

PNA Finder

A bioinformatics-based tool for rapid antisense PNA design

Overview

PNA Finder is a Python-based toolbox that combines the functions of several bioinformatics alignment programs with a series of custom scripts as a means to design antisense peptide nucleic acids (PNA). The *PNA Finder v0.1* consists of the *Get Sequences* and *Find Off-Targets* function, which can both be used in tandem to identify, design, and screen antisense sequences in a high-throughput manner.

Table of Contents

Setup.....	1
<i>Installing PNA Finder</i>	2
<i>PNA Finder start_info</i>	2
Running PNA Finder	3
<i>Get Sequences</i>	3
<i>Find Off-Targets</i>	7
PNA Finder Examples	9
References.....	10

Setup

Prerequisite programs

The *PNA Finder* toolbox is built in Python 3.7 to run on a Windows machine. It requires the following programs:

- Python (\geq version 2.7)
- Alignment programs
 - Bowtie 2
 - BEDTools
 - SAMTools

In order to run the toolbox, it is also necessary to install a Bash shell that is able to run these alignment programs. *PNA Finder* has built-in compatibility with both Cygwin and the Windows 10 Bash shell (other shell options will require slight modification to the *PNA Finder* scripts). All three alignment programs must be compiled on the Bash shell prior to using *PNA Finder* (see the programs' respective online manuals for installation instructions).

Installation should be done such that the following commands successfully display each program's installation version (calling each program should not require path specification or a different aliased command):

```
bowtie2 --version
bowtie2-build --version
samtools --version
bedtools -version
```

Installing PNA Finder

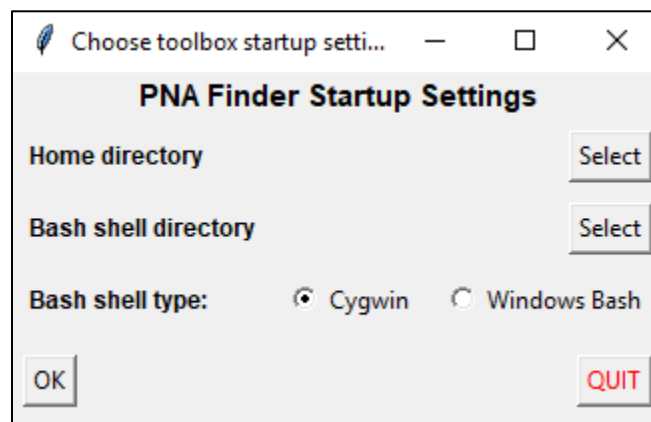
A standard command line pip command may be used for package installation:

```
pip install pna_finder
or
python -m pip install pna_finder
```

PNA Finder start_info

Before running *PNA Finder* for the first time, it is necessary to load the `start_info` module and call the `run()` function. This can be done using a Python shell or by running the script `run_start_info.py` in the package's top-level directory.

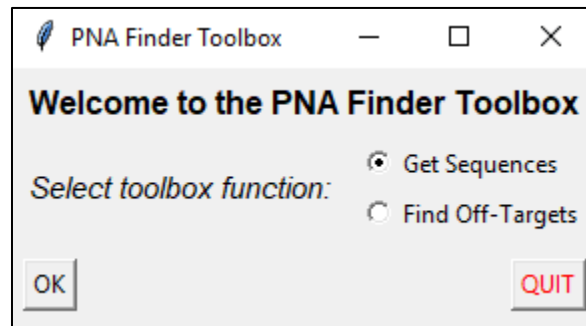
Calling this function will open the following dialog box, which allows the user to select a home directory for *PNA Finder*'s file browsing, as well as the directory where the desired `bash.exe` program is located. The dialog box also has two radio buttons for the user to select which of the supported Bash shells they have chosen to use. The user may then select OK to save these settings, and they will be saved in the module folder as `start_info.txt`. The settings can be changed by calling `start_info.run()` again.



Running PNA Finder

Similar to the `start_info.run()` command, the main functions of *PNA Finder* can be run by loading the `pna_finder` module and calling the function `run()`. This can be done using a Python shell or by running the script `run_pna_finder.py` in the package's top-level directory.

Calling this function will open the following dialog box. The user must then select either the *Get Sequences* or *Find Off-Targets* function and select OK to continue.



Get Sequences

The *Get Sequences* function is used to build a set of antisense PNA candidates based on a user-provided list of target genes. The setup dialog box for *Get Sequences* is shown below, with each individual section numbered in red:

PNA Finder - Get PNA Sequences

Job Name: 1

Sequence Window Start:

Sequence Window End: 2

PNA Length:

Annotation Record Types (comma separated):

3

Select Annotation File Option:

☒ Build new gffutils database from GFF file

☐ Select existing gffutils database file

☐ Run STRING analysis

☒ Run sequence warnings analysis 4

Species NCBI Taxon ID

1. Job Name

The results of the *Get Sequences* job will be written to a directory with this title.

2. Get Sequences parameters

The parameters “Sequence Window Start” and “Sequence Window End” determine the number of PNA sequences that will be produced for each gene target submitted. The integers refer to coordinates relative to the feature (e.g. gene, exon, mRNA) start locus. As an example, take the nucleotide sequence surrounding the start codon for the *E. coli* gene *thrL*. Using the default entries of -5 for both window parameters, as well as the default entry of 12 for the “PNA Length” parameter, we are provided with a single 12-nt PNA target sequence. This sequence is highlighted in the nucleotide sequence below:

```

E. coli MG1655 thrL
-20          -5    0                      +20
TTACAGAGTACACAA CATCCATGAAAC GCATTAGCACCAC

```

However, if instead the “Sequence Window End” parameter is changed to 0, we obtain six 12-nt PNA target sequences, as shown below:

```

E. coli MG1655 thrL
-20      -5      0      +20
TTACAGAGTACACAA CATCCATGAAAC GCATTAGCACCAC
                   ATCCATGAAACG
                   TCCATGAAACGC
                   CCATGAAACGCA
                   CATGAAACGCAT
                   ATGAAACGCATT

```

PNA Finder will raise an error message if the “Sequence Window End” parameter is less than the “Sequence Window Start” parameter. Target sequence length can be adjusted by changing the “PNA Length” parameter.

The defaults of [-5,-5] for the sequence window and 12 for the PNA length are based on observations from prior research. PNA that binds complementarily to the mRNA start codon has been found to be the most inhibitory of protein translation,¹⁻⁵ and PNA lengths between 8 and 12 bases have been shown to most effectively balance translation inhibition with intracellular delivery.⁶ Additionally, 12 bases is long enough such that the expected number of off-targets in even the largest bacterial genomes will be less than one (under the simplifying assumption of nucleotide sequence randomness).

3. File upload and annotation search options

The main sequence search function of *Get Sequences* works in two parts: (1) searching a genome annotation file for records matching the user provided gene ID list and (2) using the corresponding genome assembly to extract the desired translation start site nucleotide sequences for these records.

The “Annotation Record Types” entry will depend on the contents of the GFF/GTF annotation file that is being used. *Get Sequences* will search through the annotation file for only records that are labeled with one of the comma-separated list of feature types provided in this box. For example, the standard label for coding sequences in NCBI RefSeq annotation is “CDS,” so this is used as a default. If this entry is left blank, all feature types are included in the search.

(Note: The label “gene” is often applicable and equally useful when using prokaryotic genomes. With eukaryotic genomes, the labels “mRNA” and “exon” are often more useful, depending on the PNA application.)

The first file upload is the gene ID list. This should be formatted as a plaintext single-column file that contains a list of only gene IDs, with only one per line. Different types of gene ID can be included within the same file. The following GFF/GTF annotation identifiers are currently supported:

GFF: gene, Name, gene_synonym, protein_id, locus_tag, Dbxref

GTF: gene_id, gene_name, transcript_id, tss_id

The next upload is a genome assembly. This file should be in the FASTA file format.

Next, the user must select an option for handling the annotation file. *Get Sequences* uses the Python package gffutils to parse the annotation file, which requires creation of a gffutils database file. If the database has not yet been created, the user should select the first option and use the “Select Annotation” file below to upload a GFF/GTF file. *Get Sequences* will then automatically create the database file and use it for the gene ID search. If the required database has already been created, the user should select the second option and upload the database file.

The final selection is an output directory. This will be the location in which *Get Sequences* will create a new results directory (named for the “Job Name” entry).

4. PNA candidate analysis options

Get Sequences provides two options for gene target and sequence analysis of PNA candidates. The first is to perform an analysis of each target gene’s protein interaction network using the STRING database.⁷ The density of network connections for targeted genes/proteins has been shown to be positively correlated with PNA knockout growth inhibition of bacterial cultures. In order to use this option, the NCBI Taxon ID for the species of interest must be entered as well. The integer ID number can be looked up at the following web address: <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi> (Note: This option will only work if the STRING database can recognize the PNA candidates as gene names)

The second option is to analyze the PNA candidate sequences for potential solubility issues based on purine content, as described by Gildea et al.⁸ *Get Sequences* will produce a solubility warning if it finds, within a stretch of 10 bases, five purines in a row, four guanines in a row, or more than six total purines. Warnings are also produced for a PNA molecule longer than 30 bases or with more than 6 bases of a self-complementary subsequence.

Once all information has been entered into the *Get Sequences* dialog box, the user must click “SUBMIT” to run the *Get Sequences* function. This will produce five files in total, including two BED files, one FASTA file, one tab-delimited “.out” file of PNA sequences, and one BEDTools error record file.

The first BED file will have the same name as the gene ID list file, except with the “.bed” extension. It contains a single record for each gene ID that was matched to an entry in the genome annotation file. The second BED file is an edited version of the first, with the records replicated for as many PNA that will be designed per gene, and indices changed to match the locations of the desired sequences. The BEDTools error record file will contain error information

about the ‘bedtools getfasta’ command within *Get Sequences*. It will be empty if the job completes successfully.

The FASTA file will contain a FASTA-formatted list of all candidate PNA target sequences. This may be used as the input for the *Find Off-Targets* function. The tab-delimited “.out” file will contain a table of PNA sequences (the reverse complements of the target sequences), sequence warnings, and STRING protein network data, if the latter two options were selected in the main *Get Sequences* window.

Find Off-Targets

The *Find Off-Targets* function is used to search genomes for PNA off-targets that are likely to inhibit gene expression and/or RNA translation. The setup dialog box for *Find Off-Targets* is shown below, with each individual section numbered in red:

PNA Finder - Find Off-Targets

Job Name: 1

Off-Target Window Start:

Off-Target Window End: 2

Maximum Mismatches:

PNA Length:

Annotation Record Types (comma separated):

3

Select Index File Option:

☒ Build new index from FASTA file

☐ Select existing index file (any of the 6)

☒ Create off-target count file 4

☐ Exclude PNA self-targeting

1. Job Name

As with *Get Sequences*, the results for the *Find Off-Targets* job will be written to a

directory with this title.

2. ***Find Off-Targets* parameters**

The integer entries for “Off-Target Window Start” and “Off-Target Window End” refer to coordinates relative to the feature (e.g. gene, exon, mRNA) start locus. These coordinates are primarily used to search for off-target binding sequences that are most likely to inhibit mRNA translation. The default region for expected inhibition is set to [-20, 20] relative to the RNA start codon, based on prior work.¹

The parameter “Maximum Mismatches” allows the user to set either zero or one mismatch tolerance in the search for off-target alignments. PNA have shown very high mismatch discrimination in short sequences,^{1,9-12} and as such the default mismatch tolerance is set to zero. “PNA Length” must be specified by the user to run the appropriate Bowtie 2 alignment command. PNA of differing lengths cannot be combined into the same job.

3. **File upload and annotation search options**

As with *Get Sequences*, the “Annotation Record Types” entry will depend on the contents of the GFF/GTF annotation file. *Find Off-Targets* will search the uploaded genome for PNA alignments that are found within the off-target window for the types of features specified in this box. The default window, [-20, 20], is most applicable to prokaryotic coding sequence features (named “CDS” in NCBI RefSeq annotations), but can also knock down expression and/or translation in eukaryotic organisms as well.

The first file upload is the PNA Targets FASTA file. This can be manually constructed or taken from the output of *Get Sequences*. The FASTA file should contain a list of PNA target sequences, not the PNA sequences themselves.

Next, the user must select an option for the Bowtie 2 index. To use a genome assembly in Bowtie 2, a set of index files must first be created. To automatically create this index and use it for the *Find Off-Targets* function, the user should select the first option and use the “Select Assembly/Index” button to upload the genome FASTA file. If the index has already been created, the user should select the second option and use the “Select Assembly/Index” button to upload any of the six Bowtie 2 index files.

The next upload is the annotation file, which should be in GFF or GTF format.

As with *Get Sequences*, the final selection is an output directory. This will be the location in which *Find Off-Targets* will create a new results directory (named for the “Job Name” entry).

4. **Count file options**

For convenience, the last two options relate to an off-target count file, which tallies the

number of off-targets for each PNA candidate. The first option allows the user choose whether to create this file, and the second allows the user to specify whether they are searching for off-targets in the genome from which the PNA candidates originated. If this is the case, then the sequence that the PNA is intended to complement will be excluded from the off-target count.

Once all information has been entered into the *Find Off-Targets* dialog box, the user must click “SUBMIT” to run the *Find Off-Targets* function. This will produce 11 or 12 files in total (depending on whether the count file option was selected).

The initial Bowtie 2 alignment will produce a SAM file of all alignments that were found for the given mismatch setting, which will then be converted into a sorted BAM file through a series of SAMTools commands. Each of these commands produces at least one intermediate file, though these will likely not be useful output for the user. The Bowtie 2 and SAMTools commands will also produce individual error record files in the “error_files” subdirectory within the main job folder.

There will be a BED file produced based on these alignments, which contains all alignments within a certain distance (as defined by “Off-Target Window Start”) from any feature in the annotation file. Based on this BED file, a tab-delimited “.out” file will also be produced, containing all PNA alignments to the user-designated off-target window with the feature types specified in the *Find Off-Targets* dialog. Finally, if the user has selected the count file option, a “.count” file that tabulates total inhibitory off-target alignments for each PNA candidate will also be produced in the job directory.

PNA Finder Examples

Example files for the inputs and outputs of PNA Finder can be found in the “examples” subdirectory of the package folder.

The *Get Sequences* example uses the text file “ecoli_gene_id.txt” in the directory “gene_id”. This file lists different annotation identifiers for the first six coding sequences of the reference *E. coli* MG1655 annotation file, located in the directory “genome_files/ecoli_MG1655”. This text file, along with the annotation file and genome assembly FASTA file, can be uploaded to PNA Finder’s *Get Sequences* function to produce the output found in the “get_sequences_output/gs_ecoli_example” directory.

The *Find Off-Targets* example uses the FASTA file output found in that same “get_sequences_output/gs_ecoli_example” directory, together with the *E. coli* MG1655 annotation file and Bowtie 2 index (located in “genome_files/ecoli_MG1655”), to produce the

output found in the “find_off-targets_output” directory. Two example outputs can be found in this directory, one for a zero-mismatch screening and the other for one-mismatch screening.

References

1. Courtney, C. M. & Chatterjee, A. Sequence-Specific Peptide Nucleic Acid-Based Antisense Inhibitors of TEM-1 Beta-Lactamase and Mechanism of Adaptive Resistance. *ACS Infectious Diseases* **1**, 253–263 (2015).
2. Dryselius, R., Aswasti, S. K., Rajarao, G. K., Nielsen, P. E. & Good, L. The translation start codon region is sensitive to antisense PNA inhibition in Escherichia coli. *Oligonucleotides* **13**, 427–433 (2003).
3. Otsuka, T. *et al.* Antimicrobial activity of antisense peptide-peptide nucleic acid conjugates against non-typeable Haemophilus influenzae in planktonic and biofilm forms. *Journal of Antimicrobial Chemotherapy* **72**, 137–144 (2017).
4. Mondhe, M., Chessher, A., Goh, S., Good, L. & Stach, J. E. M. Species-selective killing of bacteria by antimicrobial Peptide-PNAs. *PLoS ONE* **9**, (2014).
5. Bai, H. *et al.* Antisense inhibition of gene expression and growth in gram-negative bacteria by cell-penetrating peptide conjugates of peptide nucleic acids targeted to rpoD gene. *Biomaterials* **33**, 659–667 (2012).
6. Good, L., Awasthi, S. K., Dryselius, R., Larsson, O. & Nielsen, P. E. Bactericidal antisense effects of peptide-PNA conjugates. *Nature Biotechnology* **19**, 360–364 (2001).
7. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic acids research* **47**, D607–D613 (2019).
8. Gildea, B. D. & Coull, J. M. Methods for Modulating the Solubility of Synthetic Polymers. (2004).
9. Doyle, D. F., Braasch, D. A., Simmons, C. G., Janowski, B. A. & Corey, D. R. Inhibition of gene expression inside cells by peptide nucleic acids: Effect of mRNA target sequence, mismatched bases, and PNA length. *Biochemistry* **40**, 53–64 (2001).
10. Good, L. & Nielsen, P. E. Inhibition of translation and bacterial growth by peptide nucleic acid targeted to ribosomal RNA. *Biochemistry* **95**, 2073–2076 (1998).
11. Weiler, J., Gausepohl, H., Hauser, N., Jensen, O. N. & Hoheisel, J. D. Hybridisation based DNA screening on peptide nucleic acid (PNA) oligomer arrays. *Nucleic Acids Research* **25**, 2792–2799 (1997).
12. Choi, J. J., Jang, M., Kim, J. & Park, H. Highly sensitive PNA array platform technology for single nucleotide mismatch discrimination. *Journal of Microbiology and Biotechnology* **20**, 287–293 (2010).