

# På genjagt i databaser

A  
B  
C



## Formål

I dette teoretiske eksperiment skal I gå på jagt i databaser på nettet. I får udleveret forskellige DNA- og aminosyre-sekvenser, som I skal forsøge at finde ud af, hvor stammer fra. Altså hvilke organismer kommer de fra, og hvilket protein kodes der for?

## Teori

Der findes mange hjemmesider med informationer om proteiner. En af disse hjemmesider er Protein Data Bank (PDB), der findes her:

**<http://www.pdb.org/pdb/home/home.do>**

På hjemmesiden kan man finde både DNA- og aminosyre-sekvenser, ligesom man kan finde 3D-modeller af alverdens proteiner.

På bilag 1 findes DNA-sekvenser opskrevet i et bestemt format, der kaldes FASTA. Dette format kan bruges i specielle programmer til at lave alignments med allerede kendte DNA-sekvenser. Derved kan man måske finde et 100 % match og dermed identificere ens ukendte DNA-sekvens.

Samme FASTA-format kan benyttes ved aminosyre-sekvenser, hvor man med de samme programmer laver alignments med kendte aminosyre-sekvenser. Også her kan man måske finde et match og derved identificere, hvilket protein der er tale om. Disse data findes på bilag 2.

Til at lave DNA-alignments bruges et program, der hedder BLAST, der står for Basic Local Alignment Search Tool. Til at lave aminosyre-alignments bruges et tilsvarende program, der hedder BLASTP. Forkortelsen står for det samme, og det sidste P står for Protein.

Begge programmer findes på denne hjemmeside:

**<http://blast.ncbi.nlm.nih.gov/Blast.cgi>**

Når man har en DNA-sekvens, kan man translaterer den helt automatisk, og dermed finde ud af hvilke aminosyrer sekvensen koder for. Det kræver dog, at man kender læserammen. Det er ikke ligegyldigt, hvor den ligger. Hjem-

mesiden Virtual Ribosome kan netop foretage sådanne translationer, hvor det undersøger alle mulige læserammer. Her er vist et eksempel på en sekvens: CATGCGA

Der er umiddelbart tre mulige læserammer:

CAT GCG A

C ATG CGA

CA TGC GA

Men faktisk kan der også læses fra den anden side, så der kommer yderligere tre muligheder:

AGC GTA C

A GCG TAC

AG CGT AC

Der findes altså teoretisk 6 mulige læserammer på en hvilken som helst DNA-sekvens. "Virtual Ribosome 2.0" kan finde alle seks og angive startkoder, stopkoder og hvilke aminosyrer, de enkelte kodons/tripletter koder for. Herved kan man finde den mest sandsynlige læseramme, som er den, der indeholder en startkode, et ukendt antal kodons og en stopkode:

**<https://services.healthtech.dtu.dk/service.php?VirtualRibosome-2.0>**

Yderligere information om bioinformatik og metoder kan findes på biotech-academy's hjemmeside:

**<https://www.biotechacademy.dk/undervisning/gymnasiale-projekter/bioinformatik-intro/>**

## Materialer

Hjemmesiderne nævnt i teori afsnittet samt data på de to bilag.

# Fremgangsmåde

## Undersøgelse af DNA-sekvenser

Åbn hjemmesiden hvor du kan udføre en BLAST-analyse.

Vælg "nucleotide blast" ovre til venstre på skærmen.

Kopier den første DNA-sekvens fra bilag 1 – sørg for kun at kopiere baserne.

Indsæt baserne som vist nedenfor (gul farve).

Tryk BLAST (blå knap nederst på skærmen).

Nu foretages et BLAST af sekvensen, hvor den sammenlignes med alle kendte DNA-sekvenser. Efter noget tid (der kan gå op til flere minutter) fås et resultat.

(se næste side)

The screenshot shows the NCBI BLAST Standard Nucleotide BLAST interface. The 'Enter Query Sequence' field contains a DNA sequence highlighted in yellow: `CCCAGGCAAAACAGGCCGACTGGCAGCCTCCTGGCTGCACTCCCACTGTTCTGAACAGCTG  
AGGGAACAGGGCAGCTGCTGCTGGGCTCACCAGGAAACTGACAGACTTACCTGAGTC  
ACCTGACACTGACTCCGACAGAATCTACTTTTCTTCATCTATCCTTGTATATTTA  
AACAACTACTACCCAAAAAATCTGGCCGCTACTTCTACTGACTTAAGCAAAAGTCATCT  
TTGATTACATAATTTTTTAATGAATAAGAAAGCTAAACAAGTTT`. The 'Choose Search Set' section shows 'Database' set to 'Nucleotide collection (nr/nr)', 'Organism' set to 'Human genomic + transcript', and 'Exclude' options for 'Models (KOMP)', 'Uncultured/environmental sample sequences', and 'Sequences from type material'. The 'Program Selection' section shows 'Optimize for' set to 'Highly similar sequences (megablast)'. The 'BLAST' button is visible at the bottom.

Hvis du scroller ned på skærmen fås en highscoreliste over de bedste alignments ud fra DNA-sekvensen. Længst ude til højre på listen ses en %-sats (se den gule markering). 100 % betyder, at der er fuld alignment med den pågældende sekvens. Hvis du trykker på den pågældende organisme kommer alignmentet frem (se blå markering).

Nu kommer en ny skærm frem, og helt ovre til højre står der "Related information" - tryk på "Gene".

Nu kommer en lang række oplysninger om genet.

Noter hvilken organisme genet kommer fra.

Hvad koder genet for?

Under overskriften "Genomic context" kan du se, hvor mange exons, genet indeholder. Her ses det også, hvilket kromosom, genet er placeret på.

Noter resultaterne i resultatskemaet.

Udfør en BLAST-analyse af DNA-sekvens nr. 2 fra bilag 1

(kopier den ind ad 2 omgange).

Noter hvilken organisme genet kommer fra.

Noter også, hvilket protein genet koder for.

Hvor mange exons består genet af?

Hvilket kromosom er det placeret på?

The image shows a BLAST search interface. At the top, there is a visualization of sequence alignments represented by red horizontal bars of varying lengths. Below this, the "Descriptions" section is visible, showing a table of sequences producing significant alignments. The table has columns for Description, Max score, Total score, Query cover, E value, Ident, and Accession. The first row is highlighted with a blue box, and the "Ident" column for the first few rows is highlighted with a yellow box.

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input checked="" type="checkbox"/> Homo sapiens Na,K-ATPase beta 2 subunit gene, complete cds	14578	14578	100%	0.0	100%	AF007876.1
<input type="checkbox"/> Homo sapiens chromosome 17, clone RP11-199F11, complete sequence	9886	17447	92%	0.0	99%	AC067388.9
<input type="checkbox"/> Homo sapiens chromosome 17, clone RP5-1030014, complete sequence	9873	16568	92%	0.0	99%	AC007421.13
<input type="checkbox"/> Homo sapiens ATPase Na+K+ transporting subunit beta 2 (ATP1B2), transcript variant 2, mRNA	3707	4585	31%	0.0	99%	NM_001303263.1
<input type="checkbox"/> Homo sapiens ATPase Na+K+ transporting subunit beta 2 (ATP1B2), transcript variant 1, mRNA	3707	6155	42%	0.0	99%	NM_001679.4
<input type="checkbox"/> PREDICTED: Pan paniscus ATPase, Na+K+ transporting, beta 2 polypeptide (ATP1B2), mRNA	3465	6892	50%	0.0	97%	XM_028652431.2
<input type="checkbox"/> PREDICTED: Nomascus leucogenys ATPase, Na+K+ transporting, beta 2 polypeptide (ATP1B2), transcript variant X2, mRNA	3286	4136	31%	0.0	96%	XM_012001454.1
<input type="checkbox"/> PREDICTED: Nomascus leucogenys ATPase, Na+K+ transporting, beta 2 polypeptide (ATP1B2), transcript variant X1, mRNA	3286	5582	42%	0.0	96%	XM_003274632.2
<input type="checkbox"/> Human sodium/potassium ATPase beta-2 subunit (atp2b2) mRNA, complete cds	2998	5052	35%	0.0	99%	M81181.1
<input type="checkbox"/> Human Na,K-ATPase beta 2 subunit (ATP1B2) mRNA, complete cds	2983	4382	30%	0.0	99%	U45945.1
<input type="checkbox"/> Homo sapiens ATPase, Na+K+ transporting, beta 2 polypeptide, mRNA (cDNA clone MGC:161453 IMAGE:8951891), complete cds	2789	4737	32%	0.0	99%	BC126175.1
<input type="checkbox"/> PREDICTED: Samiri boliviensis boliviensis ATPase, Na+K+ transporting, beta 2 polypeptide (ATP1B2), mRNA	2422	4541	42%	0.0	88%	XM_003629189.2

### Undersøgelse af aminosyre-sekvenser

Fremgangsmåden er her den samme som ovenfor, men nu skal der vælges en BLASTP-analyse i stedet for. Det gøres ved at vælge "protein blast" på forsiden. Kopier nu den første aminosyresekvens (bilag 2) ind i feltet som før og tryk BLAST. Ved hver af de fire sekvenser besvares følgende:

- hvilken organisme kommer proteinet fra
- hvilket protein er det?


## Virtual Ribosome

Nu skal vi analysere den første ukendte DNA-sekvens nærmere (nu er den ikke mere ukendt). Kopier sekvensen og åbn Virtual Ribosome. Indsæt den i det øverste felt som vist nedenfor. Midt på skærmen skal I vælge reading frame – altså fra hvilket nukleotid, læserammen begynder. Det er her, der er 6 muligheder. Vælg læseramme "1". Ved "ORF finder" vælges "start codon: Strict". Det betyder, at programmet nu leder efter den normale startkode, som er ATG. Vælg "Submit query" nederst.

### Virtual Ribosome - version 1.1

The Virtual Ribosome is a comprehensive tool for translating DNA sequences to the corresponding peptide sequences.

Besides being a strong translation tool in its own right (with an integrated ORF finder, support for all translation tables defined by the NCBI taxonomy group, and a number of options regarding START and STOP codons), the Virtual Ribosome can work directly on files containing annotation of gene structure. This makes it easy to map various aspects of Intron/Exon structure onto the translated sequence.



[Instructions](#) | [Output format](#) | [Software download](#) | [Article abstract](#)

Paste in DNA sequences in FASTA, GenBank or TAB format

```
CACTGTAACCTACTCTATGGAGGCTATCTTGGCAATACCTTGGCTTCTTGAATCTGGCTCTTGTCTCTGTGGTC  
CAGACTTAATGACCGCCCTCTACACCACTGGCCCTGCCCTCACAGCTTCTGGGGGCAAGCCGAGCCCAAGCCAC  
GCCCGACTGGCAGCTCTTGGCTGGCACTGGCACTGTTCTGACAGCTGAGGGAACTGGGGCACTGGCTGGGGC  
CTCAGCAGGGAACTGAGAGACTTACCTGACTGACTGACACTGACTTCCAGAGAACTACTCTTGGCTTCTG  
ATACTTATCTTCTTATATTTAAACAACTACTACCGAAAAATCTGGCCCTACTTCTACTGACTCTAAGGAA  
ACTCATCTTTGATTACATAATTTTTTAATCAATAAGCAAAAGCTAAACAACTTT
```

Upload DNA sequences in FASTA, GenBank or TAB format

Instructions: Basic usage - Paste in or upload one or more DNA sequences in FASTA (sequence only), GenBank (CDS sections are processed) or TAB (sequence and intron/exon annotation) format and hit submit. The Virtual Ribosome will then translate the DNA sequences using the standard genetic code (by default). Options can be customized in the section below.

#### Options

**Reading frame**  
  
Reading frame selection is ignored if Intron/Exon annotation is supplied

**ORF finder**  
Start codon:   
Scan the reading frame(s) specified above for ORFs

**Gene structure annotation**  
 Exon numbering  
 Intron position and phase  
Only for relevant for GenBank and TAB files

**Translation table:**

**Start codons**  
 First codon is start codon  
 All codons are internal  
Alternative (non-methionine) start codons will code for methionine if used as a start codon.

**Stop codons**  
 Read through  
 Terminate  
If multiple reading frames are selected (e.g. "all") this option will always be set to "Read through"

Resultaterne ser ud som vist nedenfor:  
 (ude til højre kan ses signaturforklaringer)

```

VIRTUAL RIBOSOME
-----
Translation table: Standard SGC0
>Seq1_ORF
Reading frame: 1

5' GAACCGGGGGTCCCTGTTGGCCAAAGGACCTGACCAAGACCTGTGAATTCATCGTTCACTCTGCCCCCTCAGAAGGGGAGTGGACCCCAAG 90
.....)}}***.....

5' TCCTGTGGACTTCAGGATCACTCCGGAAACCTTCAGAACCTCAAGAGAGAGCTTCACTTCCCAAATTCATTAGAGGACATCTGAA 180
.....***.....

5' CTCCTAACTAAGGCTATCAGCCAGCCCTGACAGGAGAGCTGGTGGTGGAACTCCGATGCTGCTATCCGGAGCATAGAGCTTCACTT 270
.....***}}).....

5' GGTCCGGGTGGAGACCTGTGGGTGTGCAGAAAGCTATGCACGAGATGCCACAGAGATTCAGAAATATTCAAATCGCTGATGGGACATTG 360
.....)}}).....

5' TAGAAACCTTCTGTCCCTGTACATGGTCTTCCCAAGCTGTTCACTGTCCAACTGGAGACCACCAACTTCAAAGTGGASTTTGA 450
***.....)}}).....

5' GGTCAATGTGGTGGTCCCTCTCATGCTGACCACTCATCACTGAGAACTTCCGCTGAAGCTCTGTGGACATAGTCCAGCTGCAGGAA 540
.....***.....)}}).....

5' GGGAGAGCCAGAAATGGCCATGTGGACATTTGACCCAGTCACTGCTCAGATGGGGACCTCACACAGACCCTGCTCACAGGCATCACAT 630
.....

                    H A S K L V P E R C L T S A L T S G H Q Q
5' CAGCAGCATAACATCTTCATGCTCTTCATGCCCTCAAAGCTTGTTCCTCCAGGTGCCCTTACCTCAGCCCTCACGCTCGGAATGCCAACAG 720
.....>>>.....>>>.....

D F S C V L Q D H R K R S Q W L Q M C L F L L R F I K Q C V
5' GATTTCTCCTGTGTGTGCAAGATCACAGAAAGCGTCCCAATGGCTCCAGATGTGCCCTGTTCTTCTTAGGTTTCATTAAAGCAGTGTGTC 810
.....)}}).....>>>.....)}}).....

C T L R W A A
5' TGCACATTAAGATGGCTGCCTAGTGAATCATGGAGCCAGAGAGAGCCCCACCCCTCTATCTTCCCAAGACTCCTGACAGCAAGAG 900
.....*****.....
  
```

>>>	Startkode (den ene vej)
<<<	Startkode (den anden vej)
*****	Stopkode
))) og (((	Alt. Startkoder (irrelevante)



Som det ses, er en stor del af nukleotiderne ikke kodende, mens kun en mindre del er kodende. De kodende nukleotider har vist aminosyrekoden ovenover sig.

Vælg FASTA øverst på skærmen. Herved får du opsummeret, hvilke aminosyrer genet koder for. Hvor mange kodes for?

Lav nu endnu en analyse af samme sekvens, men denne gang ændrer du læserammen til 2. Optæl antallet af aminosyrer, genet koder for, hvis læserammen er sat til 2.

Gør det samme med både læseramme 3, -1, -2 og -3.

Hvilken læseramme giver det største protein?  
Antagelsen er, at denne læseramme er den korrekte.

Foretag til sidst en analyse, hvor du vælger "Alle (6 reading frames)". Husk "ORF finder" skal stadig være "start codon: Strict".

Nu fremkommer den bedste løsning – den længste aminosyresekvens. Var det den samme reading frame, som du fandt (det står noteret øverst på skærmen)?

Hvis der er tid, så udfør lignende analyser på DNA-sekvens 2.

# RESULTATER

	Organisme	Protein	Antal Exons	Kromosom-nummer	Bedste læseramme
DNA-sekvens 1					
DNA-sekvens 2					
Aminosyresekvens 1					
Aminosyresekvens 1					
Aminosyresekvens 1					
Aminosyresekvens 1					

## Fejlkilder

### **Diskussion**

1. Forklar hvad der forstås ved begrebet læseramme, og forklar hvorfor der er 6 mulige læserammer ved ethvert stykke DNA.
2. Hvilke forskelle er der mellem prokaryoters og eukaryoters DNA?
3. Forklar kort, hvordan DNA bliver til protein hos prokaryoter og eukaryoter.
4. Hvorfor er nogle substitutionsmutationer mere alvorlige end andre?  
Giv konkrete eksempler.
5. Hvorfor har det stor betydning for proteinets struktur, hvis fx cystein udskiftes med en anden aminosyre som følge af en mutation?
6. Hvordan kan BLAST-analyser bruges til slægtskabsanalyser?

## BILAG 1: Ukendte DNA-sekvenser

### Ukendt sekvens nr. 1 (FASTA-format)

```
GAACCGGGTCACGAGTGGCCGGAAGTTGCGGCCTTTGCGGATAATCTGAGATGGGGACCACTCTGGACATCAAGATTAAGAGCGAATAAAGTGTATCACGCCGGGGAAAT-  
GCTCTCTGGTGTGGTGGTGCATCTCTAGCAAGGACTCAGTCCAGCACCAGGGAGTCTCTTTGACAATGGAAGGGACTGTAAACCTCCAGCTCAGTGCCAAAAGTGTGGGTGT-  
GTTCGAAGCCTTTTACAACCTCTGTGAAGCCGATCCAGATTATCAACAGCACCATAGACGTGCTGAAGCCAGGAAAGATTCCCAGTGGCAAGACTGAGGTCCCCTTTGAGTTC-  
CCACTGCTTGTGAAAGGCAGCAAGGTCCTGTATGAGACTTACCACGGCGTGTGTGCAACATTACAGTACACACTGCGCTGTGACATGCGGCGGTCCCTGTTGGCCAAGGAC-  
CTGACCAAGACCTGTGAATTCATCGTTCCTCTGCCCCTCAGAAGGGGAAGTTGACCCCAAGTCCTGTTGACTTCACGATCACTCCGGAAACCTTGCAGAACGTCAAAGAGA-  
GAGCTTCACTTCCCAAATTCTTCATTAGAGGACATCTGAACTCCACTAACTGCGCTATCACGCAGCCCCTGACAGGAGAGCTGGTGGTGGAACTCGGATGCTGCTATCCG-  
GAGCATAGAGCTTCAGTTGGTCCGGGTGGAGACCTGTGGGTGTGCAGAAGGCTATGCACGAGATGCCACAGAGATTGAGAATATTCAAATCGCTGATGGGGACATTTGTA-  
GAAACCTTTCTGTTCCCCTGTACATGGTCTTCCCAGACTGTTACCTGTCCAACACTGGAGACCACCAACTTCAAAGTGGAGTTTGAGGTCAATGTGGTGGTCCCTCCTTCAT-  
GCTGACCACCTCATCACTGAGAACTTTCCGCTGAAGCTCTGTGCGGACATAGTCCAGCTGCAGGAAGGGGAGAGCCAGAATGGCCATGTGGACATTGACCCAGTCATCTGC-  
TCACATGGGGACGTCACACAGACCACCTGCTCACAGGCATCACATCAGCAGCATAACATCTTCATGTCTTCTCATGGCCTCAAAGCTTGTCCCTGCAGGTGCCTTACCTCAGC-  
CCTCACGTCTGGAATGCAACAGGATTTCTCCTGTGTGTTGCAAGATCACAGAAAGCGTTCCCAATGGCTCCAGATGTGCCTGTTTCTTCTTAGGTTTATTAAGCAGTGTGTCTG-  
CACATTAAGATGGGCTGCCTAGTGAGATCATGGAGGCAGAAGAGAGCCCCCACCCTCTATCTTGCCCAGACTCCTGACAGCAAAGAGTCAGCAACCAGCAAGCTCTCTGAT-  
TGAGTCTGTAATCCCCTCCAAGGAGTGACAGTGACTTCATCTGGGAATACGGAGTTTTCCCTTCTTGACATCTTGATAGTGTTCCTAAAGCCTCTAAATGTTTTAGGGCACTG-  
GAAATAAAAACCAACAATAAGTCCTGGTTCTCACAGAGGGGAGAATGGTGAACAGTGGCTGTTGCTGGAGTGTGGTGTGTGGTACTGGCACACAGTTAATGCTGCCATAT-  
TCCCAGAATGTATTGAGAGACAAGGTGCTTGCATAAGGGGAAGAATGCACATGGACTACTCGGTAAGTACTAGGTTCTGAGTGACAGGACACCGTGCTGAGGAAGGGCAATG-  
GCCTTGGCCTGTCCCTGACCCCCTGGGTCACATGGCTTTAGCCTCACACTGTAACCTCATACTGTATGGAGGGTATGTTGCAATACTTTGCTTCCCTTGAATCTCGCTGTTGTTCT-  
GTTGGTCCAGAGTTAATGAGCGCCCTGTACACCAGCTGGCCTGCCTTACAGGTTCTGGGGGGCAGGCCGAGGCAAAGCAGGCCCGACTGGCAGCCTCCTGGCTGCACTCC-  
CACTGTTCTGAACAGCTGAGGGAAGTGGGGCACTGCTGCTGGGGCCTCACCAGGGAACTGAGAAGACTTAGCTGAGTCACCTGACAGTGAGTTCCCAGAGAATGACTTTTT-  
GCTTTGTCATACTTATCCTTGTATATTTAAACAATCTAGTACCCAAAAAATGTGGGCCGTACTTTCTAGTGAGTCTAAGGAAAAGTCATCTTTGATTACATAATTTTTTAATGAATA-  
AAGAAAGCTAAACAAGTTT
```

**Ukendt sekvens nr. 2 (FASTA-format)**

TCCCGACCCTAGAGTCCCGTCACGACCCCTGACCCTTACACCACAACTCTCCCGAAGTCCCTCTGCACTACCCTTCTACCCCTTCGGAGACATCCACCTGTCTCGGTCCAC-  
CACACCTGTCCCCGACACTCTAACCTCTTCCCTCCTCAGACCCCTGACCATCGCCGAGGTTGTCCCTTGACAGGAGACCCCCACGGGAAATAGGTAGAGACGCCCCGGACAG-  
GATCGAAGGGGAGGACCGAGGCCGCTCAGGGCCGACCAGAGAAGCGAGAGAAAGAAGTCTAGATCGATCTCTAGCTCTTTATTCTCTGACATTCTGCCCATCTGCTCCCG-  
GACTCTTCCGCCGCATCTGTATCTTAGTCTCTTTCTAACTCCCTGCCTGGGGCCCCACCTTTGAGCATATATCGGCTTCTCTCTGCCGTTGTGTTTCTCTGGGGTTATCTTTCC-  
CTCTACTCCGCCCCGACACGCCCTCTATCTTCTCCACCCCTGTCCCAGTCTAGTACCTGGGGTGGGGGTAATGGAGAGACAGCTAGGTCCTAAGAGAAGTCAGGGGG-  
GCTGGGCCAACCTCTCCATATTTACATATGTATGAGGGTCGCCTGGGCCAGTGGCGAGGAGGCGGACGTTCTGGGGGTGGGAAGGGGGCGGGCACCCCCAGAGCCGCA-  
GAGTATAAAGACCGCGCTCGGCGACCGCGGGCCCCGACTGCTGAGGAGCGGACGCTCCGCCTGGGGGGCCCCCATCCCTGGCTGTCCCCAGCTGCGCGTCCCCGCC-  
CCACCCCCGCGGCTGAGCCACCACCGGTGCAGTGGTCTCCGCTTGGCGGAGCGAGCCTTGAGCTTCGTTCCACAGCTTCTTTGCATCTTGATTTCCGGGGCGGCCCCCTCC-  
CCCACCTCTCTGCCTTTTTGTACCCCGCTTTTTTCTGCGTTCTGCTCGTTTTTGTAGCCGTCTGTTTTGCACCCCATTTCTTTTTGTTTCTAGACGGTTTGGCGGGGGGT-  
GAAGCTGCATTCATAACCCCTTCTCTTGTATTCTCCCCTGCTCTGACAGCACCCCTTTTCATCGCAGTTGGGGGGCCTAGGATCGGTGCATCTTCCGCCGCGCTGCCAGCAC-  
CCCGCAGCGCGTGGTCGTGCACCCCGGAATCTGCAGCAGCTGCATATCTGAGGGGGTCTCCTTTGCCCGCGCCGCTTCGCTCCCCGTGCTTTGGGTGTGTGGAGGGCT-  
TCAGTCGCGGCGCCCCGCTTCTCCGCAACCCCCCGCCCCGCGCCCGGACTCGCCCCGCGCCACCAAGATGGTCATCCAGAAAGAGAAGAAGAGCTGCGGGCAGGTGGTT-  
GAGGAGTGGAAGGAGTTCGTGTGGAACCCGAGGACGCACCAGTTTATGGGCCGCACCGGGACCAGCTGGGGTACGCAGGGCCGGCACGCAAGGGGGCGGGGGAAAGCCG-  
CGGGGCGACGCCTCGGGGGCGCAGGGTCCCGCCGACGCGCCCCAGCTCCCCTCCCGGGTCCCGGCGTCCAGCTCCCTGCCGGGCTCTGGGCTGGGAGGGGGCCGAATC-  
GCCAGTCTAACTCCCCGGCTGGCCGTGCGGAGGCGGAGAAAGTAGGTCACAGCCGCCTTCCGCCCCCCGCGGAGCCCCCTCGGGCGGCGGGGTGCGCAGCTCCGCCT-  
GCGTGCCGCGCCGCGGCTCACACTCCCCTCTCGGGGCTGTGCTCCACACGGGCGTCCCCACCTCCAAGAGCGCCCCCTTCCCTCCCTCCGGCTCTACTAGCTCCG-  
CAGCCCCGTCTATTTTTAGCTCGTGCCACCCCCCTGGACCCTGGGAACGTTTCATGAGGGGGCGGGTCTTGGGGGTGTGTTAGGGGGGTTCTTCACGGCGGAAGTTGTCT-  
GTATCCCACCGCCTGGCCTTGGGAGCCTTCTGGGACTGCTTTGTGGTGGGGGGCTGCTGATAGTATGAGTTTTACCGAGGCTGCAGGTTTTAGCCTCCCATGTGCGGTGACG-  
GAGGGAGGAGTGGTCGCTGTGGTGAATTTGTGTGCATCAGCCAGCCAGGTGTCTGTGACAGTCGGATGACTTGGAAGCCTCCCAGGCTGACCATGGCAGGACTCAGGGAG-  
CTGTAGTGGTCGGGGGTTGGGGGGGTGGAGGGGGTCTGGTGACCGGCACAGGTGCAGGTGAGGGGTGGAATTCATTTACATTTTCCCATAGAAACAAAGTTATAAATA-  
GTGACTGCATACTGCACCTAAAATGCCACGTATCTACAACGAAATTATTACAATGTATTTATAATATATTATTCTAAACATGTATGCATTTGAAATTATCCAATGAAATCTAAC-  
CTATTGATGTACCTATTTCTACATTAATACATTATATCATATATTAATATGCTAATAATTATATCACATTAATATATTAATAATATGTTAATTAACCAAACGAAAAGTATATACTGATCA-  
CAACACTTGAATTCCTCCAAGAGGGATGGTCAAGCTGGGACGTTGAGACACAGGGGACAGAGGACACTGTGTGACACGATTTACAATCTTTCCACACTGGGCACCGTCCCCA-  
TCAGTCCACCCATTCGGGGCCTACACGAAGTGGGTCCCATGCAATCCATTCCCTCAGGGAACTCAAACCTCCAGCCCCTGGGATGAGAAGAATCCAGCAATGCTTGGGAGAG-  
CCAGAGGACTTCATGGAAGAAGTGTCTCTGAGATGGAAGGATTGGGAGTCCAGGGTGGTGGGAACAGCCGGCCCTTGGGTCTTACTTCAGGCGGGGGAGCCATGGAGA-  
GATCCCACCAAGGGAAGGCTGTGGAGATTCTGCCTTCTCCCTGCCTCTGCCAGGGTGTGGGTGTGAACTGAGGGTGGGGTGAAGTTCTAACAAGCCGTC-  
TCTGAGAGATTTGTAGCTAGGCTAGTGTTAGGTCTTTCATTTCCAGGAACTGTGTTCAAAGTTTGGCTTCTGAAGGGCACCAAGGAGAGAGATGTTGCTATTCAAATCTGAGGGT-  
CCAGTCTCTGCGGGGTGGTATGAGGGTTTGTGTAATGGTGGCCAGTACCCGCTTTAAAAGGCACCATGCTAGCACAGCTTTAAGCATGAGTACGAATGCAGAGGTAACA-  
GATGTGTGCCCTTGTGAGGACTATGCATGGTTGAGAAGTTGGAATGTAATTGGAGGCAAATAACAGACCTCCACAAGGTCGGGCTTCACTGTGCCCTAGGACCAGGAGGGG-  
GCTGGGAGTCATGGCTAGAAGCCAGACACAACTGCCTGTTTCCAGTTTGTCTCATTTTGCCTCCAGAGGAAGGCTCTAAGACATCCCTGTGGCTCTGTGATCAGTCCCAGTGCA-  
GAACTTCAGAGTGGGTAGAGGGGTGTGTGGGGATAGTTGAGGTTATGGTGGGAACCTTGGGCCCTGCTGACCCTGTTTCTCCTCCCTAGCCTTTATCCTCCTCTTCTACCTC-  
GTTTTTTATGGGTTCTCACCGCCATGTTCCACCCTACCATGTGGGTGATGCTGCAGACTGTCTCCGACCATAACCCCAAGTACCAGGACCGACTGGCCACACCGGGTGAAGT-  
GTGGAGGCTCCCCCTGCCAGCTACTCTAACTGCTCTTGTGCCCCCAAACCTCCAGAAGGAACTCATAGTTCCTTCCAGGAGTTTGATTTTGTGATGACCCAATCCCCACGTGCT-  
TGGAAGTTCTTGAATCTGTCCACCTTCCATTTACTGCAGTTGGGAGCTGTGTGATTTGGGCATGTGGCAGATAGCCACAGGAGATCACCTCCCATGAAGACGATCTCAGA-



## BILAG 2: Ukendte aminosyre-sekvenser

### Ukendt sekvens nr. 1 (FASTA-format)

SLTKAERTIIGSMWTKISSQADTIGTETLERLFASYPQAKTYFPHFDLNPQSDQLRAHGSKVLAAVGEAVKSIDNVSAALAKLSELHAYVLRVDPVNFKFLSHCLLVTLASHFPADLTA-EAHAAWDKFLTIVSGVLTEKYR

### Ukendt sekvens nr. 2 (FASTA-format)

GIVEQCCTSICSLYQLENYCN

### Ukendt sekvens nr. 3 (FASTA-format)

MFALRAASKADKNLLPFLGQLSRSHAAKAAKAAAAANGKIVAVIGAVVDVQFDDNLPPILNALEVDNRSRPLVLEVAQH LGENTVRTIAMDGTEGLVRGQKVLDTGYPIRIPVGAET-LGRIINVIGEPIDERGPIDTDKTAAIHAEAPEFVQMSVEQEILVTGIKVV DLLAPYAKGGKIGLFGGAGVGKTVLIMELINNVAKAHGGYSVFAGVGERTREGNDLYNEMIEGGVIS-LKDKTSKVALVYGQMNEPPGARARVALTGLTVAEYFRDQEGQDVLLFIDNIFRFTQAGSEVSALLGRIPSAVGYPQLATDMGSMQERITTTT KKSITSVQAIYVPADDLTD PAPATT-FAHLDATTVLSRAIAELGIYPAVDPLDSTSRIMDPNIIGQEHYNVARGVQKILQDYKSLQDI IAILGMDELSEEDKLTVARARKIQRFLSQPFQVAEVFTGHAGKLVPLEQTIKGFSAI-LAGDYDHLPEVAFYMGPIEEVVEKADRLAKEAA

### Ukendt sekvens nr. 4 (FASTA-format)

MLSAPCCDDRRMVCPCGPRRIGIPVRSSSLPLFSDAMPAPTQLFFPLIRNCELSRIYGTACYCHHKHLCCSSSYIPQSRLRYTPHPAYATFCRPKENWWQYTQGRRYAST-PQKFYLTPPQVNSILKAN EYSFKVPEFDGKMSVLSLDLTAIKLPANAPIEDRRSAATCLQTRGM LLGVFDGHAGCAWSQAVSERLFYYIAGSLVPHETLLEIENAVESGRALL-PILQWHKHPNDYFSKEASKLYFN SLRTYWQELIDLNTGESTDIDVKEALINAFKRLDNDISLEAQVGD PNSFLNYLVLRVAFSGATA CVAHVDGVDLHVANTGDSRAMLG-VQEEDGWSAVTLSNDHNAQNERELERLKLEHPKSEAKSVVKQDRLLGLLMPFRAFGDVKFKWSIDLQKRVIESGPDQLNDNEYTKFIPP NYHTPPYLTAEPEV TYHRLRPQDKFLV-LATDGLWETMHRQDVVRIVGEYLTGMH HQPIAVGGYKVT LGQMHG LLTERRTKMSSVFEDQNAATHLIRHAVGNNEFGTV DHERLSKMLS LPEELARMYRDDITIIVVQFN SHV-VGAYQNQEK