# Dissecting the Tools of Bioinformatics: Building SpelFinder, a primer identification tool in Python

Charles Hardnett and Cynthia Bauerle
Departments of Computer Science and Biology, Spelman College, Atlanta, GA 30314

**Spelman College**
*A Choice to Change the World*

science education alliance
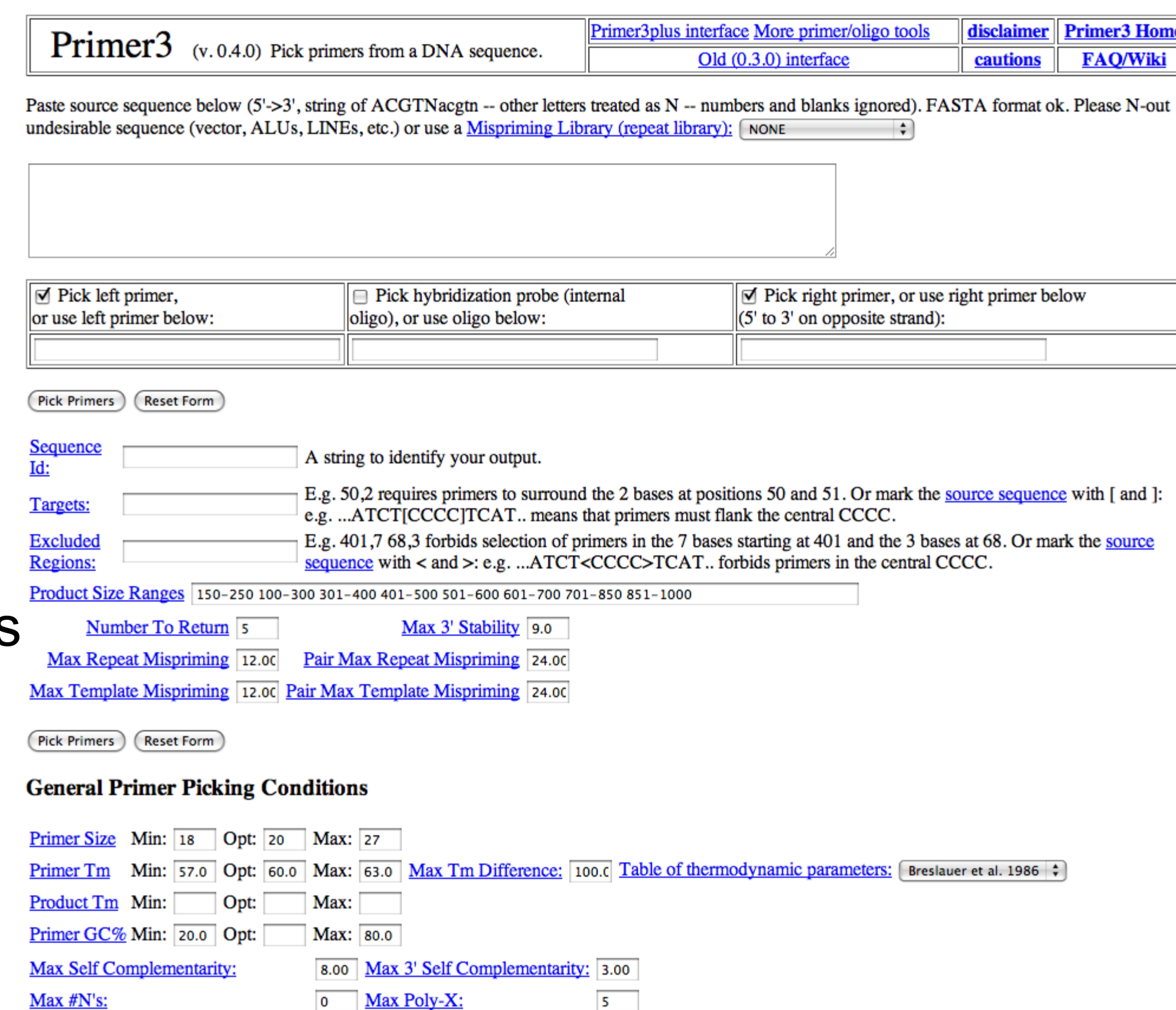HOWARD HUGHES MEDICAL INSTITUTE

## CONTEXT

Finishing the Genome: Students used Sequencher to identify low confidence/low coverage regions within the returned genomic sequence. They used Primer3 to select sites for the design of oligonucleotide primers to direct DNA sequencing across targeted regions.

Primer3: This tool contains a primer identification program readily available over the Internet, that identifies suitable primers. Users input a source sequence, set search parameters, and obtain a list of potential primers as output. Primer3 parameter variables include:

- Sequence length
- Melting T($T_m$)
- 3' stability
- GC%
- Poly-X runs
- self-complementarity
- Excludable regions



Understanding Primer3: To better understand how Primer3 operates to identify acceptable primer sequences, students built a simple sequence identification program in Python. Spelfinder searches a source sequence and identifies primers based on a simplified set of parameters:

- Sequence length: 18 ≥ 30 nucleotides
- GC%: 40 ≥ 60
- Poly-X ≤ 3 (where X = A, T, G, or C)
- $T_m$: 55°C ≤ $T_m$ ≤ 75°C
  where $T_m = [(G+C)_{primer} \times 4°C] + [(A+T)_{primer} \times 2°C]$

## GENERAL PROGRAMMING STRATEGY

Creating a program is an *iterative* process:
1. Define a search criterion
2. Design the algorithm using *pseudocode*
3. Implement the *algorithm* as a computer program
4. Test the computer program
5. Fix errors by repeating steps 1-3 until all tests pass

Program elements:
String: Nucleotide source sequence is stored as a collection of characters called a *string*
Algorithm: Each search parameter is represented by a computational procedure called an *algorithm* that sets criteria for the search
Program: All algorithms defining search parameters are assembled into a *script* that directs the function of the program

## OBJECTIVE
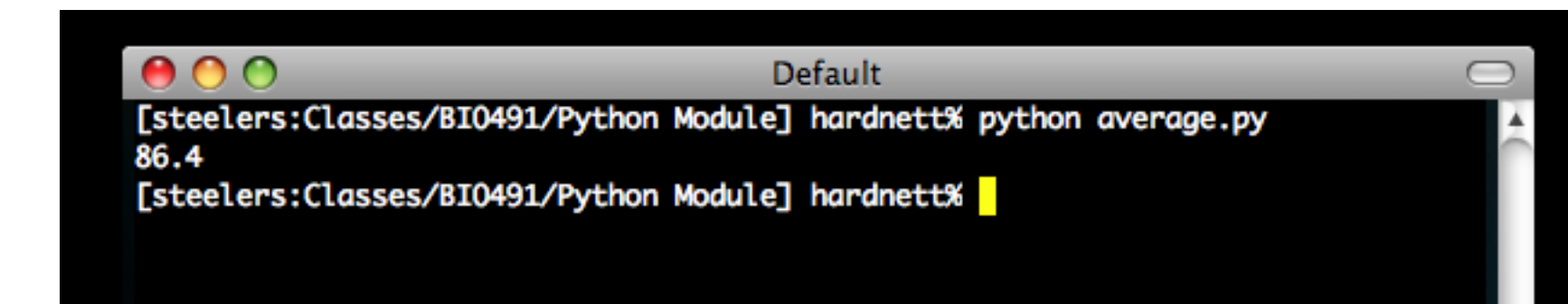### To help students better understand the construction and utility of standard sequence analysis tools

## ABSTRACT

As participant in the National Genomics Research Initiative supported by HHMI, the Biology department at Spelman College offered a course-based research opportunity for first year science students to isolate and characterize novel mycobacteriophages. In fall 2008, twenty first year students isolated novel mycophages and conducted standard microscopic and molecular characterization. In spring 2009, fourteen of these students, joined by six advanced biology majors, conducted genomic annotation of one phage isolate using standard sequence analysis approaches and bioinformatic tools available via the Internet. While prior experience *applying* bioinformatics tools varied among the cohort, no students described any prior experience in computer programming. To help students better understand the construction and utility of standard molecular analysis tools, we developed a course activity in which students used the Python programming language to build a primer identification tool that they named SpelFinder. Students were introduced to basic elements of the Python language in a hands-on group exercise, and then used these elements to build simple algorithms. Ultimately, students assembled algorithms into a script that analyzes an input nucleotide sequence and reports potential primer target sites as output. The program was based on an algorithm that reflected four relevant biological criteria. Writing the program *de novo* in Python introduced students to basic programming strategies for developing functional algorithms useful in standard bioinformatics tools. Specifically, students were challenged to understand how biological criteria may be translated into a set of rules that drives the algorithm of a sequence analysis software program.

## SPELFINDER: What the program looks like



## SPELFINDER: What the output looks like:



## PYTHON BASICS

Python is a general purpose, highly readable, high level computer programming language.

Python program may be created in a standard text editor such as Notepad or TextEdit. Once saved, the program may be run by using the python command to execute the program. The output will appear on the next line:



Python performs calculations using common arithmetic symbols:
- Addition: +
- Subtraction: -
- Multiplication: *
- Division: / and / /
- Equal: =
- Exponentiation: * *
- Modulus: % (finds remainder, e.g. 10 % 3 = 1)

## TASK 1: Compute all subsequences of lengths 18 ≥ 30 from the source sequence string

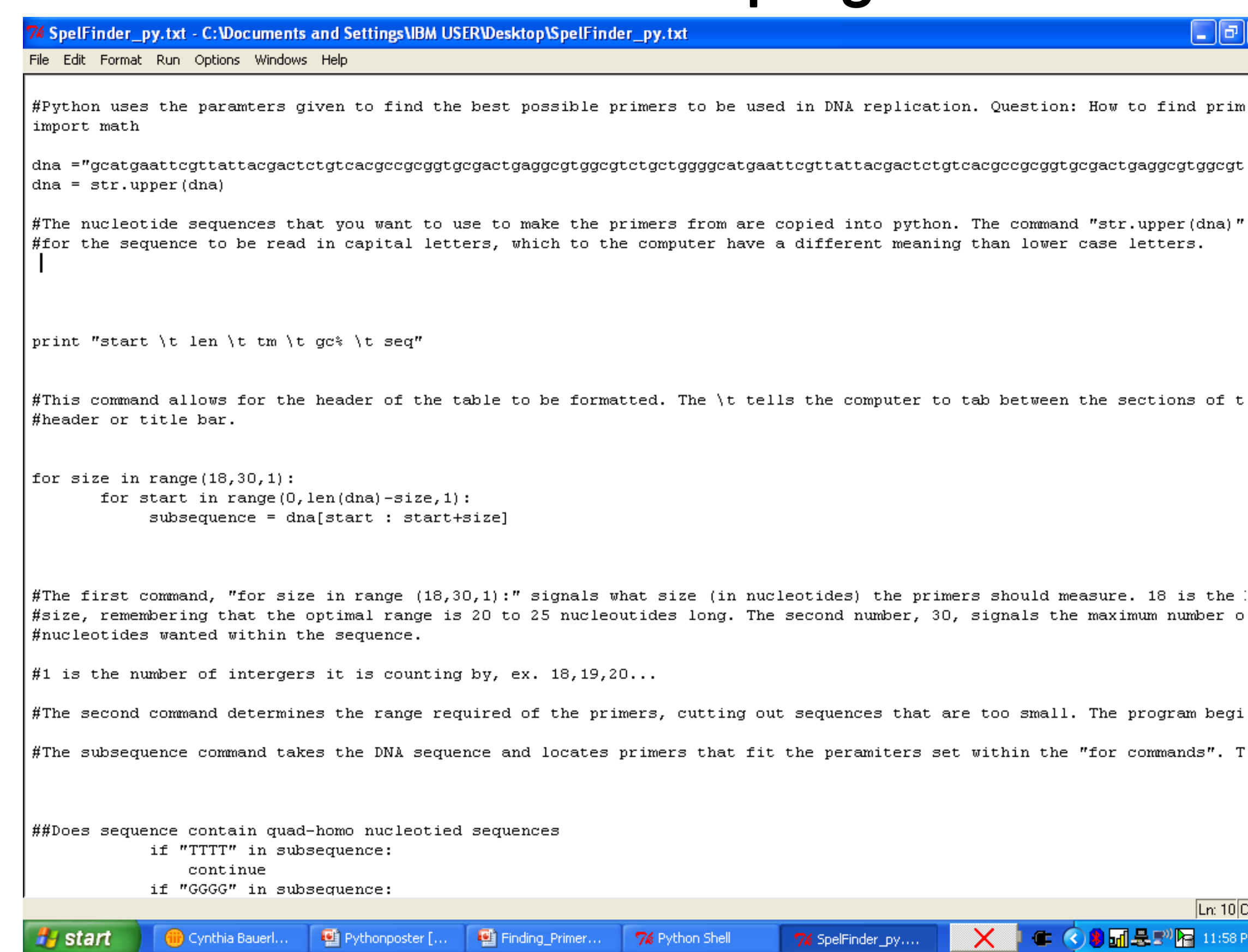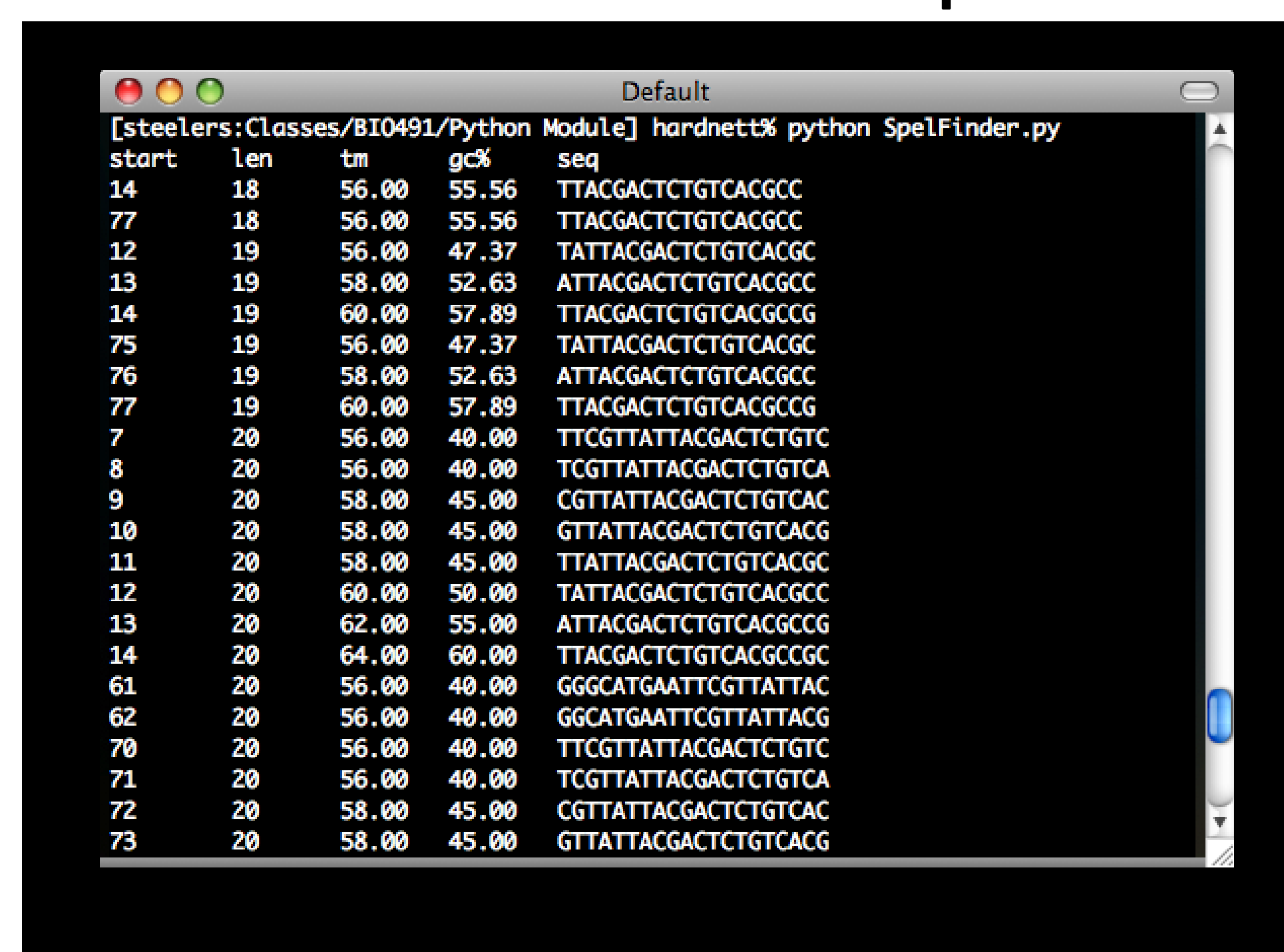Pseudocode: *Start = 1*
*While (Start < length of DNA Sequence)*
*For each extracted subsequence starting at Start with length 18, 19, 20, 21,...30*
*Analyze the subsequence using criteria*
*1. Disallowed consecutive nucleotides*
*2. GC Content range*
*3. Melting Temperature Range*
*End For*
*increment Start (create next start position)*
*End While*

Python:
```
for size in range(18,30,1):
    for start in range(0,len(dna)-size,1):
        subsequence = dna[start : start+size]
```

Output:

- The first statement assigns values 18, 19, 20, ... 30 to size on at a time

- The second statement assigns starting positions 0, 1, 2, 3, ..., end
  - Python strings start at position 0 and not 1!

- The third extracts the subsequence from the start to its position at a length of size (start+size is the end position for the subsequence)

## TASK 2: Skipping subsequences containing poly-X$_n$ where X is A, T, G, or C, and n ≥ 4

Pseudocode:
*if (subsequence contains "TTTT") or*
*(subsequence contains "AAAA") or*
*(subsequence contains "CCCC") or*
*(subsequence contains "GGGG") then*
*skip this subsequence*

Python:
```
if "TTTT" in subsequence:
    continue
if "GGGG" in subsequence:
    continue
if "AAAA" in subsequence:
    continue
if "CCCC" in subsequence:
    continue
```

Output:
- The **continue** means to skip the other criteria in the loop and go to the next loop iteration i.e. next subsequence

- The **in** operator is a Python string operator for determining the existence of a substring within a string