

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

16 de novembro de 2006

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

1 Introdução

2 Expansão do sistema EGene/CoEd

3 Geração de evidências para fins de anotação

Processamento em bioinformática

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Na área de bioinformática, realiza-se freqüentemente o processamento de seqüências de bases nitrogenadas (Adenina, Citosina, Guanina, Timina e Uracil) ou dos aminoácidos formados por essas bases.

Processamento em bioinformática

"Anotação de seqüências e expansão do sistema EGene/CoEd"

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Exemplo de dado de entrada (formato FASTA)

```
>Ac-A01-1705
CAGGCCCTTCCCAACTAAGGTAAACCCTAAGCCCTAAACCCTAAACCCTA
AACCCCTAAACCCTAAACCCTCAAACGTGTCGATGGATGGGAGTGAACCTGG
CGCACTAATGAGTAGAAAATTTGCAGACATTCTCCGCTCCCTGCGACGTA
AACGGCAGCATAATGGCGGACAGCTCCAAAATAGAAGAAGGGCTTCAACC
```

Processamento em bioinformática

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

O processamento envolve várias tarefas computacionais interconectadas, cada um com um protocolo de entrada e saída de dados distinto. Usualmente, é criado um script, chamado *pipeline* que executa as tarefas em ordem determinada pelo usuário – cada tarefa é denominada como componente.

Processamento em bioinformática

"Anotação de seqüências e expansão do sistema EGene/CoEd"

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Com o processamento das seqüências é possível obter algumas informações, chamadas evidências, a partir das quais um biólogo pode gerar anotações, ou seja, vincular partes de uma seqüência a funções reguladoras, componentes celulares, entre outros.

Processamento em bioinformática

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Entretanto, a criação de cada script torna-se onerosa no sentido de que é necessário considerar que cada par de programas que troca dados necessita de um processamento adicional para que a saída de um seja compatível com a entrada do outro.

EGene/CoEd

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Nesse cenário, existe o EGene, um sistema de geração de pipelines que torna mais fácil a implementação dos mesmos. Existe ainda o CoEd, uma ferramenta gráfica que facilita a criação dos arquivos de configuração utilizados pelo EGene.

Tela do CoEd

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

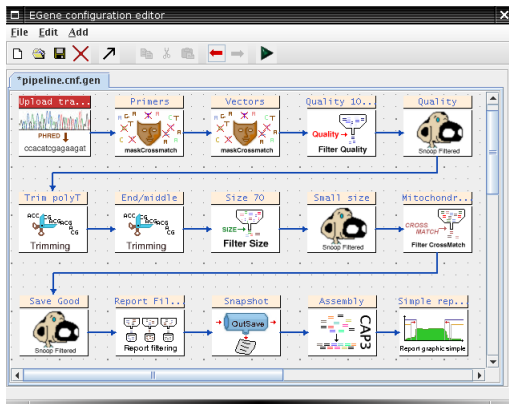


Figura: *Um pipeline gerado pelo CoEd.*

Objetivos do trabalho

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

O trabalho realizado durante a iniciação científica envolveu duas partes: a expansão do sistema EGene/Coed e geração de evidências.

O sistema original

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

A versão original do EGene permite uma única arquitetura de interligação de componentes, denominada *pipeline*. Num *pipeline*, os dados seguem um único fluxo, sendo processados seqüencialmente por cada componente.

O sistema original

"Anotação de seqüências e expansão do sistema EGene/CoEd"

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

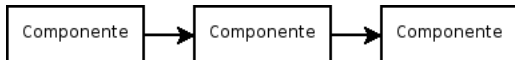


Figura: *Arquitetura disponível no EGene 1.0.*

Alterações do sistema

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Durante a iniciação científica, foram criadas duas novas estruturas para o processamento: *forks* e seletores.

Forks

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Num *fork*, um dado pode ser enviado para vários componentes diferentes, gerando novas possibilidades de fluxo.

Forks

"Anotação de seqüências e expansão do sistema EGene/CoEd"

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

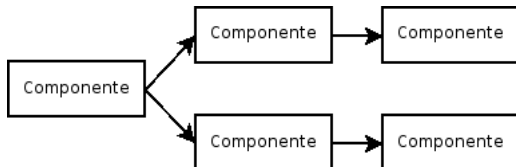


Figura: *Exemplo de arquitetura utilizando forks.*

Seletores

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Com seletores, é possível encapsular um componente e fazer com que uma seqüência só seja processada pelo mesmo se uma dada condição for satisfeita.

Forks e Seletores

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Com essas alterações, houve um ganho de desempenho de processamento. Com os seletores, é possível ignorar a execução de um ou mais componentes sobre uma dada seqüência. Com os *forks* temos ainda mais paralelismo do que o oferecido pelos *pipes* encontrados nos sistemas UNIX.

Forks e Seletores

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Com forks e seletores, o EGene é capaz de implementar o processamento de programas em praticamente toda arquitetura possível em grafos direcionados acíclicos.

Arquitetura não implementável

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Uma única estrutura não é implementável com as alterações: dois ou mais componentes enviando dados para um outro. Entretanto, essa limitação é aceitável no contexto de processamento de seqüências, dado que ela não é utilizada usualmente.

Arquitetura não implementável

"Anotação de seqüências e expansão do sistema EGene/CoEd"

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

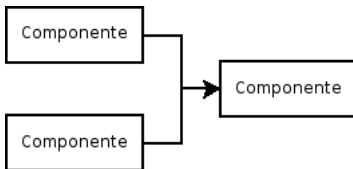


Figura: Exemplo de arquitetura não implementável no novo EGene/CoEd.

Geração de evidências

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Os componentes do EGene original permitem basicamente realizar o pré-processamento das seqüências, eliminando contaminantes e sub-seqüências de baixa qualidade.

Geração de evidências

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Em nosso grupo de bioinformática, foram desenvolvidos componentes para geração de evidências e componentes de anotação automática no formato *feature table* de submissão de seqüências anotadas.

Geração de evidências

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo Yamamoto Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

Tomando por base uma modelagem das bases de dados do sistema original, foi feita uma ampliação da mesma de modo a permitir a inclusão dos dados de evidências. Foram identificados 4 tipos de evidências: similaridade, multi-intervalo, estatística e gráfico.

“Anotação de seqüências e expansão do sistema EGene/CoEd”

Ricardo
Yamamoto
Abe

Outline

Introdução

Expansão do sistema EGene/CoEd

Geração de evidências para fins de anotação

FIM