

Biol 591 Introduction to Bioinformatics (Fall 2002): Problem Set 3

This problem set is smaller than usual (bearing in mind that it combines questions concerning bioinformatics, molecular biology, and programming) because of the heft of the study questions. Please understand that study questions are to be given the same consideration as questions on problem sets.

P3.1. What can you use to determine whether a string of characters in a Perl program is:

- | | | |
|---------------|---------------------------|-------------------------|
| a. a comment | c. a subroutine | e. a regular expression |
| b. a variable | d. a call to a subroutine | f. an array |

P3.2. Modify BlastN so that it no longer prints out a complete match but prints out instead only each initial exact match of a word.

P3.3. Examine BlastN and determine the values used for the following quantities:

- | | | |
|---------------------|--------------------------|--------------|
| a. Match reward | c. Gap open penalty | e. Word size |
| b. Mismatch penalty | d. Gap extension penalty | |

P3.4. Modify BlastN so that it prints out for each hit both the raw score and the score in bits. To do this you may need to find values for lambda and K. Do this by running ANY pairwise sequence comparison at the NCBI site, using the same parameters you use in local BlastN, and noting the values of lambda and K at the end of the output.

P3.5 Estimate how much more efficient BlastN is than a full Smith-Waterman algorithm. Proceed as follows.

- A. Presume that the total time spent by each program is proportional to the number of cells in scoring tables each has to calculate (so your job is reduced to figuring out how many cells that is in each case).
- B. Consider a specific case of a comparison of a 100-nucleotide query sequence with the E. coli genome. How big would the Smith-Waterman scoring matrix be?

Don't know how big the E. coli genome is? You have a program that can tell you! Recall that SequenceSearch reads in the genome of Nostoc in order to search for putative NtcA binding sites. Well, perhaps you can change where it reads in the Nostoc sequence and have it read in the E. coli sequence. Once the sequence is in a variable you can add the line

`print length(variable);`

and you have it! (you'll have to put in the right name in place of `variable`).

- C. OK, you got half the job done. Now you need to find out how many cells Blast would need to calculate. First of all, how many word matches would you expect Blast to find? Consider two cases: a word-size of 11 and a word-size of 7.

How do you find how many exact word matches there will be? Again, you have a program for the job – in fact the same program! Modify

