# MID-TERM  EXAMINATION

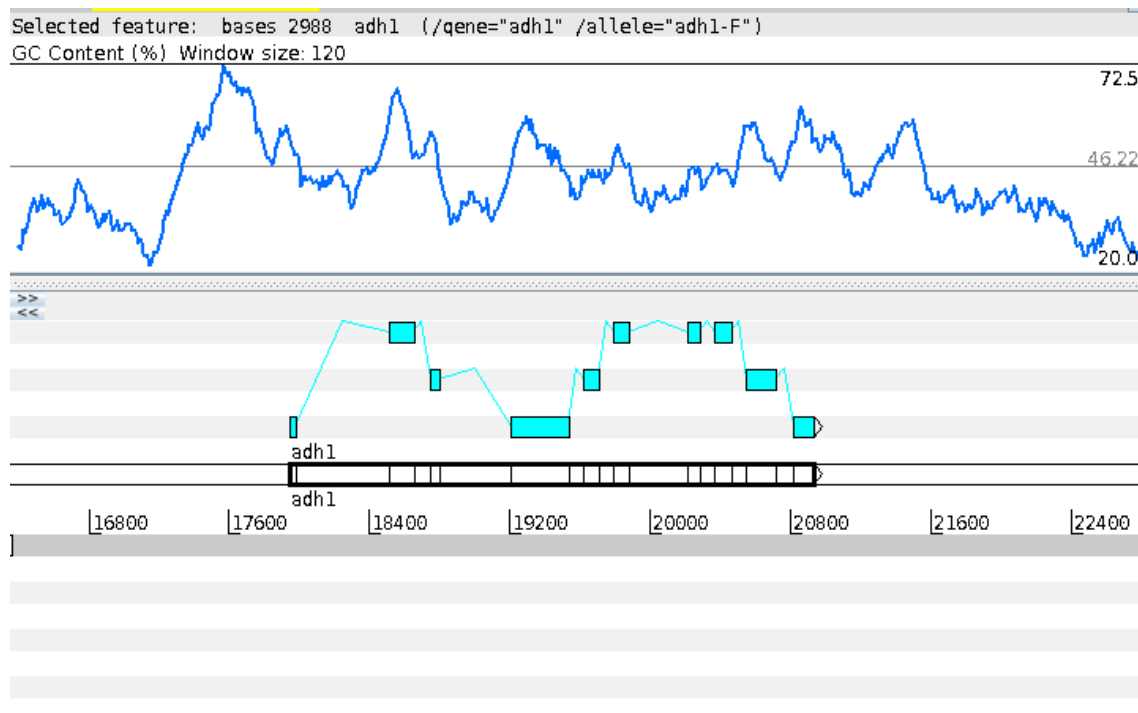08:30 -  9:45   Tuesday, October 25, 2016

Answer any combination of questions totalling to <u>exactly</u> 100 points. If you answer questions totalling more than 100 points, answers will be discarded at random until the total points are less than or equal to 100. This exam is worth 20% of the course grade.

Hand in this question sheet along with your exam book. All questions must be answered in the exam book. The exam sheets will be shreded after the exam.

---

Ways to write a readable and concise answer:
i. Just answer the question. Save time by specifically addressing what is asked. Don't give irrelevant background if it doesn't contribute to the question that was asked.
ii. Avoid stream of consciousness. Plan your answer by organizing your key points, and then write a concise, coherent answer. Make your point once, clearly, rather than repeating the same thing several times with no new information.
iii. Point form, diagrams, tables, bar graphs, figures are welcome. Often they get the point across more clearly than a long paragraph.
iv. Your writing must be legible. If I can't read it, I can't give you any credit.

---

1. (10 points) A map of one of the loci encoding alcohol dehyrogenase 1 (adh1-f) in Maize is shown.  What can you say about the structure of this gene?



2. (5 points) The BLAST database services at NCBI must process over 100,000 BLAST searches per day. Researchers at NCBI realized that the most critical bottleneck in the process was the simple matter of reading in all the sequence data when comparing a query sequence with all sequences in a database. What solution was found to solve this problem?

3. (10 points) The following features are annotated in a mouse sequence found in GenBank.

```
Key             Location/Qualifiers

source          1..1509
                /organism="Mus musculus"
                /strain="CD1"
                /mol_type="genomic DNA"
promoter        <1..9
                /gene="ubc42"
mRNA            join(10..567,789..1320)
                /gene="ubc42"
CDS             join(54..567,789..1254)
                /gene="ubc42"
                /product="ubiquitin conjugating enzyme"
                /function="cell division control"
```

a) The annotation for the promoter is expressed as "<1..9". Explain what is meant by this annotation.

b) The mRNA and CDS features imply the existence of intron and exon features. In the format of the Features Table, write the annotation for these features.

4. (15 points) Fill in the blanks. In your exam booklet, just write a term for a - e. You don't need to rewrite the entire text. Note that for e, two terms should be given.

## Protein database searches compared to DNA database searches
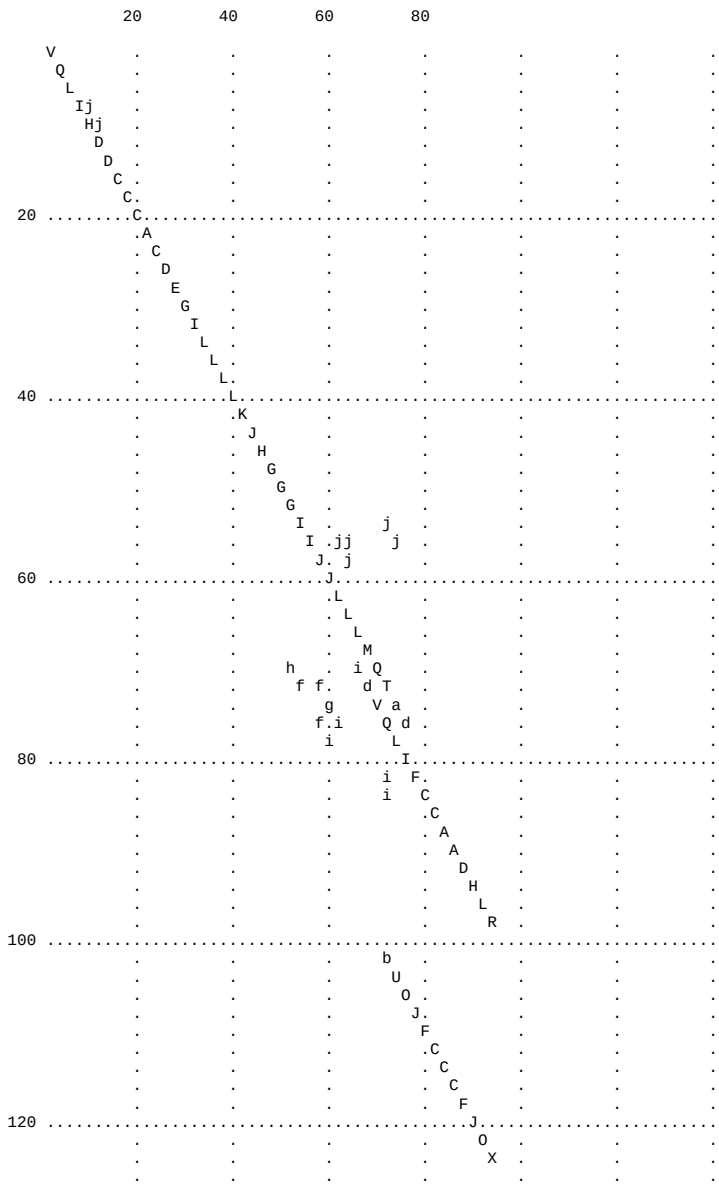
Speed

- A protein coding DNA sequence contains 3 times as many characters as the corresponding amino acid sequence
- Protein databases are much smaller than DNA databases because
    - _____a_____ are not present
    - where several DNA sequences encode identical proteins, only 1 protein database entry is usually created
- Speedup using lookup tables is more efficient with proteins. For proteins, k=2 yeilds a ___b___-fold speedup, whereas in DNA, k=4 yeilds only a 256-fold speedup.

Sensitivity
- The degeneracy of the genetic code makes it possible for DNA sequence to _____ c_____ more rapidly than amino acid sequence
- The greater complexity of _____ d _____ allows matches to be detected at very low levels of similarity
- The small alphabet size of DNA (4) compared to proteins (20) makes protein alignments more robust. That is, it is often far more obvious which _____ e-1 _____ should be aligned, compared to which _____e-2_____.
- DNA alignments tend to have more gaps, compared to protein alignments.

5. (10 points) The dot-matrix plot below shows a comparison of two SAR8  proteins. What are the most obvious differences or mutations between these two sequences?

```
P1HOM        Version  5/13/91
X-axis: XP_009786495
Y-axis: XP_009588711
SIMILARITY RANGE:  10      MIN.PERCENT SIMILARITY:  30
SCALE FACTOR:   0.90     COMPRESSION:              2

             20        40        60        80

        V      .        .        .        .        .        .        .
        Q      .        .        .        .        .        .        .
         L     .        .        .        .        .        .        .
        Ij     .        .        .        .        .        .        .
        Hj     .        .        .        .        .        .        .
        D      .        .        .        .        .        .        .
        D    .        .        .        .        .        .        .
        C    .        .        .        .        .        .        .
        C.        .        .        .        .        .        .
  20 .........C..............................................................
         .A        .        .        .        .        .        .
         . C       .        .        .        .        .        .
         .  D      .        .        .        .        .        .
         .   E     .        .        .        .        .        .
         .    G    .        .        .        .        .        .
         .     I   .        .        .        .        .        .
         .      L .        .        .        .        .        .
         .       L .        .        .        .        .        .
         .        L.        .        .        .        .        .
  40 ...................L....................................................
         .          .K      .        .        .        .        .
         .          . J     .        .        .        .        .
         .          .  H    .        .        .        .        .
         .           G      .        .        .        .        .
         .           G      .        .        .        .        .
         .           G  .   .        .        .        .        .
         .          I  .    j   .        .        .        .        .
         .          I .jj    j  .        .        .        .        .
         .          J. j        .        .        .        .        .
  60 ...........................J..........................................
         .        .        .L        .        .        .        .
         .        .        . L       .        .        .        .
         .        .         . L      .        .        .        .
         .        .          . M     .        .        .        .
         .        .        h  . i Q  .        .        .        .
         .        .         f f.  d T .        .        .        .
         .        .          g   V a .        .        .        .
         .        .         f.i  Q d .        .        .        .
         .        .          i      L .        .        .        .
  80 ...............................I.......................................
         .        .        .        i F.       .        .        .
         .        .        .        i  C       .        .        .
         .        .        .        .C          .        .        .
         .        .        .        . A         .        .        .
         .        .        .        . A         .        .        .
         .        .        .        .  D        .        .        .
         .        .        .        .   H       .        .        .
         .        .        .        .    L      .        .        .
         .        .        .        .     R .   .        .        .
 100 ....................................................................
         .        .        .        b .         .        .        .
         .        .        .        U .         .        .        .
         .        .        .        O .         .        .        .
         .        .        .        J.          .        .        .
         .        .        .        F           .        .        .
         .        .        .        .C          .        .        .
         .        .        .        . C         .        .        .
         .        .        .        .  C        .        .        .
         .        .        .        .  F        .        .        .
 120 .........................................J...........................
         .        .        .        .   O      .        .        .
         .        .        .        .   X .     .        .        .
         .        .        .        .        .        .        .
```

6. (15 points) If you wanted to design an oligonucleotide as a hybridization probe, you want to ensure that the oligo sequence is unique within the genome ie. it is not likely to occur by random chance. To help in your calculations, a table is given with some relevant information.

| n | $4^n$ | $2 \times 4^n$ |
|---|---|---|
| 10 | 1.05E+06 | 2.10E+06 |
| 11 | 4.19E+06 | 8.39E+06 |
| 12 | 1.68E+07 | 3.36E+07 |
| 13 | 6.71E+07 | 1.34E+08 |
| 14 | 2.68E+08 | 5.37E+08 |
| 15 | 1.07E+09 | 2.15E+09 |
| 16 | 4.29E+09 | 8.59E+09 |
| 17 | 1.72E+10 | 3.44E+10 |
| 18 | 6.87E+10 | 1.37E+11 |

a) How big would an oligo probe have to be for use with haploid yeast, *Saccharomyces cerevisiae*, (1N = 1.2 x 10$^7$ bp)? That is, how long does the oligo have to be to ensure that it is not likely to occur in the genome due to random chance?

b) Yeast also go through a diploid phase. If you were hybridizing to DNA extracted from diploid yeast, would you need to use a longer oligo? Explain.

c) Most eukaryotic genomes, especially for higher organisms, are largely composed of middle repetitive sequences such as the AluI family in mammals. How would this affect our estimates of the likelihood of finding a particular oligonucleotide in a eukaryotic genome?

7. (10 points) The shell script below accepts a FASTA file as input. What is written to the output? An example of a FASTA file containing 3 sequences is shown in the box.

```bash
#!/bin/bash

infile=$1
outfile=$2

cat $infile | grep '>' | cut -c2- | cut -f1 -d " " > $outfile
```

```
>BLYTHNA 137 bp
MATNKSIKSVVICVLILGLVLEQVQVEAKSCCKNTTGRNCYNACRFAGGSRPVCATACGC
KIISGPTCPRDYPKLNLLPESGEPNATEYCTIGCRTSVCDNMDNVSRGQEMKFDMGLCSN
ACARFCNDGEVIQSVEA
>TATTH20MR 131 bp
MGGGQKGLESAIVCLLVLGLVLEQVQVEGVDCGANPFKVACFNSCLLGPSTVFQCADFCA
CRLPAGLASVRSSDEPNAIEYCSLGCRSSVCDNMINTADNSTEEMKLYVKRCGVACDSFC
KGDTLLASLDD
>TGTHI13 107 bp
MMVVVILGLVVAQTQVEAKSCCRNTTARNCYNVCRLPGTPRPVCAATCDCKIISSGKCPP
GYEKLGFSDVADEALDVAEEVMKEAVERCNNACSEVCTKGSYAVVTA
```

Notes:
-c2- tells the cut command to print all columns after and including column 2.
-d " " tells the cut command to use a blank space as a field delimiter

8. (5 points) It has been demonstrated that it is impractical to extend the dynamic programming algorithm of Needleman and Wunsch/Smith Waterman to constructing multiple alignments for more than a few sequences. If there are k sequences of length n, give an equation that tells how does the problem scales, in terms of n and k?

9. (15 points) A number of alternative genetic codes have been discovered. Examples are found in mitochondria, plastids, bacteria and archaea. In all of the alternative genetic codes seen so far, most of the codons code for the same amino acids as in the Standard Genetic Code, with a few codons differing. For example in some cases, a stop codon codes for an amino acid, or a codon for an amino acid is used as a stop codon. In other cases, one or two codons are reassigned to a different amino acid.

| Type of search | NCBI | FASTA |
|---|---|---|
| a)DNA vs. DNA database | blastn | fasta3<br>ssearch3 (slow, full Smith-Waterman alignment) |
| b) protein vs. protein database | blastp | fasta3<br>ssearch3 (slow, full Smith-Waterman alignment) |
| c) protein vs. translated DNA database | tblastn | tfasta3 |
| d) translated DNA vs. translated DNA database | tblastx | tfastx3, tfasty3 |
| e) translated DNA vs. protein database | blastx | fastx3, fasty3 (especially well-suited for cDNAs, which often contain frameshift errors) |

Yeast mitochondria use a non-starndard genetic code. Suppose you had the seqeunces for a yeast mitochondrial gene, and its corresponding protein, and wished to find homolgues in other species. How would the difference in genetic codes affect each of the types of searches listed above?

10. (5 points) Sequence database search programs such as the FASTA and BLAST family of programs do not read database files that include annotation for each sequence, as would be found in GenBank or Uniprot entries. Rather, they read files in FASTA or similar formats, which include just a name and definition for each sequence, along with each sequence itself. What is the advantage, when doing a sequence database search, of eliminating the annotation?

11. (10 points) Two antifreeze proteins were aligned using both GGLSEARCH and GLSEARCH.

a) Which of the two alignments is deemed to be more statistically significant?

b) Why does the GGSEARCH alignment have a long gap, followed by a phenylalanine (F) at the end of ISP2_H? How does that gap contribute to the difference in Needleman-Wunsch (n-w) scores?


## GGSEARCH

```
Algorithm: Global/Global affine Needleman-Wunsch (SSE2, Michael Farrar 2010) (6.0 April 2007)
Parameters: BL62 matrix (11:-4), open/ext: -11/-1


>>ISP2_OSMMO 175 bp                                    (175 aa)
 n-w opt: 315  Z-score: 295.7  bits: 61.1 E(1): 1.3e-133
global/global (N-W) score: 315; 39.1% identity (65.4% similar) in 179 aa overlap (1-163:1-175)


                10        20        30        40        50
ISP2_H MLTVSLLVCAMMALTQA-NDDKILKGTATEAGPVSQRAPPNCPAGWQPLGDRCIYYETTA
         ::. .::::::.:::.: : :     ... : .    : .   .::. :. .. ::. ..
ISP2_O MLA-ALLVCAMVALTRAANGDTGKEAVMTGS---SGKNLTECPTDWKMFNGRCFLFNPLQ
                10        20        30            40        50

       60        70        80        90       100       110
ISP2_H MTWALAETNCMKLGGHLASIHSQEEHSFIQTLN-AGVV--WIGGSACLQAGAWTWSDGTP
         . :: :. .::: :..::::::: ::..:.. :. ::..  ::::: :   .   : : :.:
ISP2_O LHWAHAQISCMKDGANLASIHSLEEYAFVKELTTAGLIPAWIGGSDCHVSTYWFWMDSTS
          60        70        80        90       100       110

        120       130       140       150       160
ISP2_H MNFRSWCSTKPDDVLAACCMQMTAAADQCWDDLPCPASHKSVCAMT------------F
         :.: .::...:: .:. ::.:..... .::.: ::    : ::::             .
ISP2_O MDFTDWCAAQPDFTLTECCIQINVGVGKCWNDTPCTHLHASVCAKPATVIPEVTPPSIM
          120       130       140       150       160       170
```


## GLSEARCH

```
Algorithm: Global/Local affine Needleman-Wunsch (SSE2, Michael Farrar 2010) (6.0 April 2007)
Parameters: BL62 matrix (11:-4), open/ext: -11/-1


>>ISP2_OSMMO 175 bp                                    (175 aa)
 n-w opt: 336  Z-score: 328.6  bits: 67.2 E(1): 4e-171
global/local score: 336; 41.9% identity (69.5% similar) in 167 aa overlap (1-163:1-163)


                10        20        30        40        50
ISP2_H MLTVSLLVCAMMALTQA-NDDKILKGTATEAGPVSQRAPPNCPAGWQPLGDRCIYYETTA
         ::. .::::::.:::.: : :     ... : .    : .   .::. :. .. ::. ..
ISP2_O MLA-ALLVCAMVALTRAANGDTGKEAVMTGS---SGKNLTECPTDWKMFNGRCFLFNPLQ
                10        20        30            40        50

       60        70        80        90       100       110
ISP2_H MTWALAETNCMKLGGHLASIHSQEEHSFIQTLN-AGVV--WIGGSACLQAGAWTWSDGTP
         . :: :. .::: :..::::::: ::..:.. :. ::..  ::::: :   .   : : :.:
ISP2_O LHWAHAQISCMKDGANLASIHSLEEYAFVKELTTAGLIPAWIGGSDCHVSTYWFWMDSTS
          60        70        80        90       100       110

        120       130       140       150       160
ISP2_H MNFRSWCSTKPDDVLAACCMQMTAAADQCWDDLPCPASHKSVCAMTF
         :.: .::...:: .:. ::.:..... .::.: ::    : ::::
ISP2_O MDFTDWCAAQPDFTLTECCIQINVGVGKCWNDTPCTHLHASVCAKPATVIPEVTPPSIM
          120       130       140       150       160       170
```

12. (10 points) Draw a dot-matrix plot (eg. DXHOM) of this sequence, compared with itself. You can assume that the only repeats of any significance are the ones documented in this GenBank entry.

```
LOCUS       AY966016                 380 bp    DNA     linear   PLN 14-NOV-2005
DEFINITION  Aspergillus flavus isolate NPL GA3-3 hexose transporter-like (hexA)
            gene/telomere breakpoint junction.
ACCESSION   AY966016
VERSION     AY966016.1  GI:67944627
KEYWORDS    .
SOURCE      Aspergillus flavus
  ORGANISM  Aspergillus flavus
            Eukaryota; Fungi; Dikarya; Ascomycota; Pezizomycotina;
            Eurotiomycetes; Eurotiomycetidae; Eurotiales; Aspergillaceae;
            Aspergillus.
REFERENCE   1  (bases 1 to 380)
  AUTHORS   Chang,P.K., Horn,B.W. and Dorner,J.W.
  TITLE     Sequence breakpoints in the aflatoxin biosynthesis gene cluster and
            flanking regions in nonaflatoxigenic Aspergillus flavus isolates
  JOURNAL   Fungal Genet. Biol. 42 (11), 914-923 (2005)
   PUBMED   16154781
REFERENCE   2  (bases 1 to 380)
  AUTHORS   Chang,P.-K.
  TITLE     Direct Submission
  JOURNAL   Submitted (17-MAR-2005) Food and Feed Safety, Southern Regional
            Research Center, 1100 Robert E. Lee Boulevard, New Orleans, LA
            70124, USA
FEATURES             Location/Qualifiers
     source          1..380
                     /organism="Aspergillus flavus"
                     /mol_type="genomic DNA"
                     /isolate="NPL GA3-3"
                     /db_xref="taxon:5059"
                     /note="type: L"
     gene            62..244
                     /gene="hexA"
     misc_feature    62..244
                     /gene="hexA"
                     /note="similar to hexose transporter"
     misc_recomb     244^245
                     /gene="hexA"
                     /note="hexA-telomere breakpoint junction; recombination
                     results in deletion of aflatoxin gene cluster"
     repeat_region   245..376
                     /note="telomeric repeat"
                     /rpt_unit_seq="tcaacattaggg"
ORIGIN
        1 gtctttcccg ccaacttgaa gtccagcagt atccttaaca gtaccctttg ttactgacac
       61 catggttgct ggcggtggag ttgttccttc atccggtatg gatgcatacc gggccctgcc
      121 aaacaatacg aactcgaact ggttcaagga caagggcctc cggcgtctga atttcggcct
      181 catgcttatg tttgcatccg ctgcagcaaa tgggtatgat ggggctttga tgaatgggct
      241 cctgtcaaca ttagggtcaa cattagggtc aacattaggg tcaacattag ggtcaacatt
      301 agggtcaaca ttagggtcaa cattagggtc aacattaggg tcaacattag ggtcaacatt
      361 agggtcaaca ttagggtcaa
//
```

The IUPAC-IUB symbols for nucleotide nomenclature [Cornish-Bowden (1985)Nucl. Acids Res. 13: 3021-3030.] are shown below:

| Symbol | Meaning | Symbol | Meaning |
|--------|---------|--------|---------|
| G | Guanine | K | G or T |
| A | Adenine | S | G or C |
| C | Cytosine | W | A or T |
| T | Thymine | H | A or C or T |
| U | Uracil | B | G or T or C |
| R | Purine (A or G) | V | G or C or A |
| Y | Pyrimidine (C or T) | D | G or T or A |
| M | A or C | N | G or A or T or C |

**The Universal Genetic Code**

| | | | | | | | |
|--|--|--|--|--|--|--|--|
| UUU | phe | UCU | ser | UAU | tyr | UGU | cys |
| UUC | | UCC | | UAC | | UGC | |
| UUA | leu | UCA | | UAA | stop | UGA | stop |
| UUG | | UCG | | UAG | stop | UGG | trp |
| CUU | leu | CCU | pro | CAU | his | CGU | arg |
| CUC | | CCC | | CAC | | CGC | |
| CUA | | CCA | | CAA | gln | CGA | |
| CUG | | CCG | | CAG | | CGG | |
| AUU | ile | ACU | thr | AAU | asn | AGU | ser |
| AUC | | ACC | | AAC | | AGC | |
| AUA | | ACA | | AAA | lys | AGA | arg |
| AUG | met | ACG | | AAG | | AGG | |
| GUU | val | GCU | ala | GAU | asp | GGU | gly |
| GUC | | GCC | | GAC | | GGC | |
| GUA | | GCA | | GAA | glu | GGA | |
| GUG | | GCG | | GAG | | GGG | |

| 3-letter | 1-letter | 3-letter | 1-letter | 3-letter | 1-letter |
|----------|----------|----------|----------|----------|----------|
| Phe | F | Leu | L | Ile | I |
| Met | M | Val | V | Ser | S |
| Pro | P | Thr | T | Ala | A |
| Tyr | Y | His | H | Gln | Q |
| Asn | N | Lys | K | Asp | D |
| Glu | E | Cys | C | Trp | W |
| Arg | R | Gly | G | STOP | * |
| Asx | B | Glx | Z | UNKNOWN | X |
| Xle (Leu/Ile) | J | Pyl (pyrrolysine) | O | | |