

Genome 540

Discussion Section Week 1


Jan 4, 2021


Chengxiang (CX) Qiu

Agenda (Tu/Th)

- Introductions
- C++/general programming tips
- HW1 context/advice/questions
- Input on future topics

Introductions

- 
- Who am I?
 - 4th year Genome Sciences student
 - Shendure Lab
 - Single-cell genomics
 - Mouse embryogenesis
 - Python/R

- 
- Who are you?
 - Department
 - Programming experience

Programming style

- The more readable your code is
 - ➡ • The more I'll be able to help you if something's wrong
 - ➡ • The more useful it will be to you later
 - ➡ • Try to use Github (version control and easy sharing)

Programming style

- Tips for readability

- ➔ • Intuitive variable/function names, (a, b, c vs. seq_name, seq_length)
- ➔ • Comments
 - Outlining general structure of program/key points of implemented algorithm
 - Clarifying any tricky/unintuitive lines of code
- ➔ • Simplicity over performance optimization (until it becomes necessary)
- ➔ • Write more functions
- ➔ • Write less “for” loops

General tips for homework

- ➔ • C/C++ or Java for some homework are necessary
- ➔ • Python works for most homework
- ➔ • R and MATLAB are not suitable

General tips for homework

- ➔ • It is important to understand biological questions
- ➔ • “Gap” between understanding the algorithm and writing it out
- ➔ • Your program might work on the test data but not work on the real data (much bigger and some unexpected issues)

General tips for homework

- ➔ • Print intermediate output
- ➔ • It's not necessary to write every basic function (for example, sort)
- ➔ • Please make an effort to match the template!

For a given DNA sequence

ACCTGCACTA**AAACCGT**TACACTGGGT**TCAAGAGATTCCC**

The longest match substring?

Substring?

Longest match substring?

In assignment 1, you need to compare two DNA sequences

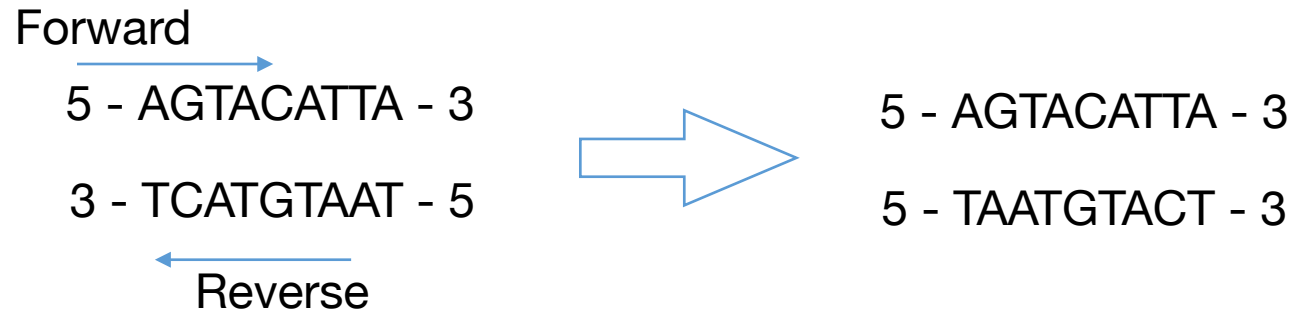
TTTTTATTATTAACCATGATAATTGTTGATAAATCTGTGGATAACTCTAAAAAAATTCCGGATTATAAA

and

AGTACATTAAAATATTTATATTTTAATGTAAATATTATATCACTTTTTCACAAAAACGTGTATTATATAT

- What's the longest match substring between them?
- Sequence 1 + Sequence 2 = A new Sequence

- Only comparing “substring” from different sequences
- Remember to consider both strands for sequence 2



- Do not store the substrings! (the real data is 10 Mb)

Difference between Python and C/C++

Python

- “interpreted”
- dynamic-typed
- data type is not required while declaring variable

```
A = 5
```

```
B = “x”
```

```
C = ['x', 'y', 'z']
```

C/C++

- “compiled”
- statically typed
- data type is required while declaring variable
- *debugging is challenging*

```
int A = 5;
```

```
char B = 'x';
```

```
char C[3] = {'x', 'y', 'z'};
```

Dynamic arrays in C++ (Vectors)

```
#include <iostream>
#include <fstream>
#include <sstream>
#include <vector> // must have this in order to use vector
#include <algorithm>
using namespace std;

int main(int argc, char* argv[] ){
    vector<string> vec; // make an empty vector that will be made up of strings
    vec.push_back("ZABC"); // adds "ZABC" to the vector, can be accessed with strings[0]
    vec.push_back("DEF"); // adds "DEF"
    for(auto x : vec){
        cout << x << endl;
    }
    vec.clear(); // empty the vec of all elements
    // documentation http://www.cplusplus.com/reference/vector/vector/
}
```

Custom data types in C++ (Structs)

```
struct complex{  
    int real;  
    int img;  
};
```

```
int main(int argc, char* argv[] ){  
    complex a = { 10 , 1 } ;  
    cout << "Real: " << a.real << "    Img: " << a.img << endl;  
  
    complex * p = &a;  
    cout << "Real: " << p->real << "    Img: " << p->img << endl;  
}
```

Pointers in C++

- Pointers are memory addresses, which point to variables
- There are two operators essential for handling pointers and memory addresses in C/C++: `*` and `&`

* as a suffix to a type is a “pointer”

`char* p;` // p is a memory location that stores a “char”

& as a unary prefix is the address-of operator and obtains the memory address of a variable

```
char x = 'h';
```

```
char* p = &x;
```

& as a suffix to a type means “pass by reference”

Example

```
void swapNums(int &x, int &y) {  
    int z = x;  
    x = y;  
    y = z;  
}  
  
int main() {  
    int firstNum = 10;  
    int secondNum = 20;  
  
    cout << "Before swap: " << "\n";  
    cout << firstNum << secondNum << "\n";  
  
    // Call the function, which will change the values of firstNum and secondNum  
    swapNums(firstNum, secondNum);  
  
    cout << "After swap: " << "\n";  
    cout << firstNum << secondNum << "\n";  
  
    return 0;  
}
```

* as a unary prefix means content of that memory location

Example

```
string food = "Pizza"; // A food variable of type string
string* ptr = &food;   // A pointer variable, with the name ptr, that stores the address of food

// Output the value of food (Pizza)
cout << food << "\n";

// Output the memory address of food (0x6dfed4)
cout << &food << "\n";

// Output the memory address of food with the pointer (0x6dfed4)
cout << ptr << "\n";
```

Difference between pointer to an array and array of pointers

- Pointer to an array - we are using the pointer to access the components of the array

```
// pointer to an array of five numbers;  
int (*ptr)[5];
```

- Array of pointers is an array of the pointer variables

```
// declaration of an array of pointers;  
int *ptr[5];
```

Arrays point to blocks of memory

- Arrays are just pointers to continuous blocks of memory

```
char word[6]; // an array of 6
characters called word
word[0] = 'a';
```

- Array indices are just pointer arithmetic and dereferencing combined

- $a[12]$ is the same as $*(a + 12)$

- $\&a[3]$ is the same as $a + 3$

	word[0]	word[1]	word[2]
	'h'	'i'	'\0'
const char *word =	0x80	0x88	0x90
"hi";	word	word+1	word+2

- Large arrays should be dynamically allocated (on the heap)

- Make sure you delete them

```
int n = some_large_number;
double * d = new double[n];
```

Sorting in C++ with sort

```
#include <iostream>
#include <fstream>
#include <sstream>
#include <vector>
#include <algorithm> // must have this in order to use sort
using namespace std;

bool my_cmp(const string & a, const string & b){
    return( a < b );
}

int main(int argc, char* argv[] ){
    vector<string> vec;
    vec.push_back("ZABC");
    vec.push_back("DEF");
    sort(vec.begin(), vec.end(), my_cmp );
    for(auto x : vec){
        cout << x << endl;
    }
}
```

Other sorting options in C++

function

qsort

<cstdlib>

```
void qsort (void* base, size_t num, size_t size,  
           int (*compar)(const void*,const void*));
```

Sort elements of array

Sorts the *num* elements of the array pointed to by *base*, each element *size* bytes long, using the *compar* function to determine the order.

The sorting algorithm used by this function compares pairs of elements by calling the specified *compar* function with pointers to them as argument.

The function does not return any value, but modifies the content of the array pointed to by *base* reordering its elements as defined by *compar*.

The order of equivalent elements is undefined.

Other sorting options in C++

function

qsort

<cstdlib>

```
void qsort (void* base, size_t num, size_t size,  
           int (*compar)(const void*,const void*));
```

Sort elements of array

Sorts the *num* elements of the array pointed to by *base*, each element *size* bytes long, using the *compar* function to determine the order.

The sorting algorithm used by this function compares pairs of elements by comparing pointers to them as argument.

The function does not return any value, but modifies the content of the array elements as defined by *compar*.

The order of equivalent elements is undefined.

<http://www.cplusplus.com/reference/>

```
1 /* qsort example */  
2 #include <stdio.h>      /* printf */  
3 #include <stdlib.h>     /* qsort */  
4  
5 int values[] = { 40, 10, 100, 90, 20, 25 };  
6  
7 int compare (const void * a, const void * b)  
8 {  
9     return ( *(int*)a - *(int*)b );  
10 }  
11  
12 int main ()  
13 {  
14     int n;  
15     qsort (values, 6, sizeof(int), compare);  
16     for (n=0; n<6; n++)  
17         printf ("%d ", values[n]);  
18     return 0;  
19 }
```


Homework 1

- Program a suffix array
- Test it on the files provided and compare results to the test output
- Run your program on the orthologous 10-megabase regions in the human and mouse genomes.
- Use the UCSC genome browser (hg38 and mm10) to figure out what biological feature they (the longest match) correspond to
- Submit your program output and your code

UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly

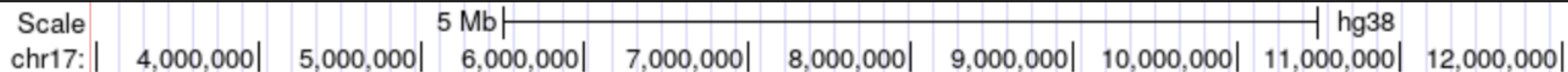
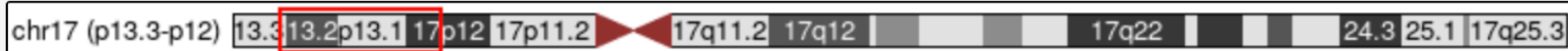
move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

multi-region

chr17:3,000,000-13,000,000 10,000,001 bp. gene, chromosome range, or other position, see examples

go

[examples](#)

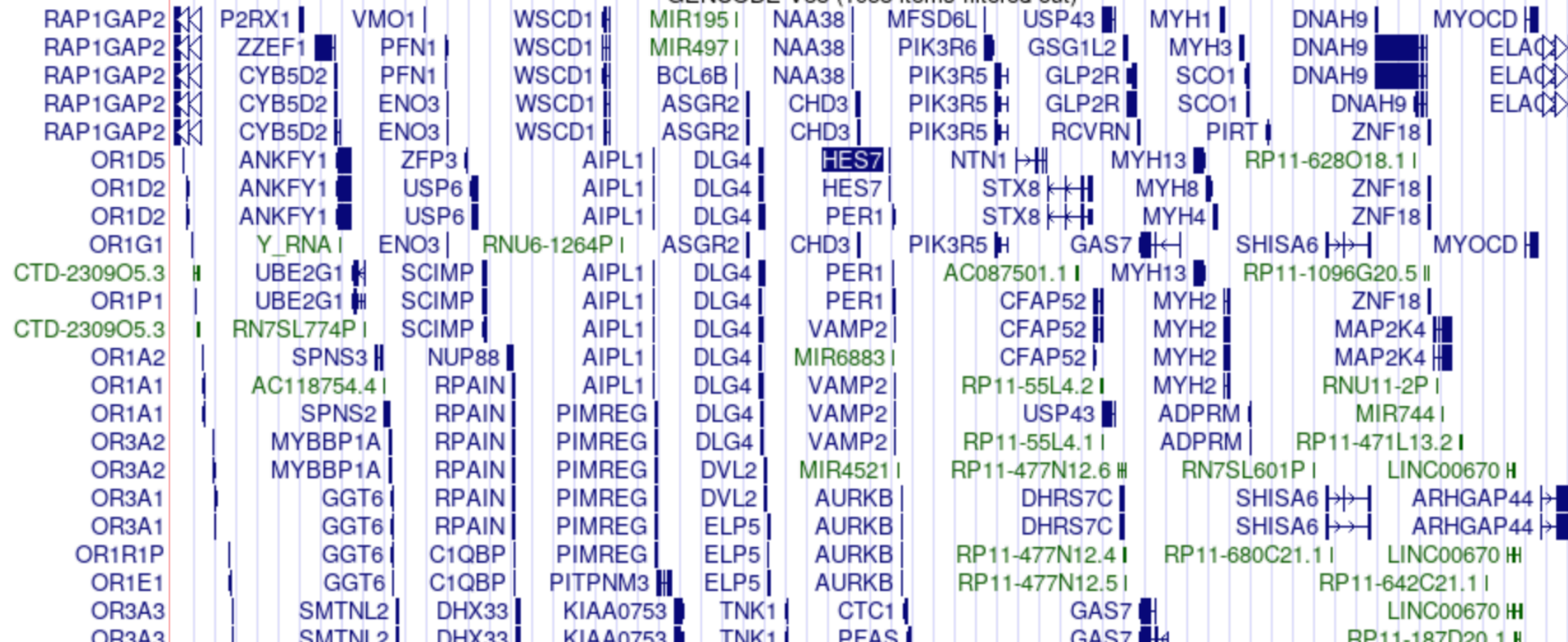


Reference Assembly Fix Patch Sequence Alignments

chr17_ML143374v1_fix
chr17_KV575245v1_fix

Reference Assembly Alternate Haplotype Sequence Alignments

GENCODE V38 (1058 items filtered out)



HW1 tips

- Start with pseudocode
- Start writing and then improve the details
- Don't hesitate, C++ is much faster than Python (at least for HW1)
- Get comfortable with pointers
- Output to template directly, also please make an effort to match the template

```

#include <stdio>
#include <iostream>
#include <fstream>
#include <string>
#include <cassert>

using namespace std;

void read_file(string filename, string& contents, int& num_lines)
{
    ifstream f;
    f.open(filename.c_str());
    string line;

    contents = "";
    num_lines = 0;
    while(getline(f, line)) {
        contents.append(line.substr(0, line.length()));
        num_lines++;
    }

    f.close();
}

int main(int argc, const char* argv[])
{
    string fn = argv[1];
    string contents;
    int num_lines;

    read_file(fn, contents, num_lines);

    cout << "Read: " << fn << "\n";
    cout << " * " << num_lines << " lines\n";
    cout << " * " << contents.length() << " characters (excluding newlines)\n";

    char * contents_cstring = (char*)contents.c_str();
    for (int i = 0; i < contents.length(); i++) {
        assert(contents_cstring[i] == *(contents_cstring + i));
        assert(contents_cstring[i] == contents.at(i));
    }
    assert(contents_cstring[contents.length()] == '\0');
}

```

```

#include <stdio>
#include <iostream>
#include <fstream>
#include <string>
#include <cassert>

using namespace std;

void read_file(string filename, string& contents, int& num_lines)
{
    ifstream f;
    f.open(filename.c_str());
    string line;

    contents = "";
    num_lines = 0;
    while(getline(f, line)) {
        contents.append(line.substr(0, line.length()));
        num_lines++;
    }

    f.close();
}

int main(int argc, const char* argv[])
{
    string fn = argv[1];
    string contents;
    int num_lines;

    read_file(fn, contents, num_lines);

    cout << "Read: " << fn << "\n";
    cout << " * " << num_lines << " lines\n";
    cout << " * " << contents.length() << " characters (excluding newlines)\n";

    char * contents_cstring = (char*)contents.c_str();
    for (int i = 0; i < contents.length(); i++) {
        assert(contents_cstring[i] == *(contents_cstring + i));
        assert(contents_cstring[i] == contents.at(i));
    }
    assert(contents_cstring[contents.length()] == '\0');
}

```

```

Annes-MacBook-Pro-2:sandbox anne$ g++ example.cpp -o example
Annes-MacBook-Pro-2:sandbox anne$ ./example example.cpp
Read: example.cpp
 * 45 lines
 * 971 characters (excluding newlines)

```

If you are using Mac
(under macOS12.1, my experience):

potentially you need to install *Xcode*,
and use clang++ instead of g++

Topics for future discussion sections?

- Scalable and reproducible bioinformatics pipelines (Snakemake)
- General programming tips
- Specific languages: Python, C++, Unix tools
- Additional applications of HMMs
- Dynamic programming
- Machine learning
- Version Control/Github
- Jupyter Notebooks/Reproducibility