WikipediA The Free Encyclopedia

Main page

Current events

Random article

About Wikipedia

Contact us

Contribute

Learn to edit

Upload file

Tools

Community portal

Recent changes

What links here

Special pages Permanent link

Related changes

Page information

Download as PDF

Printable version

0

Cite this page

Wikidata item

Print/export

Languages

العربية

Deutsch

Español

فارسى Français

한국어

Русский

中文

Português

文 4 more

Edit links

Donate

Help

Contents

Article

Talk

FASTA format From Wikipedia, the free encyclopedia

4.1 Filename extension

4.2 Compression

4.3 Encryption

5 Extensions

programs.

>SEQUENCE_1

NCBI identifiers [edit]

third-party GenBank 2

Sequence representation [edit]

U

third-party EMBL 2

third-party DDBJ 2

Nucleic Acid Code ◆

U

TrEMBL

Type

In bioinformatics and biochemistry, the **FASTA format** is a text-based format for representing either nucleotide sequences or amino acid (protein) sequences, in which nucleotides or amino acids are represented using single-letter codes. The format also allows for sequence names and comments to precede the sequences. The format originates from the

FASTA software package, but has now become a near universal standard in the field of bioinformatics.^[4]

The simplicity of FASTA format makes it easy to manipulate and parse sequences using text-processing tools and scripting languages like the R programming language, Python, Ruby, Haskell, and Perl. Contents [hide] 1 Original format & overview 2 Description line 2.1 NCBI identifiers 3 Sequence representation 4 FASTA file

FASTA format **Filename** .fasta, .fna, .ffn, extensions .faa, .frn, .fa Internet text/x-fasta media type **Uniform Type** no Identifier (UTI) **Developed by** David J. Lipman William R. Pearson^{[1][2]} **Initial release** 1985 Type of format **Bioinformatics** Extended from **ASCII for FASTA** Extended to FASTQ format[3] Website www.ncbi.nlm.nih.gov /BLAST/fasta.shtml ☑

Not logged in Talk Contributions Create account Log in

Search Wikipedia

Read

Edit | View history

Q

6 Working with FASTA files 7 See also 8 References 9 External links Original format & overview [edit]

The original FASTA/Pearson format is described in the documentation for the FASTA suite of programs. It can be downloaded with any free distribution of FASTA (see fasta20.doc, fastaVN.doc or fastaVN.me—where VN is the Version Number).

In the original format, a sequence was represented as a series of lines, each of which was no longer than 120 characters and usually did not exceed 80 characters. This probably was to allow for preallocation of fixed line sizes in software: at the time most users relied on Digital Equipment Corporation (DEC) VT220 (or compatible) terminals which could display 80 or 132 characters per line. [citation needed] Most people preferred the bigger font in 80character modes and so it became the recommended fashion to use 80 characters or less (often 70) in FASTA lines. Also, the width of a standard printed page is 70 to 80 characters (depending on the font). Hence, 80 characters became the norm. [citation needed]

The first line in a FASTA file started either with a ">" (greater-than) symbol or, less frequently, a ";" [citation needed] (semicolon) was taken as a comment. Subsequent lines starting with a semicolon would be ignored by software. Since the only comment used was the first, it quickly became used to hold a summary description of the sequence, often starting with a unique library accession number, and with time it has become commonplace to always use ">" for the first line and to not use ";" comments (which would otherwise be ignored).

Following the initial line (used for a unique description of the sequence) was the actual sequence itself in standard one-letter character string. Anything other than a valid character would be ignored (including spaces, tabulators, asterisks, etc...). It was also common to end the sequence with an "*" (asterisk) character (in analogy with use in PIR formatted sequences) and, for the same reason, to leave a blank line between the description and the sequence. Below are a few sample sequences: ;LCBO - Prolactin precursor - Bovine

; a sample sequence in FASTA format MDSKGSSQKGSRLLLLLVVSNLLLCQGVVSTPVCPNGPGNCQVSLRDLFDRAVMVSHYIHDLSS **EMFNEFDKRYAQGKGFITMALNSCHTSSLPTPEDKEQAQQTHHEVLMSLILGLLRSWNDPLYHL** VTEVRGMKGAPDAILSRAIEIEEENKRLLEGMEMIFGQVIPGAKETEPYPVWSGLPSLQTKDED ARYSAFYNLLHCLRRDSSKIDTYLKLLNCRIIYNNNC* >MCHU - Calmodulin - Human, rabbit, bovine, rat, and chicken

MADOLTEE0IAEFKEAFSLFDKDGDGTITTKELGTVMRSLGONPTEAELODMINEVDADGNGTID FPEFLTMMARKMKDTDSEEEIREAFRVFDKDGNGYISAAELRHVMTNLGEKLTDEEVDEMIREA **DIDGDGQVNYEEFVQMMTAK*** >gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus] LCLYTHIGRNIYYGSYLYSETWNTGIMLLLITMATAFMGYVLPWGQMSFWGATVITNLFSAIPYIGTNLV EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYYTIKDFLG LLILILLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLALFLSIVIL GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFLPIAGX **IENY**

A multiple sequence FASTA format would be obtained by concatenating several single sequence FASTA files in a common file (also known as multi-FASTA format). This does not imply a contradiction with the format as only

Nowadays, modern bioinformatic programs that rely on the FASTA format expect the sequence headers to be preceded by ">", and the actual sequence, while generally represented as "interleaved", i.e. on multiple lines as in

the first line in a FASTA file may start with a ";" or ">", hence forcing all subsequent sequences to start with a ">" in order to be taken as different ones (and further forcing the exclusive reservation of ">" for the sequence

the above example, may also be "sequential" when the full stretch is found on a single line. Users may often need to perform conversion between "Sequential" and "Interleaved" FASTA format to run different bioinformatic

Description line [edit] The description line (defline) or header/identifier line, which begins with '>', gives a name and/or a unique identifier for the sequence, and may also contain additional information. In a deprecated practice, the header line sometimes contained more than one header, separated by a ^A (Control-A) character. In the original Pearson FASTA format, one or more comments, distinguished by a semi-colon at the beginning of the line, may occur after the header. Some databases and bioinformatics applications do not recognize these comments and follow the NCBI FASTA specification . An example of a multiple sequence FASTA file follows:

IPQFASRKQLSDAILKEAEEKIKEELKAQGKPEKIWDNIIPGKMNSFIADNSQLDSKLTL MGQFYVMDDKKTVEQVIAEKEKEFGGKIKIVEFICFEVGEGLEKKTEDFAAEVAAQL >SEQUENCE_2 SATVSEINSETDFVAKNDQFIALTKDTTAHIQSNSLQSVEELHSSTINGVKFEEYLKSQI

The NCBI defined a standard for the unique identifier used for the sequence (SeqID) in the header line. This allows a sequence that was obtained from a database to be labelled with a reference to its database record. The

Example(s)

tpg|BK003456|

tpe|BN000123|

tpd|FAA00017|

tr|Q90RT2|Q90RT2_9HIV1

database identifier format is understood by the NCBI tools like makeblastdb and table2asn. The following list describes the NCBI FASTA defined format for sequence identifiers. [5]

Format(s)

lcl|123 lcl|*integer* local (i.e. no database reference) lcl|string lcl|hmm271 GenInfo backbone seqid bbs|integer bbs | 123 bbm|*integer* bbm | 123 GenInfo backbone moltype gim|*integer* gim|123 GenInfo import ID gb|accession|locus gb | M73307 | AGMA13GT GenBank ☑ emb|accession|locus EMBL ♂ emb|CAM43271.1| pir||G36364 PIR ♂ pir|accession|name SWISS-PROT ∠ sp|accession|name sp|P01013|OVAX_CHICK pat | US | RE33188 | 1 pat|country|patent|sequence-number patent pgp|country|application-number|sequence-number pgp|EP|0238993|7 pre-grant patent ref|NM_010450.1| ref|accession|name RefSeq ☑ gnl|taxon|9606 gnl|database|integer general database reference (a reference to a database that's not in this list) gnl|database|string gnl|PID|e1632 GenInfo integrated database gi|*integer* gi|21434723 dbj|BAC85684.1| DDBJ ♂ dbj|accession|locus PRF 🛂 prf|accession|name prf||0806162C pdb|1I4L|D PDB 🛂 pdb|*entry*|*chain*

tpg|accession|name

tpe|accession|name

tpd|accession|name

tr|accession|name

Mnemonic

Uracil

definition line). Thus, the examples above may as well be taken as a multisequence (i.e multi-FASTA) file if taken together.

MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDKAVQLLREKGLGKAAKKADRLAAEG LVSVKVSDDFTIAAMRPSYLSYEDLDMTFVENEYKALVAELEKENEERRRLKDPNKPEHK

ATIGENLVVRRFATLKAGANGVVNGYIHTNGRVGVVIAAACDSAEVASKSRDLLRQICMH

codes supported are:[6][7][8]

The vertical bars ("I") in the above list are not separators in the sense of the Backus-Naur form, but are part of the format. Multiple identifiers can be concatenated, also separated by vertical bars.

Adenine C C Cytosine G G **G**uanine **T**hymine

Following the header line, the actual sequence is represented. Sequences may be protein sequences or nucleic acid sequences, and they can contain gaps or alignment characters (see sequence alignment). Sequences are

represent a gap character; and in amino acid sequences, U and * are acceptable letters (see below). Numerical digits are not allowed but are used in some databases to indicate the position in the sequence. The nucleic acid

expected to be represented in the standard IUB/IUPAC amino acid and nucleic acid codes, with these exceptions: lower-case letters are accepted and are mapped into upper-case; a single hyphen or dash can be used to

(i)	i	inosine (non-standard)		
R	A or G (I)	pu R ine		
Υ	C, T or U	p Y rimidines		
К	G, T or U	bases which are Ketones		
М	A or C	bases with aMino groups		
S	C or G	Strong interaction		
W	A, T or U	Weak interaction		
В	not A (i.e. C, G, T or U)	B comes after A		
D	not C (i.e. A, G, T or U)	D comes after C		
Н	not G (i.e., A, C, T or U)	H comes after G		
V	neither T nor U (i.e. A, C or G)	V comes after U		
N	ACGTU	Nucleic acid		
-	gap of indeterminate length			
The amino acid codes supported (22 amino acids and 3 special codes) are:				
Amino Acid Code +	Meaning	♦		
Α	Alanine			
В	Aspartic acid (D) or Asparagine (N)			
С	Cysteine			
D	Aspartic acid			
F	Glutamic acid			

Meaning

Glutamic acid Phenylalanine

G	Glycine	
Н	Histidine	
I	Isoleucine	
J	Leucine (L) or Isoleucine (I)	
K	Lysine	
L	Leucine	
M	Methionine/Start codon	
N	Asparagine	
0	Pyrrolysine (rare)	
Р	Proline	
Q	Glutamine	
R	Arginine	
S	Serine	
Т	Threonine	
U	Selenocysteine (rare)	
V	Valine	
W	Tryptophan	
Υ	Tyrosine	
Z	Glutamic acid (E) or Glutamine (Q)	
X	any	
*	translation stop	
-	gap of indeterminate length	
FASTA fil	.e [edit]	
Filename e	xtension [edit]	
There is no sta	andard filename extension for a text file contain	ing FASTA formatted sequences. The table below shows each extension and its respective mea

Meaning

compression algorithms, see Hosseini et al., 2016, [11] and Kryukov et al., 2020. [12]

generic FASTA

FASTA nucleic acid fna FASTA nucleotide of gene regions FASTA amino acid

Extension ♦

fasta, fa^[9]

FASTA non-coding RNA Contains non-coding RNA regions for a genome, in DNA alphabet e.g. tRNA, rRNA frn Compression [edit]

Used generically to specify nucleic acids.

Contains coding regions for a genome.

Any generic fasta file. See below for other common FASTA file extensions

Notes

Contains amino acid sequences. A multiple protein fasta file can have the more specific extension mpfa.

The compression of FASTA files requires a specific compressor to handle both channels of information: identifiers and sequence. For improved compression results, these are mainly divided in two streams where the

The encryption of FASTA files has been mostly addressed with a specific encryption tool: Cryfa. [13][14] Cryfa uses AES encryption and enables to compact data besides encryption. It can also address FASTQ files.

compression is made assuming independence. For example, the algorithm MFCompress^[10] performs lossless compression of these files using context modelling and arithmetic encoding. For benchmarks of FASTA files

A2M/A3M are a family of FASTA-derived formats used for sequence alignments. In A2M/A3M sequences, lowercase characters are taken to mean insertions, which are then indicated in the other sequences as the dot (""") character. The dots can be discarded for compactness without loss of information. As with typical FASTA used in alignments, the gap ("-") is taken to mean exactly one position. [15] A3M is similar to A2M, with the added rule that gaps aligned to insertions can too be discarded. [16]

See also [edit]

Encryption [edit]

Extensions [edit]

Working with FASTA files [edit] A plethora of user-friendly scripts are available from the community to perform FASTA file manipulations. Online toolboxes are also available such as FaBox^[17] or the FASTX-Toolkit within Galaxy servers.^[18] For instance, these can be used to segregate sequence headers/identifiers, rename them, shorten them, or extract sequences of interest from large FASTA files based on a list of wanted identifiers (among other available functions). A tree-

FASTQ format is a form of FASTA format extended to indicate information related to sequencing. It is created by the Sanger Centre in Cambridge. [3]

based approach to sorting multi-FASTA files (TREE2FASTA^[19]) also exists based on the coloring and/or annotation of sequence of interest in the FigTree viewer. Additionally, Bioconductor.org's Biostrings package can be used to read and manipulate FASTA files in R. [20] Several online format converters exist to rapidly reformat multi-FASTA files to different formats (e.g. NEXUS, PHYLIP) for their use with different phylogenetic programs (e.g. such as the converter available on phylogeny.fr. [21]

The FASTQ format, used to represent DNA sequencer reads along with quality scores.

• The GVF format (Genome Variation Format), an extension based on the GFF3 format.

References [edit] 1. ^ Lipman DJ, Pearson WR (March 1985). "Rapid and sensitive protein 9. A "Alignment Fileformats" 2. 22 May 2019. Retrieved 22 May 2019. similarity searches". Science. 227 (4693): 1435-41. 10. ^ Pinho AJ, Pratas D (January 2014). "MFCompress: a compression tool for FASTA and multi-FASTA data" ∠. Bioinformatics. 30 (1): 117–8.

• The SAM format, used to represent genome sequencer reads, generally but not necessarily after they have been aligned to genome sequences. [22]

Bibcode:1985Sci...227.1435L 2. doi:10.1126/science.2983426 2. PMID 2983426 ℃.

PMID 20015970 ℃.

- Bibcode:1988PNAS...85.2444P \(\mathbb{L}\). doi:10.1073/pnas.85.8.2444 \(\frac{1}{2}\). PMC 280013 ∂. PMID 3162770 ₺. 12. A Kryukov K, Ueda MT, Nakagawa S, Imanishi T (July 2020). 3. ^ a b Cock PJ, Fields CJ, Goto N, Heuer ML, Rice PM (April 2010). "The "Sequence Compression Benchmark (SCB) database—A Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants" . Nucleic Acids Research. 38 (6):
- 4. ^ "What is FASTA Format?" ☑. zhanglab.ccmb.med.umich.edu. explains the FASTA format 5. ^ NCBI C++ Toolkit Book ☑. National Center for Biotechnology Information. Retrieved 2018-12-19.

2. ^ Pearson WR, Lipman DJ (April 1988). "Improved tools for biological

sequence comparison" . Proceedings of the National Academy of

Sciences of the United States of America. 85 (8): 2444-8.

1767–71. doi:10.1093/nar/gkp1137 \(\mathbb{Z}\). PMC 2847217 \(\mathrea\).

Learning Center]. National Center for Biotechnology Information. Retrieved 2012-03-15. 7. ^ "IUPAC code table" ☑. NIAS DNA Bank. Archived from the original ☑.

6. * Tao Tao (2011-08-24). "Single Letter Codes for Nucleotides" . [NCBI

- External links [edit] Bioconductor ☑
- on 2011-08-11. 8. ^ "anysymbol" ☑. MAFFT - a multiple sequence alignment program.
- comprehensive evaluation of reference-free compressors for FASTAformatted sequences" . GigaScience. 9 (7): giaa072. doi:10.1093/gigascience/giaa072 d. PMC 7336184 PMID 32627830 ₺. 13. A Pratas D, Hosseini M, Pinho A (2017). "Cryfa: a tool to compact and encrypt FASTA files". 11th International Conference on Practical

Applications of Computational Biology & Bioinformatics (PACBB).

Advances in Intelligent Systems and Computing. Vol. 616. Springer.

14. A Hosseini M, Pratas D, Pinho A (2018). Cryfa: a secure encryption tool

for genomic data. Bioinformatics. Vol. 35. pp. 146-148.

pp. 305–312. doi:10.1007/978-3-319-60816-7_37 \(\mathbb{Z}\). ISBN 978-3-319-

compression methods for biological sequences. Information 7(4):(2016):

doi:10.1093/bioinformatics/btt594 \(\mathbb{L}\). PMC 3866555 \(\precent{\partial}\).

11. ^ M. Hosseini, D. Pratas, and A. Pinho. 2016. A survey on data

PMID 24132931 ℃.

60815-0.

- doi:10.1093/bioinformatics/bty645 2. PMC 6298042 6. PMID 30020420 ℃.
- scientific software with Galaxy ToolShed" ☑. Genome Biology. 15 (2): 403. doi:10.1186/gb4161 2. PMC 4038738 3. PMID 25001293 2. 19. A Sauvage T, Plouviez S, Schmidt WE, Fredericq S (March 2018). "TREE2FASTA: a flexible Perl script for batch extraction of FASTA sequences from exploratory phylogenetic trees" . BMC Research Notes. 11 (1): 403. doi:10.1186/s13104-018-3268-y \(\mathbb{Z}\). PMC 5838971 \(\partial\). PMID 29506565 ₺.

15. ^ "Description of A2M alignment format" ☑. SAMtools.

doi:10.1111/j.1471-8286.2007.01821.x ₺.

16. ^ "soedinglab/hh-suite: reformat.pl" ∠. GitHub. 20 November 2022.

18. A Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N,

Galaxy Team, Taylor J, Nekrutenko A (2014). "Dissemination of

17. A Villesen P (April 2007). "FaBox: an online toolbox for fasta

sequences". Molecular Ecology Resources. 7 (6): 965-968.

20. A Pagès, H; Aboyoun, P; Gentleman, R; DebRoy, S (2018). "Biostrings: Efficient manipulation of biological strings" ☑. Bioconductor.org. R package version 2.48.0. doi:10.18129/B9.bioc.Biostrings ♂. 21. A Dereeper A, Guignon V, Blanc G, Audic S, Buffet S, Chevenet F, Dufayard JF, Guindon S, Lefort V, Lescot M, Claverie JM, Gascuel O (July 2008). "Phylogeny.fr: robust phylogenetic analysis for the non-

doi:10.1093/nar/gkn180 2. PMC 2447785 6. PMID 18424797 2.

22. ^ https://samtools.github.io/hts-specs/SAMv1.pdf [pare URL PDF]

[hide]

specialist" . Nucleic Acids Research. 36 (Web Server issue): W465-9.

- **Bioinformatics** • Sequence databases: GenBank, European Nucleotide Archive and DNA Data Bank of Japan • Secondary databases: UniProt, database of protein sequences grouping together Swiss-Prot, TrEMBL and Protein Information Resource • Other databases: Protein Data Bank, Ensembl and InterPro • Specialised genomic databases: BOLD, Saccharomyces Genome Database, FlyBase, VectorBase, WormBase, Rat Genome Database, PHI-base, Arabidopsis Information Resource and Zebrafish Information Network BLAST · Bowtie · Clustal · EMBOSS · HMMER · MUSCLE · SAMtools · SOAP suite · TopHat Software
- Broad Institute · China National GeneBank (CNGB) · Computational Biology Department (CBD) · Microsoft Research University of Trento Centre for Computational and Systems Biology (COSBI) · Database Center for Life Science (DBCLS) · DNA Data Bank of Japan (DDBJ) · European Bioinformatics Institute (EMBL-EBI) · European Molecular Biology Laboratory (EMBL) · Flatiron Institute · J. Craig Venter Institute (JCVI) · Max Planck Institute of Molecular Cell Biology and Genetics (MPI-CBG) · US National Center for Biotechnology Information (NCBI) · Japanese Institute of Genetics · Netherlands Bioinformatics Centre (NBIC) · Philippine Genome Center (PGC) · Scripps Research · Swiss Institute of Bioinformatics (SIB) · Wellcome Sanger Institute · Whitehead Institute African Society for Bioinformatics and Computational Biology (ASBCB) · Australia Bioinformatics Resource (EMBL-AR) · European Molecular Biology network (EMBnet) · International Nucleotide Sequence Database Collaboration (INSDC) · **Organizations**
- Categories: Bioinformatics | Biological sequence format

This page was last edited on 21 November 2022, at 11:29 (UTC). organization.

Text is available under the Creative Commons Attribution-ShareAlike License 3.0; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit

Databases

FASTX-Toolkit

FigTree viewer ☑

Phylogeny.fr ☑

GTO ☑

V •T •E

Server: ExPASy · Ontology: Gene Ontology · Rosalind (education platform) Institutions

Meetings Research in Computational Molecular Biology (RECOMB)

Privacy policy About Wikipedia Disclaimers Contact Wikipedia Mobile view Developers Statistics Cookie statement

CRAM format · FASTA format · FASTQ format · NeXML format · Nexus format · Pileup format · SAM format · Stockholm format · VCF format File formats

Related topics Computational biology · List of biological databases · Molecular phylogenetics · Sequencing · Sequence database · Sequence alignment (h) Category · 🚵 Commons

International Conference on Computational Intelligence Methods for Bioinformatics and Biostatistics (CIBB) · ISCB Africa ASBCB Conference on Bioinformatics · Pacific Symposium on Biocomputing (PSB) ·

International Society for Biocuration (ISB) · International Society for Computational Biology (ISCB) (Student Council (ISCB-SC)) · Institute of Genomics and Integrative Biology (CSIR-IGIB) · Japanese Society for Bioinformatics (JSBi) Basel Computational Biology Conference ([BC2]) · European Conference on Computational Biology (ECCB) · Intelligent Systems for Molecular Biology (ISMB) · International Conference on Bioinformatics (InCoB) ·