

TfReg user manual

TfReg Version 5.1 user manual

Gerald Weber gweberbh@gmail.com

Departamento de Física, Universidade Federal de Minas Gerais, Brazil

March 13, 2020

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/>.



CONTENTS

1	Introduction	3
2	Changes	4
3	Installation	6
3.1	Download and direct install	6
3.2	What will be installed?	6
3.3	Specific instructions	6
3.3.1	OpenSUSE	6
3.4	Compiling the source files	7
3.4.1	Libraries needed for compiling the source code	7
3.5	If things go wrong	7
3.5.1	License	7
4	Usage	8
4.1	Arguments	8
4.1.1	-o=<basename>	8
4.1.2	-reg=<filename>	8
4.1.3	-par=<filename or comma-separated list of filenames>	8
4.1.4	-data=<filename or comma-separated list of filenames>	8
4.1.5	-matrix=<new directory>	8
4.1.6	-res=<resulttype>	8
4.1.7	-model=<model acronym>	9
4.1.8	-duplextype=<DNA or RNA>	9
4.1.9	-dict=<list of nucleotides>	9
4.1.10	-expand=<nearest neighbours>	9
4.1.11	-cutoff=<integer number>	10
4.1.12	-pbc=<0 or 1>	10
4.1.13	-t=<temperature>	10
4.1.14	-int=<range>	10
4.1.15	-ee=<experimental error>	10
4.1.16	-rs=<seed>	10
4.1.17	-mlr=<sequence length>	11
4.1.18	-mar=<number of sequences>	11
4.1.19	-pm=<prediction method number>	11
4.1.20	-seq=<nucleotide sequence>	11
4.1.21	-cseq=<nucleotide sequence>	12
4.1.22	-salt=<salt concentration>	12
5	Models	13
5.1	Peyrard-Bishop (-model=pb)	13
5.2	The Dauxois variant (-model=dpb)	13
5.3	PB model with added solvent potential (-model=hms)	14
5.4	The Joyeux and Buyukdagli model (-model=jb)	14
5.5	The PB model including rise step h (-model=mes)	15
5.6	Build your own model (-model=test)	15
5.6.1	Changing the model potentials	15
5.6.2	Changing the Hamiltonian	16

6	Parameter specification	17
6.1	Generic * parameters	17
6.2	BP parameters	17
6.2.1	Context dependence for BP parameters	18
6.3	NN parameters	18
6.4	Parameter precedence	18
7	Formatting the sequence data file	20
7.1	DNA or RNA duplexes	20
7.2	Adding comments to your files	20
7.2.1	Special use of comments to group temperatures	20
7.3	Hybrid DNA-RNA duplexes	20
7.4	Hybrid DNA-TNA duplexes	21
7.5	RNA single bulges	21
7.6	Cy3 or Cy5 attached to DNA	21
8	Result files	22
8.1	.reg regression parameters	22
8.2	.dat melting temperatures and melting index results	22
8.3	.dat average opening if used with -res=averagey	23
8.3.1	Single sequence	23
8.3.2	Multiple sequences	23
8.4	.ver quality of the prediction	23
8.5	Matrix files if used with -matrix=	23
9	Examples	25
9.1	Task: given a set of melting temperatures, find the regression parameters	25
9.1.1	Results	25
9.1.2	Results	25
9.2	Task: prediction of DNA melting temperatures	26
9.2.1	Single sequence example	26
9.2.2	Multiple sequence example	26
9.3	Task: prediction of DNA melting temperatures with different salt concentrations	26
9.4	Task: prediction of RNA melting temperatures	27
9.5	Task: prediction of melting temperatures DNA containing inosine mismatches	27
9.6	Task: calculating the average opening $\langle y \rangle$	27
9.6.1	Calculating $\langle y \rangle$ at a given temperature	28
9.6.2	Calculating $\langle y \rangle$ for a range of temperatures	28
9.6.3	Using Perl scripts and parallel processing	28
10	Files included in the package	29
10.1	Sequences and melting temperatures	29
10.1.1	Folder /usr/share/data/	29

1 INTRODUCTION

TfReg implements the calculation of Peyrard-Bishop [1] style Hamiltonians to obtain some physical properties of DNA and RNA duplexes. The method uses the transfer matrix technique for the calculation of the classical partition function. Also, TfReg calculates the regression of experimental versus predicted melting temperatures using the equivalent melting index [2].

What will this software do for you? Given a set of experimental melting temperatures and a set of model parameters you will be able to calculate the regression parameters which will allow you predict melting temperatures of any DNA or RNA sequences. Alternatively, you may use one of the calculated regression sets which are provided and start calculating melting temperatures right away. You may chose between four different “flavours” of Hamiltonians if you wish to investigate the effect of different model parameters. If you do have basic programming skills in C++ you should be able to add new types of model Hamiltonians, as long as they fit within the 1D framework of the original Peyrard-Bishop model [1].

Evidently, this is work in progress. I hope to add further programs in the near future as well as increase the number of parameters for other types of oligonucleotides.

I would find it truly helpful indeed if you would let me know if this software is of any use to you. Showing a list of interested users to funding agencies often helps to secure the necessary resources to keeping such projects running. So, please, if you find this software useful let me know and if you use it for your scientific work please cite the appropriate papers which are listed at the end of this manual.

I wish you all the best in using TfReg

Gerald

Belo Horizonte, March 13, 2020

2 CHANGES

Version 5.1, March 2020

- New parameters and data files for mismatches [3]

Version 5, February 2020

- New parameters and data files
 - Added terminal related parameters, files of type `dna_tpb_*.par` [4]
 - Adapted sequence interpretation for TNA with new parameters for TNA-DNA hybrids [5].
- New features
 - Multiple files in `-data`, see section 4.1.4
 - New regression option `-pm=-2`, see section 4.1.19
- Corrections
 - Fixing bug that was using always `-pm=2` for predicting single sequences.
 - Fixing bug that was giving wrong trimer counts with `-res=nncheck`.

Version 4.0, September 2018

- For DNA/RNA hybrids: new melting temperature data file `rna-dna-1000.dat` from references 6–8 and parameters `rna-dna-pb-1000.par` [9].
- Now prints out CPU time after finishing calculations.
- New optimization in matrix calculations to reduce CPU time.
- Fixed error that would generate an infinite loop if matrix elements were either infinite or NaN.
- Fixed error in file `TestModel.h` that would prevent users to implement their own model.

Version 3.2, July 2017

- New parameters for Cy3 and Cy5 terminally attached to DNA [10] with new files:
`dna_pb_cy3.par dna_pb_cy5.par`
`dna_pb_cy3.reg dna_pb_cy5.reg`
`moreira15cy3.dat moreira15cy5.dat`
- A warning about missing parameters was reclassified as an error and now the program will terminate. Previously the program would continue with a value of zero for those missing parameters.
- Fixed bug for `-res=averagey` which, under certain very specific circumstances, would erroneously apply periodic boundary conditions as if `-pbc=1` had been set. Symptoms were missing end-fraying, for example for DNA terminating in AT no end-fraying would be observed.
- documentation updates.
- removed obsolete files:
`JobControl.h JobControl.cpp`

Version 3.1, April 2017

- New parameters for single type I bulges in RNA [11] with new files:

```
rna_pb_bulge_group1.par
rna_bulge_group1_adenosine_adj200.dat rna_bulge_group1_cytosine_adj200.dat
rna_bulge_group1_guanosine_adj200.dat rna_bulge_group1_uridine_adj200.dat
rna_pb_bulge_group1_adenosine.reg rna_pb_bulge_group1_cytosine.reg
rna_pb_bulge_group1_guanosine.reg rna_pb_bulge_group1_uridine.reg
```

Version 3.0, January 2016

- Fixed important bug concerning matrix multiplications which affects the average opening profiles, especially for symmetric sequences. This may also affect, although very slightly, the melting index and predicted temperatures. Differences in both cases are around 0.01% or less compared to previous versions of tfreg.
- Fixed bug from version 2.0 which was generating empty files for option `-res=prediction` with single sequences.
- Added new base pair representation of type `XY~z` which introduces context-dependent base pairs. See section 6.2.1.
- Added optimized parameters for RNA GU [12].
- Linking to the library `libboost_regex` is now required.
- Code clean up, removed `gbc_exp` from C++ name space.
- Adding version identification, see first printed line when running tfreg.
- You can now add easily your own model potentials, see section 5.6.
- Example scripts were revised for consistency, in particular they no longer accept the data folder as argument. Alternative data folders should be provided by setting the `PREFIX` command line variable.
- Add new option `-mar`.

Version 2.0, December 2014

- Corrected harmless bug which was causing the last line of a parameter to be read twice.
- Added a new model `-model=mes`, see section 5.5.
- Added new option `-res=nncheck`, see section 4.1.6, which shows how the sequences are analysed in terms of base pairs and nearest-neighbours by TfReg.
- Added optimized parameters for deoxyinosine [13], see section 9.5.
- `.ver` files now have an additional column providing ΔT_{RMS} .

Version 1.2, November 2013

- Corrected bug which would fail to predict temperatures with `-pm=-1` from reg files also generated with `-pm=-1`.
- Removed generation of `TfX` files.
- Adding new option `-dict` which allows you to add new characters to represent nucleotides.

Version 1.1, February 2013

- We corrected a software bug which recalculated unnecessarily the Gauss-Legendre quadrature weights. As a result TfReg runs considerably faster.
- We added a script for adjusting the regression coefficients to other salt concentrations. See section 9.3.

Version 1.0, November 2012

First release of TfReg.

3 INSTALLATION

3.1 Download and direct install

You can find the compiled binary package, for many Linux distributions, as well as source code, at these locations:

- <http://bioinf.fisica.ufmg.br/software>
- <https://download.opensuse.org/repositories/home:/drgweber/>

browse to the desired version of the software and then download the appropriate package for your Linux distribution. If your distribution is not covered then download the source code and compile manually, see below for instructions.

3.2 What will be installed?

Typically, there will be at least a binary executable file

`/usr/bin/tfreg`

and further files, such as model parameter files, pre-calculated regression parameters as located in

`/usr/share/TfReg/data` or `/usr/share/tfreg/data`

note that these folder may vary depending on your Linux distribution. Example scripts which help to understand how to run TfReg are to be found in

`/usr/share/TfReg/example` or `/usr/share/tfreg/example`

The documentation (which you are reading right now) should be located at

`/usr/share/doc/packages/TfReg/tfreg-user-manual.pdf` or
`/usr/share/doc/packages/tfreg/tfreg-user-manual.pdf`

3.3 Specific instructions

3.3.1 OpenSUSE

Installing via repository Using the graphical interface Yast2 or the command line zypper add the following repository URL

https://download.opensuse.org/repositories/home:/drgweber/openSUSE_Leap_15.1/

then search for the package TfReg and select install. If you have an older OpenSUSE then change the last number to your installed version accordingly. The installation via repository has the advantage that you may simply update for future versions instead of repeating the whole installation procedure.

Command line download/install Download the appropriate package for your system from

https://download.opensuse.org/repositories/home:/drgweber/openSUSE_Leap_15.1/

for example if your system is 64bits, you may download the package

http://download.opensuse.org/repositories/home:/drgweber/openSUSE_Leap_15.1/x86_64/TfReg-5.1-1.21.1.x86_64.rpm

note: version numbers may vary from this example. Then install

`zypper install TfReg-5.1-1.21.1.x86_64.rpm`

3.4 Compiling the source files

Please read this section if you are unable to find the packages for your specific Linux distribution or if you are interested in modifying the source code.

Download the source package from or from my personal webpage

<https://bioinf.fisica.ufmg.br/software/tfreg-5.1>

Typically the package is called something like `tfreg-5.1.tar.bz2` After unpacking the `tar` package

```
tar -xvjf tfreg-5.1.tar.bz2
```

change into the unpacked folder

```
cd TfReg-5.1
```

if all necessary packages are available you should try to compile using the `make` command

```
make
```

If the compilation is successful, you should see something like

```
g++ -O3 -o tfreg -Isrc src/Options.cpp src/JobControl.cpp src/Nucleotide.cpp src/tfreg.cpp -lz  
-lgsl -llapack -lgslcblas -lm -lboost_filesystem -lboost_system -lboost_regex
```

and nothing else, that is it! This generates the binary file `tfreg`, which you may copy to your main installation at `/usr/bin` (you will need root permission) or into your local folder `/home/user/bin` (replace `user` with you actual user name).

3.4.1 Libraries needed for compiling the source code

This software was developed and tested under OpenSUSE Linux 15.1 and depends on some libraries to function properly:

1. `libboost_filesystem libboost_system libboost_regex` <http://www.boost.org>
2. `gsl` <http://www.gnu.org/software/gsl/>
3. `lapack gslcblas` <http://www.netlib.org/lapack/>
4. `gcc-fortran` <http://gcc.gnu.org/>

which means that you will need at *least* these specialized packages in addition to the usual `gcc` and `g++` compiler. If you already installed pre-compiled binaries then probably you will already have all the required packages installed in your Linux system.

3.5 If things go wrong

Most problems will come from missing library packages or from erroneous usage of your system. It is not possible for me to cover everything that may go wrong, so please feel free to contact me. Please include a detailed description of error messages, which system you are using and a step by step description of what you tried to do. Please understand that I will need as much information as possible. I can do nothing with messages saying simply “TfReg is not working on Ubuntu”.

3.5.1 License

This software is published under the GNU General Public License version 3 (GPLv3), the complete text of this license can be found in the documentation folder. If you wish to use this software under a different license or wish to make changes to the software without distributing it under the GPLv3 please contact me so that we can arrange for a specific license.

4 USAGE

TfReg take all its program option from the command line, therefore you should invoke the binary `tfreg` with some of the arguments which are described below. Example shell scripts are provided and I recommend you study them as they are the best source to illustrate how to use TfReg. The following section details every available program option for TfReg.

4.1 Arguments

4.1.1 `-o=<basename>`

specifies the output basename, that is, all files which are generated start with `basename`

4.1.2 `-reg=<filename>`

specifies the input regression file name. Use this option together with `-res=prediction`, in this case the regression coefficients are taken from this file. See section 4.1.19 regarding the prediction and regression options.

4.1.3 `-par=<filename or comma-separated list of filenames>`

specifies the input parameter file name. This can be a list of files `-par=file1.par,file2.par,file3.par`. In case of multiple specifications of the same model parameter, the last one supersedes previous parameters. For example if `file1.par` has the parameter `AT:Morse.D 0.05` and `file3.par` has `AT:Morse.D 0.03`, the final value for `AT:Morse.D` will be 0.03.

4.1.4 `-data=<filename or comma-separated list of filenames>`

specifies the input files containing nucleotide sequences and melting temperatures. Multiple files in a comma-separated list are accepted since version 5.0. See section 7 for more information regarding preparing those files.

Note on memory limits: TfReg was not created with very large datasets in mind, and all sequences are loaded into memory. A safe limit is of the order of 50000 sequences which will require about 1GB of memory to run.

4.1.5 `-matrix=<new directory>`

if specified this will create a directory where the calculated matrices will be located. Once calculated these matrices will be reloaded the next time the program is run again. This speeds up the calculations, however if you do change parameters or the model from one run to the next you should not use this option as you will load matrices which are incorrect. Also, the Gauss-Legendre quadrature weights are saved into this directory if set, these weights are totally independent of any model parameters.

4.1.6 `-res=<resulttype>`

which type of result we wish to obtain.

(default) `-res=regression` given a set of parameters and a set o melting temperatures, calculates the regression parameters (stored in a file with extension `.reg`) and also calculated the prediction of temperatures.

`-res=prediction` predicts temperatures for a given set of parameters, requires a pre-calculated regression file passed through the `-reg=` option, see 4.1.2. Some calculated regression files can be found in the `data` folder, see 10.1.1.

-res=nncheck only checks consistency of the nearest-neighbour decomposition of the nucleotide sequences. This will display the sequence analysis and exit immediately without doing any calculations.

-res=averagey calculates the average opening $\langle y \rangle$, results are given in Ångstrom. For a single sequence use with **-seq**, and optionally with **-cseq**, and the result will be a file with $\langle y \rangle$ arranged column-wise. For larger number of sequences arrange all sequences in one or more files with **-data**. Note that for short sequences the temperature needs to be unrealistically low, see example 9.6.

-res=zyfy calculates the configurational part of the partition function Z_y (see for instance Eq. 28 of Ref. [14]) and its total Helmholtz free energy $-kT \ln Z_y$.

4.1.7 -model=<model acronym>

Selects the Peyrard-Bishop model which should be used. Each model requires specific parameters which should be passed via the **-par=<filename>** option.

(default) **-model=pb** original Peyrard-Bishop model [1], see section 5.1

-model=dpb the anharmonic variant [15], see section 5.2

-model=hms PB model with added solvent potential [16], see section 5.3

-model=jb finite enthalpy model [17–21], see section 5.4

-model=mes Morse-exact stacking model [22], see section 5.5

-model=test Reserved for *your* test models, please see section 5.6 for instructions. Please do not use without first implementing the changes described in 5.6 and recompiling the source code.

4.1.8 -duplextype=<DNA or RNA>

Selects the type of duplex we should expect, this is important for selecting the base pair complementarity. Note that IUPAC codes cannot be used since we need to know which nucleotide parameters to use, that is, we cannot use N for example since we would not know which potentials to use. If you need to define other types of nucleotides see option **-dict**.

(default) **-duplextype=DNA** , will expect A, C, G and T nucleotides and will consider A complementary to T and C complementary to G.

-duplextype=RNA , same as for DNA but considers A and U as complementaries, and removes T from the dictionary. [23].

4.1.9 -dict=<list of nucleotides>

When using non-canonical nucleotides, for example inosine, you must tell TfReg which characters to use in addition to its usual dictionary of A, C, G, T or U. You should also tell TfReg if there is a complementary pair to this new nucleotide. If there is none, simply repeat the character. In the following example we are adding the letter I for inosine

-dict=I:I

here we are saying: consider I and its complementary I as new letters. If you wish to add two or more new letters simply make a comma-separated list, but do not add blank spaces.

4.1.10 -expand=<nearest neighbours>

Selects which nearest neighbours to expand.

(default) **-expand=CG_CG** for DNA or RNA

`-expand=dCrG_dCrG` for DNA-RNA hybrids [9].

`-expand=dCtG_dCtG` for DNA-TNA hybrids [5].

4.1.11 `-cutoff=<integer number>`

This controls the truncation P in Eq. (22) of Ref. 14, that is, you will be using only the first P eigenvalues of the diagonalized matrix. All subsequent matrix multiplications will be $P \times P$, using a small P will make TfReg run considerably faster.

(default) `-cutoff=0` no cutoff. While this is the default, it is not recommended and not necessary. If you wish a lot of precision $P = 80$ is more than enough.

(recommended for T_m) `-cutoff=10` a cutoff $P = 10$ gives quite good results and reduces the computational cost for melting temperatures.

(recommended for $\langle y \rangle$) `-cutoff=80` a cutoff $P = 80$ gives quite good results for average openings.

4.1.12 `-pbc=<0 or 1>`

Controls the type of boundary conditions.

(default) `-pbc=0` open boundary conditions, this is what you normally would have and usually shows end fraying

`-pbc=1` periodic boundary conditions, this would be the case for a circular sequence.

4.1.13 `-t=<temperature>`

Selects the temperature in kelvin for which the calculation of the matrices is carried out. Please note: this temperature is completely unrelated to the melting temperatures.

(default) `-t=370`

4.1.14 `-int=<range>`

Range of integration, this specifies the limits of the integral shown in Eq. (14) of Ref. 14 and the size M of the matrices.

(default) `-int=-1:30/100` integrates from $y = -1$ to $y = 30 \text{ \AA}$ and uses matrices of size 100

(recommended) `-int=-1:200/400` integrates from $y = -1$ to $y = 200 \text{ \AA}$ and uses matrices of size 400. This gives very accurate results without too much computational cost.

4.1.15 `-ee=<experimental error>`

This sets the experimental error (in $^{\circ}\text{C}$) of the melting temperature set. If given, the data set (provided through the `-data` option, section 4.1.4) will be modified by small positive or negative amounts such that the standard deviation is close to the experimental error.

(example) `-ee=0.5` will modify the dataset to within $0.5 \text{ }^{\circ}\text{C}$

4.1.16 `-rs=<seed>`

Sets the seed of the random number generator (C function `srand`) which is used to modify the dataset. If you use the same seed you will get exactly the same random modifications. Only make sense to be used together with `-ee`.

4.1.17 -mlr=<sequence length>

Sets the minimal sequence length to be considered in calculating regression parameters.

(example) -mlr=6 only considers sequences with 6 bp and above

4.1.18 -mar=<number of sequences>

Sets the minimal number of sequences, within a group of length N , to be considered in calculating regression parameters.

(default) -mar=3 does not attempt regressions with less than 3 sequences for any group

(example) -mar=2 accepts at least 2 sequences for any group

4.1.19 -pm=<prediction method number>

Selects the prediction/regression method.

(default) -pm=2 this considers two equations for regression

$$T_p = a_0(N, [\text{Na}^+]) + a_1(N, [\text{Na}^+])\tau, \quad (4.1)$$

there will be one equation for each length N and for each salt concentration $[\text{Na}^+]$ and each coefficient a_k is calculated as

$$a_k(N, [\text{Na}^+]) = b_{0,k}([\text{Na}^+]) + b_{1,k}([\text{Na}^+])N^{1/2}, \quad (4.2)$$

where N is the size of the sequence and τ is the equivalence index. Note that there need to be at least 3 different sequences for each length N (this number can be changed with the option -mar). **Important:** if the dataset has non-uniform salt concentrations, there will be separate regressions for each salt concentrations. For example, say that the dataset mixes measurements at 119 mM and 220 mM, there will be one regression Eq. (4.1) for 119 mM and another for 220 mM. That is the reason why the coefficients are given as a function of N and $[\text{Na}^+]$.

-pm=3 in addition of the two previous equations, this considers a third equation for the case where there are two or more salt concentrations in the data file

$$b_{j,k}([\text{Na}^+]) = c_{0,j,k} + c_{1,j,k} \log[\text{Na}^+]. \quad (4.3)$$

Note that if there is only one salt concentration $[\text{Na}^+]$ in the data file this reverts automatically to option -pm=2. Do not use

-pm=-1 this considers only one regression

$$T_p = a_0 + a_1\tau, \quad (4.4)$$

which is useful if the dataset has only a few sequences, or if all sequences are of the same length or if you have otherwise trouble in getting good linear regression coefficients for Eq. (4.2). The last situation may happen if your model parameters are very far from the optimized values. This option was first introduced for Ref. [23].

-pm=-2 this considers separate regressions of type -pm=-1 per group, where the group type is specified with the option -pg. -pg=ct groups all T_m data by the species concentration C_t , -pg=salt groups by salt concentration, and -pg=key groups by an arbitrary key which is the first word that follows the comment sign # (see 7.2).

4.1.20 -seq=<nucleotide sequence>

Instead of providing a file with your sequences you can give them on the command line. This is useful if you want to see a melting for just one sequence. You should give the main strand from 5' to 3', the complementary sequence will be worked out automatically.

(example) -seq=ACGTTGAATT

4.1.21 -cseq=<nucleotide sequence>

You should provide a 3' → 5' sequence if your sequence is not perfectly complementary, say like in a sequence with nucleotide mismatches.

(example) -cseq=TGCTACTTAA from 3' → 5'

4.1.22 -salt=<salt concentration>

In the case where a -reg file contains several salt concentration, this will select the one that should be used for the calculations. Otherwise, this option is silently ignored. *If you need to calculate melting temperatures for different salt concentrations please see section 9.3.*

5 MODELS

There are several variants of the PB model, each of which requires different model parameters. In this section you will find which models TfReg currently supports and which model parameters are needed. Example files with model parameters are provided and typically you will find these in `/usr/share/TfReg/data` (or `url/usr/share/tfreg/data`) with file extension `.par`. Note that the parameters are specified per nucleotide type and follow a very specific notation described in section 6.

5.1 Peyrard-Bishop (`-model=pb`)

This is the original Peyrard-Bishop model proposed in Ref. [1] which uses a Morse potential for modelling the hydrogen bonds and

$$V_{\text{Morse}}(y_i) = D \left(e^{-y_i/\lambda} - 1 \right)^2, \quad (5.1)$$

and the nearest-neighbour stacking interaction as a harmonic oscillator

$$w_{\text{harm.}}(y_i, y_{i-1}) = \frac{k}{2} (y_i - y_{i-1})^2. \quad (5.2)$$

because of the divergence of the partition function [24] we modified this to

$$w_{\text{harm.}}(y_i, y_{i-1}) = \frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2), \quad (5.3)$$

Hamiltonian	model parameter	program parameter	units	type
$D (e^{-y_i/\lambda} - 1)^2$	D	<code>Morse.D</code>	eV	BP
	λ	<code>Morse.lambda</code>	Å	BP
$\frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2)$	k	<code>harmonic.k</code>	eV/Å ²	NN
	θ	<code>harmonic.theta</code>	rad	NN

5.2 The Dauxois variant (`-model=dpb`)

In 1993 Dauxois, Peyrard and Bishop introduced an anharmonicity term to account for sharp transitions in the original PB model [15],

$$w_{\text{an.}}(y_i, y_{i-1}) = \left[1 + \rho e^{-\alpha(y_i + y_{i-1})} \right] w_{\text{harm.}}(y_i, y_{i-1}), \quad (5.4)$$

Hamiltonian	model parameter	program parameter	units	type
$D (e^{-y_i/\lambda} - 1)^2$	D	<code>Morse.D</code>	eV	BP
	λ	<code>Morse.lambda</code>	Å	BP
$[1 + \rho e^{-\alpha(y_i + y_{i-1})}]$	k	<code>harmonic.k</code>	eV/Å ²	NN
$\times \frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2)$	θ	<code>harmonic.theta</code>	rad	NN
	ρ	<code>anharmonic.rho</code>	adimensional	NN
	α	<code>anharmonic.alpha</code>	adimensional	NN

data/weber06-1.par

```

1 weber06
2 AT:AU:morse.D      0.05
3 CG:morse.D         0.08
4 AT:AU:morse.lambda 0.33333
5 CG:morse.lambda    0.125
6 *:harmonic.theta   0.01
7 *:harmonic.k        0.025
8 *:anharmonic.alpha 0.35
9 *:anharmonic.rho    2.0

```

5.3 PB model with added solvent potential (-model=hms)

A solvent term was added to the harmonic PB model,

$$V(y_i) = V_{\text{Morse}}(y_i) - f_s D [\tanh(y_i/\lambda_s) + 1], \quad (5.5)$$

Hamiltonian	model parameter	program parameter	units	type
$D (e^{-y_i/\lambda} - 1)^2$	D	Morse.D	eV	BP
	λ	Morse.lambda	Å	BP
$-f_s D [\tanh[(y_i + y_e)/\lambda_s] + s]$	f_s	solvent.f_s	adimensional	BP
	λ_s	solvent.lambda	Å	BP
	y_e	solvent.eq_sol	Å	BP
	s	solvent.sign_sol	adimensional	BP
$\frac{k}{2} (y_i^2 - 2y_i y_{i-1} \cos \theta + y_{i-1}^2)$	k	harmonic.k	eV/Å ²	NN
	θ	harmonic.theta	rad	NN

data/weber06b-1.par

```

1 weber06b
2 AT:morse.D      0.05
3 CG:morse.D      0.08
4 AT:morse.lambda 0.33333
5 CG:morse.lambda 0.125
6 *:solvent.eq_sol 0.0
7 *:solvent.sign_sol 1.0
8 *:harmonic.theta 0.01
9 *:harmonic.k     0.025
10 *:solvent.lambda 1.0
11 *:solvent.f_s    0.1

```

5.4 The Joyeux and Buyukdagli model (-model=jb)

The model by Joyeux and Buyukdagli [17–20] introduces a finite stacking enthalphy

$$w_{\text{fin.}}(y_i, y_{i-1}) = \frac{\Delta H}{C} \left[1 - e^{-b(y_i - y_{i-1})^2} \right] + \frac{K_b}{2} (y_i - y_{i-1})^2. \quad (5.6)$$

data/jb.par

```

1 var_jb_owczarzy04_init2
2 AT:morse.D      0.041
3 CG:morse.D      0.054
4 AT:morse.lambda 0.1667
5 CG:morse.lambda 0.1667
6 *:finite_enthalpy.C 4.0
7 *:finite_enthalpy.DeltaH 0.409
8 *:finite_enthalpy.b 0.80
9 *:finite_enthalpy.kb 4.0e-4

```

Hamiltonian	model parameter	program parameter	units	type
$D(e^{-y_i/\lambda} - 1)^2$	D	<code>Morse.D</code>	eV	BP
	λ	<code>Morse.lambda</code>	Å	BP
$\frac{k}{2}(y_i - y_{i-1})^2 s$	k	<code>harmonic.k</code>	eV/Å ²	NN
$\frac{\Delta H}{C} \left[1 - e^{-b(y_i - y_{i-1})^2} \right]$	ΔH	<code>finite_enthalpy.DeltaH</code>	eV	NN
	C	<code>finite_enthalpy.C</code>	adimensional	NN
	b	<code>finite_enthalpy.b</code>	adimensional	NN
$\frac{K_b}{2}(y_i - y_{i-1})^2$	K_b	<code>finite_enthalpy.Kb</code>	eV/Å ²	NN

5.5 The PB model including rise step h (`-model=mes`)

This Hamiltonian [22] rewrites the torsional 3D Hamiltonian [25] in the notation of the original PB model [1] by setting the angles to zero ($\phi_i = 0$ and $\theta_0 = 0$), we obtain

$$U_{i,i-1} = D(e^{-\sqrt{2}y_i/\lambda} - 1)^2 + k \left(\sqrt{h^2 + \frac{1}{2}(y_i - y_{i-1})^2} - h \right)^2 \quad (5.7)$$

where h is stacking distance (rise) between base pairs. Unlike the PB model where the stacking factor is approximate, this model evaluates the stacking exactly. The acronym `mes` stands for "Morse Exact Stacking".

For the computational implementation the original PB Morse potential is implemented unchanged and the additional factor $\sqrt{2}$ should be embedded into the `Morse.lambda` parameter. Note that the `harmonic.theta` parameter used for the PB model is unnecessary here.

Hamiltonian	model parameter	program parameter	units	type
$D(e^{-y_i/\lambda} - 1)^2$	D	<code>Morse.D</code>	eV	BP
	λ	<code>Morse.lambda</code>	Å	BP
$k \left(\sqrt{h^2 + \frac{1}{2}(y_i - y_{i-1})^2} - h \right)^2$	k	<code>harmonic.k</code>	eV/Å ²	NN
	h	<code>exactstack.h</code>	Å	NN

5.6 Build your own model (`-model=test`)

We have been contacted by some users who wanted to modify `tfreg` for their own models. Here we are trying to make this as simple as possible, but please be warned that recompilation of the code is required. Instructions for downloading and compiling the source code are given in section 3.4.

No knowledge C++ is required to make the modifications to the code described in this section. The only changes you will make require a tiny bit of knowledge of how to write a simple equation in standard C.

If you are interested in testing different potentials see the next section. If you would like to try different combinations of existing potentials, without introducing new potentials, see section 5.6.2.

5.6.1 Changing the model potentials

Parameters We reserved 8 parameters for each potential. For the $V(x)$ potential those are

`test_vx.a test_vx.b test_vx.c test_vx.d test_vx.e test_vx.f test_vx.g test_vx.h`

and for $w(x, y)$

`test_wxy.a test_wxy.b test_wxy.c test_wxy.d test_wxy.e test_wxy.f test_wxy.g test_wxy.h`

your potentials will have to be written considering those available parameters.

Prepare your model First work out what the required parameters of your model will be. Typically you will have two potentials, one is the base-pair potential $V(x)$ and the second is the stacking potential $w(x, y)$.

Lets suppose that your $V(x)$ potential is $V = x(a + b^2)$, and your stacking potential $w(x, y) = ax + by^2$, then your parameter file should look more or less like this


```

identification
AT:test_vx.a 0.2
AT:test_vx.b 0.3
AT_AT:test_wxy.a 1.2
AT_AT:test_vxy.b 3.3

```

where the first line is an arbitrary identification string.

Change the source code Using a program file editor, open the file `TestModel.h`, and locate the line

```
class TestVx: public BasePotential<_Tp>
```

and scroll down until you find

```
return x*(a+b+c+d+e+f+g+h); // MODIFY HERE
```

for our example, $V = x(a + b^2)$ we would change this to

```
return x*(a+pow(b,2));
```

where we used simple C math formatting.

Now locate

```
class TestWxy: public BasePotential<_Tp>
```

and scroll down until you find

```
return x*y*(a+b+c+d+e+f+g+h); // MODIFY HERE
```

and for the example $w = ax + by^2$ change this to

```
return a*x+b*pow(y,2);
```

Now recompile the code, which mostly boils down to issue the `make` command, see section 3.4. To run use `-model=test`.

5.6.2 Changing the Hamiltonian

If you simply want a different combination of existing model potentials, then you only have to change the main Hamiltonian. Open the file `TestModel.h` and locate the line

```
class TestHamiltonian: public Hamiltonian<TestVx<_Tp>,TestWxy<_Tp> > // Here specify the how the....
```

suppose that you would like to use a Hamiltonian where the V potential is a simple Morse-solvent potential but with anharmonic stacking, you should then replace `TestVx` with `MorseSolvent` and `TestWxy` with `AnharmonicStacking`. The above line would then look like

```
class TestHamiltonian: public Hamiltonian<MorseSolvent<_Tp>,AnharmonicStacking<_Tp> >
```

also change in the exact same way the line

```
typedef Hamiltonian<TestVx<_Tp>,TestWxy<_Tp> > base_type; // Here specify the how the Hamiltonian
```

to

```
typedef Hamiltonian<MorseSolvent<_Tp>,AnharmonicStacking<_Tp> > base_type;
```

Now recompile the code, which mostly boils down to issue the `make` command, see section 3.4. To run use `-model=test`.

6 PARAMETER SPECIFICATION

This chapter explains how to write a parameter file.

Let's start with the most simple possible parameter file

```
1 initial
2 *:harmonic.theta 0.01
3 *:harmonic.k 0.025
4 *:morse.D 0.0324083
5 *:morse.lambda 0.350149
```

The first line holds a simple identifier, this identifier can be letters, numbers and some other characters but should contain no spaces. This identifier is used when you specify the `-matrix` option.

The second line starts with an asterisk which means that this parameter should apply to any base pair. In case of a DNA sequence this will apply to AT and CG base-pairs. TA base pairs are considered equivalent to AT, and CG are equivalent to GC, from the perspective of the models considered in this manual.

6.1 Generic * parameters

The generic base `*` specification can be used when a given parameter should be applied to any base pair (see previous example). However, if the parameter should be applied to a specific base pair you should specify either in BP or NN form (see next section).

6.2 BP parameters

Base-pair parameters are those which generally do not depend on the context, that is, it is not relevant which are the neighbours of the given base pair. For DNA we use Watson-Crick base pairs AT=TA and CG=GC, but may also use mismatched base pairs such as GT or AA. Switching the nucleotides, eg. GT to TG, makes no difference for BP parameters. The reason for this is that they usually represent properties of the hydrogen bond.

The format which needs to be specified in the parameter file is the base pair code, followed by the model parameter and its value, like in the following example

```
1 AT:Morse.D 0.5
```

You can specify multiple base pair with the same parameters as in the following example

```
1 AT:CG:Morse.D 0.5
```

which is equivalent to

```
1 AT:Morse.D 0.5
2 CG:Morse.D 0.5
```

Hybrid DNA-RNA and DNA-TNA an additional character is required, `d` for the DNA strand, `r` for RNA strand and `t` for TNA strand. For example `dCrG` means C on the DNA strand and G on the RNA strand. `TfReg` should be able to read a construct such as `dCdG` for DNA-DNA or `rCrG` for RNA-RNA, however this is generally unnecessary.

```
1 dArU:morse.D 0.027758
2 dCrG:morse.D 0.073671
3 dGrC:morse.D 0.062591
4 dTrA:morse.D 0.040300
```

6.2.1 Context dependence for BP parameters

There are situations where you may need different BP parameters for the same base pairs. Such a situation arises for example for the GU wobble pair in RNA [12]. In this case we will denote the base pair by an arbitrary superscript like GU^a or GU^b . First however we need to add context rules for these base pairs. This is all done in the parameter file as illustrated here:

```
Context BP parameter example
1 +GU^a AGU/UUG,GUA/UGU
2 +GU^b AGG/UUU,AUU/UGG,GGA/UUU,UGG/AUU
3 GU^a:morse.D 0.50
4 GU^b:morse.D 0.50
```

where we are saying that any GU in the context AGU/UUG or GUA/UGU is given the specific name GU^a , and that all further parameters use this name like GU^b :morse.D. NN parameter specification (see next section) are equally affected by this new specification.

Note that the actual parameter set for RNA GU mismatches is more complicated than the above example. You can find the complete optimized parameters in file `rna_pb_GU.par` [12].

6.3 NN parameters

Nearest neighbour (NN) parameters follow the same conventions as usually found in linear regression models [26, 27]. The convention used is of type AB_CD, where AB is the first base pair and CD the second base pair. Note that DC_BA is equivalent to AB_CD, for example, AT_AT is the same as TA_TA. When specifying a NN sequence always specify in lexical ordering, that is write AT_AT and not TA_TA. For hybrid DNA-RNA an example would be dArU_dCrG.

The following example shows all 10 irreducible N parameters for `finite_enthalpy.DeltaH`.

```
Examples of NN parameters
1 AT_AT:finite_enthalpy.DeltaH 0.42085
2 AT_CG:finite_enthalpy.DeltaH 0.416718
3 AT_GC:finite_enthalpy.DeltaH 0.400366
4 AT_TA:finite_enthalpy.DeltaH 0.330006
5 CG_AT:finite_enthalpy.DeltaH 0.445258
6 CG_CG:finite_enthalpy.DeltaH 0.385034
7 CG_GC:finite_enthalpy.DeltaH 0.411902
8 GC_AT:finite_enthalpy.DeltaH 0.44666
9 GC_CG:finite_enthalpy.DeltaH 0.497106
10 TA_AT:finite_enthalpy.DeltaH 0.431154
```

Note that if you are using context-dependent BP parameters, see section 6.2.1, you should use the superscript as normal

```
Examples of NN parameters with context dependent BP
1 GC_UG^j:harmonic.k 0.025
2 GU^a_UG^a:harmonic.k 0.025
```

For DNA-RNA or DNA-TNA hybrids it is necessary to use the d, r or t prefixes as in the following example for DNA-RNA

```
Examples of NN parameters for DNA-RNA
1 dArU_dArU:harmonic.k 0.009002
2 dGrC_dTrA:harmonic.k 0.042731
```

6.4 Parameter precedence

Since the program can read more than one parameter file, the last parameter read is the final value. There may also be several specification in the same file as well. Consider the following example

```
Precedence example 1
1 *:Morse.D 0.3
2 AT:Morse.D 0.5
```

the first line says that all base pairs should a `Morse.D` value of 0.3. The second line however says that AT base pairs should use 0.5. In this case CG base pairs for instance will use 0.3 since nothing different was specified.

However, a generic base pair `*` does not supersede a specific base pair as in the following example

Precedence example 2

```
1 AT:Morse.D 0.5  
2 *:Morse.D 0.3
```

in this case AT base pairs will continue using 0.5, not 0.3. You should understand the generic base pair `*` as: *if nothing else matches, use this value.*

7 FORMATTING THE SEQUENCE DATA FILE

A file containing all sequences can be used to make bulk calculations. The sequence file names should be passed through the option `-data`, see section 4.1.4.

7.1 DNA or RNA duplexes

The file format is composed of columns

```
data/owczarzy04-69.dat
1 temperature
2 ATCAATCATA TAGTTAGTAT 21.3 69 2
3 TTGTAGTCAT AACATCAGTA 24.7 69 2
4 GAAATGAAAG CTTACTTTC 22.1 69 2
```

The first column is the sequence (from 5' → 3') and the second column is the secondary strand (from 3' → 5'). The third column is the melting temperature in °C, the fourth column is the salt concentration (in mM). The last column the species concentration (in μM) which is currently not used. The secondary strand does not necessarily need to be the complementary of the main strand as long as there are parameters for the mismatched pairs. If TfReg fails to find parameters for your sequences it will complain loudly.

7.2 Adding comments to your files

Since version 5.0 it is possible to add comments after the # sign for all sequences, as exemplified next

```
Example data with comments
1 temperature #sequences from schoning00
2 d(AAAAAAAAAAAAAAAAAA) t(TTTTTTTTTTTTTTTT) 32 1010 10 #DNA-TNA poly-A and poly-T
3 d(AAAATTTATATTATTA) t(TTTTAAATATAATAAT) 47 1010 10 #DNA-TNA
```

7.2.1 Special use of comments to group temperatures

There might be some interest in running separate melting temperature regressions by some arbitrary key. Consider the following example

```
Example data with comments used for grouping
1 XACGATCGTV VTGCTAGCAX 48.3 1000 1 #Cy3
2 XCTGATCAGV VGACTAGTCX 44.1 1000 1 #Cy3
3 YACGATCGTV VTGCTAGCAY 48.0 1000 1 #Cy5
4 YCTGATCAGV VGACTAGTCY 44.9 1000 1 #Cy5
```

the first two sequences are for Cy3 and last two for Cy5. Normally, one would have to split this in two files and run the regression separately for each file. However, an alternative is to use the file as it is, and with `-pm=-2` together with `-pg=key` (see 4.1.19) run a single regression. What will happen is that TfReg will read the first word after # (without spaces) and use this as index. Please note that this is case-sensitive.

7.3 Hybrid DNA-RNA duplexes

For hybrid DNA-RNA, all sequences should contain the additional `d` and `r` characters to distinguish if they are of DNA or RNA type [9]. The sequence should be enclosed as show in the following example.

```
Example data files for DNA-RNA
1 temperature
2 r(UUUGUAUCCAAU) d(AAACATAGGTTA) 45.6 1000 100
3 d(GTTGGTTGGTTG) r(CAACCAACCAAC) 60.0 1000 8.3
```

7.4 Hybrid DNA-TNA duplexes

In a very similar way as for DNA-TNA, you need the addition **d** and **t** characters as in the example below

Example data files for DNA-TNA

```
1 temperature
2 t(GCCGTGAG) d(CGGCACTC) 39.9 1010 10
3 t(ACGTCATTCCTC) d(TGCAGTAAGGAG) 44.6 1010 10
```

Note on strand concentrations: TfReg makes no adjustments to strand concentration C_t . Ideally the sequences should be adjusted beforehand to the same strand concentration that was used to calculate the parameters.

7.5 RNA single bulges

Single bulges occur when one strand is shorter than the other, but otherwise they are complementary to each other. To mark the position of the bulge (if known) we mark this with an **X** so that both strands become of equal length [11], as in the following example.

Example data file for RNA single bulges

```
1 temperature
2 GGCGACUCG CCGCXGAGC 57.6588010149076 1000 200
3 GAGCAGGUC CUCGXCCAG 52.8854893430268 1000 200
```

for a full example see the files

rna_bulge_group1_adenosine_adj200.dat rna_bulge_group1_cytosine_adj200.dat
rna_bulge_group1_guanosine_adj200.dat rna_bulge_group1_uridine_adj200.dat

included with TfReg.

7.6 Cy3 or Cy5 attached to DNA

When Cy3 or Cy5 is attached through a flexible linker to the 5' end of DNA it behaves like an additional base pair [28, 29]. We use this to represent them as a XV (Cy3) or YV (Cy5) base pair, where X or Y are linked to the 5' end. The following example shows how the sequences are formatted for the case of Cy3 [10].

Example data file for Cy3 attached to DNA

```
1 temperature
2 XACGATCGTV VTGCTAGCAX 48.3 1000 1
3 XCTGATCAGV VGACTAGTCX 44.1 1000 1
```

for a full example see files `moreira15cy3.dat` and `moreira15cy5.dat` included with TfReg.

8 RESULT FILES

TfReg produces many result files, here we will make an attempt to describe them in some detail. Some result files can be used again as input file such as the file with extension `.reg`.

8.1 `.reg` regression parameters

The regression parameters will be stored in a file with extension `.reg`
`examples/verify/epl2011-69.reg`

```

1 0.5 0.5 2
2 b
3 69 -86.7907 48.7773 15.3198 -6.02399
4 a
5 69
6 10 -39.4893 30.041
7 15 -26.6021 25.2989
8 20 -16.7027 21.2912
9 25 -10.8378 18.784
10 30 -3.52253 16.0363

```

Lets start from line 4 which contains the single letter `a`, this flags the start of all regression parameters which are length and salt concentration dependent. Line 5 holds the first salt concentration 69 mM in this case. The first regression equation starts at line 6 for all sequences of length 10 bp, and the next two numbers are the a_0 and a_1 coefficients of Eq. (4.1), that is

$$T_p(N = 10, [\text{Na}^+] = 69) = -39.4819 + 30.0373\tau, \quad (8.1)$$

At line 2 we see the letter `b` which flags the start of the regression parameters which are length dependent. There will be one line for each salt concentration. At line 3 we see the coefficients for 69 mM corresponding to $b_{0,0}, b_{1,0}, b_{1,0}$ and $b_{1,1}$ of Eq. (4.2),

$$a_0([\text{Na}^+] = 69) = -86.7785 + 48.7698N^{1/2}, \quad (8.2)$$

$$a_1([\text{Na}^+] = 69) = 15.3185 - 6.02284N^{1/2}. \quad (8.3)$$

these equations are necessary to calculate the a_0 and a_1 coefficients for sequence length which are not in the regression file.

8.2 `.dat` melting temperatures and melting index results

One of the main result files has the extension `.dat` and typically contains the melting temperatures and melting index calculated for each sequence. Here is an example (output of 9.2.2).

```

Main/Complementary alpha salt_concentration species_concentration temperature.measured ...
2 TACTAACATTAAGTA/ATGATTGTAATTGAT 4 69 2 0 0 34.222 0 0 0 0 0 2 6.006 2.45071 0 ...
3 ATACTTACTGATTAG/TATGAATGACTAATC 4 69 2 0 0 35.8269 0 0 0 0 0 2 6.318 2.51356 0 ...
4 GTACACTGTCTTATA/CATGTGACAGAATAT 4 69 2 0 0 39.712 0 0 0 0 0 2 7.106 2.66571 0 ...

```

The first line (which is too long to be shown completely) identifies each column. Importantly, the 7th column is the predicted temperature and the last column is the melting index τ . The 8th column is Z_y and the 9th column is $-kT \ln Z_y$, in this specific case they are zero as those are not calculated with `-res=prediction`.

8.3 .dat average opening if used with -res=averagey

8.3.1 Single sequence

If a single sequence is specified using `-seq` and optionally with `-cseq` the file will be arranged column-wise

```
1 0 1.30847
2 1 0.875882
3 2 0.848514
4 3 0.591862
5 4 0.635131
```

Shown are the first few lines of the result from example 9.6. The first column is the base pair position and the second column is the average opening $\langle y \rangle$, results are given in Ångstrom.

8.3.2 Multiple sequences

Multiple sequences should be arranged in one or more files, and the file names should be passed through the `-data` option. In this case the average opening will be arranged row-wise.

8.4 .ver quality of the prediction

The file with extension `.ver` is a short file intended to show how close the predicted melting temperatures are when compared to the experimental melting temperatures. This file only make sense if a data file with the experimental melting temperatures was given.

```
1 average diff_deviation sqr_diff sqrt_diff2 relative_sqr_diff
2 Tm 0.811472 0.771594 115.354 1.11975 0.0806365
3 prediction method=2
```

where the first number average difference in melting temperature prediction

$$\langle \Delta T \rangle = \frac{1}{N} \sum_{i=1}^N |\Delta T_i|, \quad (8.4)$$

the second is the standard deviation of ΔT_i ,

$$\delta(\Delta T) = \sqrt{\frac{1}{N} \sum_{i=1}^N [\Delta T_i - \langle \Delta T \rangle]^2} \quad (8.5)$$

the third is

$$\chi^2 = \sum_i [T_i - T'_i(P_k)], \quad (8.6)$$

and the 4th is

$$\Delta T_{RMS} = \sqrt{\frac{1}{N} \sum_{i=1}^N [\Delta T_i]^2} = \sqrt{\frac{\chi^2}{N}}. \quad (8.7)$$

and the last the relative squared difference

$$\left\langle \frac{\Delta T}{T} \right\rangle_{\text{RMS}} = \sqrt{\frac{1}{N} \sum_{i=1}^N \frac{[\Delta T_i]^2}{T_i^2}}, \quad (8.8)$$

8.5 Matrix files if used with -matrix=

If you specify a matrix directory say, `-matrix=open-1`, the program will create files containing all matrices used in the calculation:

Matrix files

```

dna-pb-69-CG_CG-270-A.AT_AT dna-pb-69-CG_CG-270-A.AT_CG dna-pb-69-CG_CG-270-A.AT_GC
dna-pb-69-CG_CG-270-A.AT_TA dna-pb-69-CG_CG-270-A.CG_AT dna-pb-69-CG_CG-270-A.CG_CG
dna-pb-69-CG_CG-270-A.CG_GC dna-pb-69-CG_CG-270-A.GC_AT dna-pb-69-CG_CG-270-A.GC_CG
dna-pb-69-CG_CG-270-A.TA_AT dna-pb-69-CG_CG-270-C.AT_AT dna-pb-69-CG_CG-270-C.AT_CG
dna-pb-69-CG_CG-270-C.AT_GC dna-pb-69-CG_CG-270-C.AT_TA dna-pb-69-CG_CG-270-C.CG_AT
dna-pb-69-CG_CG-270-C.CG_CG dna-pb-69-CG_CG-270-C.CG_GC dna-pb-69-CG_CG-270-C.GC_AT
dna-pb-69-CG_CG-270-C.GC_CG dna-pb-69-CG_CG-270-C.TA_AT dna-pb-69-CG_CG-270-eigenvalues
dna-pb-69-CG_CG-270-eigenvectors dna-pb-69-CG_CG-270-Y

```

The first part of the matrix file name **dna-pb-69** is the parameter identifier which you will find in the first line of the file **dna-pb-69.par**. The next field **CG_CG** refers to the base pair used as expansion basis (see option **-expand**). Next comes the temperature at which the calculation was performed (in this case 270 K). The matrix type is given by the last letter or identification and can be either A [Eq. (64) of Ref. 24] or C [Eq. (56) of Ref. 24]. The extensions like **AT_AT** refer to the nearest neighbour base pairs, there will be one for each type of nearest neighbour present in the sequence. The matrices labeled **eigenvalues** and **eigenvectors** are self-explanatory and represent the result of the diagonalization. Finally the matrix ending in **Y** is that of Eq. (60) of Ref. 24.

9 EXAMPLES

The example scripts are located in `/usr/share/TfReg/examples` (or `/usr/share/tfreg/examples`). In `/usr/share/TfReg/examples/verify` (or `/usr/share/tfreg/examples/verify`) you will find the output of some of these script which will allow you to check if your installed version of TfReg is working properly.

9.1 Task: given a set of melting temperatures, find the regression parameters

This scenario appears when you have some parameters for the hydrogen bond and perhaps stacking interaction and you wish to know how close these may get to experimental melting temperatures. First you will need to calculate the regression parameters which you may later use to calculate melting temperatures for untested sequences.

What you will need:

1. A set of experimental melting temperatures (for example `data/owczarzy04-69.dat`)
2. A set of parameters (for example `dna_pb_69.par`)

9.1.1 Results of Ref. 30

In this example we take the parameters which were calculated by the minimization procedure, stored in file `dna_pb_X.par` (X being 69, 119, 220, 621 or 1020), and calculate the regression parameters which are going to be stored in file `np2009-X`.

```
1  _____ examples/np2009-1.sh _____  
2  tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=regression -t=370 -v=1 ...  
3  -par=../data/dna_pb_69.par -data=../data/owczarzy04-69.dat  
4  
5  tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=regression -t=370 -v=1 ...  
6  -par=../data/dna_pb_119.par -data=../data/owczarzy04-119.dat  
7  
8  tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=regression -t=370 -v=1 ...  
9  -par=../data/dna_pb_220.par -data=../data/owczarzy04-220.dat  
10  
11 tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=regression -t=370 -v=1 ...  
12 -par=../data/dna_pb_621.par -data=../data/owczarzy04-621.dat  
13  
14 tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=regression -t=370 -v=1 ...  
15 -par=../data/dna_pb_1020.par -data=../data/owczarzy04-1020.dat
```

9.1.2 Results of Ref. 21

In this example we take the parameters which were calculated by the minimization procedure, stored in file `stat_jb_owczarzy04-69.par`, and calculate the regression parameters which are going to be stored in file `ep12011-69`. The file `var_jb2_owczarzy04_init.par` contains the initial parameters which were used in the minimization procedure, most parameters will be superseded by the contents of `stat_jb_owczarzy04-69.par`.

```
1  _____ examples/ep12011-69.sh _____  
2  tfreg -cutoff=10 -int=-1:200/400 -m=jb -o=ep12011-69 \  
3  -par=../data/dna_jb_69.par -data=../data/owczarzy04-69.dat \  
4  -pbc=0 -pm=2 -res=regression -t=370 -v=1
```

9.2 Task: prediction of DNA melting temperatures

You would like to predict melting temperatures of some nucleotide sequence for which you have no experimental data.

What you will need:

1. A set of parameters (for example `data/dna_pb_69.par`)
2. A file with calculated regression parameters (for example `data/np2009-1-69.reg` for salt concentration of 69 mM)

9.2.1 Single sequence example

If you want to predict the melting temperature of just one sequence the easiest is to specify the sequence as a command argument `-seq` as in the following example.

```
examples/predict-1.sh
1 tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -t=370 -v=1 -o=predict-1 \
2 -par=../data/dna_pb_69.par \
3 -reg=../data/np2009-1-69.reg -salt=69 \
4 -res=prediction \
5 -seq=ACAGCGAATGGACCTACGTGGCCTT
```

9.2.2 Multiple sequence example

If you want to predict the melting temperature of many sequences it is advisable to edit a simple file like this:

```
data/example2.dat
1 temperature
2 TACTAACATTA ACTA ATGATTGTAATTGAT 0 69 0
3 ATACTTACTGATTAG TATGAATGACTAATC 0 69 0
4 GTACACTGTCTTATA CATGTGACAGAATAT 0 69 0
```

the 69 refers to the salt concentration at which you want to predict these temperatures. Note that the second column is the complementary sequence of the first column. The way to run this example is as follows:

```
examples/predict-2.sh
1 tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=prediction -t=370 -v=1 ...
2 -par=../data/dna_pb_69.par \
3 -reg=../data/np2009-1-69.reg \
4 -res=prediction \
5 -data=../data/example2.dat
```

The result file is this:

```
examples/verify/predict-2.dat
1 Main/Complementary alpha salt_concentration species_concentration temperature.measured ...
2 TACTAACATTA ACTA/ATGATTGTAATTGAT 4 69 2 0 0 34.222 0 0 0 0 0 2 6.006 2.45071 0 ...
3 ATACTTACTGATTAG/TATGAATGACTAATC 4 69 2 0 0 35.8269 0 0 0 0 0 2 6.318 2.51356 0 ...
4 GTACACTGTCTTATA/CATGTGACAGAATAT 4 69 2 0 0 39.712 0 0 0 0 0 2 7.106 2.66571 0 ...
```

where the 7th column is the predicted temperature and the 8 column is the melting index τ .

9.3 Task: prediction of DNA melting temperatures with different salt concentrations

If you need to predict salt concentrations which differ from the ones currently provided you will need to generate a new *regression file*. To ease this task we provided a Perl script which does this regression for you called `tfreg-salt-regression.pl`. To use it you will need to provide some existing regression files for different salt concentrations which will use the regression equation (4.3). Here is an practical example where we use 5 files containing different salt concentration for which we generate a new regression file for a salt concentration of 50 mM.

```
Example usage of tfreg-salt-regression.pl
1 tfreg-salt-regression.pl 50 \
2 "np2009-1-69.reg,np2009-1-119.reg,np2009-1-220.reg,np2009-1-621.reg,np2009-1-1020.reg" \
3 new50.reg
```

the first argument is the new salt concentration, then a list of files comma-separated (the order is unimportant), and the last is the name of the new file. And here is an example script showing how to use the new regression file.

```

1  tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -t=370 -v=1 -o=predict-salt50 ...
2  -par=../data/dna_pb_69.par \
3  -reg=../data/new50.reg -salt=50 \
4  -res=prediction \
5  -seq=ACAGCGAATGGACCTACGTGGCCTT

```

Select the model parameters closest to the salt concentration you need. In the example above, we used 69 mM. The model parameters do not vary much with salt concentration, in fact, the Morse potentials hardly change (see 30).

9.4 Task: prediction of RNA melting temperatures

This is very much the same as predicting the melting temperatures for DNA described in the previous sections. The main difference is that you need to use the additional command parameter `-duplextype=RNA`.

What you will need:

1. A set of parameters for RNA (for example `data/rna_pb.par` from Ref. [23])
2. A file with calculated regression parameters (for example `data/reg_pb_xia98-t1.reg` for salt concentration of 1000 mM, this was published as supplementary tables IV and V of Ref. [23])

In the following example we calculate the melting temperatures of all sequences from Ref. 31.

```

1  tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -pm=2 -res=prediction -t=370 -v=1 ...
2  -duplextype=RNA \
3  -par=../data/rna_pb.par \
4  -reg=../data/reg_pb_xia98-t1.reg \
5  -data=../data/xia98-t1.dat

```

9.5 Task: prediction of melting temperatures DNA containing inosine mismatches

Optimized parameters were calculated in Ref. 13, and pre-calculated regression parameters are given in files

`deoxyinosine_pb_ia.reg` `deoxyinosine_pb_ic.reg` `deoxyinosine_pb_ig.reg`
`deoxyinosine_pb_ii.reg` `deoxyinosine_pb_it.reg`

for IA, IC, IG, II and IT inosine mismatches. The parameter file is given in `deoxyinosine_pb.par` and the sequence data from [32] are given in files

`watkins05ia.dat` `watkins05ic.dat` `watkins05ig.dat`
`watkins05ii.dat` `watkins05it.dat`

Please note that you will need to add the letter ‘I’ to the list of recognized characters with

`-dict=I:I`

see section 4.1.9 for further information.

9.6 Task: calculating the average opening $\langle y \rangle$

Using the option `-res=averagey` you will obtain the average opening $\langle y \rangle$ as a function of nucleotide position (see Ref. 24). Please note that the temperatures used here are unrelated to the predicted melting temperatures of section 9.2, and that for short sequences those temperatures may be unrealistically small. Typically you will want to use this for a qualitative study of localized helix opening.

What you will need:

1. A set of parameters (for example `data/dna_pb_69.par`)

9.6.1 Calculating $\langle y \rangle$ at a given temperature

```
examples/open-1.sh
1 tfreg -matrix=open-1 -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -t=270 -v=1 \
2 -o=open-1 \
3 -par=../data/dna_pb_69.par -salt=69 -res=averagey \
4 -seq=ACAGCGAATGGACCTACGTGGCCTT
```

the results will be in the file with extension `.dat`.

9.6.2 Calculating $\langle y \rangle$ for a range of temperatures

Here we use a simple bash loop to loop over the temperatures, each results is stored in a file `open-2-T.dat` where `T` is the temperature.

```
examples/open-2.sh
1 for T in {250..300}
2 do
3 tfreg -cutoff=10 -int=-1:200/400 -m=pb -pbc=0 -v=1 -salt=69 -res=averagey \
4 -o=open-2-$T -t=$T \
5 -par=../data/dna_pb_69.par -seq=ACAGCGAATGGACCTACGTGGCCTT
6 done
```

This example takes about 6 min to run (intel i7-2620M 2.70GHz).

9.6.3 Using Perl scripts and parallel processing

Sometimes, doing something more complicated with shell scripts can be very cumbersome. For example, a simple `for` loop with non-integer values can be very complicated. In this case, I would recommend spending a few moments to learn the basics of Perl. In the following example we run the average opening for RNA in increments of 0.5 K. Also each instance is send into background to be processed in a sort of poor man's parallel processing.

```
examples/open-rna-3.pl
1 #!/usr/bin/perl
2 my $prefix='../data';
3 $prefix=$ARGV[0] if (exists $ARGV[0]);
4 system("mkdir -p open-rna-3");
5 my $common='-cutoff=80 -int=-1:30/400 -m=pb -pbc=0 -v=1 -res=averagey -matrix=open-rna-3';
6
7 for (my $T=280; $T < 320; $T += 0.5)
8 {
9     my $fT=sprintf("%.1f",$T); #here we format $T with 1 decimal place
10    my $com="tfreg $common -o=open-rna-3/open-rna-3-$fT -t=$fT -duplextype=RNA -par=$prefix/rna_pb.par";
11    my $seq='-seq=GUGCCCAUUUAGGGUAUAUAUGGCCGAGUGAGCGAGCAGGAUCUCCAUUUUGACCGCAAAUUUGAACG';
12    system("$com $seq");# . ' < /dev/null &> open-rna-3/open-rna-3-$T.echo &');
13 }
```

Note in this example a `-cutoff=80` is used, that is, all matrix multiplications will be 80×80 . You may reduce the processing time by using a smaller value such as `-cutoff=10` but for higher temperatures the loss of precision is considerable.

10 FILES INCLUDED IN THE PACKAGE

10.1 Sequences and melting temperatures

10.1.1 Folder `/usr/share/data/`

This folder contains the input data such as parameter files (extension `.par`) and melting temperature data (extension `.dat`)

Melting temperatures

`chen12_1000.dat` RNA sequences with GU mismatches from Ref. [33]

`moreira15cy3.dat` DNA with Cy3 from Ref. [29]

`moreira15cy5.dat` DNA with Cy5 from Ref. [29]

`tna-tm.dat` DNA-TNA hybrid sequences from Refs. 5, 34, 35

`owczarzy04-*.dat` DNA at *=69, 119, 220, 621 and 1020 mM $[\text{Na}^+]$ from Ref. [36]

`rna-dna-1000.dat` DNA-RNA from Refs. [6, 8, 37]

`watkins05*.dat` for DNA with deoxyinosine, *=ia,ic,ig,ii and it from Ref. [32]

`xia98-t1.dat` RNA from table 1 of Ref. [31]

`dna_mm_60_tm.dat` DNA mismatches at 60 mM [3]

BIBLIOGRAPHY

- [1] M. Peyrard, A. R. Bishop, Statistical mechanics of a nonlinear model for DNA denaturation, *Phys. Rev. Lett.* 62 (23) (1989) 2755–2757.
- [2] G. Weber, N. Haslam, N. Whiteford, A. Prügel-Bennett, J. W. Essex, C. Neylon, Thermal equivalence of DNA duplexes without melting temperature calculation, *Nat. Phys.* 2 (2006) 55–59. doi:10.1038/nphys189.
- [3] L. M. Oliveira, A. Long, T. Brown, K. Fox, G. Weber, Melting temperature measurement and mesoscopic evaluation of single, double and triple dna mismatches, ChemRxiv Preprint. URL <https://doi.org/10.26434/chemrxiv.11977887.v1>
- [4] I. Ferreira, T. D. Amarante, G. Weber, DNA terminal base pairs have weaker hydrogen bonds especially for AT under low salt concentration, *J. Chem. Phys.* 143 (2015) 175101. doi:10.1063/1.4934783.
- [5] M. I. Muniz, H. H. Lackey, J. M. Heemstra, G. Weber, DNA/TNA mesoscopic modeling of melting temperatures suggest weaker hydrogen bonding of CG than in DNA/RNA, ChemRxiv Preprint. URL <https://doi.org/10.26434/chemrxiv.11929821.v1>
- [6] N. Sugimoto, S.-i. Nakano, M. Katoh, A. Matsumura, H. Nakamuta, T. Ohmichi, M. Yoneyama, M. Sasaki, Thermodynamic parameters to predict stability of RNA/DNA hybrid duplexes, *Biochem.* 34 (35) (1995) 11211–11216.
- [7] P. Wu, S. Nakano, N. Sugimoto, Temperature dependence of thermodynamic properties for DNA/DNA and RNA/DNA duplex formation, *Eur. J. Biochem.* 269 (12) (2002) 2821–2830. URL <http://www.blackwell-synergy.com/doi/abs/10.1046/j.1432-1033.2002.02970.x>
- [8] B. Rauzan, E. McMichael, R. Cave, L. R. Sevcik, K. Ostrosky, E. Whitman, R. Stegemann, A. L. Sinclair, M. J. Serra, A. A. Deckert, Kinetics and thermodynamics of DNA, RNA, and hybrid duplex formation, *Biochem.* 52 (5) (2013) 765–772.
- [9] E. d. O. Martins, V. B. Barbosa, G. Weber, DNA/RNA hybrid mesoscopic model shows strong stability dependence with deoxypyrimidine content and stacking interactions similar to RNA/RNA, *Chem. Phys. Lett.* 715C (2019) 14–19. doi:10.1016/j.cplett.2018.11.015.
- [10] P. Miranda, L. M. Oliveira, G. Weber, Mesoscopic modelling of Cy3 and Cy5 dyes attached to DNA duplexes, *Biophys. Chem.* 230C (2017) 62–67. doi:10.1016/j.bpc.2017.08.007. URL <http://www.sciencedirect.com/science/article/pii/S0301462217302831>
- [11] E. d. O. Martins, G. Weber, An asymmetric mesoscopic model for single bulges in RNA, *J. Chem. Phys.* 147 (2017) 155102. doi:10.1063/1.5006948.
- [12] T. D. Amarante, G. Weber, Evaluating hydrogen bonds and base stackings of single, tandem and terminal GU in RNA mismatches with a mesoscopic model, *J. Chem. Inf. Model.* 56 (1) (2016) 101–109. arXiv: <http://dx.doi.org/10.1021/acs.jcim.5b00571>, doi:10.1021/acs.jcim.5b00571. URL <http://dx.doi.org/10.1021/acs.jcim.5b00571>
- [13] R. V. Maximiano, G. Weber, Deoxyinosine mismatch parameters calculated with a mesoscopic model result in uniform hydrogen bonding and strongly variable stacking interactions, *Chem. Phys. Lett.* 631–632 (2015) 87–91. doi:10.1016/j.cplett.2015.04.045.
- [14] G. Weber, N. Haslam, J. W. Essex, C. Neylon, Thermal equivalence of DNA duplexes for probe design, *J. Phys.: Condens. Matter* 21 (2009) 034106. doi:10.1088/0953-8984/21/3/034106.
- [15] T. Dauxois, M. Peyrard, A. R. Bishop, Entropy-driven DNA denaturation, *Phys. Rev. E* 47 (1) (1993) R44–R47.

- [16] G. Weber, Sharp DNA denaturation due to solvent interaction, *Europhys. Lett.* 73 (5) (2006) 806–811. doi:10.1209/epl/i2005-10466-6.
- [17] M. Joyeux, S. Buyukdagli, Dynamical model based on finite stacking enthalpies for homogeneous and inhomogeneous DNA thermal denaturation, *Phys. Rev. E* 72 (2005) 051902.
- [18] S. Buyukdagli, M. Joyeux, Scaling laws at the phase transition of systems with divergent order parameter and/or internal length: The example of DNA denaturation, *Phys. Rev. E* 73 (5) (2006) 51910.
- [19] S. Buyukdagli, M. Joyeux, Theoretical investigation of finite size effects at DNA melting, *Phys. Rev. E* 76 (2) (2007) 021917.
- [20] S. Buyukdagli, M. Joyeux, Statistical physics of the melting of inhomogeneous DNA, *Physical Review E* 77 (3) (2008) 031903.
- [21] G. Weber, Finite enthalpy model parameters from DNA melting temperatures, *Europhys. Lett.* 96 (2011) 68001. doi:10.1209/0295-5075/96/68001. URL <http://iopscience.iop.org/0295-5075/96/6/68001>
- [22] T. D. Amarante, G. Weber, Analysing DNA structural parameters using a mesoscopic model, *J. Phys.: Conf. Ser.* 490 (1) (2014) 012203. doi:10.1088/1742-6596/490/1/012203. URL <http://iopscience.iop.org/1742-6596/490/1/012203>
- [23] G. Weber, Mesoscopic model parametrization of hydrogen bonds and stacking interactions of RNA from melting temperatures, *Nucleic Acids Res.* 41 (2013) e30. doi:10.1093/nar/gks964. URL <http://nar.oxfordjournals.org/content/41/1/e30>
- [24] Y.-L. Zhang, W.-M. Zheng, J.-X. Liu, Y. Z. Chen, Theory of DNA melting based on the Peyrard-Bishop model, *Phys. Rev. E* 56 (6) (1997) 7100–7115.
- [25] M. Barbi, S. Lepri, M. Peyrard, N. Theodorakopoulos, Thermal denaturation of a helicoidal DNA model, *Phys. Rev. E* 68 (2003) 061909.
- [26] K. J. Breslauer, R. Frank, H. Blocker, L. A. Marky, Predicting DNA duplex stability from the base sequence, *Proc. Natl. Acad. Sci. USA* 83 (11) (1986) 3746–3750.
- [27] J. SantaLucia, Jr., A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics, *Proc. Natl. Acad. Sci. USA* 95 (4) (1998) 1460–1465. arXiv:<http://www.pnas.org/cgi/reprint/95/4/1460.pdf>. URL <http://www.pnas.org/cgi/content/abstract/95/4/1460>
- [28] D. G. Norman, R. J. Grainger, D. Uhrin, D. M. Lilley, Location of cyanine-3 on double-stranded DNA: importance for fluorescence resonance energy transfer studies, *Biochemistry* 39 (21) (2000) 6317–6324.
- [29] B. G. Moreira, Y. You, R. Owczarzy, Cy3 and Cy5 dyes attached to oligonucleotide terminus stabilize DNA duplexes: Predictive thermodynamic model, *Biophys. Chem.* 198 (2015) 36–44. doi:10.1016/j.bpc.2015.01.001.
- [30] G. Weber, J. W. Essex, C. Neylon, Probing the microscopic flexibility of DNA from melting temperatures, *Nat. Phys.* 5 (2009) 769–773. doi:10.1038/nphys1371.
- [31] T. Xia, J. SantaLucia, Jr., M. E. Burkard, R. Kierzek, S. J. Schroeder, X. Jiao, C. Cox, D. H. Turner, Thermodynamic parameters for an expanded nearest-neighbor model for formation of RNA duplexes with Watson-Crick base pairs, *Biochem.* 37 (1998) 14719–14735.
- [32] J. Watkins, Norman E., J. SantaLucia, John, Nearest-neighbor thermodynamics of deoxyinosine pairs in DNA duplexes, *Nucleic Acids Res.* 33 (19) (2005) 6258–6267. arXiv:<http://nar.oxfordjournals.org/cgi/reprint/33/19/6258.pdf>, doi:10.1093/nar/gki918. URL <http://nar.oxfordjournals.org/cgi/content/abstract/33/19/6258>
- [33] J. L. Chen, A. L. Dishler, S. D. Kennedy, I. Yildirim, B. Liu, D. H. Turner, M. J. Serra, Testing the nearest neighbor model for canonical RNA base pairs: Revision of GU parameters, *Biochem.* 51 (16) (2012) 3508–3522.
- [34] K.-U. Schöning, P. Scholz, S. Guntha, X. Wu, R. Krishnamurthy, A. Eschenmoser, Chemical etiology of nucleic acid structure: the α -threofuranosyl-(3' \rightarrow 2') oligonucleotide system, *Science* 290 (5495) (2000) 1347–1351.

- [35] H. H. Lackey, E. M. Peterson, Z. Chen, J. M. Harris, J. M. Heemstra, Thermostability trends of TNA:DNA duplexes reveal strong purine dependence, *ACS synthetic biology* 8 (5) (2019) 1144–1152.
- [36] R. Owczarzy, Y. You, B. G. Moreira, J. A. Manthey, L. Huang, M. A. Behlke, J. A. Walder, Effects of sodium ions on DNA duplex oligomers: Improved predictions of melting temperatures, *Biochem.* 43 (2004) 3537–3554.
- [37] P. Wu, Y. Kawamoto, H. Hara, N. Sugimoto, Effect of divalent cations and cytosine protonation on thermodynamic properties of intermolecular DNA double and triple helices, *J. Inorg. Biochem.* 91 (2002) 277–285.