

PubDNA Finder in a Nutshell

Searching the Life Sciences Literature with Sequences of Nucleic Acids

Biomedical researchers and clinicians working with molecular technologies in routine clinical practice often need to review the available literature to gather information regarding specific sequences of nucleic acids. This includes, for instance, finding articles related to a concrete DNA sequence, or identifying empirically-validated primer/probe sequences to evaluate the presence of different micro-organisms.

Unfortunately, these hard and time-consuming tasks often need to be manually performed by researchers themselves since no publicly available biomedical literature search engine, e.g. PubMed, PubMed Central, etc., provides the required search functionalities. In this article, we describe PubDNA Finder, a web service that enables users to perform advanced searches on PubMed Central-indexed full text articles with sequences of nucleic acids.

Searching the Life Sciences Literature

PubDNA Finder [1] is a web service we developed linking more than 180,000 full text articles available at PubMed Central (PMC) at the time

of writing, to the DNA/RNA sequences appearing in them. PubDNA Finder extends the functionality provided by the PMC search engine by enabling researchers to perform queries involving both keywords and DNA/RNA sequences. To our knowledge, PubDNA Finder is the first search engine providing such advanced search capabilities. PubDNA Finder can be accessed free of charge at <http://servet.dia.fi.upm.es:8080/pubdnafinder>

Search Functionalities provided by PubDNA Finder

Researchers using PubDNA Finder can perform three different types of queries: (1) sequence-



Dr. Miguel García-Remesal, Universidad Politécnica de Madrid

based queries, (2) keyword-based queries and (3) combined queries. A detailed description of each type of query follows.

Sequence-based Queries

Sequence-based queries (SBQs) are targeted at retrieving all articles mentioning the DNA/RNA sequences specified by the user. Users can perform two different types of SBQs: simple and complex, depending on how the target sequences are specified.

Simple SBQs involve one or more DNA/RNA sequences linked by a single logical operator. Sequences are represented as strings composed of symbols belonging to the IUPAC standard



PubDNA Finder
A web database linking PubMed Central full-text articles to sequences of nucleic acids

Sequences
Operator: **OR** Detect Sequences

Free Text
Search text:

New lines as separators
 Blanks as separators

Enviar Restaurar Clear fields

Expand all Collapse all Report false positives
Showing results 1 to 3 of 3

Pubmed Central ID	Article Title
PMC: 1065263	Competitive enzymatic reaction to control allele-specific extensions
Sequence	Context
tggggcagaggggacgggaaa	13 gggcaacattctgcttaag aggagagacacatggggc cattctgcttaagcgcaatttc tggggcagaggggacgggaaa
PMC: 1092277	The third helix of the homeodomain of paired class homeodomain proteins acts as a recognition helix both for DNA and protein interactions
Sequence	Context
TGGTCTCGAGATTTTTGCAGCAAGTCTTTCTCG	57-GAGCAGCGAAATGGGCGAGGGAGAAAAG-3? DelB 57- GTAGCTCGAGCTTCAAACCTCTTTTCA-3? DelC 57- TGGTCTCGAGATTTTTGCAGCAAGTCTTTCTCG 59 z1.pax8-HD1.5? 57- CTAGATCTGAGGCTTCAGCTTAAACGAAAC-3? z1.pax8-HD1.3? 57- TACGAGCTCGGC-TTGTCTCTTTGATTTCTAAC-3?
PMC: 1145188	Development and characterization of positively selected brain-adapted SIV
Sequence	Context
ACTTCTCGATGGCAGTGACC	Rev 2 AGAGGGTGGGGAAGAGAACTG Rev 3 AGTCTCGATGGCAGTGACC Rev 4 CCAGACATAATGAGACTGTAA Rev 5

Showing results 1 to 3 of 3

W3C HTML 4.01 W3C CSS

Proudly Powered by PrimerXtractor and Apache Lucene © 2010 Biomedical Informatics Group

Fig. 1: A screenshot showing the results of the execution of a sample simple SBQ.

nucleotide codes. To execute a simple SBQ, we would have to type all the target sequences, one per line, in the text box labeled with "Sequences", select either the AND or OR operator in the "Operator" combo box, and click on the "Submit" button. For each hit in the results set, the user would be presented with the relevant information on the manuscript. This includes the PubMed Identifier (PMCID) associated with the article, the article's title, the genetic sequences – mentioned in the paper – that match the user query, the context in which each matched sequence occurs, and a link to the full text of the article.

For instance, as shown in figure 1 if we launched the query "tggggcagaggggacgggaaa OR acttctcgatggcagtgacc OR tggctcagatTTTTGCAGCAAGTCTTTCTCG", we would be presented with all papers in the database containing at least one of the three sequences specified in the query.

On the other hand, advanced SBQs involve complex sub-searches such as wildcard searches, fuzzy searches and proximity searches. We briefly describe each complex search type below.

Wildcard searches enable users to use the single and multiple character wildcard symbols, "?" and "*" respectively, to define patterns

PubDNA Finder
A web database linking PubMed Central full-text articles to sequences of nucleic acids

Sequences
Operator: **AND** Detect Sequences

Free Text
Search text:

New lines as separators
 Blanks as separators

Enviar Restaurar Clear fields

Expand all Collapse all Report false positives
Showing results 1 to 5 of 5

Pubmed Central ID	Article Title
PMC: 1939728	STAT5 is an Ambivalent Regulator of Neutrophil Homeostasis
Sequence	Context
AAGGGACAGGAAGAGAGAAGG	[23]. The primers consisted of STAT5a: F9 primer (5'-AAGGGACAGGAAGAGAGAAGG-3'), R1 primer (5'-CCATACACACTTGCATCT-3?); herpes simplex virus thymidine kinase
CCCATACAACACTTGCATCT	consisted of STAT5a: F9 primer (5'-AAGGGACAGGAAGAGAGAAGG-3'), R1 primer (5'-CCATACACACTTGCATCT-3?) herpes simplex virus thymidine kinase (TK) cassette: TKp
GCAAACCACACTGCTCGAC	simplex virus thymidine kinase (TK) cassette: TKp primer (5'-GCAAACCACACTGCTCGAC-3?); STAT5b: R8 primer (5'-GGAGATCTGCTGGTCAAAG-3?); F11 primer (5'-TCAAACACACTCAATTAGTC-3?); A
GGAGATCTGCTGGTCAAAG	(TK) cassette: TKp primer (5'-GCAAACCACACTGCTCGAC-3?); STAT5b: R8 primer (5'-GGAGATCTGCTGGTCAAAG-3?); F11 primer (5'-TCAAACACACTCAATTAGTC-3?); A representative PCR performed on
TCAAACACACTCAATTAGTC	primer (5'-GCAAACCACACTGCTCGAC-3?); STAT5b: R8 primer (5'-GGAGATCTGCTGGTCAAAG-3?); F11 primer (5'-TCAAACACACTCAATTAGTC-3?); A representative PCR performed on tail DNA from
TGCTTAAGTCCCTGGAGCAA	nM of specific primers for either murine G-CSF (TGCTTAAGTCCCTGGAGCAA) or murine STAT5b (AGCTTGAAGTCCCTGGAGCAA) and
AGCTTGAAGTCCCTGGAGCAA	(CAGGTGGTCCCGAGTTGCA and CAGATCGAAGTCCCGATCGGTA). Real time

Showing results 1 to 5 of 5

W3C HTML 4.01 W3C CSS

Proudly Powered by PrimerXtractor and Apache Lucene © 2010 Biomedical Informatics Group

Fig. 2: A screenshot showing the results of the execution of a sample BQ.

for matching the target sequences. For instance, the sample query "cga?ttg OR tta*" would retrieve papers containing sequences such as "cgacttg" or "ttatttcc".

By contrast, fuzzy searches are aimed at performing approximate matching by retrieving manuscripts containing sequences that are "similar" to these specified in the query. The similarity between two sequences is calculated using the Levenshtein Distance [2]. These searches can be performed by appending a tilde character at the end of the target sequence. It is also possible to optionally specify a similarity threshold. The latter is a value between 0 and 1. The

greater the threshold, the more similar are the matched sequences to the target sequence. For instance, if we issued the query "cgattg~0.6", we would retrieve articles containing sequences such as "ctgatcg" or "tgcattg". Conversely, if we executed the query "cgattg~0.8", we would retrieve papers containing sequences such as "cggattg" or "cgacttg".

Proximity searches are aimed at retrieving articles that contain two specific sequences which are within a given distance, i.e. a number of words, away. Proximity searches can be performed by enclosing the target sequences between double quotes and appending the tilde

character plus the distance threshold after the last double quote character. For instance, the query „cacttggaaaacgctacttcagacgcttcattctgctgtttgtg”~3 would retrieve the article with PMID 2374257, which mentions both target sequences within a distance of two words – note that the original query requires both sequences being at a distance of at most three words.

Keyword-based Queries

Keyword-based queries (KBQs) are aimed at retrieving all DNA/RNA sequences mentioned in papers matching the search terms, a functionality that is also missing in the PMC search engine. KBQs are composed of either keywords or phrases – i.e. sequences of keywords enclosed between double quotes – linked by explicitly using the AND and OR logical operators. For instance, the KBQ ‘probe OR probe AND “E. coli”’ would retrieve all

the sequences mentioned in articles that contain the phrase “E. coli” and either the word “primer” or “probe” – or both. It is also possible to use wildcard, fuzzy and proximity modifiers in KBQs if required. For instance, to search for primer/probe sequences for the Herpes virus, we could execute the following KBQ “herpes primer”~10 OR “herpes probe”~10’. As shown in figure 2, the system would return all the sequences mentioned in articles in which the word “herpes” co-occurs either with “primer” or “probe” within distance 10.

Combined Queries

Combiner queries (CQs) combine the results of a SBQ and a KBQ by means of an AND operation, thus retrieving the records of all articles matching both queries. For each hit in the results set – i.e. papers matching the KBQ and containing any sequence specified in the SBQ

– , the system presents the user with the article’s PMID, its title, a link to its full text and a list of the sequences mentioned in the article that match the SBQ, together with the context in which they occur. Figure 3 shows the result of executing a combined query that is aimed at determining whether there are any sequences beginning with either “CTTCTAAC” or “ATAGTTC” that are somehow connected to the H1N1 virus or the swine flu disease.

Additional Features

PubDNA finder provides other additional features, such as automatically identifying and extracting all sequences mentioned in a plain-text document provided by the user, or retrieving all sequences

mentioned in a concrete article identified by its PMID.

References

- [1] García-Remesal M. et al.: Bioinformatics 26(21), 2801–2802 (2010)
- [2] Levenshtein V.I.: Soviet Physics Doklady 10, 707–10 (1966)
- [3] García-Remesal M. et al.: BMC Bioinformatics 11, 410 (2010)

Contact

Dr. Miguel García-Remesal
 Departamento de Inteligencia Artificial
 Facultad de Informática
 Universidad Politécnica de Madrid
 Campus de Montegancedo S/N
 Madrid, Spain
 mgarcia@infomed.dia.fi.upm.es

www.laboratory-journal.com

PubDNA Finder
 A web database linking PubMed Central full-text articles to sequences of nucleic acids

Sequences: Operator: LUCENE Detect Sequences
 CTTCTAAC" OR ATAGTTC"
 New lines as separators
 Blanks as separators

Free Text: Search text:
 "H1N1 primer"~3 OR "H1N1 probe"~3 OR "swine flu" OR H1N1

Enviar Restaurar Clear fields

Expand all Collapse all Report false positives Showing results 1 to 6 of 6

PubMed Central ID	Article Title
PMC: 1524785	A sensitive one-step real-time PCR for detection of avian influenza viruses using a MGB probe and an internal positive control
Sequence	Context
CTTCTAACCGAGGTCGAAACGTA	LOCATION (nt) SENSE M-Flu1 CTTCTAACCGAGGTCGAAACGTA 32-54 + M-Flu2 GGATTGGTCTTGTCTTAGCCA
PMC: 1876802	First introduction of highly pathogenic H5N1 avian influenza A viruses in wild and domestic birds in Denmark, Northern Europe
Sequence	Context
CTTCTAACCGAGGTCGAAACG	bromide. RT-PCR primers were matrix forward primer FB-A1-M52c: 5'-GCTT GTA ACC GAG GTC GAA AGG-3'; matrix reverse primer FB-A1-M253R: 5'-AGG GCA TTT TGG
PMC: 2657637	Oseltamivir-Resistant Influenza Viruses A (H1N1), Norway, 2007-08
Sequence	Context
CTTCTAACCGAGGTCGAAACG	reaction mixture containing 0.3 μM forward primer M52c (5'-CTTCTAACCGAGGTCGAAACG-3'); 0.3 μM reverse primer M149r (5'-CTT GTC TTT
PMC: 2770640	Novel Pandemic Influenza A(H1N1) Viruses Are Potently Inhibited by DAS181, a Sialidase Fusion Protein
Sequence	Context
CTTCTAACCGAGGTCGAAACG	0.5 mM of each primer. The forward primer (5'-CTTCTAACCGAGGTCGAAACG-3') and the reverse primer (5'-GGCATTGTGGACAAKCGTCTA-3') were used for
PMC: 2859049	Characterization of Quasispecies of Pandemic 2009 Influenza A Virus (A/H1N1/2009) by De Novo Sequencing Using a Next-Generation DNA Sequencer
Sequence	Context
ATAGTTCATTCTCCCTCTTGACC	common primer: pdmFlu09-HA.F, 5'-CGAACAAAGGTGTAACGGCAGCAT-3'; HA-Major-specific reverse primer: pdmFlu-HA-R_Major, 5'-ATAGTTCATTCTCCCTCTTGACC-3'; HA-Minor-specific reverse primer: 5'-ATAGTTCATTCTCCCTCTTGATT-3'), and the SuperScript III

Fig. 3: A screenshot showing the results of the execution of a sample CQ.