS  $\geq 99$  e uma taxa de verdadeiros positivos de 97%, sendo altamente provável a existência de SNPs, e, então, a variável linguística probabilidade foi definida nos termos: improvável (*P<sub>im</sub>*), pouco provável (*P<sub>pp</sub>*), medianamente provável (*P<sub>mp</sub>*), provável (*P<sub>v</sub>*), muito provável (*P<sub>mv</sub>*) e altamente provável (*P<sub>hp</sub>*);

2. na qualidade das bases (PQS), que varia entre 4 e 60, separadas, pelo limiar PQS = 20, em duas classes de valores *crisp*s e, então, a variável linguística qualidade foi definida nos termos: ruim (*Q<sub>R</sub>*), boa (*Q<sub>B</sub>*) e ótima (*Q<sub>O</sub>*).

Assim, no modelo de inferência aqui proposto, os valores discretos de entrada - os PPSs, encontrado pelo Polyphred e pelo seu equivalente no Polybayes, e o PQS - têm seus graus de pertinência estabelecidos pelas expressões 2, 3 e 4, que "disparam" regras difusas, cujo resultado é discretizado pelo método do "Centro do Máximo" (*Middle-of-Maxima* - MoM), visto que esse considera a ocorrência de múltiplos disparos de regras sobre uma mesma saída, "valorizando" essa saída. Desse modo, como resultado, determina-se um novo valor, mais apurado, indicativo da existência de polimorfismo, para cada SNP anteriormente identificado, onde foram considerados os valores iniciais dos PPSs e da PQS no ponto.

**Palavras-chave:** *Supporte à decisão com inferência difusa, Polimorfismo e Polimorfismo de base única.*

### Referências

[1] Nickerson, D. A., Toke, V. O. and Taylor, S. L., PolyPhred: automating the detection and genotyping of single nucleotide substitutions using fluorescence-based resequencing. *Nucleic Acids Research*, 25 (14): 2745-2751, 1997.

[2] Marth, G. T., Korn, L., Yandell, M. D., Teh, R. T., Gu, Z., Zakari, H., Stitzel, N. O., Hillier, L., Kwok, P. Y. and Gish, W. R., A general approach to single-nucleotide polymorphism discovery. *Nature Genetics*, 23 (4): 452-456, 1999.

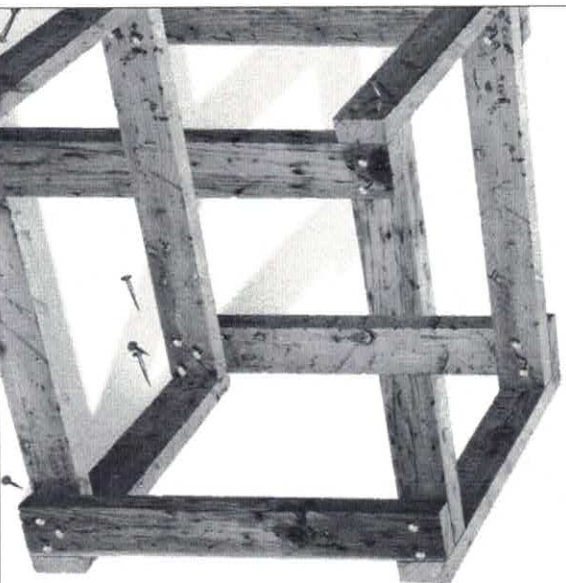
[3] Ewing, B., Hillier, L., Wendt, M. C. and Green, P., Basecalling of automated sequencer traces using Phred (I). *Genome Research*, 8 175-185, 1998.

**Sessão Técnica:  
Melhores Trabalhos  
Selecionados**

II ENCONTRO ACADÊMICO  
**MODELAGEM COMPUTACIONAL**

## Laboratório Nacional de Computação Científica

13 e 14 de Janeiro de 2009  
Petrópolis - RJ



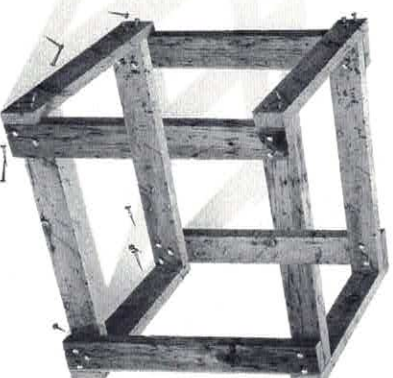
# RESUMOS

II Encontro Acadêmico de Modelagem Computacional  
do Laboratório Nacional de Computação Científica

IIEAMC-LNCC/MCT

Petrópolis, Rio de Janeiro, Brasil  
13 – 14 de Janeiro de 2008

# RESUMOS



### Comitê Organizador

Priscila V. Z. Capriles Gollatt  
Jonas Joacir Radtke  
Sicilia Ferreira Judice  
Raphael T. R. de Oliveira  
Raquel Lopes Costa  
Diego Augusto T. O. Leite

### Comitê Científico

Laurent Emmanuel Dardenne  
Abimael F. Dourado Loula  
Hello José Corrêa Barbosa  
Jauvane C. de Oliveira  
Rauli Antonino Feijóo  
Sandra Mara C. Malta



## Programação

### 13 de Janeiro de 2008 – Auditório A

HORÁRIO	EVENTO
09:00 - 10:00	<b>P1 – Palestra de Abertura</b> Eugenius Kashnericz (FINEP)
10:00 - 10:15	<b>CO1 – Modelagem Multiescala - Técnica de Homogeneização</b> Riedson Baptista (LNCC)
10:15 - 10:30	<b>CO2 – Dimensão Dinâmica e Fractalidade</b> Marcelo Miranda Barros (LNCC)
10:30 - 10:50	<b>ST-P – Apresentação de Pôsteres</b>
11:00 - 11:30	<b>P2 – Modelagem Computacional e Verificações Experimentais em Dinâmicas das Estruturas</b> Flavio de Souza Barbosa (UFEP)
11:30 - 11:45	<b>CO3 – A Celular Automata Framework for Real Time Fluid Animation</b> Sicilia Ferreira Ponce Pasini Judice (LNCC)
11:45 - 12:00	<b>CO4 – Evolução Diferencial para Otimização Restrita</b> Eduardo Krempser da Silva (LNCC)
12:00 - 12:30	<b>P3 – Jogos Educativos e Inovadores</b> Tulio Sorra (UNESP)
12:30 - 13:30	Almooço
13:30 - 14:00	<b>P4 – Criação e Uso de Bases de Dados Secundárias</b> José Miguel Ortega (UFMG)
14:00 - 14:15	<b>CO5 – The Importance of Different Temperatures and Electrostatic Treatments on Molecular Dynamics Simulations</b> Priscila Vanessa Zabelia Capriles Gollart (LNCC)
<b>P – Palestra</b> <b>CO – Comunicação Oral</b> <b>ST-P – Sessão Técnica-Poster</b>	

### 13 de Janeiro de 2008 – Auditório A

HORÁRIO	EVENTO
14:15 - 14:30	<b>CO6 – Por que Predadores não Têm Efeitos Positivos sobre Presas?</b> José Carlos Lisboa Recarey Eiras (LNCC)
14:30 - 14:45	<b>CO7 – A Multidisciplinaridade da Computação Quântica - Desenvolvimento e Perspectivas</b> Amanda Castro Oliveira (LNCC)
14:45 - 15:00	<b>CO8 – Assimilação de Dados de Limitado do Solo Aplicado ao Modelo Hidrológico NOAH</b> Claudia Adam Ramos (LNCC)
15:00 - 15:30	<b>P5 – Modelagem Integrada da Atmosfera-Biosfera-Hidrosfera</b> Marcos Heil Costa (UFV)
15:30 - 15:50	<b>ST-P – Apresentação de Pôsteres</b>
16:00 - 16:15	<b>ST-CO1 – Processo de Busca do Parâmetro de Regularização Ótimo em um Problema de Restauração de Imagens</b> Dalmo Stutz (UERJ/PPRJ)
16:15 - 16:30	<b>ST-CO2 – A Performance Evaluation of a Parallel Implementation of the Inexact Jacobian-Free Newton-Krylov Method Using the PETSc Numerical Framework</b> Rafael Santos Coelho (UFES)
16:30 - 16:45	<b>ST-CO3 – A Nonlinear Subgrid Discontinuous Galerkin Method for Transport Problems</b> Natalia Arruda (LNCC)
16:45 - 17:00	<b>ST-CO4 – Modeling Heat Treatment of Steels by the Calorimeter Method</b> Cecilia Grisel Arone (KB Engineering S.R.L.)
17:00 - 17:15	<b>ST-CO5 – Interferência Difusa como Suporte à Descoberta de Possíveis SNPs em Sequências de DNA</b> Wagner Arbx (EMBRAPA)
<b>P – Palestra</b> <b>CO – Comunicação Oral</b> <b>ST-P – Sessão Técnica-Poster</b> <b>ST-CO – Sessão Técnica-Comunicação Oral</b>	

## Sumário

MODELAGEM COMPUTACIONAL E VERIFICAÇÕES EXPERIMENTAIS EM DINÂMICAS DAS ESTRUTURAS	2
JOGOS EDUCATIVOS E INOVADORES	3
MODELAGEM INTEGRADA DA ATMOSFERA-BIOSFERA-HIDROSFERA	4
MODELAGEM MULTIESCALA - TÉCNICA DE HOMOGENEIZAÇÃO	5
DIMENSÃO DINÂMICA E FRACTALIDADE	6
A CELLULAR AUTOMATA FRAMEWORK FOR REAL TIME FLUID ANIMATION	7
EVOLUÇÃO DIFERENCIAL PARA OTIMIZAÇÃO RESTRITA	9
THE IMPORTANCE OF DIFFERENT TEMPERATURES AND ELECTROSTATIC TREATMENTS ON MOLECULAR DYNAMICS SIMULATIONS	10
COMPUTAÇÃO, INFORMAÇÃO QUÂNTICA E MULTIDISCIPLINARIDADE	11
ASSIMILAÇÃO DE DADOS DE UMIDADE DO SOLO NO NORDESTE BRASILEIRO	12
PROCESSO DE BUSCA DO PARÂMETRO DE REGULARIZAÇÃO ÓTIMO EM UM PROBLEMA DE RESTAURAÇÃO DE IMAGENS	14
A PERFORMANCE EVALUATION OF A PARALLEL IMPLEMENTATION OF THE INEXACT-JACOBIAN-FREE NEWTON-KRYLOV METHOD USING THE PETSC-NUMERICAL FRAMEWORK	15
A NONLINEAR SUBGRID DISCONTINUOUS GALERKIN METHOD FOR TRANSPORT PROBLEMS	17
MODELING HEAT TREATMENT OF STEELS BY THE CALORIMETER METHOD	18
INFERÊNCIA DIFUSA COMO SUPORTE À DESCOBERTA DE POSSÍVEIS SNPs EM SEQUÊNCIAS DE DNA	19
SEQUENCE ANALYSIS OF BOVINE PAPILLOMAVIRUS FOR THE ESTABLISHMENT OF DEGENERATE AND TYPE-SPECIFIC PRIMERS	22
THE ADAPTIVE MESH REFINEMENT AND COARSENING (AMRC) SCHEME IN THE LIBMESH	23
MOLECULAR DYNAMICS SIMULATIONS OF CALMODULIN: A COMPARATIVE STUDY OF REACTION FIELD AND PARTICLE-MESH Ewald ELECTROSTATIC TREATMENTS	24
HIBERNATE: UM FRAMEWORK DE Mapeamento Objeto Relacional	25
THE USE OF NEURAL NETWORKS IN CONTROL SYSTEMS FOR DETECTION AND FAILURES IN SENSORS	27
MOJLine: PORTAL PARA MODELAGEM COMPARATIVA EM GRANDE ESCALA USANDO WORKFLOW	28
DETERMINAÇÃO DO NÚMERO REPRODUTIVO BÁSICO EM MODELOS NÃO AUTÔNOMOS	29
NOVOS MÉTODOS DE ELEMENTOS FINITOS ENRIQUECIDOS E ESTABILIZADOS APLICADOS A EQUAÇÃO DE REAÇÃO-DIFUSÃO	31

CONTROLE DE VELOCIDADE DE UM MOTOR CC UTILIZANDO UM MICROCONTROLADOR PIC COM PROGRAMAÇÃO EM ASSEMBLY	33
IMPACTO DA QUEBRA DE ONDAS OCEÂNICAS NA ESTRUTURA DA CAMADA LIMITE PLANETÁRIA	34
ALGORITMOS DE BIOINFORMÁTICA PARA DETECÇÃO DE NULLOMERS NO TRANSCRIPTOMA HUMANO	35
BIFURCATIONS AND CHAOTIC DYNAMICS ANALYSIS OF THE ELASTIC PENDULUM	37
SOLUÇÃO ANALÍTICA DAS EQUAÇÕES DA CINÉTICA PONTOAL PARA VARIAÇÃO LINEAR DA REATIVIDADE DURANTE A PARTIDA DE UM REATOR NUCLEAR	39
MÉTODO ANALÍTICO DE IDENTIFICAÇÃO DE SISTEMAS MULTIVARIÁVEIS NO DOMÍNIO DA FREQUÊNCIA	41
MODELOS DE PROGRAMAÇÃO LINEAR INTEIRA PARA REDES	43
UMA ARQUITETURA DE MONITORAMENTO AMBIENTAL	45
CONSTRUÇÃO DE BIBLIOTECAS DE FRAGMENTOS PARA A PREDIÇÃO DE ESTRUTURAS DE PROTEÍNAS	47
IMPLEMENTAÇÃO DA ANÁLISE DE BIOLOGIA DE SISTEMAS DO FUNGO CAUSADOR DA YASSOURA DE BRUXA DO CACAUEIRO	49
APLICAÇÃO DE SISTEMAS BASEADOS EM REGRAS FUZZY PARA O ROTEAMENTO EM REDES ÓPTICAS	51

## LISTA DE AUTORES