



## Ferramentas de Bioinformática: Dos Cromatogramas ao Agrupamento

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

## Estrutura

- Cromatogramas: máquinas seqüenciadores;
- Phred: sintaxe, diretórios e arquivos;
- Phd2fasta: sintaxe e arquivos;
- Cross\_match: sintaxe e vector;
- Phrap: sintaxe e arquivos;
- Cap3: sintaxe e arquivos;
- Contig Viewer: visualização;
- Consed: visualização;
- Pipeline.

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes






## Seqüenciador

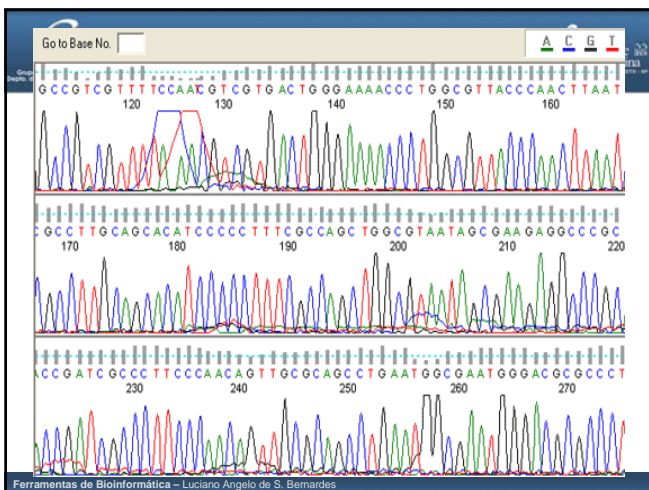
**1986 - Seqüenciador Automático de DNA (Leroy Hood)**



Lloyd M. Smith et al. Fluorescence detection in automated DNA sequence analysis. *Nature* 321, 674-679 (June 12, 1986).

T. Hunkapiller, R. J. Kaiser, B. F. Koop, L. Hood. Large-Scale and automated DNA sequence determination. *Science* 254, 59-67 (October 4, 1991).




Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes



## Phred

- Lê dados de sinais de fluorescência vindos de seqüenciadores automáticos de DNA
  - SCF (standard chromatogram format);
  - ABI (applied biosystems);
  - ESD (MegaBACE) e
  - LI-COR.
- Define as bases (base calling)
- Atribui valores de qualidade às bases
  - valor baseado na estimação do erro calculado para cada base individualmente.
- Gera arquivo de saída
  - bases definidas e valores de qualidade.

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

Phred

RESEARCH

# Base-Calling of Automated Sequencer Traces Using *Phred*. I. Accuracy Assessment

Brent Ewing,<sup>1</sup> LaDeana Hillier,<sup>2</sup> Michael C. Wendl,<sup>2</sup> and Phil Green<sup>1,3</sup>

<sup>1</sup>Department of Molecular Biotechnology, University of Washington, Seattle, Washington 98195-7730 USA;  
<sup>2</sup>Genome Sequencing Center, Washington University School of Medicine, Saint Louis, Missouri 63108 USA

8.175-185 ©1998 by Cold Spring Harbor Laboratory Press ISSN 1054-9803/98 \$5.00; www.genome.org GENOME RESEARCH 4:175

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

Phred

RESEARCH

# Base-Calling of Automated Sequencer Traces Using *Phred*. II. Error Probabilities

Brent Ewing and Phil Green<sup>1</sup>

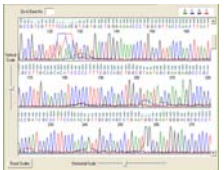
Department of Molecular Biotechnology, University of Washington, Seattle, Washington 98195-7730 USA

186 GENOME RESEARCH 8.186-194 ©1998 by Cold Spring Harbor Laboratory Press ISSN 1054-9803/98 \$5.00; www.genome.org

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

Phred

**Linha de comando:**  
 phred -id chromat\_dir -pd phd\_dir -sd seq\_dir -qd qual\_dir



↓

```
>sequence1_name
ccccagctcggatgcaggcacaanaacgctatggcaataggcactcttaccagggtg
aatcaggaacaatagtcataaaagtgaaggtccctgaacagatagocctctctttacag
aaacaggcagtgtaataggogaaggtgtgtaattgctaccaaagocctatccacagagc
agaagccatattcgttatgaatagcgttacccatcctgtacagaactggctcctaaaa
gtgccttcgctaccacaaaagtgaagtcgctgccggaacgagcctgggtattacagaaa
ggcgcagtgccctctactcttgcgagaaacttggaactacagacaagaacagaacta
tcatccttaagaagaccgcagtcgggaatacaaaccttcttttaactctc
```

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

.seq

```
>sequence001
TGGGNNAGNNNNNAGGGNNNGAGGAACTGAGCCCCTCGAGAGCCTTTTCG
GAGACACTATAGAACATGTTTTGTACAAAAAAGCAGGCTGGTACCGGTCGC
GAATTCGCCGGGATCGGAAACACTTAATATATATTTTCAGTCTAAATGCCAGC
GAACAGCGTGAGCTGGCCAATCGGATGGAAGGAGCAGATGAAGGAGTT
CATGACGGTATACTCGAACAATCGTCCAGCGCTGCTTGAAGACTGCGTCA
ATGATTTCAACAAGAAATCCCTCATCAGCCGCAACAAGGCTGCGTCAAC
CGCTGCTTCGACAAAGTTCATGAAGGCTCCGAGCGCATCAACCAGCGCTT
CCAGGAGCAGAAGCCTCAGATGATGAGTCCGGCCAACCTGGGAAAT
AATATCATTGAAAGAAAACGGCCATTTAGACTTGGGGTTTGAAGCAGTCA
TTTGAAAGGGGGTTCTGATGGGACGATTGTCATCGCCTTTTGGTTCCCTT
ATTGCTCGTCAGGGAACCTTTTTGTTGTCACAGGCCTTTCTGCAATGATT
GTGGATTTCAATTTCCATGAACAAGAATTTTCATTGGCAGCTAGTGGG
```

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

.qual

```
>sequence001
6 6 9 9 4 0 4 4 0 0 0 0 4 7 9 9 4 0 0 0 4 8 6 0 0
6 6 6 6 13 13 7 7 6 6 6 6 6 8 9 8 8 9 9 6 7 6 6 8 10
12 17 24 20 24 9 7 7 11 23 29 24 19 13 8 6 6 8 16 21
25 28 30 33 37 40 37 35 29 29 29 29 35 37 42 42
43 42 37 40 37 37 37 45 45 45 51 40 40 40 40 40
40 56 56 51 45 45 45 45 45 46 44 27 22 22 22 27 31
46 45 45 51 51 51 56 51 51 45 45 45 45 40 40 42
56 56 56 56 45 40 40 40 40 40 40 40 40 40 40 40
51 51 51 56 56 44 56 40 37 35 35 35 35 39 51 56
56 51 51 51 51 56 56 56 56 56 56 56 56 44 42 42
46 56 51 51 51 51 51 56 56 56 56 56 51 51 51 51
51 51 51 43 45 51 45 45 40 51 51 51 51 51 51 51
```

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes



.phd

```
BEGIN_SEQUENCE sequence001
BEGIN_COMMENT
CHROMAT_FILE: sequence001
ABI_THUMBPRINT:
067052220276000004012201162200
PHRED_VERSION: 0.020425.c
CALL_METHOD: phred
QUALITY_LEVELS: 99
TIME: Wed Jun 27 16:22:27 2007
TRACE_ARRAY_MIN_INDEX: 0
TRACE_ARRAY_MAX_INDEX: 9281
TRIM: 69 613 0.0500
TRACE_PEAK_AREA_RATIO: 0.0090
CHEM: term
DYE: big
END_COMMENT
BEGIN_DNA
t 6 9
g 6 11
g 9 13
g 9 15
a 45 1424
a 45 1437
t 51 1447
a 51 1459
t 51 1471
a 56 1483
t 51 1495
a 51 1506
t 45 1518
t 45 1532
c 20 9021
c 25 9035
a 27 9047
t 15 9059
t 11 9072
c 15 9086
a 13 9094
c 9 9124
a 9 9131
t 9 9149
c 10 9165
a 10 9171
t 10 9185
g 12 9201
a 7 9215
g 7 9222
g 7 9239
a 9 9253
t 9 9266
g 7 9280
END_DNA
END_SEQUENCE
```

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes







 **Phrap** 

Grupo de Bioinformática  
Depart. de Genética - FINEP/UFPA

**Linha de comando:**  
`phrap seqs_fasta.screen -new_ace > phrap.out`

- Programa de “montagem” das seqüências
  - [arquivo].ace
  - [arquivo].contigs
  - [arquivo].contigs.qual
  - [arquivo].log
  - [arquivo].problems
  - [arquivo].problems.qual
  - [arquivo].singlets
  - Phrap.out



Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

 **CAP3** 

Grupo de Bioinformática  
Depart. de Genética - FINEP/UFPA

- Usa partes de alta qualidade das “reads”, para conter o erro no final das seqüências consenso;
- Reconhece informações do usuário;
- Fornece extensas informações;
- Restrições de “forward-reverse”;
- Produção de ordem parcial de contigs por restrição, tornando mais fácil e rápido sua construção;
- Menos erros internos no consenso;
- Contigs mais curtos, com menos erros

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

 **CAP3** 

Grupo de Bioinformática  
Depart. de Genética - FINEP/UFPA



**CAP3: A DNA Sequence Assembly Program**

Xiaojin Huang<sup>1,2</sup> and Anup Madan<sup>3</sup>

<sup>1</sup>Department of Computer Science, Michigan Technological University, Houghton, Michigan 49931 USA; <sup>2</sup>Department of Molecular Biotechnology, University of Washington, School of Medicine, Seattle, Washington 98195 USA

We describe the third generation of the CAP sequence assembly program. The CAP3 program includes a number of improvements and new features. The program has a capability to clip 5' and 3' low-quality regions of reads. It uses base quality values in computation of overlaps between reads, construction of multiple sequence alignments of reads, and generation of consensus sequences. The program also uses forward-reverse constraints to correct assembly errors and link contigs. Results of CAP3 on four BAC data sets are presented. The performance of CAP3 was compared with that of PHRAP on a number of BAC data sets. PHRAP often produces longer contigs than CAP3 whereas CAP3 often produces fewer errors in consensus sequences than PHRAP. It is easier to construct scaffolds with CAP3 than with PHRAP on low-pass data with forward-reverse constraints.

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes



 **CAP3** 

Grupo de Bioinformática  
Depart. de Genética - FINEP/UFPA

**Linha de comando:**  
`cap3 seqs_fasta.screen > cap3.out`

- Programa de “montagem” das seqüências
  - [arquivo].cap.ace
  - [arquivo].cap.contigs
  - [arquivo].cap.contigs.qual
  - [arquivo].cap.info
  - [arquivo].cap.singlets
  - Cap3.out



Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

 **Visualizadores** 

Grupo de Bioinformática  
Depart. de Genética - FINEP/UFPA

- Permitem visão geral da distribuição das seqüências;
- Permitem algumas interações com as visualizações;
- Permitem comparações e buscas;
- Permitem certas modificações;

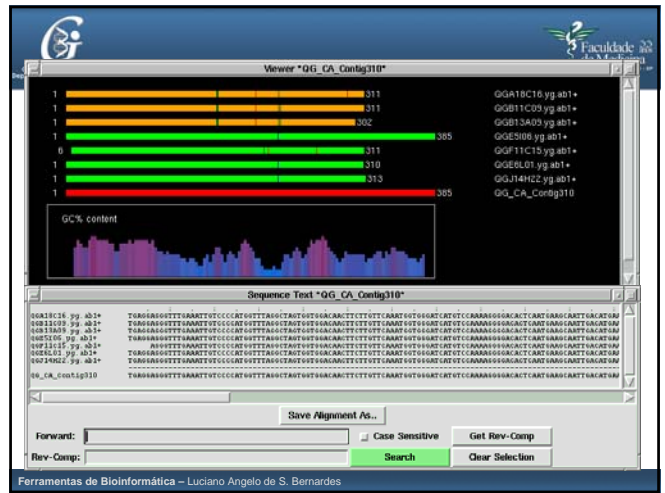
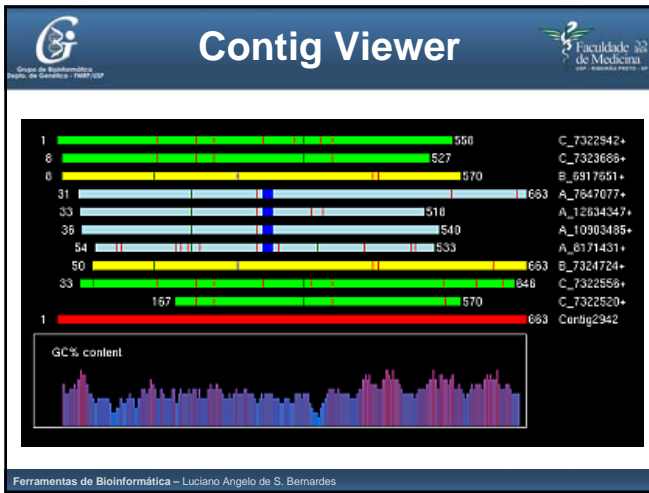
Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

 **Contig Viewer** 

Grupo de Bioinformática  
Depart. de Genética - FINEP/UFPA

- Desenvolvido em linguagem Python;
- Sistemas Operacionais:
  - Windows
  - MAC OS X
  - Linux/UNIX
- Cap3

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes



# CONSED

- Visualização, edição e finalização das montagens;
- Facilita a checagem da exatidão das montagens;
- Permite visualização dos valores de qualidade e cromatogramas ;
- Estrutura de diretórios: chromat\_dir, edit\_dir e phd\_dir;
- Phrap;

Ferramentas de Bioinformática - Luciano Angelo de S. Bernardes

# Consed

RESEARCH

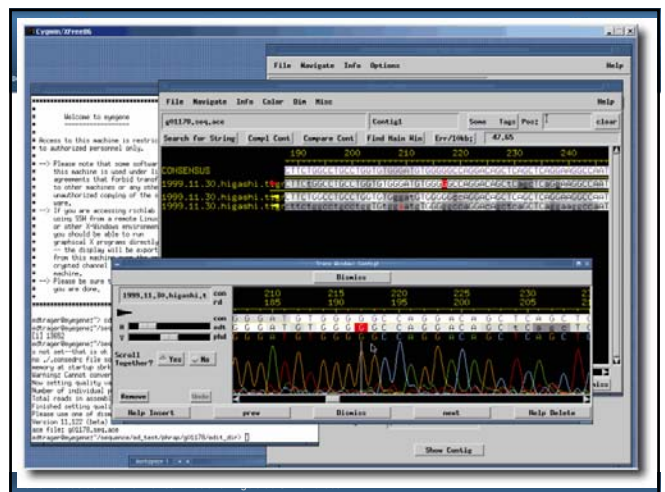
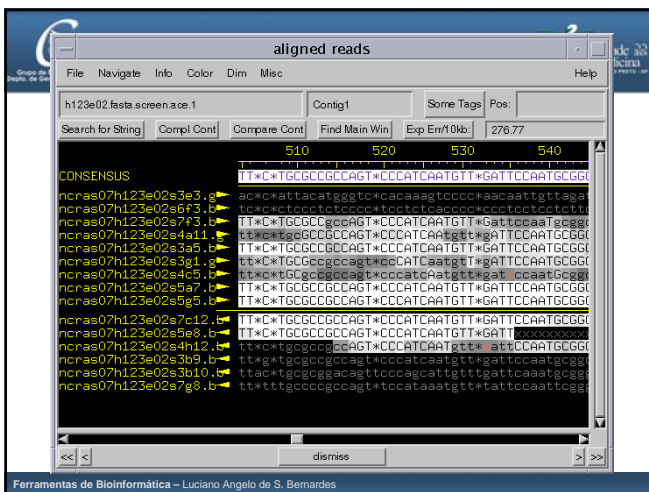
## Consed: A Graphical Tool for Sequence Finishing

David Gordon,<sup>2</sup> Chris Abajian,<sup>1</sup> and Phil Green<sup>2</sup>

Department of Molecular Biotechnology, University of Washington, Seattle, Washington 98195-7730 USA


8:195-202 ©1998 by Cold Spring Harbor Laboratory Press ISSN 1054-9803/98 \$5.00. www.genome.org GENOME RESEARCH # 195

Ferramentas de Bioinformática - Luciano Angelo de S. Bernardes

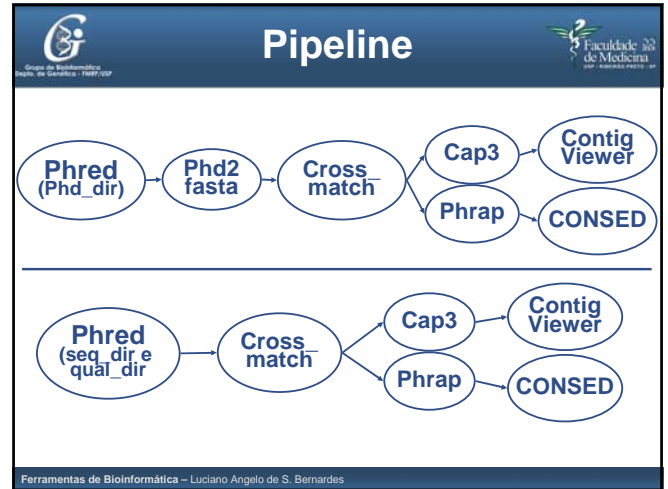


**Pipeline**

- Um conjunto de processos encadeados através das suas saídas padrão, de forma que a cada uma dessas é utilizada como entrada do processo seguinte.



Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes



**PhredPhrap**

- Script em linguagem perl;
- PhredPhrap roda:
  - phred;
  - phd2fasta;
  - cross\_match;
  - phrap;
- Execução:
  - **phredphrap**

Ferramentas de Bioinformática – Luciano Angelo de S. Bernardes

