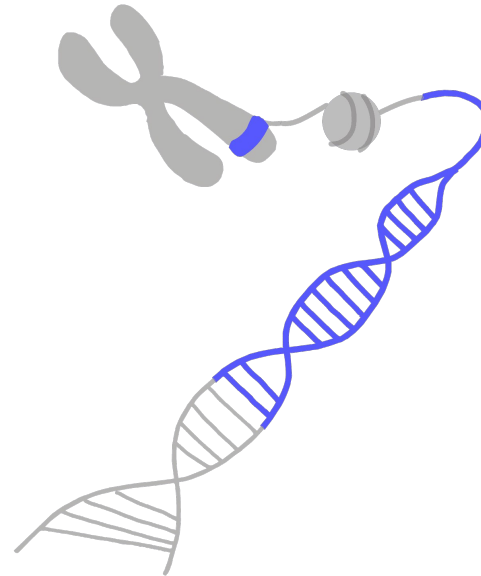


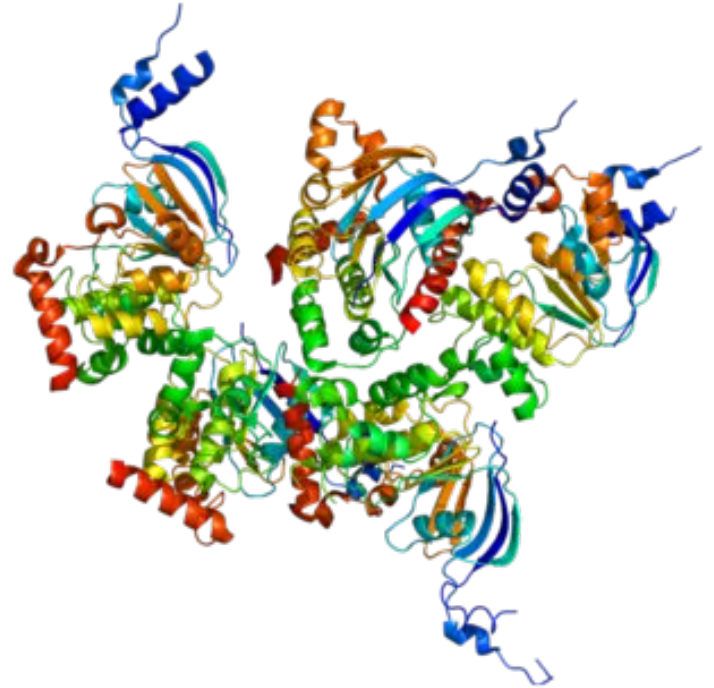
Análisis de la información de los archivos FASTA y GenBank

- Gen CFTR (Cystic Fibrosis Transmembrane Conductance)
- Proteína 4INS (Insulina Humana)



Gen CFTR

- Fundamental en el estudio de la fibrosis quística
- Gen específico





Datos en formato FASTA

```
fasta_seq = SeqIO.read(fasta_file_path, "fasta")
```

```
fasta_info = {  
    "ID": fasta_seq.id,  
    "Description": fasta_seq.description,  
    "Sequence Length": len(fasta_seq.seq),  
    "Composition": {  
        "A": fasta_seq.seq.count("A"),  
        "T": fasta_seq.seq.count("T"),  
        "C": fasta_seq.seq.count("C"),  
        "G": fasta_seq.seq.count("G")  
    }  
}
```

- **ID:** 'NM_000492.3',
- **Description:** 'NM_000492.3 Homo sapiens CF transmembrane conductance regulator (CFTR) mRNA'
- **Sequence Length:** 6132
- **Composition:**
 - 'A': 1887
 - 'T': 1731
 - 'C': 1182
 - 'G': 1332

Datos en formato GenBank

```
genbank_seq = SeqIO.read(genbank_file_path, "genbank")
```

```
genbank_info = {
    "ID": genbank_seq.id,
    "Description": genbank_seq.description,
    "Sequence Length": len(genbank_seq.seq),
    "Composition": {
        "A": genbank_seq.seq.count("A"),
        "T": genbank_seq.seq.count("T"),
        "C": genbank_seq.seq.count("C"),
        "G": genbank_seq.seq.count("G")
    },
    "Annotations": genbank_seq.annotations,
    "Features": genbank_seq.features,
    "Origin": genbank_seq.seq
}
```

- **Referencias bibliográficas** importantes de la secuencia: el gen CFTR, de gran relevancia, ej:
 O [Reference \(title='Domain-interface dynamics of CFTR protein stabilizes nanobodies', ...\)](#),
 O [Reference \(title='Unusual Cystic Fibrosis Transmembrane Conductance Regulator Mutations and Effects on the Protein and Review of the Literature'\)](#)
- **Accesiones y versión:** Versión 3 de la secuencia que ha sido reemplazada por versión 4
- **Comentarios** que han sido reemplazados por versión 4
- **Palabras clave:** RefSeq y RefSeq completa (del gen por NCBI)
- **Resumen funcional:** regulación del transporte de iones, relación con Vertebrata, Mammalian Homiidae, Metazoa, Chordata, Craniata, etc.

Proteína 4INS

- Estructura de la insulina de cerdo
- Formada por 2 cadenas A y B
- Cadenas específicas de la proteína de la insulina





Datos en formato FASTA

```
fasta_seq = SeqIO.read(fasta_file_path, "fasta")
```

```
for seq_record in fasta_sequences:
    fasta_info = {
        "ID": seq_record.id,
        "Description": seq_record.description,
        "Sequence Length": len(seq_record.seq),
        "Composition": {
            "A": seq_record.seq.count("A"),
            "C": seq_record.seq.count("C"),
            "D": seq_record.seq.count("D"),
            "E": seq_record.seq.count("E"),
            "F": seq_record.seq.count("F"),
            "G": seq_record.seq.count("G"),
            "H": seq_record.seq.count("H"),
            "I": seq_record.seq.count("I"),
            "K": seq_record.seq.count("K"),
            "L": seq_record.seq.count("L"),
            "M": seq_record.seq.count("M"),
            "N": seq_record.seq.count("N"),
            "P": seq_record.seq.count("P"),
            "Q": seq_record.seq.count("Q"),
            "R": seq_record.seq.count("R"),
            "S": seq_record.seq.count("S"),
            "T": seq_record.seq.count("T"),
            "V": seq_record.seq.count("V"),
            "W": seq_record.seq.count("W"),
            "Y": seq_record.seq.count("Y")
        }
    }
```

- **ID:** pdb|4INS|A,
- **Description:** pdb|4INS|A Chain A, INSULIN (CHAIN A)
- **Sequence Length:** 30
- **Composition:** 'A': 0, 'C': 2, 'D': 0, 'E': 2, 'F': 0, 'G': 3, 'H': 0, 'I': 0, 'K': 0, 'L': 2, 'M': 0, 'N': 2, 'P': 0, 'Q': 2, 'R': 0, 'S': 2, 'T': 1, 'V': 3, 'W': 0, 'Y': 2



Datos en formato GenPept

```
genbank_seq = SeqIO.read(genbank_file_path, "genbank")
```

```
for seq_record in genpept_sequences:
    genpept_info = {
        "ID": seq_record.id,
        "Description": seq_record.description,
        "Sequence Length": len(seq_record.seq),
        "Amino Acid Composition": {
            "A": seq_record.seq.count("A"),
            "C": seq_record.seq.count("C"),
            "D": seq_record.seq.count("D"),
            "E": seq_record.seq.count("E"),
            "F": seq_record.seq.count("F"),
            "G": seq_record.seq.count("G"),
            "H": seq_record.seq.count("H"),
            "I": seq_record.seq.count("I"),
            "K": seq_record.seq.count("K"),
            "L": seq_record.seq.count("L"),
            "M": seq_record.seq.count("M"),
            "N": seq_record.seq.count("N"),
            "P": seq_record.seq.count("P"),
            "Q": seq_record.seq.count("Q"),
            "R": seq_record.seq.count("R"),
            "S": seq_record.seq.count("S"),
            "T": seq_record.seq.count("T"),
            "V": seq_record.seq.count("V"),
            "W": seq_record.seq.count("W"),
            "Y": seq_record.seq.count("Y"),
        },
        "Annotations": seq_record.annotations,
        "Features": seq_record.features,
        "Origin": seq_record.seq
    }
```

- **Annotation base y comentario:**
 - Palabras clave: Insulina (estructural)
 - Topología: Linear
 - Comentario: Estructura cristalina de insulina 2ZN a 1.5 Å de resolución
 - Clasificación: MAM
- **Referencias destacadas:**
 - Fecha: 25 octubre 2024
 - Estructura de la insulina 2ZN: análisis cristalográfico
- **Accesiones y fuente:**
 - Métodos de refinamiento: evaluación comparativa en Modelistic4INS_A
 - Relaciones estructurales: hexámero de insulina
 - Fuente de datos: pdb: molécula 4INS, cadena A coordinado con zinc
- **Metodo experimental:** Difracción de rayos X
 - Ubicación fuente: posición 0-21
 - Organismo: Sus scrofa (cerdo)
 - Estructura secundaria: posiciones 0-9
 - Taxonomía: Eukaryota, Metazoa, Chordata, Enlace: posiciones 9-6 y 10-11
- **Origen:** Craniata, Vertebrata, Mammalia, Laurasiatheria,
 - Secuencia proteica: GIVEQCCTSIICSLYQLENYCN Suina, Sus.

INTERPRETACIÓN 2º ENLACE

NIH National Library of Medicine
National Center for Biotechnology Information

Nucleotide

GenBank

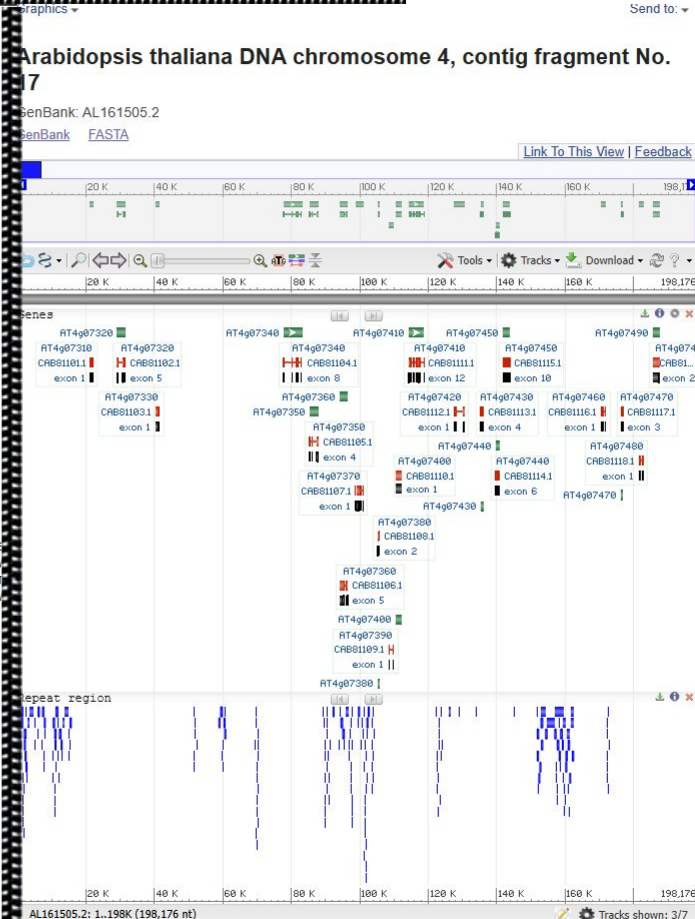
Arabidopsis thaliana DNA chromosome 4, c

GenBank: AL161505.2
[FASTA](#) [Graphics](#)

Go to:

LOCUS	AL161505	198176 bp	DNA	linear
DEFINITION	Arabidopsis thaliana DNA chromosome 4, contig fragme			
ACCESSION	AL161505			
VERSION	AL161505.2			
KEYWORDS	.			
SOURCE	Arabidopsis thaliana (thale cress)			
ORGANISM	Arabidopsis thaliana Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Spermatophyta; Magnoliopsida; eudicotyledons; Gunner Pentapetales; rosids; malvids; Brassicales; Brassica Camelineae; Arabidopsis.			
REFERENCE	1 (bases 1 to 1; 63093 to 175721)			
AUTHORS	Wilson,R., Lamar,B., Stoneking,T., Stumpf,J., Mewes, and Mayer,K.F.X.			
JOURNAL	Unpublished			
REFERENCE	2 (bases 1 to 198176)			
AUTHORS	EU Arabidopsis sequencing,project.			
TITLE	Direct Submission			
JOURNAL	Submitted (10-MAR-2000) MIPS, at the Max-Planck-Inst Biochemie, Am Klopferspitz 18a, D-82152 Martinsried, lemcke@mips.biochem.mpg.de,mayer@mips.biochem.mpg.de Coordinator: Mike Bevan, Molecular Genetics Departme Laboratory, John Innes Centre, Colney Lane, NR4 7UJ E-mail: michael.bevan@bbsrc.ac.uk			
COMMENT	Information on performance of analysis and a more de annotation of this entry and other sequences of chro and 5 can be viewed at: http://www.mips.biochem.mpg.de this fragment has an overlap with ATCHRIV16 at the 5 overlap with ATCHRIV18 at the 3' end.			
FEATURES	Location/Qualifiers			
source	1..198176 /organism="Arabidopsis thaliana" /mol_type="genomic DNA" /db_xref="taxon:3702" /chromosome="4" /ecotype="Columbia"			

misc feature 12310..15735
/note="pseudogene, similar to retrovirus-related
polyproteins"
repeat region 12735..13052
/note="T32N15_del_retroTn"
repeat region 13053..13086
/note="T32N15_del_retroTn"
repeat region 13087..13862
/note="T32N15_del_retroTn"
repeat region 13864..14927
/note="T32N15_del_retroTn"
repeat region 14928..15027
/note="T32N15_del_retroTn"
repeat region 15024..15226
/note="T32N15_del_retroTn"
repeat region 15224..15350
/note="T32N15_del_retroTn"
repeat region 15340..15739
/note="T32N15_del_retroTn"
21183..22195
/gene="AT4g07310"
join(21183..21320,21398..21662,21795..21897,21973..22195)
/gene="AT4g07310"
/note="contains similarity to Arabidopsis thaliana
hypothetical proteins, see GB:AF058826;
similarity to"
/codon_start=1
/protein_id="CAB81101.1"
/db_xref="InterPro:IPR003871"
/db_xref="InterPro:IPR012340"
/db_xref="InterPro:IPR016027"
/translation="MESFHLSSLKPISIRGWCIRGRVRTLFLVLPSSKVMGLIAE
EHGHTIATVGYKMSOHYKDFINEGEWITITFGVVENSGSVRAATHSFKIGFVDIV
VRLTSLPAIPHYRLASFSSIIODEIDKSVLDVLGAIDYDVGELINTRPKQNNVDLT
LFKISDNENRVLECLATKKEALDFDHNRYRGVGVIVAVLGMKIDRYFDGPKNVR
CTAGTITVFPDIPESNEIHEI"
21183..21320
/gene="AT4g07310"
/number=1
21321..21397
/gene="AT4g07310"
/number=1
21398..21662
/gene="AT4g07310"
/number=2
21663..21794
/gene="AT4g07310"
/number=2
21795..21897
/gene="AT4g07310"
/number=3
21898..21972
/gene="AT4g07310"
/number=3



Send to:

[Analyze this sequence](#)

[Run BLAST](#)

[Pick Primers](#)

Related information

[Protein](#)

[Taxonomy](#)

[Component Of](#)

[Full text in PMC](#)

[Gene](#)

[PubMed \(Weighted\)](#)

[GEO Profiles](#)

Recent activity

[Turn Off](#) [Clear](#)

[Arabidopsis thaliana DNA chromosome 4, contig frag](#) Nucleotide

[1A3N \(10\)](#) Nucleotide

[4INS \(5\)](#) Protein

[Chain A. GREEN FLUORESCENT PROTEIN](#) Protein

[1EMA \(3\)](#) Protein

[See more...](#)

INTERPRETACIÓN 1º ENLACE

intron complement(186935..187163)
/number=1
exon complement(187164..187579)
/gene="AT4g07490"
/number=2
misc feature 188252..191155
/note="contains similarity to athila retroelement
ORF1-like proteins"
misc feature 192885..195541
/note="pseudogene, contains similarity to athila
retroelement ORF1-like proteins"
misc feature 196274..198176
/note="pseudogene"

ORIGIN

```
1 gaattcgaaa ttaggttaaa tattctaaat ttatttaaat aactggaaaa ataattttaa
61 ctaaaattatt ttattgttt gataaataa atagaaaaa tatgttttaa catttaaaaa
121 taagattaat tacatttttg gttatatta atttttgta atatatata ttgaattcct
181 ttgcagcaa cgtaatcata ttaattttaa atattcaatg cgtatatccc ttcttcaaat
241 ttaaataata cacttctcca taattatagc taatcaataa ttcaattgga tattcatttt
301 tgaagttaac ataaagtctt cgacaattt ccataattaa ataaactact ataaagataat
361 tgcataaacat tctcttataa ttatgagatt tattataatc attaaacttg acttccacga
421 ttttaattatt tattcatatt tttaaccata aattttatata ctataaataa ctgtacgtcc
481 catttcaaga actcaactcag cagttctcat ttctttttac tctctgcctt ctttttttat
541 tcaagcgaat aaaaacgaga ccatgaagaa aactgaggtt gtggatcgac tagaccatgt
601 agaaagcaat atatcgaggt gctctctcac actctctgtt ttctcttatt ttttctatgt
661 ttcttccaca tctttttttt ataattatta tgcaggtaat gtggaagatt aaaaagttag
721 gcacacatgt attttgttaa agtatgtctt gctttctaat tatttatctt ttatgcaaa
781 taatagttaa accttttttt accctctcat atatttataa gatgagtaat tctgattgaa
841 ttttttgata ggggatacgg tctttttcct gttaaattgc agatacacga tgaatgatta
901 taatttgact actatcttcc aatgaacaac aaaaattataa ggtatgcaat attttttctt
961 aaactcgatt gataaattct ccaattattt tccaatttat aaagtatgaa gaagaacatt
1021 ttccaattat tatttgtcaa attttataat tcactacat gttaattgaca tacacgaaaa
1081 tgtacaaatc tctatgtatg acaactacat atacaaaata tacaattatg tatcacaagg
1141 tacataccaa ctatagacat atatttttta tgcataaac cgaggaaact acatatata
1201 aacaaaatat ggtttgacca gactaaatat atttagttgg taccaataat acagataaaa
1261 cactacatgt tttaaagata atgaatatat agtagatttt ctgtggaata cacagaactg
1321 agatttgtatg attgattaaa ccaaagatac cgtgcgtagc acgggtactg acctagtccc
1381 atattataac caaactattt ccatagttaa attctaatct ataaaccaatc agatctttca
1441 atcgaaatcat aaacttgttt tcaggaaaaa gctattgtat cagcaataaaa atcttcaatt
1501 taaactttag caatacaaga aaaaagaagt tgcataattc ttttttgttt acacttgaa
1561 ataatgatatc ttatttataa attatgtaat ttttttttgt tctaatactt ttggacactt
1621 gaatacaaaa acaaaattcaa aatgaccaaa tagcagaata ataaagagagg aaaaacaca
1681 gaatctaactc cagcaagctt actatgtacc tagcggacac ttctctattt gatgccggaa
1741 ccaccacctc tctgtcattg cttatctgaa atcgaaaaaa ctacaaaaaa ataaagtcat
1801 attgtctatt gagtatgcta caattacaaa atcaaaattca aatcaaaata tatgtttatc
1861 tcgagaaaaa gataatttag agtttttatc gtccaatgta ttatatagta tctttttata
1921 ttttttaaat taataaaaaa ggtttttttg ttttaccat gtcaatatat gatttttagt
1981 tctcaaaatt tgccataaac ataattatgg tacatttga taattttctt atatatcat
2041 attatcccat atattaaaaa atcaaaattt cttatttaa tgaagaata tttttttctt
2101 ataaacataa ttgcatgta tttttttaa gactttttt aatatatact tagtatacat
2161 ttatatatgc aaaaagtctt ttaaaattct aactataag aaaaactaac aaaaagtctaa
2221 aatatgatat taatggttat tattttgtta aaaaataat aatcaatac aaaaataaaa
2281 taaatatta ttgttatagt tttaaagtgt aagatatcgt gcaactcttt ttcaaaaaa
2341 gaatgtataa tttcaagtaa cacataataa gaatgaacaa aaaaattgttc atgcaggag
2401 tctagtctta ttttcacca cttccaagta attttgaatt tgaaaaattt gttttttatt
```



National Library of Medicine
National Center for Biotechnology Information

Nucleotide

Nucleotide

Advanced

FASTA

Send to:

Arabidopsis thaliana DNA chromosome 4, contig fragment No. 17

GenBank: AL161505.2

[GenBank](#) [Graphics](#)

>AL161505.2 Arabidopsis thaliana DNA chromosome 4, contig fragment No. 17

```
GAATTGCGAAATAGGTTAAATATTTCTAAATTTATTAAATAACTGGAAAAATAAATTTAACTAAATTTT
TTTTATGTTTGATAAATTAATAAGAAAAATATATGTTTAAACATTTAAAAATAAGATTATACATTTTGT
GTTTATATTAATTTTGTAAATATATTAATGAATCCCTTTTGCAGCACTGAATCATATTAATTTAAA
ATATTCAATGCGTATATTCCTTTCTTAAATTTAAATATACACTCTCCATAATATAGCTAATCAAAATA
TCGATTTGGATATTCATTTTGAAGTTAAACATAAAGTTTCCGGACAAATTTCCATAATTAATAACTACTT
ATAGATAAATGCTAAACATCTCTTATAATATGAGATTATTATAAATCATTAATAACTTGACTTCCACGA
TTTAATATTTTATCATATTTTACCTAAAAATTTATATCTATAAATACTGTACGTCCTTCAAGTTCACGA
ACTCAGTCAGCAGTTCTTATTTCTTTACTCTCTGCGCTCTTTTATTCAAGCAGAAAAAACGACGAG
CATGAAGAAAAATCGAGGTTGTGGATCGACTAGACCATGGAGAAAGCAATATATCGAGGTGTCTCTCA
ACTCTCTGTTTTCTCTAATTTTTCTCATGTTTTCTCACATCTTTTTTTATATAATATTATGACAGGTAAT
GTGGAGATTAATAATGTAGACGCACACATGATTTTGTAAAGTATGTCTGCTTCTAATTAATTTATCTT
TTATGCAAAAGTAATGTGTTAACTTTTTTACCTCTCTATATATTATAAGATGAGTAATCTGATTGAA
TTTTTTGATAGGGATACGGCTCTTTCTCTGTAAATTTGCAGATACACGATGATGATTAATAATTGACT
ACTACTCTTCCAATGAACAACAAAAATATAAGGTATGCAATTTTTTCTTAAATCTGATTGATAAATCTT
CCAATTTTCTTCAATTTATAAAGTATGAAGAAAGACATTTTCCAATTTATTTTGTGCAAAATTTTATAAT
TCACTACATGTTAATGACATACACGAAATGTACAAATCTCTATGTATACACACTACATATACAAAAATA
TACAAATATGTATCCAAAGGTACATACCAACTATAGACATATATTTTTATGCTCAAAACGAGGAAACCT
ACATATACTAAACAAAAATGTTTTGACCAGACATAAATATTAGTTGGTACCAATATACATGATAAAA
CAACTACATGTTTAAAGATAATGAAAAATATAGTAGAGTTTCTGTGGAATACACAGAACTGAGATTGTATG
ATTGATTAAACCAAGATACCGTGTAGCAGGGGTACTGACCTAGTTCCATATTATAACCAACATATTT
CATAGTAAAAATTTCTAATCTATAACCAATCAGATCTTCAATCGAATCATAAAACTTGTTTCAGGAAAAA
GCTATTGTATCAGCAATAAAATCTTCAATTTAACTTAGCAATACAGAAAAAGCAAGTTGTCAAAATCTT
TTTTTGTGTACACTGAAATATAAGATATCTTATTTAAAAATATGTAATTTTTTTTTTTGTTCTAAATCTT
TGACGACATGGAATACAAAAACAAATTTCAAAATGACCAAAATAGCAGAAATATAAGAGAGGAAAAACCAA
GAATCTAATCCAGCAAGCTTACTATGTACTAGCGGACACTTTCTCTATTGATGCGGAGAACCCACCTCT
TTCGTCTTGGCTTATCTGAATCGAAAAAATACAAAAATAAGCTGATATTGCTCATTTGATATGCTA
CTACTCAAAATCAAAATCAAAATCAAAATATATGTTTATCTCGAGAAAAATGATAATTTAGAGTTTTCATC
GTTCAATGATTTATAGATTATCTTTTATATTTTAAAAATTAATAAAAAAGGTTTTTTTGTGTTTACACAT
GTCAATATGATGATTTAGTTTCTCAAAATTTGGCAATTAACATAAATATGGTACATTTTGATTTATTTCTC
ATATTATCATATTATCCCATATATTAAAAAATCAAAATTTCTTATTATAAGGAAGAATATTTTTTTCTC
ATAAACTAATTTGCATGATTTTTTTTAAAGCATTTTTTAAATATATACTTAGTACATATTATATATCG
AAAGGTTTTTTTAAAAATCTTAACTATAAGAAAAAACTAACAAATGTCTAAAAATGATGATTAATGGTTAT
TTTTTTGTTAAAAAATAATATACTAATCAATCAAAATAAAAATAAAATATTATTTGTTATAGTTTAAAGGTG
TAGATTCTGTGCAACCTCTTTTCAAAATAGAATGATAAATTTCAAGTAAACATAAATAAGATAGAACCA
AAATGGTTGATGATGAGGATCTAGTCTTATTTTCCACCACTCCAAGTAATTTGAAATTTGAAATTTT
```