

Bioinformatyka 2 - kurs mały

Adrian Kania¹

¹Zakład Biofizyki Obliczeniowej i Bioinformatyki

2025/2026

Dopasowanie sekwencji

Dopasowanie pomiędzy sekwencjami określa które pozycje odpowiadają sobie.

acgtctag

||

actctag-

acgtctag

|||||

-actctag

acgtctag

|| | | | |

ac-tctag

2 matches

5 mismatches

1 not aligned

5 matches

2 mismatches

1 not aligned

7 matches

0 mismatches

1 not aligned

Takie porównania mogą być użyte do:

- szukania relacji ewolucyjnych,
- zidentyfikowania konserwatywnych miejsc,
- zidentyfikowania odpowiadających sobie genów pomiędzy różnymi modelami (np. ludzkimi czy mysimi).



Które dopasowanie jest najlepsze?

Na początku trzeba ustalić punktację za dopasowanie, niedopasowanie i przerwę. Poniżej przykładowe punkty.
Łączny wynik dopasowania to suma punktów za każdą pozycję.

acgtctag

||

actctag-

2 matches

5 mismatches

1 not aligned

$$S = 2*1 + 5*(-1) + 1*(-2) = -5$$

acgtctag

|||||

-actctag

5 matches

2 mismatches

1 not aligned

$$S = 5*1 + 2*(-1) + 1*(-2) = 1$$

acgtctag

|| | | | |

ac-tctag

7 matches

0 mismatches

1 not aligned

$$S = 7*1 + 0*(-1) + 1*(-2) = 5$$

Sekwencje aminokwasowe

Do porównywania sekwencji aminokwasowych używamy odpowiednich macierzy podobieństwa (np. BLOSUM)

A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4																		
R	-1	5																	
N	-2	0	6																
D	-2	-2	1	6															
C	0	-3	-3	-3	9														
Q	-1	1	0	0	-3	5													
E	-1	0	0	2	-4	2	5												
G	0	-2	0	-1	-3	-2	-2	6											
H	-2	0	1	-1	-3	0	0	-2	8										
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4									
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4								
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5							
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5						
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6					
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7				
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4			
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5		
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11	
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1

AABCDA...BBCDA
DABCDAA.A.BBCBB
BBCDABABA.BCCAA
AAACDAC.DCDBCDB
CCBADAB.DBBDCC
AAACAA...BBCCC

Dopasowanie wielu sekwencji (MSA)

ATGGCGAC 8

ATGCCGA- 7

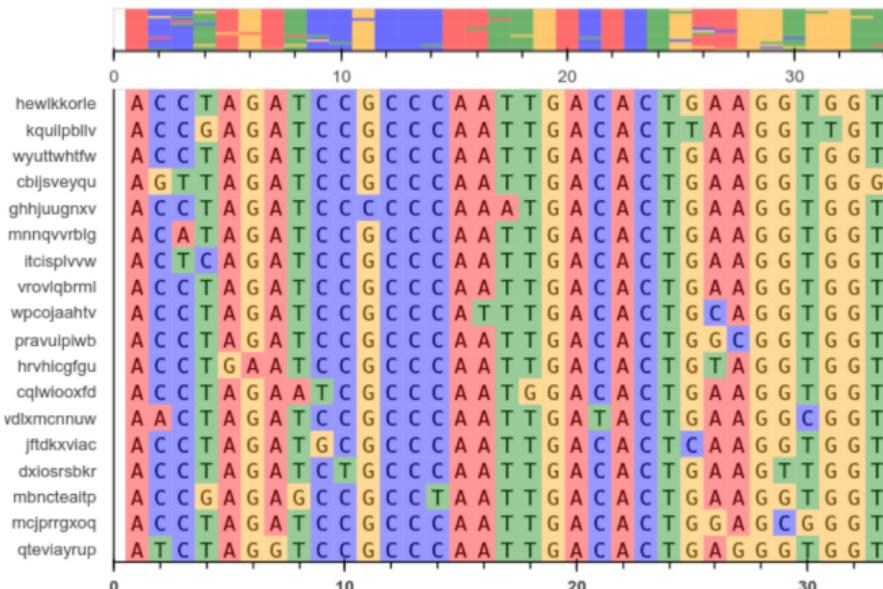
*** ***

ATGCCGA- 7

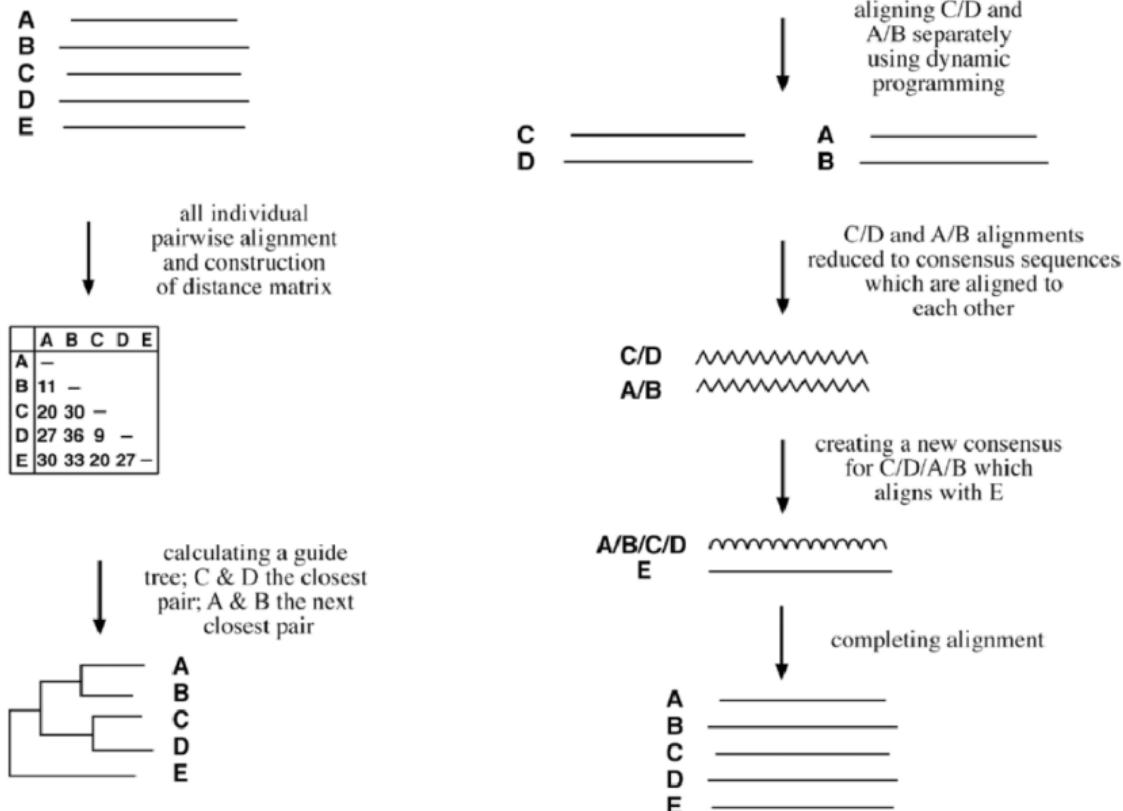
ATGGCGAC 8

TTGGCGAC 8

** ***



Schemat podejścia progresywnego (np. Clustal)



REVIEW

Open Access



CrossMark

Alignment-free sequence comparison: benefits, applications, and tools

Andrzej Zielezinski¹, Susana Vinga², Jonas Almeida³ and Wojciech M. Karlowski^{1*}

Abstract

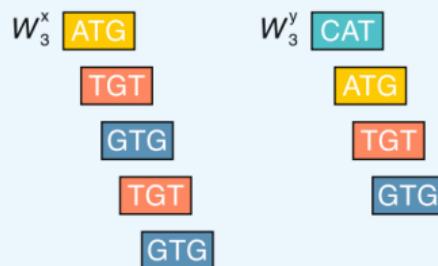
Alignment-free sequence analyses have been applied to problems ranging from whole-genome phylogeny to the classification of protein families, identification of horizontally transferred genes, and detection of recombined sequences. The strength of these methods makes them particularly useful for next-generation sequencing data processing and analysis. However, many researchers are unclear about how these methods work, how they compare to alignment-based methods, and what their potential is for use for their research. We address these questions and provide a guide to the currently available alignment-free sequence analysis tools.

Pięć przypadków, w których analiza sekwencji oparta na dopasowaniu może być kłopotliwa:

- Programy tworzące dopasowanie zakładają, że sekwencje homologiczne składają się z szeregu ułożonych liniowo i mniej lub bardziej konserwatywnych odcinków sekwencji.
- Dokładność dopasowania sekwencji gwałtownie spada w przypadkach, gdy identyczność sekwencji spada poniżej pewnego punktu krytycznego.
- Podejścia oparte na dopasowaniu są na ogół czasochłonne i pamięciowe, a zatem mają ograniczone zastosowanie w przypadku danych sekwencji w skali wielogenomowej.
- Obliczenie dokładnego dopasowania wielu sekwencji jest problemem NP trudnym, co oznacza, że trudno o rozwiązanie w akceptowalnych ramach czasowych.
- Dopasowanie sekwencji zależy od wielu założeń apriorycznych dotyczących ewolucji porównywanych sekwencji (np. macierzy podstawień, kar za przerwy i wartości progowych parametrów statystycznych).

Query sequences x ATGTGTG y CATGTG

Word size: 3



Union of two sets

$$W_3 = W_3^x \cup W_3^y \quad \text{CAT} \quad \text{ATG} \quad \text{TGT} \quad \text{GTG}$$

Word counts

$$c_3^x \quad \begin{matrix} 0 & 1 & 2 & 2 \end{matrix} \quad c_3^y \quad \begin{matrix} 1 & 1 & 1 & 1 \end{matrix}$$

Euclidean distance

$$\|c_3^x - c_3^y\| = \sqrt{(0-1)^2 + (1-1)^2 + (2-1)^2 + (2-1)^2} = \sqrt{3} = 1.73$$

Query sequences

x ATGTGTG

y CATGTG

xy ATGTGTGCATGTG

Lempel-Ziv complexity

 ATGTGTG

$$c(x)=4$$

 CATGTG

$$c(y)=5$$

 ATGTGTGCATGTG

$$c(xy)=7$$

Normalized compression distance

$$\frac{C(xy) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}}$$

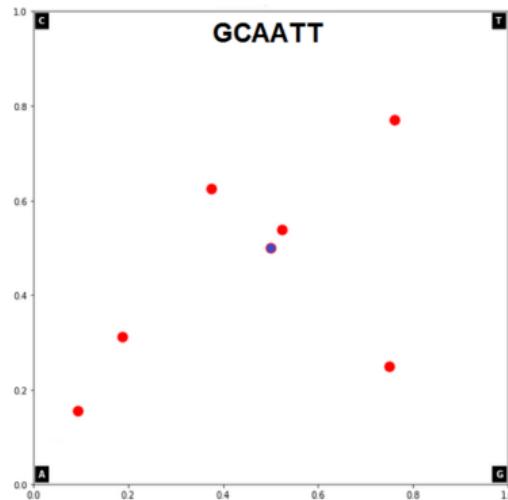
$$\frac{7-4}{5} = 0.6$$

Reprezentacja gry chaosu (CGR)

Jest to metoda służąca do reprezentacji numerycznej sekwencji biologicznych. Reprezentacja gry chaosu jest generowana wg

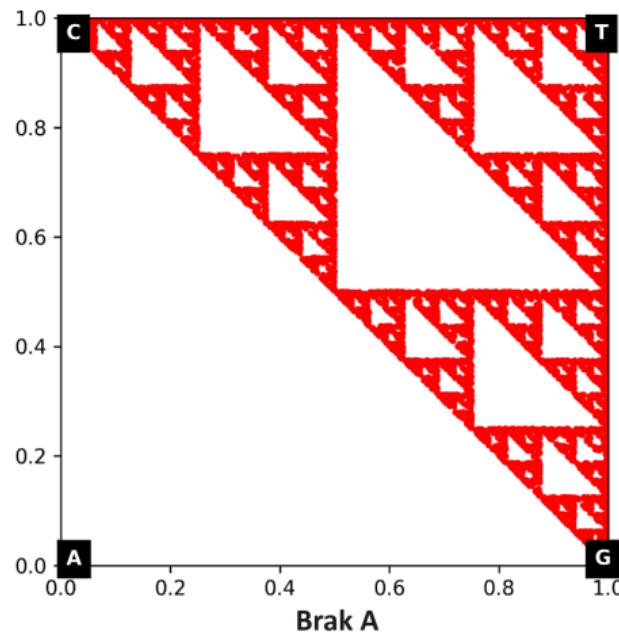
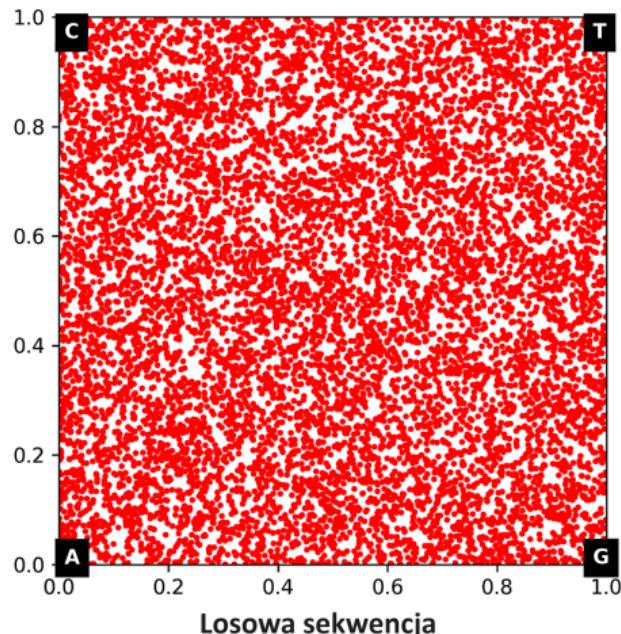
$$(x_{n+1}, y_{n+1}) = (x_n, y_n) - 0.5 \cdot ((x_n, y_n) - N_{n+1})$$

gdzie $N_A = (0, 0)$, $N_T = (1, 1)$, $N_C = (0, 1)$, $N_G = (1, 0)$, a $(x_0, y_0) = (0.5, 0.5)$ jest punktem początkowym.



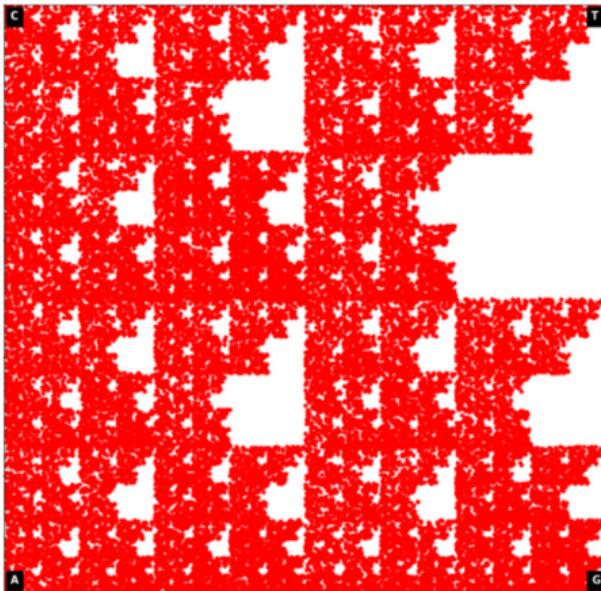
Reprezentacja gry chaosu

Za pomocą takiej reprezentacji możemy łatwo wskazać na pewne wzorce występujące w sekwencjach.

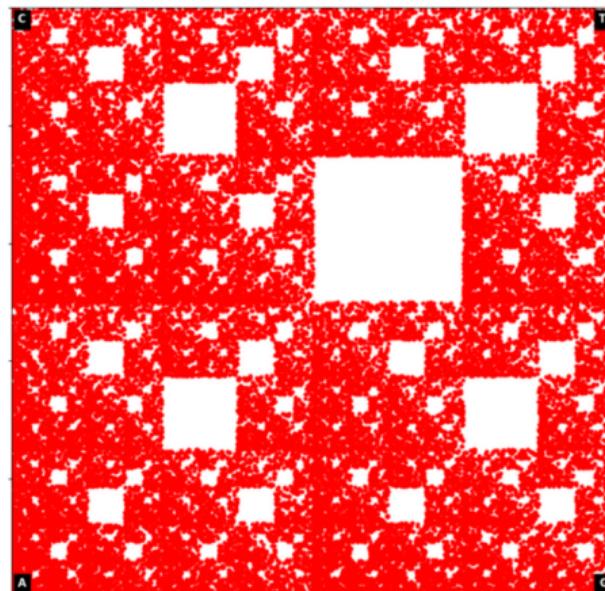


Reprezentacja gry chaosu

Za pomocą takiej reprezentacji możemy łatwo wskazać na pewne wzorce występujące w sekwencjach.



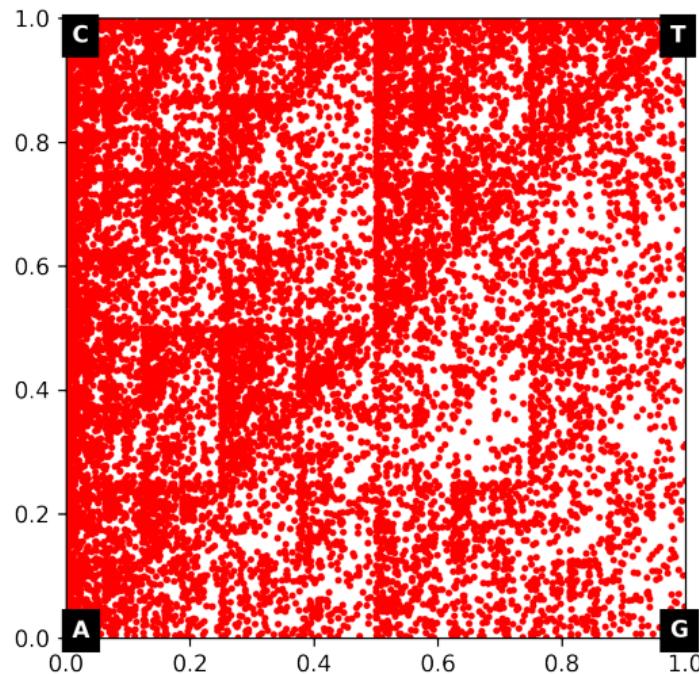
Brak GT w sekwencji



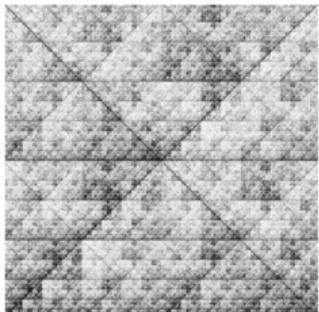
Brak AT w sekwencji

Reprezentacja gry chaosu

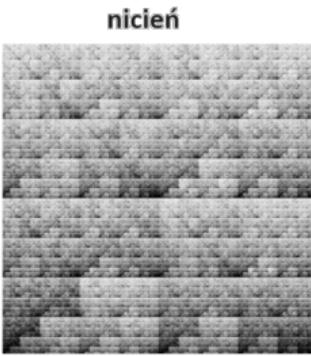
Genom mitochondrialny (NCBI ID: FJ986465.1)



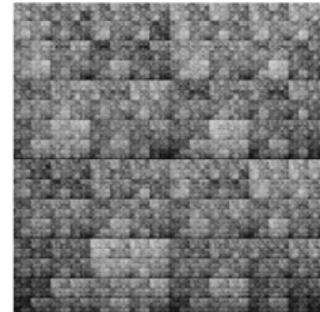
Reprezentacja gry chaosu



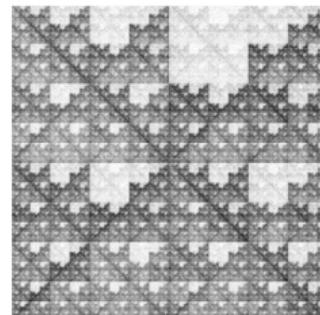
mysz



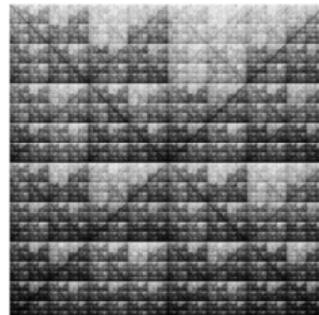
nicień



muszka owocowa



pszczoła



topola
kalifornijska

Deschavanne, Patrick, Alain Giron, Joseph Vilain, Guillaume Fagot, and Bernard Fertil. "Genomic Signature: Characterization and Classification of Species Assessed by Chaos Game Representation of Sequences". Mol. Biol. Evol. 16(10): 1391 - 1399. 1999.



ABOUT

BIOLOGICAL DATABASES

- NCBI – National Center for Biotechnology Information
- NCBI Entrez – zintegrowany system ponad 30 baz danych gromadzących informację z nauk biomedycznych, głównie w postaci sekwencji i literatury naukowej. Obejmuje m.in Pubmed (streszczenia publikacji), Nucleotide (sekwencje nukleotydowe), Protein (sekwencje aminokwasowe).

NCBI Resources How To

All Databases Search

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [NCBI News & Blog](#)

Submit Deposit data or manuscripts into NCBI databases 

Download Transfer NCBI data to your computer 

Learn Find help documents, attend a class or watch a tutorial 

Popular Resources

[PubMed](#)
[Bookshelf](#)
[PubMed Central](#)
[BLAST](#)
[Nucleotide](#)
[Genome](#)
[SNP](#)
[Gene](#)
[Protein](#)
[PubChem](#)

Develop Use NCBI APIs and code libraries to build applications

Analyze Identify an NCBI tool for your data analysis task

Research Explore NCBI research and collaborative projects

Adrian Kania (ZBOiB)

Bioinformatyka 2 - kurs mały

2025/2026

17 / 26

NCBI Resources How To

NCBI National Center for Biotechnology Information

All Databases

- All Databases
- Assembly
- Biocollections
- BioProject
- BioSample
- BioSystems
- Books
- ClinVar
- Conserved Domains
- dbGaP
- dbVar
- Gene
- Genome
- Genomes & Maps
- Homology
- Literature
- Proteins
- Sequence Analysis
- Taxonomy
- Training & Tutorials
- Variation

NCBI

National Center for Biotechnology Information advances science and health by providing access to genomic information.

[Mission](#) | [Organization](#) | [NCBI News & Blog](#)

Search

Submit

Download

Learn

Transfer NCBI data to your computer

Find help documents, attend a class or watch a tutorial

Drop

Analyze

Research

Identify an NCBI tool for your data analysis task

Explore NCBI research and collaborative projects

Popular Resources

- PubMed
- Bookshelf
- PubMed Central
- BLAST
- Nucleotide
- Genome
- SNP
- Gene
- Protein
- PubChem

TABLE 2.2. Several Selected PubMed Tags and Their Brief Descriptions

Tag	Name	Description
AB	Abstract	Abstract
AD	Affiliation	Institutional affiliation and address of the first author and grant numbers
AID	Article identifier	Article ID values may include the PII (controlled publisher identifier) or doi (digital object identifier)
AU	Author	Authors
DP	Publication date	The date the article was published
JID	Journal ID	Unique journal ID in the National Library of Medicine's catalog of books, journals, and audiovisuals
LA	Language	The language in which the article was published
PL	Place of publication	Journal's country of publication
PT	Publication type	The type of material the article represents
RN	EC/RN number	Number assigned by the Enzyme Commission to designate a particular enzyme or by the Chemical Abstracts Service for Registry Numbers
SO	Source	Composite field containing bibliographic information
TA	Journal title abbreviation	Standard journal title abbreviation
TI	Title	The title of the article
VI	Volume	Journal volume

Source: www.ncbi.nlm.nih.gov/entrez/query/static/help/pmhelp.html.

Format sekwencji

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide TDF | Search Help

Species Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

Animals (251)
Plants (1,588)
Fungi (48)
Protists (148)
Viruses (848)
Customize ...

Molecule types Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

genomic DNA/RNA (2,554)
mRNA (1,779)
Customize ...

Source databases Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

INSDC (GenBank) (5,345)
RefSeq (30)
Customize ...

Sequence Type Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

Nucleotide (3,677)
EST (1,098)

Genetic compartments Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

Mitochondrion (19)
Plastid (1)

Sequence length Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

Custom range...

Release date Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

Custom range...

Revision date Summary ▾ 20 per page ▾ Sort by Default order ▾ Send to: ▾ Filters: Manage Filters

Custom range...

See [SRY \(TDF\) sex determining region Y](#) in the Gene database
tdf reference sequences [Genomic \(1\)](#) [Transcript \(1\)](#) [Protein \(1\)](#)

Items: 1 to 20 of 5375 << First < Prev Page 1 of 269 Next > Last >>

[SRY \(SRY109S\)=testis-determining factor \[human, XY female, mRNA Partial Mutant, 424 nt\]](#)
1. 424 bp linear mRNA
Accession: S53156.1 GI: 263858
[Protein](#) [PubMed](#) [Taxonomy](#)
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens sex determining region Y \(SRY\).mRNA](#)
2. 828 bp linear mRNA
Accession: NM_003140.3 GI: 1653960407
[Protein](#) [PubMed](#) [Taxonomy](#)
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens sex determining region Y \(SRY\).RefSeqGene on chromosome Y](#)
3. 7,897 bp linear DNA
Accession: NG_011751.1 GI: 226053418
[Protein](#) [PubMed](#) [Taxonomy](#)
[GenBank](#) [FASTA](#) [Graphics](#)

[Mus musculus sex determining region of Chr Y \(Sry\).mRNA](#)
4. 1,188 bp linear mRNA
Accession: NM_011564.1 GI: 6755760
[Protein](#) [PubMed](#) [Taxonomy](#)

Results by taxon Top Organisms [Tree]
synthetic construct (1526)
Human immunodeficiency virus 1 (786)
Populus simonii x Populus nigra (314)
Diospyros kaki (283)
unidentified (242)
All other taxa (2224)
More...

Find related data Database: Select

Find items

Search details TDF [All Fields]

Adrian Kania (ZBOiB) Bioinformatyka 2 - kurs mały 2025/2026 20 / 26

FASTA format

*Record name;
starts with a “>”*

Fasta file:

```
>sequence_1  
TTTCCGGGGCACATAATCTTCAGCCGGGCGC  
>sequence_1  
TCAGCCGGGCCTTCAGCCGGGCACATAATA
```

DNA sequence

```
>Seq1 [organism=Carpodacus mexicanus] [clone=6b] actin (act) mRNA, partial cds  
CCTTATCTAACATTGGAGCATGAGCTGGCATAGTTGGAAACCGCCCTCAGCCTCTCATCGTGAGAA  
TAATAATTCTTATAGTAATAACCAATCATGATCGGTGGTTCGGAAGCTGACTAGTCCCACTCATAAT
```

```
>Seq2 [organism=uncultured bacillus sp.] [isolate=A2] corticotropin (CT) gene, complete cds  
GGTAGGTACCGCCCTAACGCTCTAACCGAGCAGAACTANGCCAACCCGGAGCCCTCTGGAGACGAC  
TCAACACCACCTCTTGACCCAGCAGCAGGAGGAGACCCAGTACTATACCGACCTATTCTGATTCTT
```

GenBank

Header

```
LOCUS      g92069          440 am      linear  BCT 15-JUN-2002
DEFINITION Light-independent protochlorophyllide reductase subunit N (LI-PCR
SUBUNIT N)
ACCESSION  g92069
VERSION    g92069 GI:18203671
SOURCE     bacterial isolate BACTN_HELJO, accession Q95059;
CLASS: protein
CREATED: created; Oct 16, 2001.
SEQUENCE_UPDT: sequence updated; Oct 16, 2001.
ANNOT_UPDT: annotation updated; Jun 15, 2002.
XREFS: gi: 18203671, gi: 3820556
XREF (mso-sequence databases); InterProIPI2000510, PfamPF00148
KEYWORDS: Photosynthesis; Bacteriochlorophyll biosynthesis; Oxidoreductase;
SOURCE: Helobacillus mobilis
ORGANISM: Helobacillus mobilis
BACTERIA; Firmicutes; Clostridia; Clostridiales; Helobacteriaceae;
Helobacillaceae.
REFERENCE: 1 (residues 1 to 440)
AUTHORS: Xiong,J., Itoue,J. and Bauer,C.B.
TITLE: Tracing molecular evolution of photosynthesis by characterization
of a major photosynthesis gene cluster from Helobacillus mobilis
JOURNAL: Proc. Natl. Acad. Sci. U.S.A. 95 (25), 14851-14856 (1998)
MEDLINE: 98611957
PUBLISHED: 08/1998
REMARK: SEQUENCE FROM N.D.
COMMENT: -----
This SWISS-PROT entry is copyright. It is produced through a
collaboration between the Swiss Institute of Bioinformatics and
the EMBL outstation - the European Bioinformatics Institute.
The original entry is available from http://www.ebi.ac.uk/swissprot
and http://www.sib.ch/swissprot.
```

Features

```
[FUNCTION] Uses Mg-ATP and reduced ferredoxin to reduce ring D of
protochlorophyllide (Pchlide) to form chlorophyllide a (Chlide) (By
similarity). This reaction is light-independent.
[PATHWAY] Light-independent bacteriochlorophyll biosynthesis.
[SUBUNIT] Protochlorophyllide reductase is thought to be composed
of three subunits: bchH, bchM and bchN. Could form a heterotetramer
of two bchH and two bchN subunits.
[SIMILARITY] BELONGS TO THE BCHN / CHLN FAMILY.
```

Sequence

```
FEATURES
  SOURCE
  /organism='Helobacillus mobilis'
  /db_xref='taxon:28064'
  gene
  /gene='BCHN'
  Protein
  /gene='BCHN'
  /product='Light-independent protochlorophyllide reductase
  subunit N'
  /EC_number='1.18.1.-'
ORIGIN
  1 msvvseggc ftticplasw awlhkrkhd fflivgthtc ahfigtaldv myvahsrifg
  61 avlvesdilv ampteslegh vggvvdewhp krlifvlstns vdilkmlnay skclklstrifg
  121 fpvlpastag idrftqpdq avhllfslpv pkeapavgv eekpkpwaf qkessckcak
  181 parhkhvlgg vcttqigq wlkqkqjpk vdtfpqgqk kmpvngqvg vqplqppqnd
  241 tlaatirzr svhlytvpq gpqdtatfls aelcfqglfd wrkshenqa wtlleplqql
  301 lighkhkqfz dnnellip:at Ellscudvqvq esqtpylnsh dlquqellk wivvilevsp
  361 dftkqkqmgc sykpldkwag lqicpleam qfttawmies tfqaqihgva aidliklfth
  421 pllkrgalne hgwaaagwle
  /f
```

Figure 2.3: NCBI GenBank/GenPept format showing the three major components of a sequence file.



LOCUS	X56734 1859 bp mRNA linear PLN 26-NOV-2005	FEATURES	Location/Qualifiers	ORIGIN
DEFINITION	Trifolium repens mRNA for non-cyanogenic beta-glucosidase.	source	1..1859	1 aaacaacca aataaggatt ttatgttagc catatgtgc ctgtttgtta ttatcttatt
ACCESSION	X56734 S46826		/organism="Trifolium repens"	61 cacaattact tcacaaagc cagtgtggc ttctactctt ctgcacatcg gtaaccgtgg
VERSION	X56734_1 GI:21954		/mol_type="mRNA"	121 tcggagcagt ttccctcggt gtttcacatcc tggtgtcgga ttccacgat accaaatgg
KEYWORDS	beta-glucosidase.		/db_xref="taxon:3899"	181 aggtgcgtta aacgaaaggg giaggagacc aagtatgtgg gataccatcca ccataaaata
SOURCE	Trifolium repens (white clover)		/clone="TRE361"	241 tccagaaaaa ataaggatg gaacgttg agacatcacg gtggcaatc atcaccgcta
ORGANISM	Trifolium repens		/tissue_type="leaves"	...
	Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; eudicotyledons; core eudicotyledons; rosids; eurosids I; Fabales; Fabaceae; Papilionoideae; Trifolieae; Trifolium.	mRNA	/clone_lib="lambda gt10"	1741 agaaacctatcataacta taggttgtatc cttcatgtat cagttttggat ttgaaatcac
REFERENCE	1 (bases 1 to 1859)		1..1859	1801 ttgttaattaa aagtctttttt ttatttttt aaaaaaaaaaaa aaaaaaaaaaaa aaaaaaaaaaaa
AUTHORS	Oxtoby,E., Dunn,M.A., Panciro,A. and Hughes,M.A.		/experiment="experimental evidence, no additional details recorded"	
TITLE	Nucleotide and derived amino acid sequence of the cyanogenic beta-glucosidase (linamarase) from white clover (<i>Trifolium repens</i> L.)	CDS	14..1495	
JOURNAL	Plant Mol. Biol. 17 (2), 209-219 (1991)		/EC_number="3.2.1.21"	
PUBLMED	1907511		/note="non-cyanogenic"	
REFERENCE	2 (bases 1 to 1859)		/codon_start="1"	
AUTHORS	Hughes,M.A.		/product="beta-glucosidase"	
TITLE	Direct Submission		/protein_id="CA4A0058.1"	
JOURNAL	Submitted (19-NOV-1990) Hughes M.A., University of Newcastle Upon Tyne, Medical School, Newcastle Upon Tyne, NE2 4HH, UK		/db_xref="GI:21954"	
COMMENT	On Jun 10, 2005 this sequence version replaced gi:233395.		/db_xref="P26204"	
			/db_xref="InterPro:IPR001360"	
			/db_xref="UniprotKB/Swiss-Prot:P26204"	
			/translation="#MFIVIAIFALFVISSFTITSTIAVEASTLIDGNLSRSSFPRGF	
			IFPGASSAYOFEGAVNEDGCRGPSPSUDWTFTHYPEKIRDGNSADITVTDQHRYKEDVGI	
			M6QDQNMGYNPJ313WPA1LPGOKL9Q1HNG1HNG1KYNNH1LNELANG1COPFVTLRHD	
			LPQVLEDEYFGCPFLNSGV1NDFRDYTDLCFKEFODWRVYKSTLNEPWFVFSNGSGALCTN	
			APGRCSASNVAKPQD867GPIYV7HNQNL1AHABAVHRYKKYQYAOYKOKICITLWSNV	
			LMPLEDDNSIPDKAERSLDFQQLPMEQLTGDYSSKSNR1AVKNALPKPSKFESSLV	
			NGSFDFG1INYYSSSYISNAPSHQAKPSYSTNPMTN1SFEKH1CPLGPRAAS1IYV	
			YPYPMFIQEDFIEPCYI1KIN1TILQPSITTECNNEFHNDATL1VVEALLHITYRIDYVYR	
			HPYYTIRSAIRAGSNVKGFYAWSFLDCNWEFGTIVRFGLNFVD"	

- ➊ **Entrez Programming Utilities** – oprogramowanie umożliwiające dostęp do danych zintegrowanych w ramach systemu Entrez bez konieczności obsługi formularzy na stronach WWW
 - ➌ **ESearch** – wyszukiwanie identyfikatorów/kodów dostępu rekordów wybranej bazy danych,
 - ➍ **EFetch** – pobieranie rekordów o wskazanych identyfikatorach lub wyników wskazanego wyszukiwania

EFetch

EFetch (efetch.fcgi) returns full data records for a list of unique identifiers (UIDs) in a format specified in the parameters. The list of UIDs is either provided in the parameters, or is retrieved from the [History server](#).

EFetch Parameters

EFetch Required Parameters

- **db** (required): Database containing the unique identifiers (UIDs) for which you wish to retrieve records. You can see NCBI's [table of Entrez Unique Identifiers \(UIDs\)](#) for a complete list of allowable database names, but some example values include:
 - [pubmed](#): PubMed
 - [pmc](#): PubMed Central
 - [nlmcatalog](#): NLM Catalog
- **id** (required): Either a single unique identifier (UID) or a comma-delimited list of UIDs. All of the UIDs must be from the database specified by the **db** parameter.

EFetch Optional Parameters

- **restart** (optional): Setting this parameter helps limit which records will be shown in the output, as it determines whether the record for the first input unique identifier (UID) is retrieved, or whether to skip to a later UID in the input list. For example, if **restart** is set to [10](#), the output will begin with the record for the tenth UID. The default of this parameter is [1](#), corresponding to the first UID in the input list. This parameter can be used in conjunction with **retmax** to download an arbitrary subset of records.
- **retmax** (optional): Total number of records to be shown in the output, up to a maximum of 10,000. If the set of records you are trying to retrieve is larger than 10,000, you can submit multiple EFetch requests, and increase the **restart** parameter each time.
- **retmode**/**rettype**: These two parameters determine how your results will be displayed. **retmode** determines the data format your records will be returned in (e.g. XML, plain text, etc.). **rettype** determines the specific view your records will be returned in (e.g. MEDLINE, Abstract, list of PMIDs, etc.). Different databases have different allowable data formats and record views, and not all **retmode** data formats are compatible with all **rettype** record views, and vice versa.

The table below shows the allowable combinations of **retmode** and **rettype** for some of the databases. **Bold** **retmode** values are the default data format for the specified database. **Bold** **rettype** parameters are the default record view for the specified data format and database.

EFetch Examples

- Retrieve the abstract view (text format) of two PubMed records.
 - <https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=pubmed&id=17284678,9997&retmode=text&rettype=abstract>
- Retrieve two PubMed records in XML format.
 - <https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=pubmed&id=11748933,11700088&retmode=xml>

ESearch Formatting Parameters

- `restart` (optional): Setting this parameter helps limit which of the unique identifiers (UIDs) in the results set will be shown in the output, as it determines whether the output begins at the first retrieved UID, or with a UID that is later in the results set. For example, if `restart` is set to `10`, the first ten UIDs in the results set will be skipped, and the output will begin with the eleventh UID. The default of this parameter is `0`, corresponding to the first record in the entire set. This parameter can be used in conjunction with `retmax` to download an arbitrary subset of UIDs retrieved from a search.
- `retmax` (optional): Total number of unique identifiers (UIDs) from the retrieved set to be shown in the output (default=20). Increasing `retmax` allows more of the retrieved UIDs to be included in the output, up to a maximum of 100,000 UIDs. If you need to retrieve more than 100,000 UIDs, you can submit multiple ESearch requests, and increase the `restart` parameter each time. For example:
 - <https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=pubmed&term=cancer&retmax=100>: This URL will return results 1 through 100 of a search for "Cancer".
 - <https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi?db=pubmed&term=cancer&restart=100&retmax=100>: This URL will return results 101 through 200 of the same search for "Cancer".
- `rettype` (optional): Retrieval type. There are two supported values:
 - `uillist` (default): Displays the standard XML output, including a list of unique identifiers (UIDs), the total number of results, and the query translation for the search.
 - `count`: Displays only the total number of results, without the list of UIDs or query translation.
- `retmode` (optional): Determines the format of the returned output. The default value is `xml`, but `json` is also supported.