

UFRN, Instituto Metr pole Digital
Bacharelado em Tecnologia da Informa  o
DIM501 – Boas Pr ticas de Programa  o, 2016.1
Prof. Handerson Bezerra Medeiros

Relat rio:

**Aplica  o de Boas Pr ticas de Programa  o para o aprimoramento do
projeto *Sequencer***

Aluno:

Pit goras de Azevedo Alves Sobrinho

Natal-RN, 22 de maio de 2016.

Sumário:

<u>Lista de figuras</u>	3
<u>1. Introdução</u>	4
<u>2. Remoção da paralelização</u>	5
<u>3. Aplicação da STL do C++11 para simplificação dos algoritmos</u>	6
<u>4. Modularização e revisão para nomes mais significativos</u>	9
<u>5. Programação defensiva</u>	11
<u>6. Testes</u>	13
<u>7. Documentação</u>	14
<u>8. Métricas</u>	15
<u>9. Referências</u>	16

Lista de figuras:

Figura 1: Código inicial de entrada de dados na versão não modificada de <i>Sequencer</i>	6
Figura 2: Modificações na entrada/saída	7
Figura 3: Método de fundir os segmentos de texto aprimorado	8
Figura 4: Diagrama de funções/classes da versão ainda não aprimorada de <i>Sequencer</i>	9
Figura 5: Diagrama da versão aprimorada de <i>Sequencer</i> , sem os módulos de testes.	10
Figura 6: asserts sendo aplicados na fusão de um par de strings	11
Figura 7: try-catch no construtor do “SegmentSet”	11
Figura 8: try-catch sendo utilizados para tentar abrir arquivos de forma segura	12
Figura 09: Aplicação de testes para a classe “StringMerge”	13
Figura 10: Aplicação de testes para a classe “SegmentSet”	13

1. Introdução

Shotgun Sequencing^[1] é um método de sequenciamento genético em que são amostradas múltiplas subsequências de uma sequência de DNA ou RNA e sequência completa final é inferida a partir deste conjunto de subsequências.

Sequencer^[2] é um software para realizar o sequenciamento de um texto a partir de um conjunto de pequenas passagens (segmentos do texto). O conjunto de passagens cobre toda a extensão do texto final sendo que sempre existirá uma sobreposição entre passagens subsequentes. O algoritmo foi desenvolvido para ter bom desempenho em máquinas com muitos núcleos, através do uso de paralelismo. Isso foi feito utilizando a API de paralelismo genérica OpenMP^[3].

O objetivo era cumprir a tarefa do Bioinformatics Hackaton^[4], por isso o algoritmo foi desenvolvido de forma rápida sem haver qualquer preocupação com boas práticas, a única prioridade era o desempenho e a boa aplicação das técnicas de paralelização com OpenMP.

2. Remoção da paralelização

Os algoritmos do projeto *Sequencer* utilizavam OpenMP, uma API para programação paralela multiplataforma.

A utilização desta API se baseia na aplicação de clausulas “#pragma” para definir formas de paralelização para determinadas regiões do código. Tais clausulas podem ser de difícil interpretação para um leitor do código que não conheça bem a API em questão.

Sequencer foi desenvolvido originalmente para o supercomputador do NPAD (Núcleo de Processamento de Alto Desempenho)^[5], onde ele teria que processar uma grande quantidade de dados. Porém foi verificado durante o desenvolvimento que, numa maquina com poucos núcleos e com a entrada sendo constituída por poucos fragmentos de texto, o desempenho não é melhorado pelo OpenMP. Na verdade, muitas vezes, ele é piorado. Isso ocorre porque o *Sequencer* acaba passando muito tempo criando *threads* desnecessariamente.

A especificação do OpenMP só é adequada oficialmente ao padrão ISO/IEC 14882:1998^[6], ou seja, até o C++98 e também não é recomendado o uso de templates. Assim, é inseguro e uma péssima pratica utilizar capacidades implementadas na *Standard Template Library* após 1998. Em decorrência dessa e das outras desvantagens citadas anteriormente nesta sessão, foi decidido que (no contexto deste trabalho e da disciplina “Boas Praticas de Programação”) é melhor abandonar o uso do OpenMP e fazer uso do (não tão novo) padrão C++11^[9].

O OpenMP foi removido facilmente retirando do código as suas clausulas “#pragma”.

3. Aplicação da STL do C++11 para simplificação dos algoritmos

A STL do C++^[7] tem diversas classes e rotinas que podem melhorar a legibilidade do algoritmo, tornando tudo mais simples. Nesta sessão será descrito como foi feita a aplicação deste padrão da linguagem no projeto *Sequencer*.

3.1 Modificação do processo de I/O

A entrada de dados vem diretamente da entrada padrão dos sistema e a saída é feita diretamente na saída padrão.

Para possibilitar a inserção de uma quantidade indefinida de segmentos de texto, os segmentos são armazenados inicialmente numa estrutura de dados “std::vector”, um vetor genérico no qual espaço é alocado dinamicamente conforme novos elementos são inseridos. Porém, para evitar incompatibilidades com o OpenMP, antes de entrar no processo de paralelização, estas estruturas de dados são convertidas para vetores comuns. Esse processo confuso pode ser visto no seguinte fragmento de código:

```
int segmentsPerBucket = 9;
int totalSegments = 0;
string line;
vector<string> tempSegments;
vector<vector<string> > tempBuckets;
cout << "vars created\n";
while(getline(cin, line)){
    //Caso o texto use quebras de linha do tipo "\n\r", retira as \r
    if(line[line.size()-1] == '\r'){
        line = line.substr(0, line.size()-1);

        //cout << "quebra de linha barra r retirada  \n";
    }
    tempSegments.push_back(line);
    totalSegments++;
    if(tempSegments.size() == segmentsPerBucket){
        tempBuckets.push_back(tempSegments);
        tempSegments.clear();
    }
}
if(!tempSegments.empty()){
    tempBuckets.push_back(tempSegments);
    tempSegments.clear();
}
cout << "temp buckets instantiated\n";
int nBuckets = tempBuckets.size();
Bucket* buckets = new Bucket[nBuckets];

for(int i = 0; i < nBuckets; i++){
    buckets[i].nSegments = tempBuckets[i].size();
    buckets[i].segments = new string[buckets[i].nSegments];
    for(int j = 0; j < buckets[i].nSegments; j++){
        buckets[i].segments[j] = tempBuckets[i][j];
    }
    tempBuckets[i].clear();
}
tempBuckets.clear();
cout << "buckets instantiated\n";
```

Figura 1: Código inicial de entrada de dados na versão não modificada de *Sequencer*

Um “Bucket” (balde, nome de classe nada significativo) é uma classe que, basicamente, armazena um vetor de segmentos de texto e o número destes segmentos. Seu único método é a rotina “process()”, que “funde” todos os segmentos presentes numa instancia do “Bucket” em um único texto final. Ele não possui qualquer tipo de modularização e ou construtor.

Como, sem o OpenMP, não há mais motivos para não fazer uso total das estruturas de dados da STL, o código foi simplificado para fazer amplo uso da estrutura “std::vector”.

3.1.1. Passagem de arquivo de segmentos de texto por argumentos para o programa:

A entrada de dados é feita, dentro da lógica do programa, através da entrada padrão do sistema, ao invés de ler os dados vindo de um arquivo. Porém na prática é diferente. Os operadores de direcionamento de entrada e saída “>” e “<” do UNIX são utilizados para ler os segmentos de texto vindo de um arquivo e escrever toda a saída em outro arquivo. Por isso o script para rodar o programa tem essa aparência:

```
$ ./bin/main-normal < inputs/A.input > outputs/A_normal.output
```

Esse comando, basicamente, executa o *Sequencer* de forma que os segmentos de texto sejam lidos do arquivo “inputs/A.input” e o resultado do sequenciamento seja escrito em “outputs/A_normal.output”. A entrada/saída foi feita dessa forma porque esse é o padrão para competições de programação, porém caso um usuário precise utilizar o programa, ele teria que seguir exatamente essa sintaxe, própria de sistemas UNIX. É desejável que o próprio *Sequencer* possa acessar os arquivos e ler/escrever neles, através dos seguintes argumentos:

```
$ ./sequencer -i “arquivo_de_entrada” -o “arquivo_de_saída”
```

Foi utilizada uma biblioteca de interpretação de argumentos de linha de comando, a AnyOption^[8]. A biblioteca não era compatível com C++11, então seu código foi “modernizado” para corrigir isso. Sem a necessidade de paralelização OpenMP, utilizando o padrão c++11 e com a saída e entrada do programa sendo definida pelos argumentos passados pelo usuário, o código inicial do programa ficou assim:

```
23
24 int main(int argCount, char* argVector[]){
25     Arguments arguments(argCount, argVector);
26     if(!arguments.helpFlag){
27         istream* inputStream;
28         ostream* outputStream;
29
30         //seleciona a entrada de dados
31         if(arguments.inputFileDefined){
32             ifstream *inputFileStream = new ifstream;
33             inputStream->open(arguments.inputFile.c_str());
34             inputStream = inputStream;
35         }else{
36             inputStream = &cin;
37         }
38         //seleciona a saída de dados
39         if(arguments.outputFileDefined){
40             ofstream* outputFileStream = new ofstream;
41             outputFileStream->open(arguments.outputFile.c_str(), ios::trunc);
42             outputStream = outputFileStream;
43         }else{
44             outputStream = &cout;
45         }
46         //loop de leitura
47         Bucket bucket;
48         string line;
49         while(getline(*inputStream, line)){
50             //Caso o texto use quebras de linha do tipo "\n\r", retira as \r"
51             if(line[line.size()-1] == '\r'){
52                 line = line.substr(0, line.size()-1);
53                 //cout << "quebra de linha barra r retirada \n";
54             }
55             bucket.segments.push_back(line);
56         }
57     }
```

Figura 2: Modificações na entrada/saída

3.2. Trocar matriz de semelhança com cópias de vetores por armazenamento em *hashmap*

Todo o processo de transformar os diversos segmentos de texto inseridos num “Bucket” se baseia em combinar o par de segmentos que se combinam melhor e então repetir o processo, até que não haja mais o que combinar ou só reste um único segmento. Isso é feito pelo método `Bucket::process()`. É um enorme método com 155 linhas de código e nenhuma modularização. Esse método faz uso de uma matriz de semelhanças para decidir qual a melhor opção de combinação entre segmentos de texto. Isso tem um grande custo de memória e o cálculo da matriz, normalmente, teria complexidade quadrática ($O(n^2)$) e seria feito a cada iteração do algoritmo, resultando numa complexidade cúbica ($O(n^3)$). Porém, inspirado na programação dinâmica, foi adicionado a este método uma técnica que, através do armazenamento da versão da matriz na iteração anterior, consegue fazer com que só os elementos da matriz de semelhança que seriam diferentes sejam calculados. Isso diminuiu a complexidade do método de $O(n^3)$ para $O(n^2)$. Essa técnica apenas torna o método mais complexo.

Agora, usando o padrão *C++11* e tendo toda a liberdade para usar *templates*, foi possível substituir o `Bucket::process()` antigo por uma nova versão, que armazena as semelhanças em uma estrutura de *hashmap*, o `std::map`. A complexidade *big O* ainda é a mesma (quadrática), porém o comprimento do método foi reduzido de 155 linhas para apenas 35 e agora o algoritmo é muito mais fácil de ser entendido por um leitor:

```
14 void Bucket::process(bool forceMerge){
15     map<StringPair, CompareResult> hashmap;
16     bool merged = false;
17     int iteration = 0;
18     do{
19         StringPair* keyMax = NULL;
20
21         for(int i = 1; i < segments.size(); i++){
22             for(int j = 0; j < i; j++){
23                 StringPair key = StringPair(segments[i], segments[j]);
24                 if(hashmap.find(StringPair(segments[i], segments[j])) == hashmap.end()){
25                     hashmap[key] = CompareResult(segments[i], segments[j], true);
26                 }
27                 if(hashmap[key].haveResult()){
28                     if(keyMax == NULL){
29                         keyMax = new StringPair();
30                         *keyMax = key;
31                     }else{
32                         if(hashmap[key].module > hashmap[*keyMax].module){
33                             *keyMax = key;
34                         }
35                     }
36                 }
37             }
38         }
39
40         if(keyMax != NULL){
41             segments.erase(find(segments.begin(), segments.end(), keyMax->x));
42             segments.erase(find(segments.begin(), segments.end(), keyMax->y));
43             segments.push_back(hashmap[*keyMax].result);
44             segments = segments;
45             merged = true;
46         }
47         iteration++;
48     }while(segments.size() > 1 && merged);
49 }
```

Figura 3: Método de fundir os segmentos de texto aprimorado

4. Modularização e revisão para nomes mais significativos

O diagrama a seguir mostra os componentes do software *Sequencer* e como eles se relacionavam antes das modificações feitas neste trabalho:

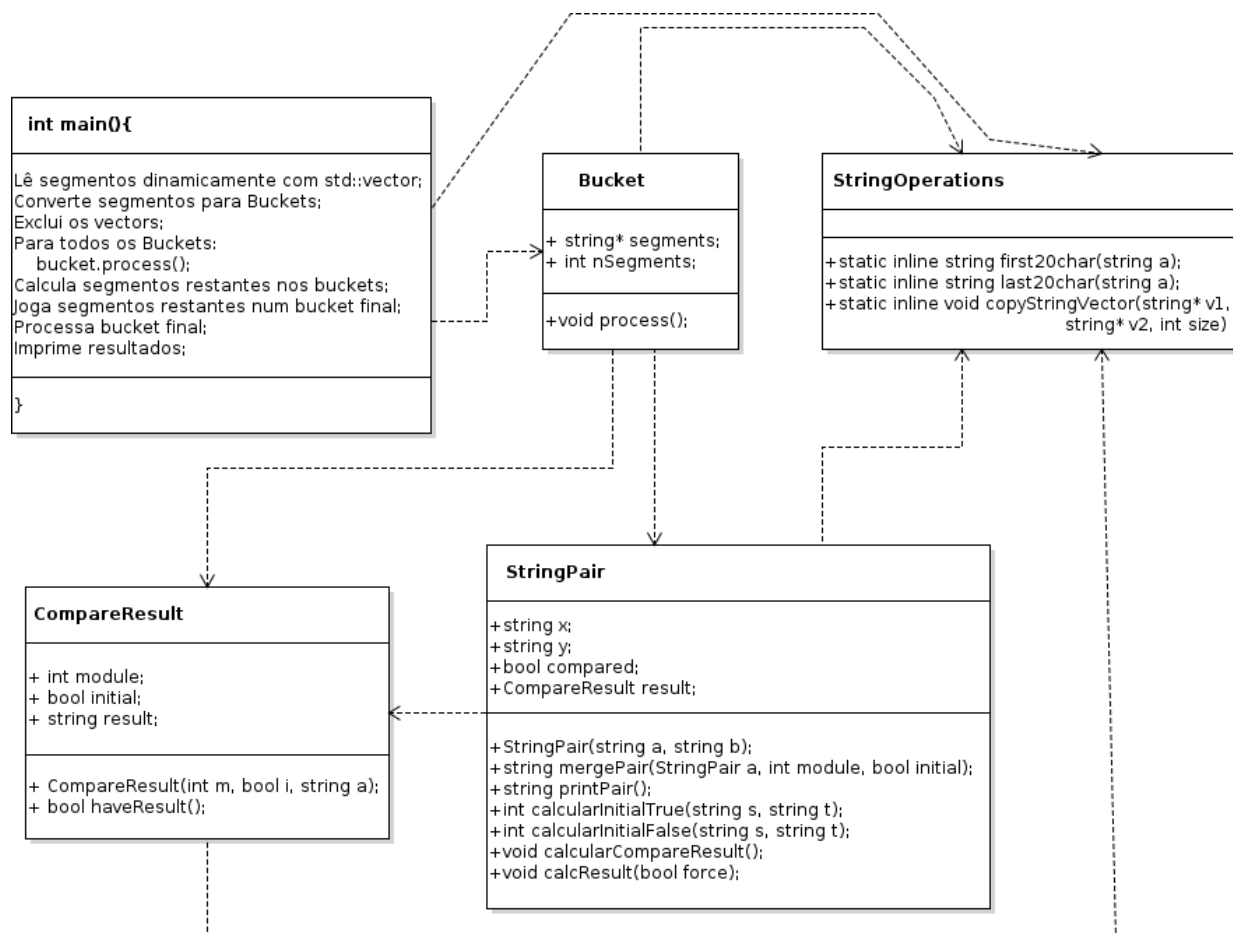


Figura 4: Diagrama de funções/classes da versão ainda não aprimorada de *Sequencer* (sim, não é uma boa prática de UML descrever algoritmos onde deveriam estar descritas classes, mas cheguei a conclusão de que neste caso seria uma melhor forma de descrever o sistema como um todo)

Há um evidente excesso de tarefas sendo feito unicamente na função “*int main()*” e no método `Bucket.process()`. Além disso os nomes em geral são pouco significativos e os próprios relacionamentos entre as funções e classes.

Serão feitas diversas mudanças, a principal delas é a mudança da classe “*Bucket*”. *Bucket* (balde) não é um nome significativo, nem chega perto de ser. O nome da classe será substituído por *SegmentSet* (conjunto de segmentos). Essa classe também deixa de ser construída dentro da função `int main()`, ganhando um construtor adequado. Outra classe que teve seu nome alterado foi “*CompareResult*”, que cuida da comparação de pares de strings e guarda o resultado dessas comparações. Ela foi renomeada para “*StringMerge*”, para evidenciar melhor seu funcionamento. Diversas mudanças nos nomes de variáveis e separações de funções grandes em diferentes funções menores foram feitas.

O diagrama a seguir retrata o estado atual do software, após as mudanças efetivadas ao longo do trabalho:

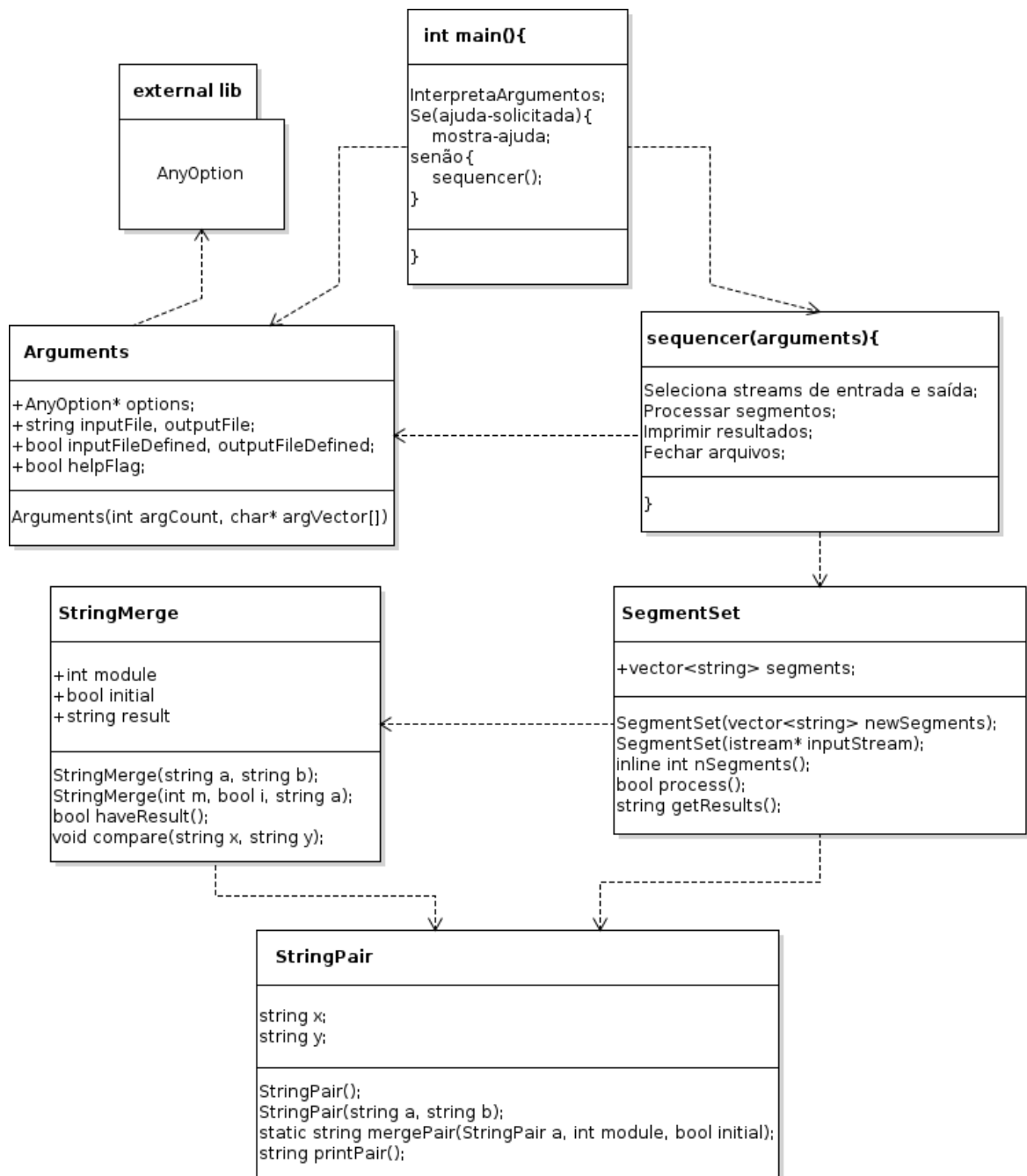


Figura 5: Diagrama da versão aprimorada de *Sequencer*, sem os módulos de testes.

5. Programação defensiva

Foram aplicadas algumas técnicas de programação defensiva no código das classes e funções com o objetivo de prevenir erros na lógica e na entrada/saída.

Em quase todos os métodos e funções foram utilizados *asserts* para se certificar de que os parâmetros recebidos estão adequados. Um exemplo é o método “StringPair::mergePair”, que combina duas strings em uma só. Nele, os *asserts* servem para ter certeza de que as strings são válidas e que o módulo da combinação (quantos caracteres de uma *string* já existem na outra) é maior ou igual a zero.

```
20 string StringPair::mergePair(StringPair a, int module, bool initial){
21     assert(module >= 0);
22     assert(a.x.size() >= 1 && a.y.size() >= 1);|
23
24     string result;
25     string cuttedY;
26     if(initial){
27         cuttedY = a.y.substr(0, a.y.size() - module);
28         result = cuttedY;
29         result.append(a.x);
30     }else{
31         cuttedY = a.y.substr(module, a.y.size() - module);
32         result = a.x;
33         result.append(cuttedY);
34     }
35
36     return result;
37 }
```

Figura 6: asserts sendo aplicados na fusão de um par de strings

Também foram utilizadas as estruturas de decisão *try-catch*. Diferentemente dos procedimentos *assert*, os *try-catch* são utilizados aqui com a finalidade de tentar detectar erros relacionados com entrada/saída e então tentar corrigir eles.

Durante a construção do “SegmentSet” isso foi utilizado para detectar a chegada do final do arquivo de leitura:

```
15 SegmentSet::SegmentSet(istream* inputStream){
16     assert(inputStream != NULL);
17
18     string line;
19     bool endOfFile = false;
20     do{
21         try{
22             getline(*inputStream, line);
23             //Caso o texto use quebras de linha do tipo \r
24             if(line[line.size()-1] == '\r'){
25                 line = line.substr(0, line.size()-1);
26                 //cout << "quebra de linha barra r ret
27             }
28             segments.push_back(line);
29         }catch(...){
30             endOfFile = true;
31         }
32     }while(!endOfFile);
33 }
```

Figura 7: try-catch no construtor do SegmentSet

Já durante a abertura dos arquivos de entrada/saída, no procedimento “sequencer()”, foram utilizados para verificar se foi possível abrir o arquivo com sucesso:

```
36     ostream* outputStream;
37     //seleciona a entrada de dados
38     ifstream *inputFileStream;
39     if(arguments->inputFileDefined){
40         inputFileStream = new ifstream;
41         inputFileStream->exceptions ( ifstream::failbit);
42         try{
43             inputFileStream->open(arguments->inputFile.c_str());
44             inputStream = inputFileStream;
45         }catch(ifstream::failure e){
46             cout << "ERROR OPENING A INPUT FILE:" << endl;
47             cout << arguments->inputFile;
48             cout << "USING CIN AS DEFAULT INPUT" << endl;
49             inputStream = &cin;
50         }
51     }
52     }else{
53         inputStream = &cin;
54     }
55
56     //seleciona a saída de dados
57     ofstream* outputFileStream;
58     if(arguments->outputFileDefined){
59         outputFileStream = new ofstream;
60         outputFileStream->exceptions ( ofstream::failbit);
61         try{
62             outputFileStream->open(arguments->outputFile.c_str(), ios::trunc);
63             outputStream = outputFileStream;
64         }catch(ifstream::failure e){
65             cout << "ERROR OPENING A OUTPUT FILE:" << endl;
66             cout << arguments->outputFile;
67             cout << "USING COUT AS DEFAULT OUTPUT" << endl;
68             outputStream = &cout;
69         }
70     }
71     }else{
72         outputStream = &cout;
73     }
```

Figura 8: try-catch sendo utilizados para tentar abrir arquivos de forma segura.

6. Testes

Alguns testes foram desenvolvidos com o objetivo de testar, separadamente, duas partes essenciais do sistema. Estes testes são construídos separadamente do resto do sistema. Cada um testa um determinado componente do sistema fornecendo entradas e verificando se as saídas são adequadas ao que é esperado.

```
9  struct MergeTest{
10     string name;
11     string a, b;
12     bool forcemerge;
13     string expectedResult;
14 };
15
16 vector<MergeTest> getTests();
17 void doTests();
18
19 int main(){
20     doTests();
21     return 0;
22 }
23
24 void doTests(){
25     vector<MergeTest> tests;
26     tests = getTests();
27
28     for(auto test : tests){
29         StringMerge merge(test.a, test.b, test.forcemerge);
30         bool success = (merge.result.compare(test.expectedResult) == 0);
31         if(!success){
32             cout << "StringMerge não passou no teste " + test.name;
33         } else {
34             cout << "StringMerge passou no teste " + test.name;
35         }
36         cout << endl;
37     }
38 }
```

Figura 09: Aplicação de testes para a classe “StringMerge”

Também foram feitos testes para a classe “SegmentTest”. Como esta classe faz uso de “StringMerge”, então por agregação “StringMerge” também está sendo testada neste teste.

```
9  struct ProcessTest{
10     vector<string> testSegments;
11     string expectedResult;
12     string name;
13 };
14
15 vector<ProcessTest> getTests();
16 void doTests();
17
18 int main(){
19     doTests();
20     return 0;
21 }
22
23 void doTests(){
24     vector<ProcessTest> tests;
25     tests = getTests();
26
27     for(auto test : tests){
28         SegmentSet testSet(test.testSegments);
29         testSet.process();
30         bool success = (testSet.getResults().compare(test.expectedResult) == 0);
31         if(!success){
32             cout << "SegmentSet não passou no teste " + test.name;
33         } else {
34             cout << "SegmentSet passou no teste " + test.name;
35         }
36         cout << endl;
37     }
38 }
```

Figura 10: Aplicação de testes para a classe “SegmentSet”

7. Documentação

8. Métricas

9. Referências:

1. <http://www.nature.com/scitable/topicpage/complex-genomes-shotgun-sequencing-609>;
2. <https://github.com/pentalpha/sequencer>;
3. <http://openmp.org/wp/>;
4. <http://bioinformatica.imd.ufrn.br/snnb/>;
5. OpenMP C and C++ Application Program Interface, página 5. Disponível em: <http://www.openmp.org/mp-documents/cspec20.pdf>;
6. <http://npad.imd.ufrn.br/>;
7. https://en.wikipedia.org/wiki/Standard_Template_Library;
8. <https://github.com/pentalpha/AnyOption>;
9. <http://www.cprogramming.com/c++11/what-is-c++0x.html>;