

Modules

numpy

os

random

Functions

```
main()
    Generates random sequence data files.

    datadir/datatype.fasta
    ...
    >rseq_{i}_{energy}

    sequence

    structure

    >rseq_{i+1}_{energy}

    sequence

    structure

    ...
```

Modules

absl

numpy

os

Functions

```
main()
```

Modules

absl

argparse

os

tensorflow

Functions

```
main()
    RNAdeep training interface.

parse_rnadeep_args(p)
    Arguments that are used by RNAdeep.

training(datatag, dbn_dir, ali_dir, spotmodel=None, basemodel=None, savedir='.', epochs=50, epoch0=0, batch_size=4)
```

Modules

absl

argparse

os

Functions

```
main()
    RNAdeep training interface.
    python train.py -d l30 -t train_data/fixlen30_n100000.fa-train -v train_data/fixlen30_n100000.fa-valid -m 4 -b 250 -e 20 --model-
    log-dir intermediate_models -l intermediate_models/sm4_l30_010/ --epoch0 10 2>> sm4_l30.err >> sm4_l30.out &

parse_rnadeep_args(p)
    Arguments that are used by RNAdeep.

training(datatag, ftrain, fvalid, spotmodel=None, basemodel=None, savedir='.', epochs=50, epoch0=0, batch_size=4)
```

Modules

absl

numpy

os

Functions

main()

rfam\_converter

[index](#)

/home/julian-zim/Desktop/proj/rnadeep/rnaconv/rfam\_converter.py

Modules

RNA  
numpy

os  
shutil

sys  
textdistance

Functions

**convert\_rfam\_data**(seed\_filepath, tree\_dirpath, outpath)

Calls the necessary functions to convert the whole rfam database into single files, preparing them to be used by SISSI.

Parameters:

- seed\_filepath (str): Path to the Rfam.seed file in STOCKHOLM format, containing the families
- tree\_dirpath (str): Path to the directory in which Rfam tree files in newick format (.seed\_tree) files are located
- outpath (str): Path to the directory in which to save the converted data

**fix\_newick\_strings**(tree\_dirpath, outpath)

Fixes newick strings by replacing every control character (e.g. '(', ')', ',', '.', ':') within a node name with an underscore.  
Additionally, multifurcations are resolved and non-leaf node labels are removed.  
(These three steps are nessecary for SISSI to be able to parse the Rfam tree files.)

Parameters:

- tree\_dirpath (str): path to the directory containing the tree files in newick string format
- outpath (str): path to the directory in which to save the trees in the fixed newick string format.

**main**()

**obtain\_equilibrium\_frequencies**(ali\_dirpath, neigh\_dirpath, outpath)

Extracts the equilibrium frequencies for unpaired single nucleotides and nucleotide pairs from an alignment, by counting the occurences of single nucleotides in unpaired site and saving them in a 4-vector, counting the occurences of nucleotide pairs in paired sites and saving them in a 16-vector, adding pseudocounts to both (+1 for each element) and normalizing in the end.

Parameters:

- ali\_dirpath (str): Path to the directory containing the alignment files in CLUSTAL format
- neigh\_dirpath (str): Path to the directory containing the alignment consensus structure files in dot bracket notation format
- outpath (str): Path to the directory in which to save the extracted unpaired single and paired nucleotide equilibrium frequencies

**rescale\_newick\_strings**(tree\_dirpath, ali\_dirpath, outpath)

Rescales the tree branch lengths for trees which corresponding sequence alignments sequences are over 95% similar with respect to their mean pairwise hamming distance, in order to increase the evolution rate when using the tree for evolutionary simulation.  
The rescale factor is 2.

Parameters:

- tree\_dirpath (str): path to the directory containing the tree files in newick string format
- ali\_dirpath (str): path to the directory containing the alignment files in CLUSTAL format
- outpath (str): path to the directory in which to save the rescaled trees in the newick string format.

**stockholm\_to\_alignments**(filepath, outpath)

Converts the alignments contained in the STOCKHOLM input file into CLUSTAL files.

Parameters:

- filepath (str): Path to the Rfam.seed file in STOCKHOLM format, containing the families
- outpath (str): Path to the directory in which to save the extracted alignments in the CLUSTAL format

**stockholm\_to\_neighbourhoods**(filepath, outpath)

Calls the necessary functions to convert the consensus structures contained in the STOCKHOLM input file into single files in the wuss and dbn formats, respectively.

Parameters:

- filepath (str): Path to the Rfam.seed file in STOCKHOLM format, containing the families
- outpath (str): Path to the directory in which to save the extracted consensus structures

**stockholm\_to\_wuss**(filepath, outpath)

Converts the consensus structures contained in the STOCKHOLM input file into single files in the washington university secondary structure (wuss) format.

Parameters:

- filepath (str): Path to the Rfam.seed file in STOCKHOLM format, containing the families
- outpath (str): Path to the directory in which to save the resulting wuss file

## wuss\_to\_db(filepath, outpath)

Converts the consensus structures contained in the wuss input file into the dot bracket notation format

Parameters:

filepath (str): Path to the file in wuss format, containing the secondary structure

outpath (str): Path to the directory in which to save the resulting dot bracket notation file

## data\_filter

[index](#)

[/home/julian-zim/Desktop/proj/rnadeep/rnaconv/data\\_filter.py](/home/julian-zim/Desktop/proj/rnadeep/rnaconv/data_filter.py)

## Modules

[RNA](#)

[numpy](#)

[os](#)

[sys](#)

## Functions

### filter\_data(ali\_dirpath, seq\_dirpath, neigh\_dbn\_dirpath, neigh\_ct\_dirpath, max\_dbrs\_deviation=20)

Filters the data generated by the data\_generator: In each alignment, every sequence is removed that deviates by over <max\_dbrs\_deviation> from the consensus structure that was used by SISSI to generate it. This is achieved by using RNAfold to predict the secondary structure and the base pair distance to compare it to the desired consensus structure. If this results in leaving less than 2 sequences in the alignment, essentially deleting it or turning it into a single sequence, the whole alignment file with the corresponding consensus structure and sequence is removed.

Note: If families were generated, the consensus structures used for the alignment generation were generated by RNAfold and saved into the neigh\_dirpath.

If only alignments were generated, the consensus structures used for the generation were provided by the user, most likely from a converted Rfam database, but then copied to the neigh\_dirpath anyway, for the sake of integrity.

Therefore, in both cases, neigh\_dirpath can be used to retrieve the desired consensus structures to compare the sequences with.

Parameters:

ali\_dirpath (str): Path to the directory of the generated alignments.

seq\_dirpath (str): Path to the directory of the generated or copied sequences.

neigh\_dbn\_dirpath (str): Path to the directory of the generated or copied dbn files.

neigh\_ct\_dirpath (str): Path to the directory of the generated ct files.

max\_dbrs\_deviation (int, None, optional): maximum allowed base pair distance deviation from the consensus structure in percent. Default is 20.

### main()

### obtain\_and\_compare\_equilibrium\_frequencies(ali\_dirpath, neigh\_dirpath, orig\_single\_freq\_dirpath, orig\_doublet\_freq\_dirpath, outpath)

Extracts the equilibrium frequencies for unpaired single nucleotides and nucleotide pairs from the generated alignments and forms the differences to the already extracted equilibrium frequencies of the original alignments.

Parameters:

ali\_dirpath (str): Path to the directory of the generated alignment files in CLUSTAL format

neigh\_dirpath (str): Path to the directory containing the alignment consensus structure files in dot bracket notation format

orig\_single\_freq\_dirpath (str): Path to the directory of the extracted single frequency files of the original alignments

orig\_doublet\_freq\_dirpath (str): Path to the directory of the extracted doublet frequency files of the original alignments

outpath (str): Path to the directory in which to save the extracted unpaired single and paired nucleotide equilibrium frequencies and frequency differences

## rfam\_filter

[index](#)

[/home/julian-zim/Desktop/proj/rnadeep/rnaconv/rfam\\_filter.py](/home/julian-zim/Desktop/proj/rnadeep/rnaconv/rfam_filter.py)

## Modules

[os](#)

[sys](#)

## Functions

### filter\_rfam\_data(ali\_dirpath, single\_freq\_dirpath, doublet\_freq\_dirpath, neigh\_wuss\_dirpath, neigh\_dbn\_dirpath, tree\_fixed\_dirpath, tree\_rescaled\_dirpath, max\_length=700)

Filters the alignments, consensus structures, frequencies and trees of a converted rfam database (not touching the original tree files and original Rfam.seed file):

Every data point (consisting of the four filetypes mentioned above) which alignment exceeds a certain length is removed.

Parameters:

ali\_dirpath (str): Path to the directory of the converted alignments.

single\_freq\_dirpath (str): Path to the directory of the extracted single frequency files.

doublet\_freq\_dirpath (str): Path to the directory of the extracted doublet frequency files.

neigh\_wuss\_dirpath (str): Path to the directory of the extracted wuss files.

neigh\_dbn\_dirpath (str): Path to the directory of the extracted dbn files.

tree\_fixed\_dirpath (str): Path to the directory of the fixed newick string tree files.

tree\_rescaled\_dirpath (str): Path to the directory of the rescaled tree files.

max\_length (int, None, optional): Maximum allowed length of an alignment. Default is 700.

### main()

[index](#)

## Modules

[RNA](#)  
[argparse](#)

[os](#)  
[random](#)

[subprocess](#)

## Functions

**db\_to\_ct**(dbn, seq)

Converts the consensus structures contained in the dot bracket notation input file into the connect table format.

Parameters:

dbn (str): Secondary structure in dot bracket notation  
seq (str): Sequence

**generate\_alignment\_set**(sissi\_filepath, number, tree\_dirpath, neigh\_dirpath, sfreq\_dirpath, dfreq\_dirpath, ali\_dirpath, outpath)

Generates <number> alternative alignments for each tree file in the given tree-directory, searching in the respectively given consensus-structure-, single- & doublet-frequencies- and, optionally for readding ndels, alignment-directories for files of the same name to use.

Parameters:

sissi\_filepath (str): Path to the compiled sissi099 file  
number (int): The number of alignments to generate  
tree\_dirpath (str): Path to a directory containing tree files in the newick string format ('.seed\_tree')  
neigh\_dirpath (str): Path to a directory containing neighbourhood files in the dot-bracket notation format ('.dbn')  
sfreq\_dirpath (str): Path to a directory containing files storing a single frequency vector ('.sfreq')  
dfreq\_dirpath (str): Path to a directory containing files storing a doublet frequency vector ('.dfreq')  
ali\_dirpath (str, None): Path to a directory containing alignment files in the clustal format ('.aln')  
outpath (str): The Path to which to write the generated sequences, alignments & copied consensus structures

**generate\_alignments**(sissi\_filepath, number, tree\_filepath, neigh\_filepath, sfreq\_filepath, dfreq\_filepath, ali\_filepath, outpath)

Generates <number> alternative alignments for the given tree-, consensus-structure-, single- & doublet-frequencies- and, optionally for readding indels, alignment-file, using:

- RNAinverse to generate an ancestral sequence for the provided consensus structure
- SISSI simulate homologous sequence alignments (taking the generated ancestral sequence, provided tree, provided consensus structure and provided equilibrium frequencies as input).

Note:

The provided consensus structure will also be copied into an additional file per generated alignment, in order to create pairs of samples and tags to be used for training (during the process, the dbn files are converted to ct files, which are also saved).

The generated ancestral sequence will also be saved to maintain integrity.

Parameters:

sissi\_filepath (str): Path to the compiled sissi099 file  
number (int): The number of alignments to generate  
tree\_filepath (str): Path to a directory containing tree files in the newick string format ('.seed\_tree')  
neigh\_filepath (str): Path to a directory containing neighbourhood files in the dot-bracket notation format ('.dbn')  
sfreq\_filepath (str): Path to a directory containing files storing a single frequency vector ('.sfreq')  
dfreq\_filepath (str): Path to a directory containing files storing a doublet frequency vector ('.dfreq')  
ali\_filepath (str, None): Path to a directory containing alignment files in the clustal format ('.aln')  
outpath (str): The Path to which to write the generated sequences, alignments & copied consensus structures

**generate\_families**(sissi\_filepath, number, min\_length, max\_length, tree\_filepath, sfreq\_filepath, dfreq\_filepath, outpath)

Generates <number> families for the given tree- and single- & doublet-frequencies-file, using:

- random ancestral sequences of uniformly distributed lengths up to <maxlength>
- RNAfold to predict secondary structures for these sequences to be used as consensus structures for the alignment generation
- SISSI simulate corresponding homologous sequence alignments (taking the random ancestral sequences, provided tree, predicted consensus structures, and provided equilibrium frequencies as input).

Note:

During the process, the generated dbn files are converted to ct files, which are also saved.

Parameters:

sissi\_filepath (str): Path to the compiled sissi099 file  
number (int): The number of families to generate  
min\_length (int): Minimum allowed length of the ancestral sequences used to generate the families  
max\_length (int): Maximum allowed length of the ancestral sequences used to generate the families  
tree\_filepath (str): Path to a tree file in the newick string format ('.seed\_tree')  
sfreq\_filepath (str): Path to a file containing a single frequency vector ('.sfreq')  
dfreq\_filepath (str): Path to a file containing a doublet frequency vector ('.dfreq')  
outpath (str): The path to which to write the generated families

**generate\_family\_set**(sissi\_filepath, number, min\_length, max\_length, tree\_dirpath, sfreq\_dirpath, dfreq\_dirpath, outpath)

Generates <number> RNA families of uniformaly distributed lengths up to <maxlength> for each tree file in the given tree-directory, searching in the respectively given single- & doublet-frequencies-directories for files of the same name to use.

For more information, refer to the [generate\\_families\(\)](#) function.

Parameters:

sissi\_filepath (str): Path to the compiled sissi099 file  
number (int): The number of families to generate  
min\_length (int): Minimum allowed length of the ancestral sequences used to generate the families  
max\_length (int): Maximum allowed length of the ancestral sequences used to generate the families  
tree\_dirpath (str): Path to a directory containing tree files in the newick string format ('.seed\_tree')  
sfreq\_dirpath (str): Path to a directory containing files storing a single frequency vector ('.sfreq')  
dfreq\_dirpath (str): Path to a directory containing files storing a doublet frequency vector ('.dfreq')  
outpath (str): Path to which to write the generated families

**generate\_sequence\_structure\_pair**(length=85, min\_paired\_sites\_percent=20)

Repeatedly generates a random sequence and predicts its secondary structure using RNAfold, until the structure has at least min\_paired\_sites paired sites.

Parameters:  
length (int, optional): Length of the random sequence  
min\_paired\_sites\_percent (int, optional): Minimal required sites to be paired in percent

**main()**

**setup\_args(parser)**

## sampling\_ali

[index](#)  
[/home/julian-zim/Desktop/proj/rnadeep/rnadeep/sampling\\_ali.py](/home/julian-zim/Desktop/proj/rnadeep/rnadeep/sampling_ali.py)

### Modules

[numpy](#)

[os](#)

### Functions

**draw\_ali\_sets**(ali\_directory, dbn\_directory, splits=None)

**parse\_families**(ali\_dirpath, dbn\_dirpath)

Combines pairs of the same name of alignment CLUSTAL files and neighbourhood Dot Bracket String files found in the respective directories to be used for training.

Parameters:  
ali\_dirpath (str): Path to the directory containing the alignment CLUSTAL files  
dbn\_dirpath (str): Path to the directory containing the neighbourhood Dot Bracket String files

**parse\_family**(ali\_filepath, dbn\_filepath)

Reads an alignment CLUSTAL file and neighbourhood Dot Bracket String file and combines them into a pair to be used for training.

Parameters:  
ali\_filepath (str): Path to the alignment CLUSTAL file  
dbn\_filepath (str): Path to the neighbourhood Dot Bracket String file

## models

[index](#)  
</home/julian-zim/Desktop/proj/rnadeep/rnadeep/models.py>

### Modules

[keras.api.v2.keras.layers](#)

[tensorflow](#)

### Functions

**spotrna\_alignment\_models**(model=1, use\_mask=True)

Some modifications to Julia's SPOT-RNA implementations.

Supposed to be a reimplementaion of the models in the SPOT-RNA paper. If you find mistakes, please let us know!

Overview:

- Initial 3x3 convolution layer
- ResNet blocks
- Act./Norm.
- 2D-BLSTM
- Fully Connected blocks
- Output masking layer (optional)
- Output layer

Args:

model: select the model (0-4)  
use\_mask: for padded input/output (defaults to True!)

**spotrna\_models**(model=1, use\_mask=True)

Some modifications to Julia's SPOT-RNA implementations.

Supposed to be a reimplementaion of the models in the SPOT-RNA paper. If you find mistakes, please let us know!

Overview:

- Initial 3x3 convolution layer
- ResNet blocks
- Act./Norm.
- 2D-BLSTM
- Fully Connected blocks
- Output masking layer (optional)
- Output layer

Args:

model: select the model (0-4)  
use\_mask: for padded input/output (defaults to True!)

## sampling

[index](#)  
</home/julian-zim/Desktop/proj/rnadeep/rnadeep/sampling.py>

Modules

[numpy](#) [os](#) [random](#)

Functions

**draw\_sets**(fname, splits=None)

**generate\_random\_structures**(lengths)

**rseq**(l)

**write\_data\_file**(data, fname, mode='w')  
Save sequence/structure pairs for the given lengths.

**write\_fixed\_len\_data\_file**(seqlen, num, root='')

**write\_normal\_len\_data\_file**(central, std, num, root='')

**write\_uniform\_len\_data\_file**(minlen, maxlen, num, root='')

encoding\_utils

[index](#)  
[/home/julian-zim/Desktop/proj/rnadeep/rnadeep/encoding\\_utils.py](#)

Modules

[numpy](#)

Functions

**base\_pair\_matrix**(ss)

**binary\_encode**(structure)

**create\_windows**(sequences, window\_size)

**encode\_padded\_alignment\_matrix**(alignments, max\_length=None)

**encode\_padded\_sequence\_matrix**(sequences, max\_length=None)

**encode\_padded\_structure\_matrix**(structures, max\_length=None)

**encode\_sequence**(sequences)

**encode\_sequence\_matrix**(sequences)  
Make a BP probability matrix with one-hot encoding of basepairs.  
NOTE: This only works if all sequences have the same length, otherwise you need to use: encode\_padded\_sequence\_matrix

**encode\_sequence\_windows**(sequences, window\_size)

**encode\_structure**(structures)

**encode\_structure\_matrix**(structures)  
Make a BP probability matrix with one-hot encoding of basepairs.  
NOTE: This only works if all sequences have the same length!

**make\_pair\_table**(ss, base=0, chars=['.'])  
Return a secondary struture in form of pair table.

Args:  
ss (str): secondary structure in dot-bracket format  
base (int, optional): choose between a pair-table with base 0 or 1  
chars (list, optional): a list of characters to be are ignored, default: ['.']

**Example:**  
base=0: ((..)). => [5,4,-1,-1,1,0,-1]  
i.e. start counting from 0, unpaired = -1  
base=1: ((..)). => [7,6,5,0,0,2,1,0]  
i.e. start counting from 1, unpaired = 0, pt[0]=len(ss)

Returns:  
[list]: A pair-table

**one\_hot\_encode**(char)

**one\_hot\_matrix**(seq)

**profile\_vec\_matrix**(ali)  
Creates a profile matrix for the given alignment: For each cell a\_ij, the columns i and j or the alignment are combined by forming the outer product of two profile vectors for the two respective symbols at the current row index of the two columns and summing them all up. The two respective profile vectors created by the sheme defined in the base\_to\_ids dictionary variable.

