

MCB 5472 Midterm

Your Name

In case of multiple choice questions, more than one answer might be correct.

1) Why is %identity not a good parameter to assess the significance of a match in a blast search?
% identity does not take the lengths of a match into consideration (a 100% match of three residues frequently occurs by chance)

2) Two sequences found in two different organisms are homologous if and only if they evolved from the same ancestral sequence that existed in **some** organism in the past. **Correct - Incorrect**

3) Homologous sequences almost always show significant similarity. **Correct - Incorrect**
(most homologs have diverged beyond recognition)

4) Proteins with moderately complex amino acid composition that show significant similarity in a BLAST search almost always are homologous. **Correct - Incorrect**

5) In a BLASTP search an E-value for a match is reported as 3×10^{-62} . Which of the following statements is correct.

A) A value of this magnitude is not sufficient to suggest homology. Convergent evolution frequently gives rise to sequence matches with very small E-values.

B) This means that a match of this quality due to chance would be obtained 3×10^{-62} times. A match of this quality demonstrates homology beyond reasonable doubts.

C) This means that a match of this quality due to chance would be obtained 3×10^{-62} times (e denotes Euler's number, about 2.72).

6) Give a short definition of E-value and P-value:

E-value: Number of matches expected due to chance with a score equal or better than the one obtained.

P-value: probability to obtain a match with a score equal or better than the one obtained due to chance alone.

7) Which of the following are correctly formed fasta sequences? (new line symbols are given as \n; tabulator symbols as \t)

```
>gi|12643370|sp|Q9P997.2|VATA_THEAC this is a test sequence RecName: Full=Tac atpA intein \n MGKIIRISGPVVVAEDVEDAKMYDVVKVKGEMGLIGEIIKIEGNRSTIQVYEDTAGIRPDEKVENTRRPLS \n
```

```
>gi|12643370|sp|Q9P997.2|VATA_THEAC this is a test sequence RecName: Full=Tac atpA intein \n MGKIIRISGPVVVAEDVEDAKMYDVVKVKGEM \n GLIGEIIKIEGNRSTIQVYEDTAGIRPDEKVENTRRPLS \n
```

```
\n >gi|12643370|sp|Q9P997.2|VATA_THEAC this is a test sequence RecName: Full=Tac atpA intein \n MGKIIRISGPVVVAEDVEDAKMYDVVKVKGEMGLIGEIIKIEGNRSTIQVYEDTAGIRPDEKVENTRRPLS \n  
space following the first \n would be a problem
```

```
>gi|12643370|sp|Q9P997.2|VATA_THEAC \t Full=Tac atpA intein \n MGKIIRISGPVVVAEDVEDAKMYDVVKVKGEMGLIGEIIKIEGNRSTIQVYEDTAGIRPDEKVENTRRPLS \n
```

```
>gi|12643370|sp|Q9P997.2|VATA_THEAC \n Full=Tac atpA intein \n MGKIIRISGPVVVAEDVEDAKMYDVVKVKGEMGLIGEIIKIEGNRSTIQVYEDTAGIRPDEKVENTRRPLS \n
```

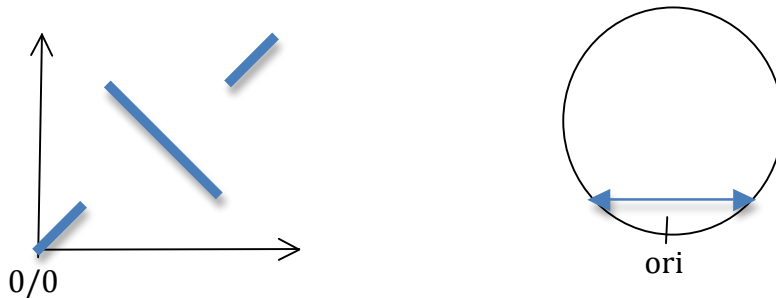
8) What steps does clustal performs to calculate a multiple sequence alignment?

A) calculate all possible pairwise alignments, create bootstrap samples, calculate guide tree, calculate cluster alignment

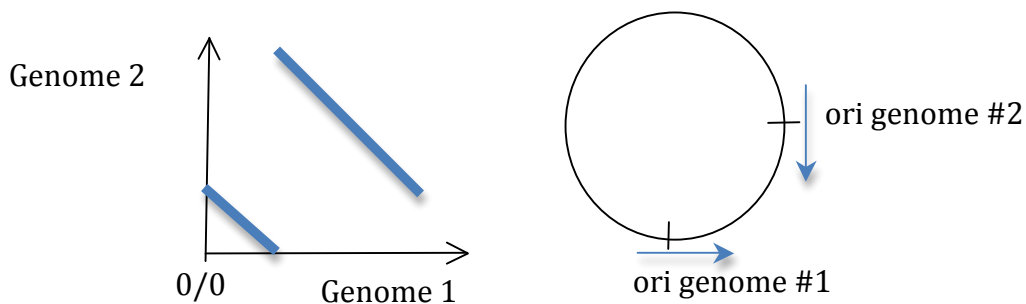
B) calculate all possible pairwise alignments, calculate master slave alignment using the most conserved sequence as master

C) calculate all possible pairwise alignments, calculate guide tree, calculate cluster alignment

9) Assume that the origin of replications are at (0/0) in the coordinate system on the left (below), and that you compare two closely related genomes. The organisms in question have circular genomes. The sketch below on the right indicates where a rearrangement event took place in one of the genomes. Draw a sketch of the expected gene plot, assuming that except for the single recombination event the two genomes are syntenic.



10) Assume that you compare two syntenic genomes in a gene plot comparison. One of the genomes has the origin of replication assigned to the wrong place. The direction in which the ORFs are listed in the genomes is given by the arrow. What result do you expect for the gene plot?



11) Regarding types of error in a normal BLAST search, which is correct:

- A) False positives (i.e. non-homologous proteins in the databank that are identified as matches) occur frequently
- B) False negatives (i.e. homologous proteins that are not identified as matches) occur frequently.**
- C) Most false negatives can be identified, if one turns on the filter for low complexity

12) What are the names given to the three domains of life.

- A) Prokaryotes, Eukaryotes, Bacteria
- B) Prokaryota, Eukaryota, Archaea
- C) Archaea, Bacteria, Eukaryotes**
- D) Plants, Animals, Protozoa

13) The leg-hemoglobin in plants is a homolog to the myo- and hemoglobin molecules found in animals. If all these proteins evolved from the same ancestral molecule that was present in the most recent common ancestor of plants and animals, then, according to Walter Fitch's terminology,

- A) the plant globin is an ortholog to the animals myoglobin genes and an ortholog to animals hemoglobin genes.**
- B) the plant homolog is a paralog to the animal hemoglobin genes
- C) hemoglobin genes are parlogs to the myoglobin genes**

14) In perl, which beginning symbols denote variables that contain Scalars, Arrays, Hashes?

\$, @, %

For part A:

```
#! \usr\bin\perl
use warnings;
while(defined($file=glob("*.faa"))){
    system("clustalw -align -infile=$file -type=protein");
}
```

For part B:

```
#!\usr\bin\perl
#MCB 5472 Exam Part 2
use warnings;
my %amino_acids;
#We will open each .faa file in the directory
#We will groom the file to remove GI numbers and new lines
#then we will set up an array that goes through the file character by character

while (defined($file=glob("*.faa")))
{
    open(IN, "< $file") or die "cannot open $file:$!";
    my $inseq = "";
    while(defined(my $line=<IN>)){
        chomp($line);
        $line =~ s/\s//g; #removes blank spaces
        if ($line=~/^>/) {} #ignores annotation lines
        else {$inseq .= $line}; #transfers file to $inseq
    }
    close(IN);
    my @seqarray = split(//,$inseq); #makes an array from $inseq
    foreach (@seqarray){
        #Now we will compare each character to all the keys %amino_acids
        $amino_acids{$_} += 1;
    }

}

# Now we must print out the hash

print "\n\nThe number of times each amino acid occurs in $file is:\n";

my @keys = keys (%amino_acids);
my @vals = values (%amino_acids);

foreach my $key (@keys){
    print"\t$key\t$amino_acids{$key}\n";
    $amino_acids{$key}=0 #Sets value back to 0 for next pass thru
}

};
```

```
#MCB 5472 Exam Part 2
use warnings;
my %amino_acids;
#We will open each .faa file in the directory
#We will groom the file to remove GI numbers and new lines
#then we will set up an array that goes through the file character by character
while (defined($file=glob("*.faa")))
{
    open(IN, "< $file") or die "cannot open $file:$!";
    my $inseq = "";
    while(defined(my $line=<IN>)){
        chomp($line);
        $line =~ s/\s//g; #removes blank spaces
        if ($line=~/^>/) {} #ignores annotation lines
        else {$inseq .= $line}; #transfers file to $inseq
    }
    close(IN);
    my @seqarray = split(//,$inseq); #makes an array from $inseq
    foreach (@seqarray){
#Now we will compare each character to all the keys %amino_acids
        $amino_acids{$_} += 1;
    }
}

# Now we must print out the hash

print "\n\nThe number of times each amino acid occurs in $file is:\n";

my @keys = keys (%amino_acids);
my @vals = values (%amino_acids);

foreach my $key (@keys){
    print "\t$key\t$amino_acids{$key}\n";
    $amino_acids{$key}=0 #Sets value back to 0 for next pass thru
}

};
```