

Genome 540 Discussion

January 8th, 2024

Clifford Rostomily

Agenda

- Assignment #1
- Questions
- Random Stuff



Assignment #1



Assignment Overview

- Read in two fasta files
 - Track # non-alpha characters
 - Also track base counts
- Combine the 3 sequences and store subseqs as a pointer array/vector
 - Forward of seq1, forward and reverse of seq2 = 3
- Implement and run the suffix array algorithm
 - Returns a sorted list of pointers
- Iterate through results
 - Track longest match length of each seq1 suffix to either the forward or reverse strand of seq2
 - Track the longest overall match

Fasta format

- File format for storing sequences
- Can store multiple sequences
- Sequences are preceded by a single header line denoted by a “>”

```
> my_sequence1
```

```
TCGATCGATGGCTTCGGATGCGCTTAG
```

```
> my_sequence2
```

```
GCTGCGCGAGAACAAACGTAGCATGGC
```

- Can be used to store protein and nucleic acid sequences
 - The extension can tell you more about what is being stored

Extension ↕	Meaning ↕	Notes ↕
fasta, fas, fa ^[9]	generic FASTA	Any generic FASTA file
fna	FASTA nucleic acid	Used generically to specify nucleic acids
ffn	FASTA nucleotide of gene regions	Contains coding regions for a genome
faa	FASTA amino acid	Contains amino acid sequences
mpfa	FASTA amino acids	Contains multiple protein sequences
frn	FASTA non-coding RNA	Contains non-coding RNA regions for a genome, e.g. tRNA, rRNA

https://en.wikipedia.org/wiki/FASTA_format

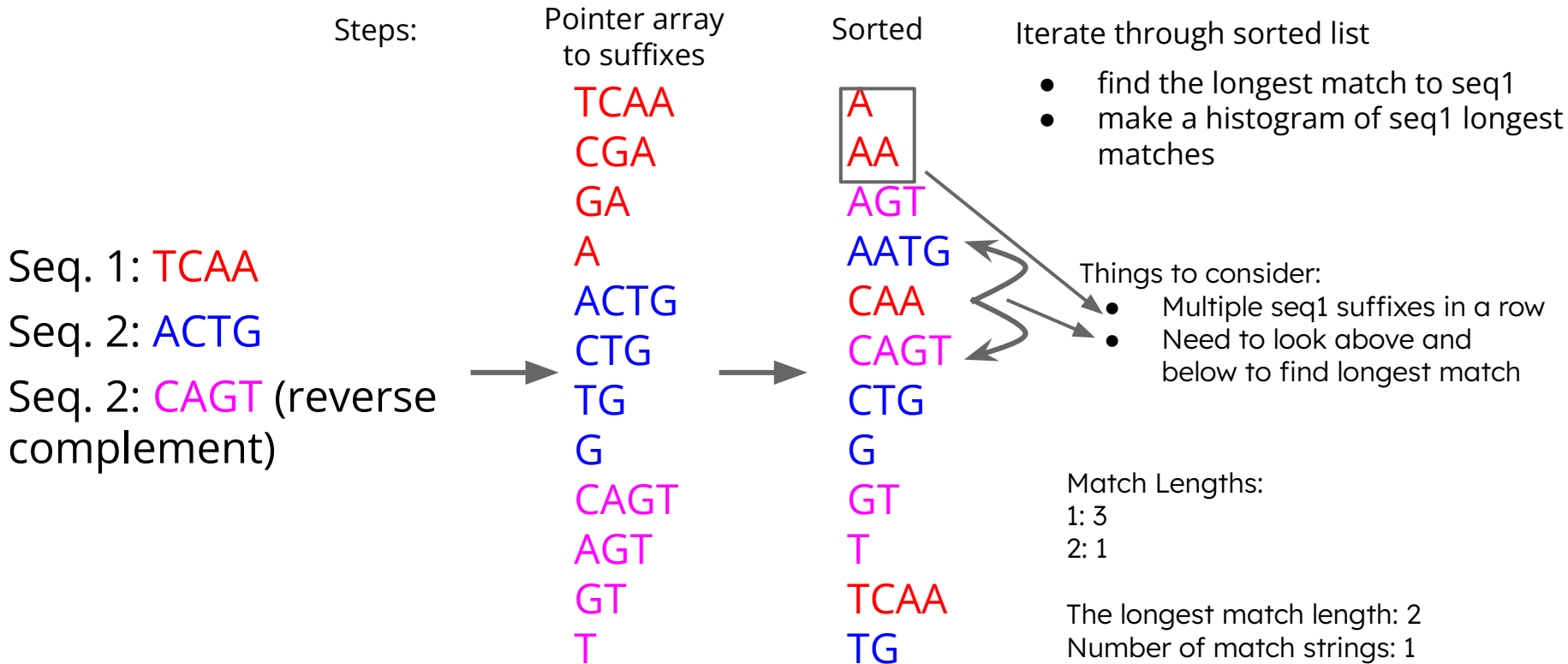
Fasta IUPAC ambiguity codes

Nucleic Acid Code ⇅	Meaning ⇅	Mnemonic ⇅
A	A	A denine
C	C	C ytosine
G	G	G uanine
T	T	T hymine
U	U	U racil
(i)	i	i nosine (non-standard)
R	A or G (I)	puR ine
Y	C, T or U	pY rimidines
K	G, T or U	bases which are K etones
M	A or C	bases with aM ino groups
S	C or G	S trong interaction
W	A, T or U	W eak interaction
B	not A (i.e. C, G, T or U)	B comes after A
D	not C (i.e. A, G, T or U)	D comes after C
H	not G (i.e., A, C, T or U)	H comes after G
V	neither T nor U (i.e. A, C or G)	V comes after U
N	A C G T U	N ucleic acid
-	gap of indeterminate length	

Non-alphabetic characters

- Exclude the header line
- Exclude white space (e.g. spaces)
- Include only digits from seq position numbers

Small Example



Questions?

- Unless you specifically ask me not to bring it up I will try to cover common questions asked on slack during the next class discussion.



Random Stuff



The Burrows Wheeler Transform

Transformation				
1. Input	2. All rotations	3. Sort into lexical order	4. Take the last column	5. Output
<div style="border: 1px solid gray; padding: 5px; width: fit-content;"> <code>^BANANA\$</code> </div>	<div style="border: 1px solid gray; padding: 5px;"> <code>^BANANA\$</code> <code>\$^BANANA</code> <code>A\$^BANAN</code> <code>NA\$^BANA</code> <code>ANA\$^BAN</code> <code>NANA\$^BA</code> <code>ANANA\$^B</code> <code>BANANA\$^</code> </div>	<div style="border: 1px solid gray; padding: 5px;"> <code>ANANA\$^B</code> <code>ANA\$^BAN</code> <code>A\$^BANAN</code> <code>BANANA\$^</code> <code>NANA\$^BA</code> <code>NA\$^BANA</code> <code>^BANANA\$</code> <code>\$^BANANA</code> </div>	<div style="border: 1px solid gray; padding: 5px;"> <code>ANANA\$^B</code> <code>ANA\$^BAN</code> <code>A\$^BANAN</code> <code>BANANA\$^</code> <code>NANA\$^BA</code> <code>NA\$^BANA</code> <code>^BANANA\$</code> <code>\$^BANANA</code> </div>	<div style="border: 1px solid gray; padding: 5px; width: fit-content;"> <code>BNN^AA\$A</code> </div>

- Used originally for compression
- The Bowtie aligner uses it for compression and indexing

Bowtie

