

INFERNAL User's Guide

Sequence analysis using profiles of RNA secondary structure consensus

<http://infern.janelia.org/>
Version 0.81; May 2007

Sean Eddy
HHMI Janelia Farm
19700 Helix Drive
Ashburn VA 20147
<http://selab.janelia.org/>

Copyright (C) 2001-2007 HHMI Janelia Farm.

Permission is granted to make and distribute verbatim copies of this manual provided the copyright notice and this permission notice are retained on all copies.

The free version of the Infernal software package is a copyrighted work that may be freely distributed and modified under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version. Alternative license terms may be obtained (for instance, for commercialization purposes) from the Office of Technology Management at Washington University. See the files COPYING and LICENSE that came with your copy of the Infernal software for details.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.

For a copy of the full text of the GNU General Public License, see www.gnu.org/licenses.

Contents

1	Introduction	4
2	Installation	5
	Quick installation instructions	5
	More detailed installation notes	5
	setting installation targets	5
	setting compiler and compiler flags	6
	turning on Large File Support (LFS)	6
	installing rigorous filters	7
	installing Message Passing Interface (MPI) programs	7
	Example configuration	7
3	Getting started	8
	Format of a simple input RNA alignment file	8
	Building a model with cmbuild	9
	Searching a sequence database with cmsearch	10
	Accelerating cmsearch	11
	Creating new multiple alignments with cmalign	12
	Accelerating cmalign	13
	Using optional annotation to completely specify model architecture to cmbuild	14
	Using local alignment in cmsearch and cmalign	14
	Parallelizing search and alignment with mpi-cmsearch and mpi-cmalign	16
	Getting more information	16
4	Profile SCFG construction: the cmbuild program	18
	Technical description of a covariance model	18
	Definition of a stochastic context free grammar	18
	SCFG productions allowed in CMs	18
	From consensus structural alignment to guide tree	19
	From guide tree to covariance model	21
	Parameterization	22
	Comparison to profile HMMs	22
	The cmbuild program, step by step	22
	Alignment input file	24
	Parsing secondary structure annotation	24
	Sequence weighting	25
	Architecture construction	26
	Parameterization	26
	Naming the model	26
	Saving the model	27

5	File and output formats	28
	RNA secondary structures: WUSS notation	28
	Full (output) WUSS notation	28
	Shorthand (input) WUSS notation	29
	Multiple alignments: Stockholm format	32
	A minimal Stockholm file	32
	Syntax of Stockholm markup	32
	Semantics of Stockholm markup	33
	Recognized #=GF annotations	33
	Recognized #=GS annotations	34
	Recognized #=GC annotations	34
	Recognized #=GR annotations	34
	Sequence files: FASTA format	34
	CM file format	35
	Null model file format	35
	Dirichlet prior files	35
6	Manual pages	39
	cmalign - use a CM to make a structured RNA multiple alignment	39
	Synopsis	39
	Description	39
	Options	39
	Expert Options	39
	cmbuild - construct a CM from an RNA multiple sequence alignment	43
	Synopsis	43
	Description	43
	Options	43
	Expert Options	43
	cmemit - generate sequences from a covariance model	47
	Synopsis	47
	Description	47
	Options	47
	Expert Options	47
	cmscore - align and score one or more sequences to a CM	49
	Synopsis	49
	Description	49
	Options	50
	Expert Options	50
	cmsearch - search a sequence database for RNAs homologous to a CM	52
	Synopsis	52
	Description	52
	Options	53
	Expert Options	53

1 Introduction

INFERNAL is a software package that allows you to make consensus RNA secondary structure profiles, and use them to search nucleic acid sequence databases for homologous RNAs, or to create new structure-based multiple sequence alignments.

To make a profile, you need to have a multiple sequence alignment of an RNA sequence family, and the alignment must be annotated with a consensus RNA secondary structure. The program **cmbuild** takes an annotated multiple alignment as input, and outputs a profile.

You can then use that profile to search a sequence database for homologs, using the program **cmsearch**.

You can also use the profile to align a set of unaligned sequences to the profile, producing a structural alignment, using the program **cmalign**. This allows you to build hand-curated representative alignments of RNA sequence families, then use a profile to automatically align any number of sequences to that profile. This seed alignment/full alignment strategy combines the strength of stable, carefully human-curated alignments with the power of automated updating of complete alignments as sequence databases grow. This is the strategy used to maintain the Rfam database of RNA multiple alignments and profiles.

INFERNAL is comparable to HMMER (hmmerr.janelia.org). The HMMER software package builds profile hidden Markov models (profile HMMs) of multiple sequence alignments. Profile HMMs capture only primary sequence consensus features. INFERNAL models are profile stochastic context-free grammars (profile SCFGs). Profile SCFGs include both sequence and RNA secondary structure consensus information.

Currently INFERNAL is really just an algorithm testbed. Output is rudimentary, and some desired features are missing. Most importantly, INFERNAL is very slow and CPU-intensive. You will probably need a large number of CPUs in order to use it for serious work. Planned algorithmic improvements should make it more practical in the future. We are making it available as a fully documented package now, only because INFERNAL has been pressed prematurely into service as the basis for constructing and maintaining the Rfam database of structurally annotated RNA multiple alignments (Griffiths-Jones et al., 2003). When we assign a 1.0 release number, that's when we'll think INFERNAL is ready for prime time. Until then, please bear with us.

2 Installation

Quick installation instructions

Download the source tarball (**infernald.tar.gz**) from <ftp://selab.janelia.org/pub/software/infernal/> or <http://infernald.janelia.org>

Unpack the software:

```
> tar xvf infernald.tar.gz
```

Go into the newly created top-level directory (named either **infernald**, or **infernald-xx** where **xx** is a release number:

```
> cd infernald
```

Configure for your system, and build the programs:

```
> ./configure
```

```
> make
```

Run the automated testsuite. This is optional. All these tests should pass:

```
> make check
```

The programs are now in the **src/** subdirectory. The user's guide (this document) is in the **documentation/userguide** subdirectory. The man pages are in the **documentation/manpages** subdirectory. You can manually move or copy all of these to appropriate locations if you want. You will want the programs to be in your **\$PATH**.

Optionally, you can install the man pages and programs in system-wide directories. If you are happy with the default (programs in **/usr/local/bin/** and man pages in **/usr/local/man/man1**), do:

```
> make install
```

That's all. More complete instructions follow, including how to change the default installation directories for **make install**.

More detailed installation notes

INFERNAL is distributed as ANSI C source code. It is designed to be built and used on UNIX platforms. It is developed on Intel GNU/Linux systems, and intermittently tested on a variety of vendor-donated UNIX platforms including Sun/Solaris, HP/UX, Digital Tru64, Silicon Graphics IRIX, IBM/AIX, and Intel/FreeBSD. It is not currently tested on either Microsoft Windows or Apple OS/X. It should be possible to build it on any platform with an ANSI C compiler. The software itself is vanilla POSIX-compliant ANSI C. You may need to work around the configuration scripts and Makefiles to get it built on a non-UNIX platform.

The GNU configure script that comes with INFERNAL has a number of options. You can see them all by doing:

```
> ./configure --help
```

All customizations can and should be done at the **./configure** command line, unless you're a guru delving into the details of the source code.

setting installation targets

The most important options are those that let you set the installation directories for **make install** to be appropriate to your system. What you need to know is that INFERNAL installs only two types of files: programs and man pages. It installs the programs in **--bindir** (which defaults to **/usr/local/bin**), and the man pages in the **man1** subdirectory of **--mandir** (default **/usr/local/man**). Thus, say you want **make**

install to install programs in `/usr/bioprog/bioprogs/bin/` and man pages in `/usr/share/man/man1`; you would configure with:

```
> ./configure --mandir=/usr/share/man --bindir=/usr/bioprog/bioprogs/bin
```

That's really all you need to know, since INFERNAL installs so few files. But just so you know; GNU configure is very flexible, and has shortcuts that accomodates several standard conventions for where programs get installed. One common strategy is to install all files under one directory, like the default `/usr/local`. To change this prefix to something else, say `/usr/mylocal/` (so that programs go in `/usr/mylocal/bin` and man pages in `/usr/mylocal/man/man1`, you can use the `--prefix` option:

```
> ./configure --prefix=/usr/mylocal
```

Another common strategy (especially in multiplatform environments) is to put programs in an architecture-specific directory like `/usr/share/Linux/bin` while keeping man pages in a shared, architecture-independent directory like `/usr/share/man/man1`. GNU configure uses `--exec-prefix` to set the path to architecture dependent files; normally it defaults to being the same as `--prefix`. You could change this, for example, by:

```
> ./configure --prefix=/usr/share --exec-prefix=/usr/share/Linux/
```

In summary, a complete list of the `./configure` installation options that affect INFERNAL:

Option	Meaning	Default
<code>--prefix=PREFIX</code>	architecture independent files	<code>/usr/local/</code>
<code>--exec-prefix=EPREFIX</code>	architecture dependent files	<code>PREFIX</code>
<code>--bindir=DIR</code>	programs	<code>EPREFIX/bin/</code>
<code>--mandir=DIR</code>	man pages	<code>PREFIX/man/</code>

setting compiler and compiler flags

By default, **configure** searches first for the GNU C compiler `gcc`, and if that is not found, for a compiler called `cc`. This can be overridden by specifying your compiler with the `CC` environment variable.

By default, the compiler's optimization flags are set to `-g -O2` for `gcc`, or `-g` for other compilers. This can be overridden by specifying optimization flags with the `CFLAGS` environment variable.

For example, to use an Intel C compiler in `/usr/intel/ia32/bin/icc` with optimization flags `-O3 -ipo`, you would do:

```
> env CC=/usr/intel/ia32/bin/icc CFLAGS="-O3 -ipo" ./configure
```

which is the one-line shorthand for:

```
> setenv CC /usr/intel/ia32/bin/icc
```

```
> setenv CFLAGS "-O3 -ipo"
```

```
> ./configure
```

If you are using a non-GNU compiler, you will almost certainly want to set `CFLAGS` to some sensible optimization flags for your platform and compiler. The `-g` default generated unoptimized code. At a minimum, turn on your compiler's default optimizations with `CFLAGS=-O`.

turning on Large File Support (LFS)

INFERNAL has one optional feature: support for Large File System (LFS) extensions that allow programs to access files larger than 2 GB. LFS is rapidly becoming standard, but not yet standard enough to be default. If you do anything with Genbank files or large genome files (like the 3 GB human genome), you will need LFS support. LFS is enabled with the `--enable-lfs` option:

```
> ./configure --enable-lfs
```

installing rigorous filters

INFERNAL includes programs by Zasha Weinberg that implement rigorous filtering. This software requires a C++ compiler, and also relies on an external library, CFSQP, that is not included in this distribution, but can be obtained by request from <http://www.aemdesign.com/>. To build the executables, include these two options to configure:

```
> ./configure --with-rigfilters --with-cfsqp=/path/to/cfsqp
```

installing Message Passing Interface (MPI) programs

INFERNAL includes two programs `mpi-cmsearch` and `mpi-cmalign` that use MPI parallelization. You have the option of compiling these two executables, if you have MPI installed. (We use LAM MPI here, but alternative MPI libraries should also work.) To enable MPI and compile these two additional executables, add `--enable-mpi` to the configuration command:

```
> ./configure --enable-mpi
```

A `make install` following a `./configure --enable-mpi` will install the non-MPI programs (`cmalign`, `cmbuild`, `cmemit`, `cmscore`, and `cmsearch`) as well as `mpi-cmsearch` and `mpi-cmalign`.

`mpi-cmsearch` and `mpi-cmalign` have the same command line arguments as `cmsearch` and `cmalign`, but you must run them in an MPI environment with `mpirun` or `mpiexec`, for instance (in our LAM environment) with:

```
> mpirun C mpi-cmsearch query.cm target.fa
```

Example configuration

The Intel GNU/Linux version installed at Janelia Farm is configured as follows:

```
> env CFLAGS="-O3" ./configure --enable-mpi --enable-lfs --prefix=/usr/local/infernal
```


Here's a quick walk-through of the package. Here, we will a) build a model of an RNA multiple alignment using **cmbuild**; b) use that model to search for new homologs using **cmsearch**; and c) use that model to align new sequences, and create a new multiple alignment, using **cmalign**.

tutorial.sto A multiple alignment of five tRNA sequences. This file is a simple example of *Stockholm format* that INFERNAL uses for structurally-annotated alignments.

tutorial.big.db A larger 300,000 nt sequence database; with the same tRNA as in **tutorial.db**.
Used to demonstrate HMM filtered search.

tutorial.fa The same sequences as in **tutorial.sto**, plus one more tRNA with an internal deletion (to demonstrate local alignment), in unaligned FASTA format.

Look at the alignment file **tutorial.sto** in the **intro/** subdirectory of the INFERNAL distribution. It is shown below, with a secondary structure of the first sequence shown to the right for reference (yeast Phe tRNA, labeled as “tRNA1” in the file):

For now, what you need to know about the key features of the input file is:

- The alignment is in an interleaved format, like other common alignment file formats such as CLUSTALW. Lines consist of a name, followed by an aligned sequence; long alignments are split into blocks separated by blank lines.
- Each sequence must have a unique name. (This is important!)

- For residues, any one-letter IUPAC nucleotide code is accepted, including ambiguous nucleotides. Case is ignored; residues may be either upper or lower case.
- Gaps are indicated by the characters `.`, `-`, `-`, or `~`. (Blank space is not allowed.)
- A special line starting with `#=GC SS_cons` indicates the secondary structure consensus. Gap characters annotate unpaired (single-stranded) columns. Base pairs are indicated by any of the following pairs: `<>`, `()`, `[]`, or `{}` . No pseudoknots are allowed; the open/close-brackets notation is only unambiguous for strictly nested base-pairing interactions.
- The file begins with the special tag line `# STOCKHOLM 1.0`, and ends with `//`.

Building a model with `cmbuild`

To build a model from this alignment, do:

```
> cmbuild my.cm tutorial.sto
```

Almost instantly, `cmbuild` reads in the alignment, constructs a model, and saves that model to the new file `my.cm`. It is a convention to use the `.cm` suffix for model files; CM stands for “covariance model”, another name for the profile SCFG architecture used by INFERNAL (Eddy and Durbin, 1994).

The output `cmbuild` contains information about the size of your input alignment (in aligned columns and # of sequences), and about the size of the resulting model. You don’t need to understand this to use the model, so for now we’ll skip describing the output.

The result, the model file in `my.cm` is a text file. You can look at it (e.g. `> more my.cm`) if you like, but it isn’t really designed to be human-interpretable. You can treat `.cm` files as compiled models of your RNA alignment.

Searching a sequence database with `cmsearch`

You can use your model to search for new homologues of your RNA family. The file `tutorial.db` contains an example sequence “database”: one 300 nt sequence, with yeast tRNA-Phe embedded at position 101...173. To search it, do:

```
> cmsearch my.cm tutorial.db
```

`cmsearch` now searches both strands of each sequence in the target database, and returns alignments for high scoring hits. In this case, it returns one hit:

```
CM 1: tutorial.1
CM lambda and K undefined -- no statistics
Using CM score cutoff of 0.00
>example

Plus strand results:

Query = 1 - 72, Target = 101 - 173
Score = 65.96, GC = 53

      ((((((, <<<<_. ____>>>>, <<<<_____>>>>, , , , <<<<_____
1  gccgacaUaGcgCAGu.GGuAgcgCgccagcuUgaaaagcuggAGguccgggGUUCgAuu 59
   GC:+A::UAGC:CAGU GG AG:GCGCCAG:+UGAA A:CUGGAGGUCC:G:GUUCGAU
101 GCGGAUUUAGCUCAGUuGGGAGAGCGCCAGACUGAAGAUCUGGAGGUCCUGUGUUCGAUC 160

      >>>>))))))):
60 Ccccgugucggca 72
   C:C:G::U+:GCA
161 CACAGAAUUCGCA 173

//
Fin
```

The first line gives the name of the CM (this can be defined in the input Stockholm alignment file or as an option to `cmbuild`, as described later). The next two lines give information on E-value statistics, which were turned off in our example search (more on this later). Next comes the results section, the name of each target sequence in the target database is given starting with a `>`, in this case there is only one: **example**. Next, all the hits to the top (Watson) strand of **example** are given, in this example there is a single hit from position 101 to 173 with a score of 65.96 bits. As discussed next, E-value statistics can give an estimate of the significance of bit scores. We have not rigorously tested our implementation of E-values, so by default E-values are not calculated. Larger bit scores are better. As a rough guide, scores greater than the log (base two) of the target database size are significant. Here, given a 600 nt target (300 nt \times 2 strands), scores over 9-10 bits are significant - so a score of 65.96 is a good hit.

The alignment is shown in a BLAST-like format, augmented by secondary structure annotation.

The top line shows the predicted secondary structure of the target sequence. The format is a little fancier and more informative than the simple least-common-denominator format we used in the input alignment file. It's designed to make it easier to “see” the secondary structure by eye. The format is described in detail later; for now, here's all you need to know. Base pairs in simple stem loops are annotated with `<>` characters. Base pairs enclosing multifurcations (multiple stem loops) are annotated with `()`, such as the tRNA acceptor stem in this example. In more complicated structures, `[]` and `{ }` annotations also show up, to reflect deeper nestings of multifurcations. For single stranded residues, `_` characters mark hairpin loops; `^` characters mark interior loops and bulges; `,` characters mark single-stranded residues in multifurcation loops; and `:` characters mark single stranded residues external to any secondary structure. Insertions relative to this consensus are annotated by a `.` character.

The second line shows that consensus of the query model. The highest scoring residue sequence is shown. Upper case residues are highly conserved. Lower case residues are weakly conserved or unconserved.

The third line shows where the alignment score is coming from. For a consensus base pair, if the observed pair is the highest-scoring possible pair according to the consensus, both residues are shown in upper case; if a pair has a score of ≥ 0 , both residues are annotated by : characters (indicating an acceptable compensatory base pair); else, there is a space, indicating that a negative contribution of this pair to the alignment score. For a single-stranded consensus residue, if the observed residue is the highest scoring possibility, the residue is shown in upper case; if the observed residue has a score of ≥ 0 , a + character is shown; else there is a space, indicating a negative contribution to the alignment score.

Finally, the fourth line is the target sequence.

Let's repeat this example search, now with E-values. When run with the **-E** option, **cmsearch** estimates the E-value of each hit found. This is computed by sampling 1000 random sequences and determining the bit score of the best hit within each sequence, these scores are fit to a Gumbel distribution. The search is then carried out, and the Gumbel distribution is used to estimate the significance of the bit scores of the hits found. This procedure used by **cmsearch** is very similar to that described for the RSEARCH program in (Klein and Eddy, 2003). Because using E-values requires searching 1000 random sequences, it takes quite a bit longer than a regular search. To repeat the search saving all hits with E-values of 10 or less, do:

```
> cmsearch -E 10 my.cm tutorial.db

CM 1: tutorial.1
CM statistics calculated with simulation of 1000 samples of length 180
Random seed: 1177592728
No partition points
Using CM E cutoff of 10.00

>example

Plus strand results:

Query = 1 - 72, Target = 101 - 173
Score = 65.96, E = 6.27e-12, P = 6.27e-12, GC = 53

(((((((, <<<<_.>>>>, <<<<____>>>>, , , , <<<<_____
1 gccgacaUaGcgCagcGGuAgcgCgcccagcuUgaaaagcuggAGguccgggGUUCgAuu 59
GC:+A::UAGC:CAGU GG AG:GCGCCAG:+UGAA A:CUGGAGGUCC:G:GUUCGAU
101 GCGGAUUUAGCUCAGUuGGGAGAGCGCCAGACUGAAGAUCUGGAGGUCCUGUGUUCGAUC 160

>>>>))))))):
60 Ccccgugucggca 72
C:C:G::U+:GCA
161 CACAGAAUUCGCA 173

//
Fin
```

The output tells us that the E-value is $6.27e - 12$, this is the number of hits we expect to find with a bit score of 65.96 or better if we were searching a database of random sequence of length 300.

Accelerating **cmsearch**

The **cmsearch** program is very slow, and several techniques have been developed for acceleration. We'll briefly discuss two here (for more information see the **cmsearch** man pages at the end of this guide).

The first acceleration strategy is called query-dependent banding (QDB). Briefly, this technique precalculates regions of the dynamic programming matrix that have very low probability before the search, and skips these regions during the search (Nawrocki and Eddy, 2007). By default, QDB is on in **cmsearch**, you can turn it off with the **--noqdb** option. QDB offers about a four-fold speedup for an average RNA family, with very little impact on sensitivity. In general, the longer the average sequence length of an RNA family, the greater the acceleration achieved by using QDB.

The second acceleration technique is filtering the target database with a profile HMM, and using the expensive CM based search on only subsequences that receive high HMM scores. This was pioneered by Zasha Weinberg and Larry Ruzzo at the University of Washington (Weinberg and Ruzzo, 2004a; Weinberg and Ruzzo, 2004b; Weinberg and Ruzzo, 2006). Here we consider two of their techniques for constructing a profile HMM. The first builds a “rigorous filter” that is guaranteed to find all CM hits above a certain threshold. We have incorporated Weinberg’s implementation of “rigorous filters” into INFERNAL but not within the **cmsearch** program. (For details see the file **rigfilters.doc.html** in the **documentation/userguide** subdirectory of the top level **infernald** directory.) The second HMM filtering method involves a maximum likelihood (ML) HMM that is not guaranteed to find all hits, but is more practical than rigorous filters for many RNA families. We have implemented a version of ML HMMs in INFERNAL, which we call these CM Plan 9 HMMs (CP9 HMMs), to distinguish them from the Plan 7 HMMs of HMMER. CP9 HMMs can be used to accelerate **cmsearch** using the **--hmmfilter** option, as well as to accelerate **cmalign** as discussed later in this guide. We still consider our implementation of the CP9 HMM filtering technique experimental, so it is not on by default. When the **--hmmfilter** option is enabled, E-values are calculated for the CP9 HMM in the same manner that CM E-values are calculated, except this time each of the 1000 random sequences is searched with the HMM, so it is much quicker. The database is then searched with the CP9 HMM and high scoring HMM hits with E-values of 500 or less (along with some flanking sequence) are searched again with the CM. The default E-value threshold of 500 can be changed with the **--hmmE** option. Also, a bit score threshold can be set instead of an E-value threshold with the **--hmmT** option. Here’s an example searching a larger sample database called **tutorial.big.db**, which is 300,000 nt and still has the tRNA from position 101 to 173:

```
> cmsearch --hmmfilter my.cm tutorial.big.db
```

This search finds the real tRNA and takes about 45 seconds on my machine. The same search without HMM filtering and with the **--time** option to print timings can be run with:

```
> cmsearch --time my.cm tutorial.big.db
```

On my machine, this non-filtered search takes about 8 minutes, meaning the HMM gave us about a 10-fold speedup. This speedup will increase for RNA families with longer sequence lengths, or if you lower the E-value cutoff with **--hmmE**. The speedup comes at a cost to sensitivity though, and we haven’t rigorously tested this to our satisfaction. This is the main reason this option is not on by default. There are more options related to HMM filtering for **cmsearch** discussed in the man pages at the end of this guide.

Creating new multiple alignments with **cmalign**

You can also use a model to structurally align any number of new RNA sequences to your consensus structure. This is how the RFAM database is constructed: we start with “seed” alignment, build a CM of it, and use that CM to align all known members of the sequence family and create a “full” alignment. This allows us to maintain representative seed alignments that are stable and small enough to be human-curated, while still being able to automatically incorporate and align all homologues detected in the rapidly growing public sequence databases.

An example of some unaligned tRNA sequences are in the file **tutorial.fa**. (In fact, these are the same sequences that are in **tutorial.sto**, reformatted into unaligned FASTA format; plus a new sequence, tRNA6, which was created by deleting some residues out of the middle of tRNA1. tRNA6 will be used a little later to demonstrate local alignment.)

To align these sequences to the model we made in **my.cm**, do:

```
> cmlalign my.cm tutorial.fa
```

This results in an alignment that looks like:

```
# STOCKHOLM 1.0
#=GF AU      Infernal 0.72

tRNA1          GCGGAUUUAGCUCAGUuGGG.AGAGCGCCAGACUGAAGAUCUGGAGGUCC
tRNA2          UCCGAUAUAGUGUAAC.GGCuAUCACAUCACGCUUUCACCGUGGAGA-CC
tRNA3          UCCGUGAUAGUUUAU.GGUcAGAAUGGGCGCUUGUCGCGUGCCAGA-UC
tRNA4          GCUCGU AUGGCGCAGU.GGU.AGCGCAGCAGAUUGCAAAUCUGUUGGUCC
tRNA5          GGGCACAU GGGCGCAGUuGGU.AGCGCGCUUCCUUGCAAGGAAGAGGUCA
tRNA6          GCGGAUUUAGCUCAGUuGGG.AGAGCGC-----CAGAC----GAGGUCC
#=GC SS_cons   ((((((, ,<<<<____.____.>>>>, <<<<_____>>>>, , , , , <
#=GC RF        gccgacaUaGcgAgu.GGu.AgcgCgccagcuUgaaaagcuggAGgucc

tRNA1          UGUGUUCGAUCCACAGAAUUCGCA
tRNA2          GGGGUUCGACUCCCCGUAUCGGAG
tRNA3          GGGGUUCAAUUCCCCGUCGCGAG
tRNA4          UUAGUUCGAUCCUGAGUGCGAGCU
tRNA5          UCGGUUCGAUUCGGUUGCGUCCA
tRNA6          UGUGUUCGAUCCACAGAAUUCGCA
#=GC SS_cons   <<<<_____>>>>))))):
#=GC RF        gggGUUCgAuuCcccgugucgca
//
```

In the aligned sequences, a . character indicates an inserted column relative to consensus; the - character is an alignment pad. A - character is a deletion relative to consensus.

The symbols in the consensus secondary structure annotation line have the same meaning that they did in a pairwise alignment from **cmsearch**.

The **#=GC RF** line is *reference annotation*. Non-gap characters in this line mark consensus columns; **cmlalign** uses the residues of the consensus sequence here, with upper case denoting strongly conserved residues, and lower case denoting weakly conserved residues. Gap characters (specifically, the . pads) mark insertions relative to consensus. As described below, **cmbuild** is capable of reading these RF lines, so you can specify which columns are consensus and which are inserts (otherwise, **cmbuild** makes an automated guess, based on the frequency of gaps in each column).

If you want to save the alignment to a file, you can use the **-o** option:

```
> cmlalign -o my.sto my.cm tutorial.fa
```

We'll use this **my.sto** alignment file in an upcoming section.

Accelerating **cmlalign**

Earlier in this guide, we described the use of CM Plan 9 HMMs (CP9 HMMs) as filters to accelerate **cmsearch** using the **--hmmfilter** option. CP9 HMMs can also be used to accelerate CM alignment. The basic idea is to first align each sequence to the CP9 HMM and to use that alignment to derive constraints, which we call bands, for the more computationally expensive CM alignment. These bands are then applied during CM alignment, reducing the time necessary for alignment. This technique can be enabled with the **--hbanded** option to **cmlalign**. Here's an example:

```
> cmalign --hbanded my.cm tutorial.fa
```

You may notice that this time the program runs a bit quicker than before. It's almost imperceptible here, but for large RNAs the difference is significant. For example, alignment of small subunit ribosomal RNA (SSU rRNA) sequences (which are roughly 1500 nucleotides) using the **--hbanded** option takes close to one second for one sequence, whereas non-banded alignment takes several minutes. Importantly, using the HMM banded technique sacrifices the guarantee that the optimal alignment will be found, but our tests suggest this happens very rarely for non-local alignment (these tests were done using the **cmscore** program). However, we have not rigorously tested the **--hbanded** option for local alignment. More information on HMM banded alignment can be found in the **cmalign** and **cmscore** man pages at the end of this guide.

Using optional annotation to completely specify model architecture to **cmbuild**

cmbuild needs to know two things to convert your alignment into a profile SCFG.

First, it needs to know the consensus secondary structure. It reads this from the **#=GC SS_cons** line, as described above. This annotation is mandatory.

It also needs to know which columns are consensus, and which columns are insertions relative to consensus. By default, it will determine this by a simple rule: if a column contains more than a certain fraction of gap characters (default >50%), the column is called an insertion. This may not be what you want; for instance, maybe you are trying to iteratively build models based on larger and larger numbers of sequences (based on an RFAM seed, say), but you don't want the curated consensus model architecture to change just because you added some new sequences to the alignment.

You can optionally override that default and specify the complete architecture of the model, using both a **#=GC SS_cons** structure annotation line and a **#=GC RF** reference column annotation line. To do this, you use the **--rf** flag to **cmbuild**.

For example, to build a model called **second.cm** from **my.sto** that has the same architecture as **my.cm**, you would do:

```
> cmbuild --rf second.cm my.sto
```

Since **cmalign** leaves an RF line on the alignments it generates, the **--rf** option allows you to propagate your consensus structure into new, larger alignments. The RF line is also handy when you want the model's coordinate system to be the same as a canonical, well-studied single sequence: you can simply use that sequence as the RF line, or manually create any consensus coordinate system you like. (This is the origin of RF as the "reference line", e.g. giving a reference coordinate system.) The only thing that matters in the RF line is nongap versus gap characters: the line can be as simple as x's marking consensus columns, .'s for insert columns.

Using local alignment in **cmsearch** and **cmalign**

The examples above required the entire model to match a subsequence of the target: so-called *glocal* alignment (global with respect to the query model, local with respect to the target sequence). But in many cases, a homologous RNA structure has undergone enough changes that parts of its structure cannot be aligned to the consensus. *Local* alignment, in which only part of the query model needs to match the target to detect a hit, can be a more sensitive searching strategy.

In primary sequence alignment, local alignment means an alignment of two subsequences of the query and target. In aligning a query RNA structure to a target sequence, local alignment means starting and ending at points inside the query structure – which, when you map that idea onto linear sequence, means an alignment that may consist of more than one discontinuous subsequence. We'll demonstrate this by example

for now, and describe local alignment is described in detail later. For the purposes of the tutorial, all you really need to know is how to activate it. It is not the default behavior for either **cmsearch** or **cmalign**.

Local alignment is activated for **cmsearch** by using the **--local** option. For example:

```
> cmsearch --local my.cm tutorial.fa
```

Look at the first alignment for the target sequence tRNA6 (the second to last alignment in the output):

```
>tRNA6

Plus strand results:

Query = 1 - 72, Target = 1 - 63
Score = 44.94, GC = 55

      ((((((, <<<<_.____>>>>, <~~~~~>, , , , <<<<_____>>>>))))))
1 gccgacaUaGcgCagu.GGuAgcgCgc*[15]*gAGguccgggGUUCgAuuCcccguguc 68
  GC:+A::UAGC:CAGU GG AG:GCGC      GAGGUCC:G:GUUCGAU C:C:G::U+
1 GCGGAUUUAGCUCAGUuGGGAGAGCGC*[ 5]*GAGGUCCUGUGUUCGAUCCACAGAAUU 59

      ))):
69 ggca 72
      :GCA
60 CGCA 63
```

The ***[15]*** and ***[5]*** in the query and target, respectively, indicate that 15 consensus residues and 5 target residues were left unaligned; the target does not appear to have the consensus structure in this region. (No kidding, since I made the tRNA6 example sequence by deleting part of the anticodon stem.) The structure annotation line is marked with **~~~~~** to indicate the gap in the alignment, and to distinguish local alignment induced gaps from normal insertions (which are marked with **.** characters).

You can activate local alignment in **cmalign** with the **-l** option:

```
> cmalign -l my.cm tutorial.fa
```


This results in the following alignment: ¹

```
# STOCKHOLM 1.0
#=GF AU      Infernal 0.72

tRNA1          GCGGAUUUAGCUCAGUuGGG.AGAGCGCCAGACUGAAGA....UCUGGA
tRNA2          UCCGAUAUAGUGUAAC.GGCuAUCACAUCACGCUUUCAC....CGUGGA
tRNA3          UCCGUGAUAGUUUAU.GGUcAGAAUGGGCGCUUGUCGC....GUGCCA
tRNA4          GCUCGU AUGGCGCAGU.GGU.AGCGCAGCAGAUUGCAA....UCUGUU
tRNA5          GGGCACAU GGGCGCAGUuGGU.AGCGCGCUUCCCUUGCAA....GGAAGA
tRNA6          GCGGAUUUAGCUCAGUuGGG.AGAGCGC-----cagac----GA
#=GC SS_cons   ((((((, ,<<<<____.____._>>>>, <<<<_____~~~~~>>>>,
#=GC RF        gccgacaUaGcgcAgu.GGu.AgcgCgccagcuUgaaaa~~~~~gcuggA

tRNA1          GGUCCUGUGUUCGAUCCACAGAAUUCGCA
tRNA2          GA-CCGGGGUUCGACUCCCCGUAUCGGAG
tRNA3          GA-UCGGGGUCAAUCCCCGUCGCGGAG
tRNA4          GGUCCUUAAGUUCGAUCCUGAGUGCGAGCU
tRNA5          GGUCAUCGGUUCGAUUCGGUUGCGUCCA
tRNA6          GGUCCUGUGUUCGAUCCACAGAAUUCGCA
#=GC SS_cons   , , , , <<<<_____>>>>))))))):
#=GC RF        GguccgggGUUCgAuuCcccgugucggca
//
```

Note how the local alignment is represented for tRNA6. The deleted consensus columns are marked by - characters. The unaligned “insertion” is shown in its own columns; those columns are again marked with ~ characters in the consensus secondary structure annotation and the reference (RF) annotation lines.

Parallelizing search and alignment with `mpi-cmsearch` and `mpi-cmalign`

As mentioned in the Installation section, INFERNAL contains two MPI programs: `mpi-cmsearch` and `mpi-cmalign`. These programs must be run using `mpirun`. These MPI programs are under current development, and we have only tested them using the LAM implementation of MPI. Here are example runs using LAM:

```
> mpirun C mpi-cmsearch query.cm target.fa
> mpirun C mpi-cmalign query.cm target.fa
```

Getting more information

For a quick refresher on the command line usage of any program and its commonly used options, just type the name of the program with no other arguments: e.g.

```
> cmbuild
```

and you’ll get a brief help:

```
FATAL: Incorrect number of arguments.
Usage: cmbuild [-options] <cmfile output> <alignment file>
The alignment file is expected to be in Stockholm format.
Available options are:
```

¹The discontinuity of structural local alignment presents a quandary for representing multiple alignments. On the one hand, you might not want to even show the unaligned target residues in the gap (e.g., cagac) – they aren’t aligned to the model. On the other hand, you sort of expect that if you pull an RNA sequence out of a multiple alignment, it represents a true subsequence of a larger sequence, not a concatenation of disjoint subsequences – you’d at least like some indication of where some residues have gone missing. One option would be to leave a *[5]* in the gap, as in the pairwise representation; but one of the nice properties of Stockholm format is that it’s easy to interconvert it to other alignment formats just by stripping off everything by the name/sequence part of the alignment, and sticking non-sequence characters like *[5]* in the alignment would prevent that.

```

-h      : help; print brief help on version and usage
-n <s>  : name this CM <s>
-A      : append; append this CM to <cmfile>
-F      : force; allow overwriting of <cmfile>

```

For version information and a complete listing of options, use the **-h** option with any program, e.g.

```
> cmscore -h
```

and you'll see something like:

```

cmscore - score RNA covariance model against sequences
Infernal 0.72 (January 2007)
Copyright (C) 2001-2007 HHMI Janelia Farm
Freely distributed under the GNU General Public License (GPL)
-----
Usage: cmscore [-options] <cmfile> <sequence file>
Most commonly used options are:
-h      : help; print brief help on version and usage
-i      : print individual timings & score comparisons, not just summary

Expert options
--local      : align locally w.r.t the model
--sub        : build sub CM for columns b/t HMM predicted start/end points
--regress <f> : save regression test data to file <f>
--stringent  : require the two parse trees to be identical
--trees      : print parsetrees

Expert stage 2 alignment options, to compare to stage 1 (D&C non-banded)
--std        : compare divide and conquer versus standard CYK [default]
--qdb        : compare non-banded d&c versus QDB standard CYK
--qdbsmall   : compare non-banded d&c versus QDB d&c
--qdbbboth   : compare          QDB d&c versus QDB standard CYK
--beta <x>   : set tail loss prob for QDB to <x> [default:1E-7]
--hbanded    : compare non-banded d&c versus HMM banded CYK
--tau <x>    : set tail loss prob for HMM bands to <x> [default: 1E-7]
--hsafe      : realign (non-banded) seqs with HMM banded CYK score < 0 bits
--hmmonly    : align with the CM Plan 9 HMM (only gives timings)
--scoreonly   : for standard CYK stage, do only score, save memory

Expert stage 2-N alignment options, to compare to stage 1 (D&C non-banded)
For --hbanded or --qdb, try multiple tau or beta values, all will = 10^-n
--betas <n>   : set initial (stage 2) tail loss prob to 10^-(<x>) for qdb
--betae <n>   : set final   (stage N) tail loss prob to 10^-(<x>) for qdb
--taus <n>    : set initial (stage 2) tail loss prob to 10^-(<x>) for hmm
--taue <n>    : set final   (stage N) tail loss prob to 10^-(<x>) for hmm

```

More detailed information on usage and command line options is available in UNIX manual pages. If they have been installed for your system, you can see this information with, e.g.:

```
> man cmalign
```

Copies of the man pages are also provided at the end of this guide.

4 Profile SCFG construction: the `cmbuild` program

INFERNAL builds a model of consensus RNA secondary structure using a formalism called a *covariance model* (CM), which is a type of *profile stochastic context-free grammar* (profile SCFG) (Eddy and Durbin, 1994; Durbin et al., 1998; Eddy, 2002).

What follows is a technical description of what CM is, how it corresponds to a known RNA secondary structure, and how it is built and parameterized.² You certainly don't have to understand the technical details of CMs to understand `cmbuild` or INFERNAL, but it will probably help to at least skim this part. After that is a description of what the `cmbuild` program does to build a CM from an input RNA multiple alignment, and how to control the behavior of the program.

Technical description of a covariance model

Definition of a stochastic context free grammar

A stochastic context free grammar (SCFG) consists of the following:

- M different nonterminals (here called *states*). I will use capital letters to refer to specific nonterminals; V and Y will be used to refer generically to unspecified nonterminals.
- K different terminal symbols (e.g. the observable alphabet, a,c,g,u for RNA). I will use small letters a, b to refer generically to terminal symbols.
- a number of *production rules* of the form: $V \rightarrow \gamma$, where γ can be any string of nonterminal and/or terminal symbols, including (as a special case) the empty string ϵ .
- Each production rule is associated with a probability, such that the sum of the production probabilities for any given nonterminal V is equal to 1.

SCFG productions allowed in CMs

A CM is a specific, repetitive SCFG architecture consisting of groups of model states that are associated with base pairs and single-stranded positions in an RNA secondary structure consensus. A CM has seven types of states and production rules:

State type	Description	Production	Emission	Transition
P	(pair emitting)	$P \rightarrow aYb$	$e_v(a, b)$	$t_v(Y)$
L	(left emitting)	$L \rightarrow aY$	$e_v(a)$	$t_v(Y)$
R	(right emitting)	$R \rightarrow Ya$	$e_v(a)$	$t_v(Y)$
B	(bifurcation)	$B \rightarrow SS$	1	1
D	(delete)	$D \rightarrow Y$	1	$t_v(Y)$
S	(start)	$S \rightarrow Y$	1	$t_v(Y)$
E	(end)	$E \rightarrow \epsilon$	1	1

Each overall production probability is the independent product of an emission probability e_v and a transition probability t_v , both of which are position-dependent parameters that depend on the state v (analogous

²Much of this text is taken from (Eddy, 2002).

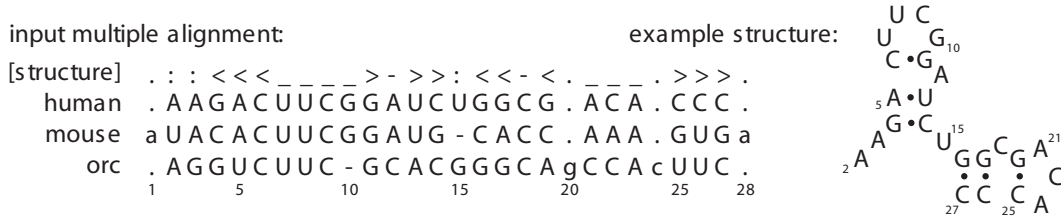


Figure 1: **An example RNA sequence family.** Left: a toy multiple alignment of three sequences, with 28 total columns, 24 of which will be modeled as consensus positions. The [structure] line annotates the consensus secondary structure in WUSS notation. Right: the secondary structure of the “human” sequence.

to hidden Markov models). For example, a particular pair (P) state v produces two correlated letters a and b (e.g. one of 16 possible base pairs) with probability $e_v(a, b)$ and transits to one of several possible new states Y of various types with probability $t_v(Y)$. A bifurcation (B) state splits into two new start (S) states with probability 1. The E state is a special case ϵ production that terminates a derivation.

A CM consists of many states of these seven basic types, each with its own emission and transition probability distributions, and its own set of states that it can transition to. Consensus base pairs will be modeled by P states, consensus single stranded residues by L and R states, insertions relative to the consensus by more L and R states, deletions relative to consensus by D states, and the branching topology of the RNA secondary structure by B, S, and E states. The procedure for starting from an input multiple alignment and determining how many states, what types of states, and how they are interconnected by transition probabilities is described next.

From consensus structural alignment to guide tree

Figure 1 shows an example input file: a multiple sequence alignment of homologous RNAs, with a line in WUSS notation that describes the consensus RNA secondary structure. The first step of building a CM is to produce a binary *guide tree* of *nodes* representing the consensus secondary structure. The guide tree is a parse tree for the consensus structure, with nodes as nonterminals and alignment columns as terminals.

The guide tree has eight types of nodes:

Node	Description	Main state type
MATP	(pair)	P
MATL	(single strand, left)	L
MATR	(single strand, right)	R
BIF	(bifurcation)	B
ROOT	(root)	S
BEGL	(begin, left)	S
BEGR	(begin, right)	S
END	(end)	E

These consensus node types correspond closely with the CM’s final state types. Each node will eventually contain one or more states. The guide tree deals with the consensus structure. For individual sequences, we will need to deal with insertions and deletions with respect to this consensus. The guide tree is the skeleton on which we will organize the CM. For example, a MATP node will contain a P-type state to model a

$i..k$ and $k + 1..j$, there will be more than one possible choice of k if $i..j$ is a multifurcation loop containing three or more stems. The choice of k impacts the performance of the divide and conquer algorithm; for optimal time performance, we will want bifurcations to split into roughly equal sized alignment problems, so I choose the k that makes $i..k$ and $k + 1..j$ as close to the same length as possible.

The result of this procedure is the guide tree. The nodes of the guide tree are numbered in preorder traversal (e.g. a recursion of “number the current node, visit its left child, visit its right child”: thus parent nodes always have lower indices than their children). The guide tree corresponding to the input multiple alignment in Figure 1 is shown in Figure 2.

From guide tree to covariance model

A CM must deal with insertions and deletions in individual sequences relative to the consensus structure. For example, for a consensus base pair, either partner may be deleted leaving a single unpaired residue, or the pair may be entirely deleted; additionally, there may be inserted nonconsensus residues between this pair and the next pair in the stem. Accordingly, each node in the master tree is expanded into one or more *states* in the CM as follows:

Node	States	total # states	# of split states	# of insert states
MATP	[MP ML MR D] IL IR	6	4	2
MATL	[ML D] IL	3	2	1
MATR	[MR D] IR	3	2	1
BIF	[B]	1	1	0
ROOT	[S] IL IR	3	1	2
B EGL	[S]	1	1	0
B EG R	[S] IL	2	1	1
END	[E]	1	1	0

Here we distinguish between consensus (“M”, for “match”) states and insert (“I”) states. ML and IL, for example, are both L type states with L type productions, but they will have slightly different properties, as described below.

The states are grouped into a *split set* of 1-4 states (shown in brackets above) and an *insert set* of 0-2 insert states. The split set includes the main consensus state, which by convention is first. One and only one of the states in the split set must be visited in every parse tree (and this fact will be exploited by the divide and conquer algorithm). The insert state(s) are not obligately visited, and they have self-transitions, so they will be visited zero or more times in any given parse tree.

State transitions are then assigned as follows. For bifurcation nodes, the B state makes obligate transitions to the S states of the child BEGL and BEGR nodes. For other nodes, each state in a split set has a possible transition to every insert state in the *same* node, and to every state in the split set of the *next* node. An IL state makes a transition to itself, to the IR state in the same node (if present), and to every state in the split set of the next node. An IR state makes a transition to itself and to every state in the split set of the next node.

This arrangement of transitions guarantees that (given the guide tree) there is unambiguously one and only one parse tree for any given individual structure. This is important. The algorithm will find a maximum likelihood parse tree for a given sequence, and we wish to interpret this result as a maximum likelihood

structure, so there must be a one to one relationship between parse trees and secondary structures (Giegerich, 2000).

The final CM is an array of M states, connected as a directed graph by transitions $t_v(y)$ (or probability 1 transitions $v \rightarrow (y, z)$ for bifurcations) with the states numbered such that $(y, z) \geq v$. There are no cycles in the directed graph other than cycles of length one (e.g. the self-transitions of the insert states). We can think of the CM as an array of states in which all transition dependencies run in one direction; we can do an iterative dynamic programming calculation through the model states starting with the last numbered end state M and ending in the root state 1. An example CM, corresponding to the input alignment of Figure 1, is shown in Figure 3.

As a convenient side effect of the construction procedure, it is guaranteed that the transitions from any state are to a *contiguous* set of child states, so the transitions for state v may be kept as an offset and a count. For example, in Figure 3, state 12 (an MP) connects to states 16, 17, 18, 19, 20, and 21. We can store this as an offset of 4 to the first connected state, and a total count of 6 connected states. We know that the offset is the distance to the next non-split state in the current node; we also know that the count is equal to the number of insert states in the current node, plus the number of split set states in the next node. These properties make establishing the connectivity of the CM trivial. Similarly, all the parents of any given state are also contiguously numbered, and can be determined analogously. We are also guaranteed that the states in a split set are numbered contiguously. This contiguity is exploited by the divide and conquer implementation.

Parameterization

Using the guide tree and the final CM, each individual sequence in the input multiple alignment can be converted unambiguously to a CM parse tree, as shown in Figure 4. Weighted counts for observed state transitions and singlet/pair emissions are then collected from these parse trees. These counts are converted to transition and emission probabilities, as maximum *a posteriori* estimates using mixture Dirichlet priors.

Comparison to profile HMMs

The relationship between an SCFG and a covariance model is analogous to the relationship of hidden Markov models (HMMs) and profile HMMs for modeling multiple sequence alignments (Krogh et al., 1994; Durbin et al., 1998; Eddy, 1998). A comparison may be instructive to readers familiar with profile HMMs. A profile HMM is a repetitive HMM architecture that associates each consensus column of a multiple alignment with a single type of model node – a MATL node, in the above notation. Each node contains a “match”, “delete”, and “insert” HMM state – ML, IL, and D states, in the above notation. The profile HMM also has special begin and end states. Profile HMMs could therefore be thought of as a special case of CMs. An unstructured RNA multiple alignment would be modeled by a guide tree of all MATL nodes, and converted to an unbifurcated CM that would essentially be identical to a profile HMM. (The only difference is trivial; the CM root node includes a IR state, whereas the start node of a profile HMM does not.) All the other node types (especially MATP, MATR, and BIF) and state types (e.g. MP, MR, IR, and B) are SCFG augmentations necessary to extend profile HMMs to deal with RNA secondary structure.

The cmbuild program, step by step

The `cmbuild` command line syntax is:

```
> cmbuild <options> [cmfile] [alifile]
```

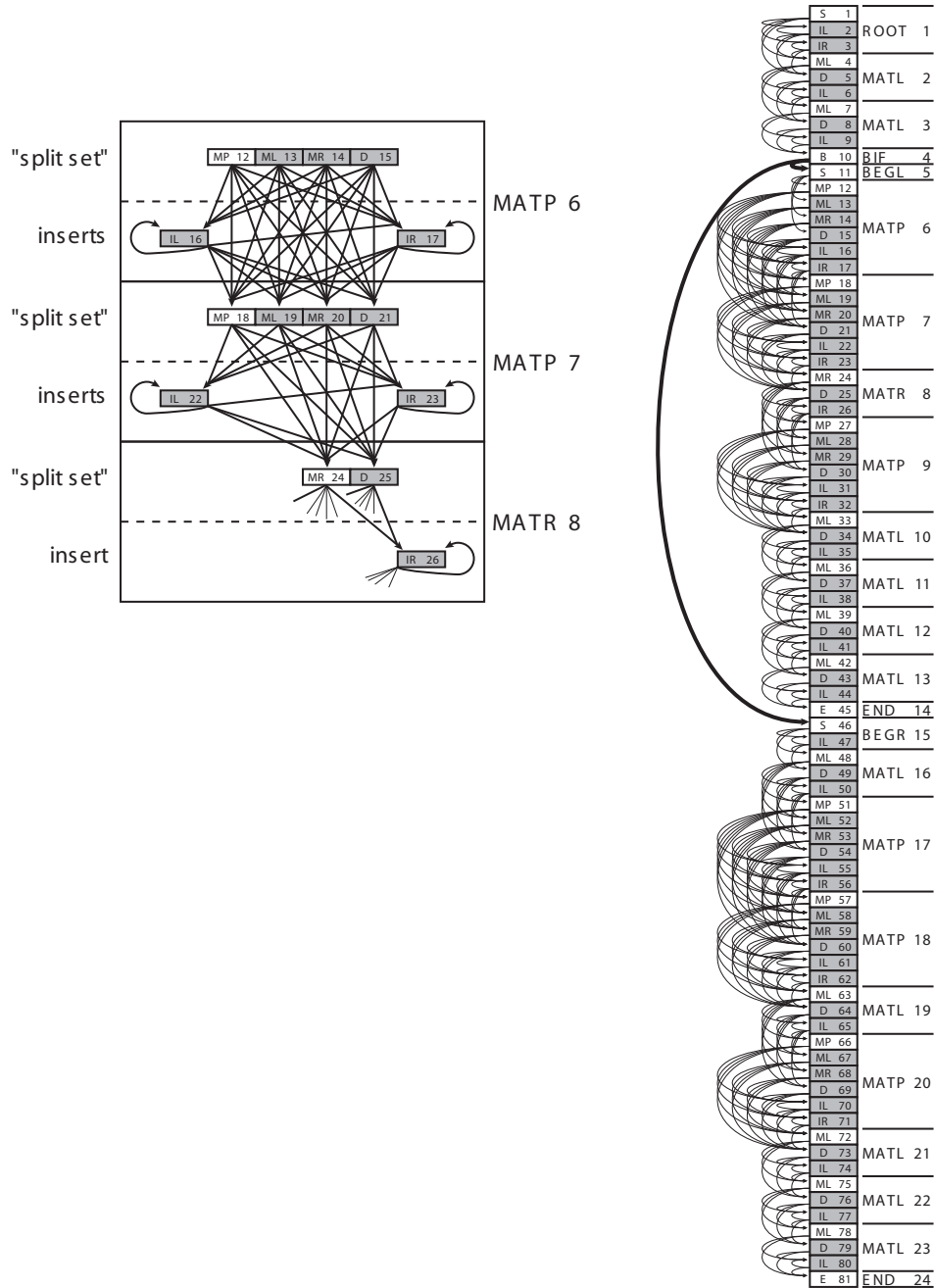


Figure 3: **A complete covariance model.** Right: the CM corresponding to the alignment in Figure 1. The model has 81 states (boxes, stacked in a vertical array). Each state is associated with one of the 24 nodes of the guide tree (text to the right of the state array). States corresponding to the consensus are in white. States responsible for insertions and deletions are gray. The transitions from bifurcation state B10 to start states S11 and S46 are in bold because they are special: they are an obligate (probability 1) bifurcation. All other transitions (thin arrows) are associated with transition probabilities. Emission probability distributions are not represented in the figure. Left: the states are also arranged according to the guide tree. A blow up of part of the model corresponding to nodes 6, 7, and 8 shows more clearly the logic of the connectivity of transition probabilities (see main text), and also shows why any parse tree must transit through one and only one state in each “split set”.

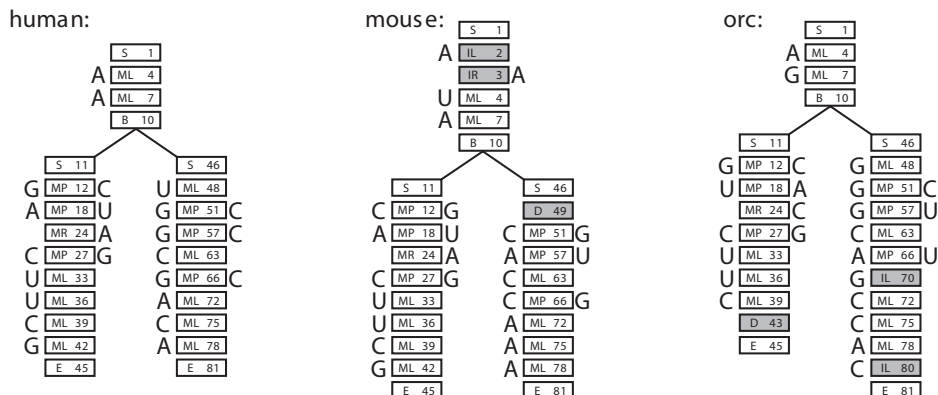


Figure 4: **Example parse trees.** Parse trees are shown for the three sequences/structures from Figure 1, given the CM in Figure 3. For each sequence, each residue must be associated with a state in the parse tree. (The sequences can be read off its parse tree by starting at the upper left and reading counterclockwise around the edge of parse tree.) Each parse tree corresponds directly to a secondary structure – base pairs are pairs of residues aligned to MP states. A collection of parse trees also corresponds to a multiple alignment, by aligning residues that are associated with the same state – for example, all three trees have a residue aligned to state ML4, so these three residues would be aligned together. Insertions and deletions relative to the consensus use nonconsensus states, shown in gray.

where **[alifile]** is the name of the input alignment file, and **[cmfile]** is the name of the output CM file. What follows describes the steps that **cmbuild** goes through, and the most important options that can be chosen to affect its behavior.

Alignment input file

The input alignment file must be in Stockholm format, and it must have a consensus secondary structure annotation line (**#=GC SS_cons**).

The program is actually capable of reading many common multiple alignment formats (ClustalW, PHYLIP, GCG MSF, and others) but no other format currently supports consensus RNA secondary structure annotation. This may change in the future, either when other formats allow structure annotation, or when **cmbuild** is capable of inferring consensus structure from the alignment by automated comparative analysis, as the earlier COVE suite was capable of (Eddy and Durbin, 1994).

If the file does not exist, is not readable, or is not in a recognized format, the program exits with a “could not be opened for reading” error. If the file does not have consensus secondary structure annotation, the program exits with a “no consensus structure annotation” error. This includes all non-Stockholm alignment files.

▷ **Why does *cmbuild* have a *--informat* option, if it only accepts Stockholm?** If you don’t specify *--informat*, the software has to autodetect the file format. Autodetection of file formats doesn’t work in certain advanced/nonstandard cases, for instance if you’re reading the alignment from standard input instead of from a file. The *--informat* allows you to override autodetection; e.g. **cat my.sto | cmbuild --informat Stockholm my.cm** - is an example of reading the alignment from piped standard input.

Parsing secondary structure annotation

The structure annotation line only needs to indicate which columns are base paired to which. It does not have to be in full WUSS notation. Even if it is, the details of the notation are largely ignored. Nested pairs of `<>`, `()`, `[]`, or `{ }` symbols are interpreted as base paired columns. All other columns marked with the symbols `:`, `_`, `-`, `~` are interpreted as single stranded columns.

A simple minimal annotation is therefore to use `<>` symbols to mark base pairs and `.` for single stranded columns.

If a secondary structure annotation line is in WUSS notation and it contains valid pseudoknot annotation (e.g. additional non-nested stems marked with AAA,aaa or BBB,bbb, etc.), this annotation is removed and a warning is printed. INFERNAL cannot handle pseudoknots. Internally, these columns are treated as if they were marked with `.` symbols.

▷ **How should I choose to annotate pseudoknots?** INFERNAL can only deal with nested base pairs. If there is a pseudoknot, you have to make a choice of which stem to annotate as normal nested structure (thus including it in the model) and which stem to call additional “pseudoknotted” structure (thus ignoring it in the model). For example, for a simple two-stem pseudoknot, should you annotate it as AAAA.<<<<aaaa...>>>>, or <<<<.AAAA>>>>...aaaa? From an RNA structure viewpoint, which stem I label as the pseudoknotted one is an arbitrary choice; but since one of the stems in the pseudoknot will have to be modeled as a single stranded region by INFERNAL, the choice makes a slight difference in the performance of your model. You want your model to capture as much information content as possible. Thus, since the information content of the model is a sum of the sequence conservation plus the additional information contributed by pairwise correlations in base-paired positions, you should tend to annotate the shorter stem as the “pseudoknot” (modeling as many base pairs as possible), and you should also annotate the stem with the more conserved primary sequence as the “pseudoknot” (if one stem is more conserved at the sequence level, you won’t lose as much by modeling that one as primary sequence consensus only).

If (aside from any ignored pseudoknot annotation) the structure annotation line contains characters other than `<>()[]{}: _ - . ~` then those characters are ignored (treated as `.`) and a warning is printed.

If, after this “data cleaning”, the structure annotation is inconsistent with a secondary structure (for example, if the number of `<` and `>` characters isn’t the same), then the program exits with a “failed to parse consensus structure annotation” error.

Sequence weighting

By default, the input sequences are weighted in two ways to compensate for biased sampling (phylogenetic correlations). Relative sequence weights are calculated by the Gerstein/Chothia/Sonnhammer method (Gerstein et al., 1994). (The `--wgsc` option forces GSC weights, but is redundant since that’s the default.) To turn relative weighting off (e.g. set all weights to 1.0), use the `--wnone` option.

Some alignment file formats allow relative sequence weights to be given in the file. This includes Stockholm format, which has `#=GS WT` weight annotations. Normally `cmbuild` ignores any such input weights. The `--wgiven` option tells `cmbuild` to use them. This lets you set the weights with any external procedure you like; for example, the `weight` utility program in SQUID implements some common weighting algorithms, including the fast $O(N)$ Henikoff position-based weights (Henikoff and Henikoff, 1994).

If for some reason you put more than one relative weighting option on the command line, the last one you give is used.

▷ **Why is cmbuild taking so much time?** The GSC weighting algorithm scales as $O(N^2)$ with the number of sequences N . Weighting may become rate-limiting for `cmbuild` if your alignment contains

many sequences. Model construction itself is fast. You might want to turn weighting off, or pre-calculate the weights by a faster algorithm.

Absolute weights (the “effective sequence number”) is calculate by “entropy weighting” (Karplus et al., 1998). This sets the balance between the prior and the data, and affects the information content of the model. Entropy weighting reduces the effective sequence number (the total sum of the weights) and increases the entropy (degrading the information content) of the model until a threshold is reached. The default entropy is 1.46 bits per position (roughing 0.54 bits of information, relative to uniform base composition). This threshold can be changed with the `--etarget <x>` option. Entropy weighting may be turned off entirely with the `--effnone` option.

Architecture construction

The CM architecture is now constructed from your input alignment and your secondary structure annotation, as described in the previous section.

The program needs to determine which columns are consensus (match) columns, and which are insert columns. (Remember that although WUSS notation allows insertions to be annotated in the secondary structure line, `cmbuild` is only paying attention to annotated base pairs.) By default, it does this by a simple rule based on the frequency of gaps in a column. If the frequency of gaps is greater than a threshold, the column is considered to be an insertion.

The threshold defaults to 0.5. It can be changed to another number `<x>` (from 0 to 1.0) by the `--gapthresh <x>` option. The higher the number, the more columns are included in the model. At `--gapthresh 1.0`, all the columns are considered to be part of the consensus. At `--gapthresh 0.0`, only columns with no gaps are.

You can also manually specify which columns are consensus versus insert by including reference coordinate annotation (e.g. a `#=GC RF` line, in Stockholm format) and using the `--rf` option. Any columns marked by non-gap symbols become consensus columns. (The simplest thing to do is mark consensus columns with `x`'s, and insert columns with `.`'s. Remember that spaces aren't allowed in alignments in Stockholm format.) If you set the `--rf` option but your file doesn't have reference coordinate annotation, the program exits with an error.

Parameterization

Weighted observed emission and transition counts are then collected from the alignment data. These count vectors c are then converted to estimated probabilities p using mixture Dirichlet priors. The default mixture priors are described in (Nawrocki and Eddy, 2007). You can provide your own prior as a file, using the `--priorfile <f>` option.

Naming the model

Each CM gets a name. Stockholm format allows the alignment to have a name, provided in the `#=GF ID` tag. If this name is provided, it is used as the CM name.

If a name is not provided, the name is the input filename, without any extension – for example, if you build a model from the alignment file `RNaseP.sto`, the model will be named `RNaseP`.

You can override this and provide your own name with the `-n <s>` option, where `<s>` is any string.

Stockholm format also allows more than one alignment per file, and **cmbuild** supports this: CM files can contain more than one model, and if you say e.g. **cmbuild Rfam Rfam.sto** where **Rfam.sto** contains a whole database of alignments, **cmbuild** will create a database of CMs in the **Rfam** file, one per alignment. But in this case, obviously you don't want them all to have the same name! Therefore when running **cmbuild** on a multi-multiple alignment database file, the alignment database file *must* provide **#=GF ID** tags with names for each alignment. If any alignment is found to not have one, the program exits at that point with an error. Attempting to set the **--n** option for an alignment database also results in an error.

Saving the model

The model is now saved to a file, according to the filename specified on the command line. By default, a new file is created, and the model is saved in a portable ASCII text format.

If the cmfile already exists, the program exits with an error. The **--F** option causes the new model to overwrite an existing cmfile. The **--A** option causes the new model to be appended to an existing cmfile (creating a growing CM database, perhaps).

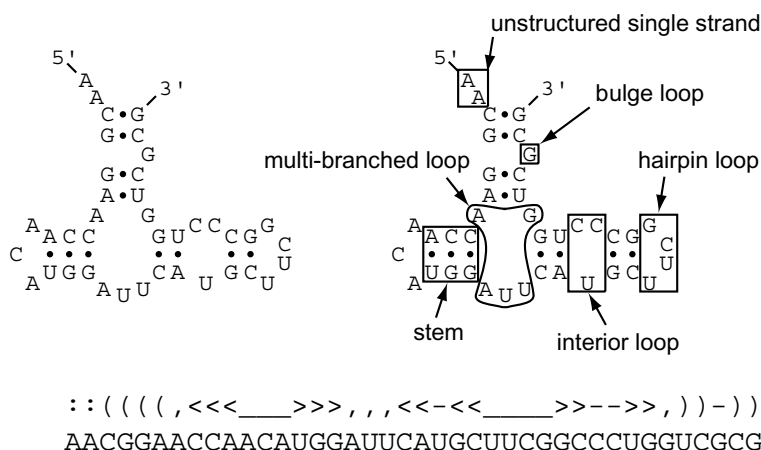
5 File and output formats

RNA secondary structures: WUSS notation

INFERNAL annotates RNA secondary structures using a linear string representation called “WUSS notation” (Washington University Secondary Structure notation).

The symbology is extended from the common bracket notation for RNA secondary structures, where open- and close-bracket symbols (or parentheses) are used to annotate base pairing partners: for example, $((((\dots)))$ indicates a four-base stem with a three-base loop. Bracket notation is difficult for humans to interpret, for anything much larger than a simple stem-loop. WUSS notation makes it somewhat easier to interpret the annotation for larger structures.

The following figure shows an example with the key elements of WUSS notation. At the top left is an example RNA structure. At the top right is the same structure, with different RNA structural elements marked. Below both structure pictures: the WUSS notation string for the structure.



Full (output) WUSS notation

In detail, symbols used by WUSS notation in *output* structure annotation strings are as follows:

Base pairs Base pairs are annotated by nested matching pairs of symbols $<>$, $()$, $[]$, or $\{\}$. The different symbols indicate the “depth” of the helix in the RNA structure as follows: $<>$ are used for simple terminal stems; $()$ are used for “internal” helices enclosing a multifurcation of all terminal stems; $[]$ are used for internal helices enclosing a multifurcation that includes at least one annotated $()$ stem already; and $\{\}$ are used for all internal helices enclosing deeper multifurcations.

Hairpin loops Hairpin loop residues are indicated by underscores, $_$. Simple stem loops stand out as, e.g. $<<<<____>>>>$.

Bulge, interior loops Bulge and interior loop residues are indicated by dashes, $-$.

Multifurcation loops Multifurcation loop residues are indicated by commas, $,$. The mnemonic is “stem 1, stem 2”, e.g. $<<<____>>>, , <<<____>>>$.

External residues Unstructured single stranded residues completely outside the structure (unenclosed by any base pairs) are annotated by colons, `:`.

Insertions Insertions relative to a known structure are indicated by periods, `.`. Regions where local structural alignment was invoked, leaving regions of both target and query sequence unaligned, are indicated by tildes, `~`. These symbols only appear in alignments of a known (query) structure annotation to a target sequence of unknown structure.

Pseudoknots WUSS notation allows pseudoknots to be annotated as pairs of upper case/lower case letters: for example, `<<<<_AAAA_____>>>>aaaa` annotates a simple pseudoknot; additional pseudoknotted stems could be annotated by `Bb`, `Cc`, etc. INFERNAL cannot handle pseudoknots, however; pseudoknot notation never appears in INFERNAL output; it is accepted in input files, but ignored.

An example of WUSS notation for a complicated structure (*E. coli* RNase P) is shown in Figure 5. An example of WUSS notation for a local INFERNAL alignment of *B. subtilis* RNase P to *E. coli* RNase P, illustrating the use of local alignment annotation symbols, is in Figure 6.

Shorthand (input) WUSS notation

While WUSS notation makes it easier to visually interpret INFERNAL *output* structural annotation, it would be painful to be required to *input* all structures in full WUSS notation. Therefore when INFERNAL reads input secondary structure annotation, it uses simpler rules:

Base pairs Any matching nested pair of `()`, `[]`, `{ }` symbols indicates a base pair; the exact choice of symbol has no meaning, so long as the left and right partners match up.

Single stranded residues All other symbols `_`, `-`, `:`, `.`, `~` indicate single stranded residues. The choice of symbol has no special meaning. Annotated pseudoknots (nested matched pairs of upper/lower case alphabetic characters) are also interpreted as single stranded residue in INFERNAL input.

Thus, for instance, `<<<< >>>>` and `((((_____)))` and `< ((. _ . _) >) >` all indicate a four base stem with a four base loop (the last example is legal but weird).

Remember that the key property of canonical (nonpseudoknotted) RNA secondary structure is that the pairs are *nested*. `((<) >>` is not a legal annotation string: the pair symbols don't match up properly. INFERNAL will reject such an annotation and report an input format error, suspecting a problem with your annotation. If you want to annotate pseudoknots, WUSS notation allows alphabetic symbols `Aa`, `Bb`, etc. see above; but remember that INFERNAL ignores pseudoknotted stems and treats them as single stranded residues.

Because many other RNA secondary structure analysis programs use a simple bracket notation for annotating structure, INFERNAL's ability to input this format makes it easier to use data generated by other RNA software packages. Conversely, converting INFERNAL output WUSS notation to simple bracket notation is a matter of a simple Perl or sed script, substituting the symbols appropriately.

Multiple alignments: Stockholm format

The Pfam consortium developed an annotated alignment format called “Stockholm format”, and this format has been adopted as the standard alignment format in HMMER and INFERNAL, and by the Rfam consortium. The reasons for inventing a new alignment format were two-fold. First, there really is no standard accepted format for multiple sequence alignment files, so we don’t feel guilty about inventing a new one. Second, the formats of popular multiple alignment software (e.g. CLUSTAL, GCG MSF, PHYLIP) do not support rich documentation and markup of the alignment. Stockholm format was developed to support extensible markup of multiple sequence alignments, and we use this capability extensively in both RNA work (with structural markup) and the Pfam database (with extensive use of both annotation and markup).

A minimal Stockholm file

```
# STOCKHOLM 1.0

seq1  ACDEF...GHIKL
seq2  ACDEF...GHIKL
seq3  ...EFMNRGHIKL

seq1  MNPQTVWY
seq2  MNPQTVWY
seq3  MNPQT...
```

The simplest Stockholm file is pretty intuitive, easily generated in a text editor. It is usually easy to convert alignment formats into a “least common denominator” Stockholm format. For instance, SELEX, GCG’s MSF format, and the output of the CLUSTAL multiple alignment programs are all similar interleaved formats.

The first line in the file must be `# STOCKHOLM 1.x`, where `x` is a minor version number for the format specification (and which currently has no effect on my parsers, other than identifying the file as Stockholm format). This line allows a parser to instantly identify the file format.

In the alignment, each line contains a name, followed by the aligned sequence. A dash or period denotes a gap. If the alignment is too long to fit on one line, the alignment may be split into multiple blocks, with blocks separated by blank lines. The number of sequences, their order, and their names must be the same in every block. Within a given block, each (sub)sequence (and any associated `#=GR` and `#=GC` markup, see below) is of equal length, called the *block length*. Block lengths may differ from block to block; the block length must be at least one residue, and there is no maximum.

The sequence names must be unique. (They are used to associate markup tags with the sequences.)

Other blank lines are ignored. You can add comments to the file on lines starting with a `#`.

All other annotation is added using a tag/value comment style. The tag/value format is inherently extensible, and readily made backwards-compatible; unrecognized tags will simply be ignored. Extra annotation includes consensus and individual RNA or protein secondary structure, sequence weights, a reference coordinate system for the columns, and database source information including name, accession number, and coordinates (for subsequences extracted from a longer source sequence) See below for details.

Syntax of Stockholm markup

There are four types of Stockholm markup annotation, for per-file, per-sequence, per-column, and per-residue annotation:

- #=GF** **<tag>** **<s>** Per-file annotation. **<s>** is a free format text line of annotation type **<tag>**. For example, **#=GF DATE April 1, 2000**. Can occur anywhere in the file, but usually all the **#=GF** markups occur in a header.
- #=GS** **<seqname>** **<tag>** **<s>** Per-sequence annotation. **<s>** is a free format text line of annotation type **tag** associated with the sequence named **<seqname>**. For example, **#=GS seq1 SPECIES.SOURCE Caenorhabditis elegans**. Can occur anywhere in the file, but in single-block formats (e.g. the Pfam distribution) will typically follow on the line after the sequence itself, and in multi-block formats (e.g. HMMER output), will typically occur in the header preceding the alignment but following the **#=GF** annotation.
- #=GC** **<tag>** **<s>** Per-column annotation. **<s>** is an aligned text line of annotation type **<tag>**. **#=GC** lines are associated with a sequence alignment block; **<s>** is aligned to the residues in the alignment block, and has the same length as the rest of the block. Typically **#=GC** lines are placed at the end of each block.
- #=GR** **<seqname>** **<tag>** **<s>** Per-residue annotation. **<s>** is an aligned text line of annotation type **<tag>**, associated with the sequence named **<seqname>**. **#=GR** lines are associated with one sequence in a sequence alignment block; **<s>** is aligned to the residues in that sequence, and has the same length as the rest of the block. Typically **#=GR** lines are placed immediately following the aligned sequence they annotate.

Semantics of Stockholm markup

Any Stockholm parser will accept syntactically correct files, but is not obligated to do anything with the markup lines. It is up to the application whether it will attempt to interpret the meaning (the semantics) of the markup in a useful way. At the two extremes are the Belvu alignment viewer and the HMMER profile hidden Markov model software package.

Belvu simply reads Stockholm markup and displays it, without trying to interpret it at all. The tag types (**#=GF**, etc.) are sufficient to tell Belvu how to display the markup: whether it is attached to the whole file, sequences, columns, or residues.

HMMER and INFERNAL use Stockholm markup to pick up a variety of information from the multiple alignment files. The Pfam and Rfam consortiums therefore agree on additional syntax for certain tag types, so software can parse some markups for useful (or necessary) information. This additional syntax is imposed by Pfam, HMMER, INFERNAL, and other software of mine, not by Stockholm format per se. You can think of Stockholm as akin to XML, and what my software reads as akin to an XML DTD, if you're into that sort of structured data format lingo.

The Stockholm markup tags that are parsed semantically by my software are as follows:

Recognized **#=GF** annotations

ID **<s>** Identifier. **<s>** is a name for the alignment; e.g. "RNaseP. Mandatory, if the file is an alignment database used as input for **cmbuild**, because each CM must get a unique name. One word. Unique in file.

AC **<s>** Accession. **<s>** is a unique accession number for the alignment; e.g. “PF00001”. Used by the Rfam database, for instance. Often a alphabetical prefix indicating the database (e.g. “RF”) followed by a unique numerical accession. One word. Unique in file.

DE **<s>** Description. **<s>** is a free format line giving a description of the alignment; e.g. “Ribonuclease P RNA”. One line. Unique in file.

AU **<s>** Author. **<s>** is a free format line listing the authors responsible for an alignment; e.g. “Bateman A”. One line. Unique in file.

Recognized #=GS annotations

WT **<f>** Sequence weight. **<f>** is a positive real number giving the relative weight for a sequence, usually used to compensate for biased representation by downweighting similar sequences. Usually the weights average 1.0 (e.g. the weights sum to the number of sequences in the alignment) but this is not required. Either every sequence must have a weight annotated, or none of them can.

AC **<s>** Accession. **<s>** is a database accession number for this sequence. (Compare the **#=GF AC** markup, which gives an accession for the whole alignment.) One word.

DE **<s>** Description. **<s>** is one line giving a description for this sequence. (Compare the **#=GF DE** markup, which gives a description for the whole alignment.)

Recognized #=GC annotations

RF Reference line. Any character is accepted as a markup for a column. The intent is to allow labeling the columns with some sort of mark. **cmbuild** uses this annotation to determine which columns are consensus versus insertion; insertion columns are annotated by a gap symbol, and consensus columns by any non-gap symbol.

ss_cons Secondary structure consensus. When this line is generated by INFERNAL, it is generated in full WUSS notation. When it is read by **cmbuild**, it is interpreted more loosely, in shorthand (input) WUSS notation: pairs of symbols <>, (), [], or [] mark consensus base pairs, and symbols : _ - , . ~ mark single stranded columns.

Recognized #=GR annotations

ss Secondary structure for this sequence. See **#=GC ss_cons** above.

Sequence files: FASTA format

FASTA is probably the simplest of formats for unaligned sequences. FASTA files are easily created in a text editor. Each sequence is preceded by a line starting with >. The first word on this line is the name of the sequence. The rest of the line is a description of the sequence (free format). The remaining lines contain the sequence itself. You can put as many letters on a sequence line as you want. For example:

```
>seq1 This is the description of my first sequence.
AGTACGTAGTAGCTGCTGCTACGTGCGCTAGCTAGTACGTCA CGACGTAGATGCTAGCTGACTCGATGC
>seq2 This is a description of my second sequence.
CGATCGATCGTACGTCGACTGATCGTAGCTACGTCGTACGTAG CATCGTCAGTTACTGCATGCTCG
CATCAGGCATGCTGCTGACTGATCGTACG
```

For better or worse, FASTA is not a documented standard. Minor (and major) variants are in widespread use in the bioinformatics community, all of which are called “FASTA format”. My software attempts to cater to all of them, and is tolerant of common deviations in FASTA format. Certainly anything that is accepted by the database formatting programs in NCBI BLAST or WU-BLAST (e.g. setdb, pressdb, xdformat) will also be accepted by my software. Blank lines in a FASTA file are ignored, and so are spaces or other gap symbols (dashes, underscores, periods) in a sequence. Other non-amino or non-nucleic acid symbols in the sequence are also silently ignored, mostly because some people seem to think that “*” or “.” should be added to protein sequences to (redundantly) indicate the end of the sequence. The parser will also accept unlimited line lengths, which allows it to accomodate the enormous description lines in the NCBI NR databases.

(On the other hand, any FASTA files *generated* by my software adhere closely to community standards, and should be usable by other software packages (BLAST, FASTA, etc.) that are more picky about parsing their input files. That means you can run a sloppy FASTA file thru the **sreformat** utility program to clean it up.)

Partly because of this tolerance, the software may have a difficult time dealing with files that are *not* in FASTA format, especially if you’re relying on file format autodetection (the “Babelfish”). Some (now mercifully uncommon) file formats are so similar to FASTA format that they be erroneously called FASTA by the Babelfish and then quietly and lethally misparsed. An example is the old NBRF file format. If you’re afraid of this, you can use the **--informat fasta** option to bypass the Babelfish and improve robustness. However, it is still possible to construct files perversely similar to FASTA that will still confuse the parser. (The gist of these caveats applies to all formats, not just FASTA.)

CM file format

The default CM file format is a simple, extensible tag-value format. The format being used right now is tentative and likely to change. Therefore, it is not currently documented here. If you absolutely need to interpret it, see the file `cmio.c` in the source code.

Null model file format

The Infernal source distribution includes an example prior file, **rna.null**. This null model is identical to the hardcoded default prior used by Infernal, all four RNA nucleotides are equiprobable in the null, background model.

A null model file must contain exactly four non-comment lines. A comment line begins with a “#”, that is a # followed by a single space. Each of the four non-comment lines must contain a single floating point number, the four of which sum to 1.0. The first non-comment line is interpreted as the background probability of an “A” residue, the second, third, and fourth non-comment lines are interpreted as the background probabilities of a “C”, “G” and “U” respectively.

Dirichlet prior files

A prior file is parsed into a number of whitespace-delimited, non-comment fields. These fields are then interpreted in order. The order and number of the fields is important. This is not a robust, tag-value save file format.

All whitespace is ignored, including newlines. The number of fields per line is unimportant.

Comments begin with a # character. The remainder of any line following a # is ignored.

The Infernal source distribution includes an example prior file, **default.pri**. This prior is identical to the hardcoded default prior used by Infernal. The following text may only make sense if you’re looking at that example while you read.

The order of the fields in the prior file is as follows:

Strategy. The first field is the keyword **Dirichlet**. Currently Dirichlet priors (mixture or not) are the only prior strategy used by Infernal.

Transition prior section. The next field is the number **74**, the number of different types of transition distributions. (See Figure 7 for an explanation of where the number 74 comes from.) Then, for each of these 74 distributions:

<from-uniqstate> <to-node>: Two fields give the transition type: from a unique state identifier, to a node identifier. Example: **MATP_MP MATP**.

<n>: One field gives the number of transition probabilities for this transition type; that is, the number of Dirichlet parameter vector $\alpha_1^q \dots \alpha_n^q$ for each mixture component q .

<nq>: One field gives the number of mixture Dirichlet components for this transition type’s prior. Then, for each of these **nq** Dirichlet components:

p(q): One field gives the mixture coefficient $p(q)$, the prior probability of this component q . For a single-component “mixture”, this is always 1.0.

$\alpha_1^q \dots \alpha_n^q$: The next n fields give the Dirichlet parameter vector for this mixture component q .

Base pair emission prior section. This next section is the prior for MATP_MP emissions. One field gives **<K>**, the “alphabet size” – the number of base pair emission probabilities – which is always 16 (4x4), for RNA. The next field gives **<nq>**, the number of mixture components. Then, for each of these **nq** Dirichlet components:

p(q): One field gives the mixture coefficient $p(q)$, the prior probability of this component q . For a single-component “mixture”, this is always 1.0.

$\alpha_{AA}^q \dots \alpha_{UU}^q$: The next 16 fields give the Dirichlet parameter vector for this mixture component, in alphabetical order (AA, AC, AG, AU, CA ... GU, UA, UC, UG, UU).

Consensus singlet base emission prior section. This next section is the prior for MATL_ML and MATR_MR emissions. One field gives **<K>**, the “alphabet size” – the number of singlet emission probabilities – which is always 4, for RNA. The next field gives **<nq>**, the number of mixture components. Then, for each of these **nq** Dirichlet components:

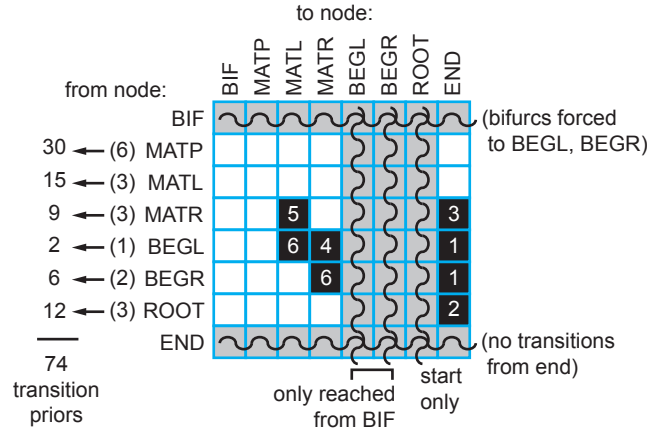
p(q): One field gives the mixture coefficient $p(q)$, the prior probability of this component q . For a single-component “mixture”, this is always 1.0.

$\alpha_A^q \dots \alpha_U^q$: The next 4 fields give the Dirichlet parameter vector for this mixture component, in alphabetical order (A, C, G, U).

Nonconsensus singlet base emission prior section. This next section is the prior for insertions (MATP_IL, MATP_IR, MATL_IL, MATR_IR, ROOT_IL, ROOT_IR, BEGR_IL) as well as nonconsensus singlets (MATP_ML, MATP_MR). One field gives **<K>**, the “alphabet size” – the number of singlet emission probabilities – which is always 4, for RNA. The next field gives **<nq>**, the number of mixture components. Then, for each of these **nq** Dirichlet components:

$\mathbf{p}(\mathbf{q})$: One field gives the mixture coefficient $p(q)$, the prior probability of this component q . For a single-component “mixture”, this is always 1.0.

$\alpha_{\mathbf{A}}^q \dots \alpha_{\mathbf{U}}^q$: The next 4 fields give the Dirichlet parameter vector for this mixture component, in alphabetical order (A, C, G, U).



STL9/63

Figure 7: Where does the magic number of 74 transition distribution types come from? The transition distributions are indexed in a 2D array, from a unique statetype (20 possible) to a downstream node (8 possible), so the total conceivable number of different distributions is $20 \times 8 = 160$. The grid represents these possibilities by showing the 8×8 array of all node types to all node types; each starting node contains 1 or more unique states (number in parentheses to the left). Two rows are impossible (gray): bifurcations automatically transit to determined BEGL, BEGR states with probability 1, and end nodes have no transitions. Three columns are impossible (gray): BEGL and BEGR can only be reached by probability 1 transitions from a bifurcation, and the ROOT node is special and can only start a model. Eight individual cells of the grid are unused (black) because of the way **cmbuild** (almost) unambiguously constructs a guide tree from a consensus structure. These cases are numbered as follows. (1) BEGL and BEGR never transit to END; this would imply an empty substructure. A bifurcation is only used if both sides of the split contain at least one consensus pair (MATP). (2) ROOT never transits to END; this would imply an alignment with zero consensus columns. Infernal models assume ≥ 1 consensus columns. (3) MATR never transits to END. Infernal always uses MATL for unpaired columns whenever possible. MATR is only used for internal loops, multifurcation loops, and 3' bulges, so MATR must always be followed by a BIF, MATP, or another MATR. (4) BEGL never transits to MATR. The single stranded region between two bifurcated stems is unambiguously assigned to MATL nodes on the right side of the split, not to MATR nodes on the left. (5) MATR never transits to MATL. The only place where this could arise (given that we already specified that MATL is used whenever possible) is in an interior loop; there, by unambiguous convention, MATL nodes precede MATR nodes. (6) BEGL nodes never transit to MATL, and BEGR nodes never transit to MATR. By convention, at any bifurcated subsequence i, j , i and j are paired but not to each other. That is, the smallest possible subsequence is bifurcated, so that any single stranded stretches to the left and right are assigned to MATL and MATR nodes above the bifurcation, instead of MATL nodes below the BEGL and MATR nodes below the BEGR. Thus, the total number 74 comes from multiplying, for each row, the number of unique states in each starting node by the number of possible downstream nodes (white), and summing these up, as shown to the left of the grid.

6 Manual pages

cmalign - use a CM to make a structured RNA multiple alignment

Synopsis

cmalign [*options*] *cmfile seqfile*

Description

cmalign aligns the RNA sequences in *seqfile* to the covariance model (CM) in *cmfile*, and outputs a multiple sequence alignment.

Currently, the sequence file must be in FASTA format.

CM files are profiles of RNA consensus secondary structure. A CM file is produced by the **cmbuild** program, from a given RNA sequence alignment of known consensus structure.

The alignment that **cmalign** makes is written in Stockholm format. It can be redirected to a file using the *-o* option. By default, the CYK algorithm is used for CM alignment. An HMM banded CYK alignment algorithm that offers a significant speed-up is enabled with the *--hbanded* option (see below).

Options

- h** Print brief help; includes version number and summary of all options, including expert options.
- l** Turn on the local alignment algorithm, which allows the alignment to span two or more subsequences if necessary (e.g. if the structures of the query model and target sequence are only partially shared), allowing certain large insertions and deletions in the structure to be penalized differently than normal indels. The default is to globally align the query model to the target sequences.
- o** *<f>* Save the alignment in Stockholm format to a file *<f>*. The default is to write it to standard output.
- q** Quiet; suppress the verbose banner, and only print the resulting alignment to stdout. This allows piping the alignment to the input of other programs, for example.

Expert Options

- informat** *<s>* Assert that the input *seqfile* is in format *<s>*. Do not run Babelfish format autodetection. This increases the reliability of the program somewhat, because the Babelfish can make mistakes; particularly recommended for unattended, high-throughput runs of Infernal. *<s>* is case-insensitive. This option is a bit forward-looking; **cmalign** currently only accepts FASTA format, but this will not be true in the future.

- nosmall** Use the normal CYK alignment algorithm. The default is to use the divide and conquer algorithm described in SR Eddy, BMC Bioinformatics 3:18, 2002. This is useful for debugging, and checking that the two algorithms give identical results. The "normal" algorithm requires too much memory for most uses.
- regress** <*f*> Save regression test information to a file <*f*>. This is part of the automated testing procedure at each release.
- full** Include all match columns in the output alignment. By default only consensus columns in which at least one sequence has a non-gap character are included.
- tfile** <*f*> Dump tabular sequence tracebacks for each individual sequence to a file <*f*>. Primarily useful for debugging.
- banddump** <*n*> Set verbosity level for debugging print statements related to query-dependent bands (QDB) (the **--qdb** option) to <*n*>. Where <*n*> is 1, 2 or 3. By default debugging print statements are turned off.
- dlev** <*n*> Set verbosity level for general debugging print statements <*n*>. Where <*n*> is 1, 2 or 3. By default debugging print statements are turned off.
- time** Print timings for band calculation and alignment for each sequence in *seqfile*.
- inside** Instead of aligning the sequences in *seqfile* to the CM with CYK, use the Inside algorithm to score each sequence against the CM, and print out scores. Each Inside score is the summed score of the all possible alignments of the sequence to the CM. When enabled, the output does not include an alignment.
- outside** Instead of aligning the sequences in *seqfile* to the CM with CYK, use the Outside algorithm to score each sequence against the CM, and print out scores. Each Outside score is the summed score of the all possible alignments of the sequence to the CM. When enabled, the output does not include an alignment.
- post** Calculate "confidence estimates" for each residue of the alignment and include them in the alignment. These estimates are based on posterior probabilities that each residue aligns at each position of the alignment. The posterior probabilities can be determined only after running the Inside and Outside algorithms. The alignment is not changed when **--post** is used, it is still the CYK optimal alignment. The confidence estimates appear as "#=GR <seq name> POST" annotation in the output Stockholm alignment for each sequence. This annotation includes the digits "0-9", "*" and "." characters. A "8" indicates that the confidence estimate for the corresponding residue being aligned to the corresponding column is between 80 and 90%. A "7" indicates the same, but between 70 and 80%, and so on for the other digits. A "*" indicates that the confidence estimate is "very nearly" 100% (it's hard to be exact here due to numerical precision issues) A "." indicates that that column aligns to a gap in the corresponding sequence.
- checkpost** Perform a check that the result of the Inside and Outside algorithms make sense by ensuring they give the same score to each sequence (while allowing for slight differences due to numerical precision issues). This is primarily useful for debugging.

- zeroinserts** Set all insert emission scores to 0.0 bits. Default behavior is to read them from the CM file. Zeroing inserts is default behavior for **cmsearch**.
- sub** Turn on the sub model construction and alignment procedure. For each sequence, an HMM is first used to predict the model start and end consensus columns, and a new sub CM is constructed that only models consensus columns from start to end. The sequence is then aligned to this sub CM. This option is useful for aligning sequences that are known to truncated, non-full length sequences. This "sub CM" procedure is not the same as the "sub CMs" described by Weinberg and Ruzzo.
- elsilent** Disallow emissions from the local end (EL) state.
- enfstart** $\langle n \rangle$ Used in combination with **--enfseq** $\langle s \rangle$ to attempt to enforce the subsequence $\langle s \rangle$ align beginning at consensus column $\langle n \rangle$. A major limitation to this option is that the consensus columns from ($\langle n \rangle - 1$) to ($\langle n \rangle + \text{length}(\langle s \rangle) - 1$) must all be modelled by MATL nodes. This may change in future versions.
- enfseq** $\langle s \rangle$ Used in combination with **--enfstart** $\langle n \rangle$ to enforce the subsequence $\langle s \rangle$ occur beginning at consensus column $\langle n \rangle$ as described above.
- hbanded** Score each sequence to an HMM derived from the CM in *cmfile* using the Forward and Backward HMM algorithms and calculate posterior probabilities each residue aligns to each state of the HMM. Use these posterior probabilities to derive constraints (bands) for the CM, and use them to accelerate CM alignment. This method sacrifices the guarantee that the optimal alignment will be found. When not run in local mode with **-l**, this option yields a significant speed-up (around 25-fold for the average family), while rarely missing the optimal alignment. It has not been rigorously tested in local mode yet, so use at your own risk with **-l**.
- tau** $\langle f \rangle$ Set the tail loss probability used during HMM band calculation to $\langle f \rangle$. This is the amount of probability mass within the HMM posterior probabilities that is considered negligible. The default value is 1E-7. In general, higher values will result in greater acceleration, but increase the chance of missing the optimal alignment due to the HMM bands.
- hsafe** Must be used in combination with **--hbanded**. In HMM banded mode, realign any sequences with a negative alignment score using non-banded CYK to guarantee finding the optimal alignment. Based on empirical tests, the fraction of HMM banded alignments that are non-optimal and have negative scores is much higher than for those with positive scores.
- hmmonly** Align each sequence to an HMM derived from the CM in *cmfile* using the Viterbi HMM alignment algorithm. Currently, no alignment is printed when this option is enabled, only scores are printed.
- qdb** Accelerate alignment using query-dependent bands (QDB) as described in (Nawrocki and Eddy, PLoS Computational Biology 3(3): e56, 2007) to constrain the CYK alignment. In practice, acceleration due to QDB seems to be significantly less than

acceleration due to the **--hbanded** option, and the chance of missing the optimal alignment is slightly more with QDB.

- beta** *<x>* Set the tail loss probability of QDB to *<x>*. The default value is 1E-7. In general, higher values will result in greater acceleration, but increase the chance of missing the optimal alignment due to the bands.
- withali** *<f>* Reads an alignment from file *<f>* and aligns it as a single object to the CM; e.g. the alignment in *<f>* is held fixed. This allows you to align sequences to a model with **cmalign** and view them in the context of an existing trusted multiple alignment. The alignment in the file *<f>* must be exactly the alignment that the CM was built from, or a subset of it with the following special property: the definition of consensus columns and consensus secondary structure must be identical between *<f>* and the alignment the CM was built from. One easy way to achieve this is to use the **--rf** option to **cmbuild** (see man page for **cmbuild**) and to maintain the **"#=GC RF"** annotation in the alignment when removing sequences to create the subset alignment *<f>*. To specify that the **--rf** option to **cmbuild** was used enable the **--rf** option to **cmalign** (described below).
- rf** Must be used in combination with **--withali** *<f>*. Specify that the alignment in *<f>* has the same **"#=GC RF"** annotation as the alignment file the CM was built from using **cmbuild** and further that the **--rf** option was supplied to **cmbuild** when the CM was constructed.
- gapthresh** *<x>* Must be used in combination with **--withali** *<f>*. Specify that the **--gapthresh** *<x>* option was supplied to **cmbuild** when the CM was constructed from the alignment file *<f>*.

cmbuild - construct a CM from an RNA multiple sequence alignment

Synopsis

cmbuild [*options*] *cmfile* *alifile*

Description

cmbuild reads an RNA multiple sequence alignment from *alifile*, constructs a covariance model (CM), and saves the CM to *cmfile*.

The alignment file must be in Stockholm format, and must contain consensus secondary structure annotation. **cmbuild** uses the consensus structure to determine the architecture of the CM.

The alignment file may be a database containing more than one alignment. If it does, the resulting *cmfile* will be a database of CMs, one per alignment.

The experimental expert options **--ctarget**, **--cmindiff**, and **--call** result in multiple CMs being built from a single alignment as described below.

Options

- h** Print brief help; includes version number and summary of all options, including expert options.
- n** <*s*> Name the covariance model <*s*>. (Does not work if *alifile* contains more than one alignment or in combination with **--call**, **--ctarget**, or **--cmindiff** as described below.) The default is to use the name of the alignment (given by the #=GF ID tag, in Stockholm format), or if that is not present, to use the name of the alignment file minus any file type extension plus a positive integer indicating the position of that alignment in the file (that is, the first alignment in a file "myrnas.sto" would give a CM named "myrnas.1", the second alignment would give a CM named "myrnas.2").
- A** Append the CM to *cmfile*, if *cmfile* already exists.
- F** Allow *cmfile* to be overwritten. Normally, if *cmfile* already exists, **cmbuild** exits with an error unless the **-A** or **-F** option is set.

Expert Options

- binary** Save the model in a compact binary format. The default is a more readable ASCII text format.
- rf** Use reference coordinate annotation (#=GC RF line, in Stockholm) to determine which columns are consensus, and which are inserts. Any non-gap character indicates a consensus column. (For example, mark consensus columns with "x", and

insert columns with ”.”.) The default is to determine this automatically; if the frequency of gap characters in a column is greater than a threshold, `gapthresh` (default 0.5), the column is called an insertion.

- gapthresh** *<x>* Set the gap threshold (used for determining which columns are insertions versus consensus; see above) to *<x>*. The default is 0.5.
- informat** *<s>* Assert that the input *alifile* is in format *<s>*. Do not run Babelfish format autodetection. This increases the reliability of the program somewhat, because the Babelfish can make mistakes; particularly recommended for unattended, high-throughput runs of Infernal. *<s>* is case-insensitive. This option is a bit forward-looking; **cmbuild** currently only accepts Stockholm format, but this may not be true in the future.
- beta** *<x>* Set the probability tail loss beta parameter for query-dependent banding (QDB) to *<x>*. QDB is used by default to set the maximum hit length of the model, *W*, that is stored in *cmfile*. Lower values of beta will yield higher values for *W*. The default beta used is 1E-7. For details on QDB see (Nawrocki and Eddy, PLoS Computational Biology 3(3): e56, 2007). Other options for setting *W* are **--window** and **--rsw**.
- window** *<n>* Set the maximum hit length *W* parameter to *<n>* where *<n>* is a positive integer. By default, *W* is calculated using QDB.
- rsw** *<n>* Set the maximum hit length *W* parameter as twice the average sequence length in *alifile*. This is how RSEARCH sets *W*, hence the "rs". By default, *W* is calculated using QDB.
- nodetach** Do not detach one of two insert states that model insertions at exactly the same position in the model. By default, one of these states is detached, making it unreachable in the model. This default behavior was adopted to address an ambiguity in the CM grammar due to a design flaw. Used primarily for debugging as earlier versions of **cmbuild** did not detach these insert states.
- wgiven** Use sequence weights as given in annotation in the input alignment file. If no weights were given, assume they are all 1.0. The default is to determine new sequence weights by the Gerstein/Sonnhammer/Chothia algorithm, ignoring any annotated weights.
- wnone** Turn sequence weighting off; e.g. explicitly set all sequence weights to 1.0.
- wgsc** Use the Gerstein/Sonnhammer/Chothia weighting algorithm. This is the default, so this option is probably useless.
- effent** Use the entropy weighting strategy to determine the effective sequence number that gives a target mean match state entropy. This option is the default, and can be turned off with **--effnone**. The default target mean match state entropy is 1.46 bits but can be changed with **--etarget**.

- etarget** *<x>* Set the target mean match state entropy as *<x>*. By default the target entropy 1.46 bits. This default value was empirically determined as optimal in a RFAM based benchmark described in (Nawrocki and Eddy, PLoS Computational Biology 3(3): e56, 2007).
- effnone** Turn off the entropy weighting strategy. The effective sequence number is just the number of sequences in *alifile*.
- cfile** *<f>* Save a file containing observed count vectors (both emissions and transitions) to a counts file *<f>*. One use for this file is as the starting point for estimating Dirichlet priors from observed RNA structure data.
- cmtbl** *<f>* Save a tabular description of the CM's topology to a file *<f>*. Primarily useful for debugging CM architecture construction.
- emap** *<f>* Save a consensus emission map to a file *<f>*. This file relates the numbering system of states in the CM's tree-like directed graph to the linear numbering of consensus columns. Primarily useful for debugging.
- gtree** *<f>* Save an ASCII picture of the high level structure of the CM's guide tree to a file *<f>*. Primarily useful for debugging.
- gtbl** *<f>* Save a tabular description of the nodes in CM's guide tree to a file *<f>*. Primarily useful for debugging.
- tfile** *<f>* Dump tabular inferred sequence tracebacks for each individual training sequence to a file *<f>*. Primarily useful for debugging.
- bfile** *<f>* Dump the query-dependent bands to a file *<f>*. This file can be read by **cmsearch**. Primarily useful for debugging.
- bdfile** *<f>* Dump the band distributions to a file *<f>*. Primarily useful for debugging.
- nobalance** Turn off the architecture "rebalancing" algorithm. The nodes in a CM are initially numbered in standard preorder traversal. The rebalancing algorithm is an optimizer that reorders the numbering of the CM in order to absolutely guarantee certain algorithmic performance bounds. However, it is a stylistic riff that has almost no real empirical impact on performance, and is a tricky algorithm to get right. This option was inserted for debugging purposes. It is sometimes also useful to obtain a simple preorder traversal numbering system in the CM architecture (for illustrative purposes, for example).
- regress** *<f>* Save regression test information to a file *<f>*. This is part of the automated testing procedure at each release.
- treeforce** After building the model, score the first sequence in the alignment using its inferred parsetree, and show both the score and the parsetree. This is a debugging tool, used to specify and score a particular desired parsetree.

- ignorant** Strip all base pair secondary structure information from *alifile* before building the model. The resulting model will be all single stranded MATL nodes, with 0 bifurcations.
- null** *<f>* Read a null model from *<f>*. The default is to use 0.25 for each RNA nucleotide. For more information on the format of the null model file, see the User's Guide.
- prior** *<f>* Read a Dirichlet prior from *<f>*, replacing the default mixture Dirichlet. The format of prior files is documented in the User's Guide.
- ctarget** *<n>* Cluster the sequence alignment in *alifile* by percent identity. Find a cutoff percent id threshold that gives exactly *<n>* clusters and build a separate CM from each cluster. If *<n>* is greater than the number of sequences in *alifile*, the program will not complain, and each sequence in *alifile* will be its own cluster. Each CM will have a positive integer appended to its name indicating the order in which it was built. For example, if **cmbuild --ctarget 3** is called with *alifile* "myrnas.sto", and "myrnas.sto" has no #=GF ID tag annotation and exactly one Stockholm alignment in it, three CMs will be built, the first will be named "myrnas.1.1", the second, "myrnas.1.2", and the third "myrnas.1.3". (As explained above, the first number "1" after "myrnas" indicates the CM was built from the first alignment in "myrnas.sto".)
- cmindiff** *<x>* Cluster the sequence alignment in *alifile* by percent identity. Define clusters at the cutoff fractional id difference of *<x>* and build a separate CM from each cluster. The CMs are named as described above for **--ctarget**.
- call** Build a separate CM from each sequence in *alifile*. Naming of CMs takes place as described above for **--ctarget**.
- corig** After building multiple CMs using **--ctarget**, **--cmindiff** or **--call** as described above, build a final CM using the complete original alignment in *alifile*. The CMs are named as described above for **--ctarget** with the exception of the final CM built from the original alignment which is named in the default manner, without an appended integer.

cmemit - generate sequences from a covariance model

Synopsis

cmemit [*options*] *cmfile seqfile*

Description

cmemit reads the first covariance model (CM) in *cmfile* and generates a number of sequences from the CM; or if the **-c** option is selected, generates a single majority-rule consensus. This can be useful for various application in which one needs a simulation of sequences consistent with a sequence family consensus. By default, **cmemit** generates 10 sequences and outputs them in FASTA (unaligned) format.

Options

- a** Write the generated sequences in an aligned format (STOCKHOLM) with consensus structure annotation rather than FASTA.
- c** Predict a single majority-rule consensus sequence instead of sampling sequences from the CMs probability distribution. Highly conserved residues are shown in upper case; others are shown in lower case.
- h** Print brief help; includes version number and summary of all options, including expert options.
- n** *<n>* Generate *<n>* sequences. Default is 10.
- o** *<f>* Save the synthetic sequences to file *<f>* rather than writing them to stdout.
- q** Quiet; suppress all output except for the sequences themselves. Useful for piping or directing the output.

Expert Options

- seed** *<n>* Set the random seed to *<n>*, where *<n>* is a positive integer. The default is to use `time()` to generate a different seed for each run, which means that two different runs of **cmemit** on the same CM will give slightly different results. You can use this option to generate reproducible results.
- begin** *<n>* Truncate the resulting alignment by beginning at consensus column *<n>*, where *<n>* is a positive integer no greater than the consensus length of the CM. Must be used in combination with **--end** and either **-a** or **--cp9**.
- end** *<n>* Truncate the resulting alignment by ending at consensus column *<n>*, where *<n>* is a positive integer no greater than the consensus length of the CM. Must be used in combination with **--begin** and either **-a** or **--cp9**.

--cp9 Use the generated sequences to train a CM Plan 9 maximum-likelihood HMM without using pseudo-counts. Print the counts and normalized probabilities of the HMM. Primarily useful for debugging.

cmscore - align and score one or more sequences to a CM

Synopsis

cmscore [*options*] *cmfile seqfile*

Description

cmscore uses the covariance model (CM) in *cmfile* to align and score the sequences in *seqfile*, and output summary statistics on timings and scores. **cmscore** is a testbed for new CM alignment algorithms, and it is also used by the testsuite. It is not intended to be particularly useful in the real world. Documentation is provided for completeness, and to aid our own memories.

cmscore aligns the sequence(s) using two alignment algorithms, and compares the scores and timings of each algorithm. By default the two algorithms compared are the full CYK algorithm and the divide and conquer CYK variant (SR Eddy, BMC Bioinformatics 3:18, 2002). The expert options allow different algorithms to be compared as explained below. Among the algorithms that can be tested are two banded variants of CYK, query-dependent banded (QDB) CYK (Nawrocki and Eddy, PLoS Computational Biology 3(3): e56, 2007) with the **--qdb** option and HMM banded CYK with the **--hbanded** option. While non-banded CYK variants are guaranteed to find the optimal alignment and score of each sequence, the banded variants sacrifice this guarantee for acceleration. The level of acceleration can be controlled by the **--beta** *<x>* and **--tau** *<x>* options for QDB and HMM banding respectively. In short, *<x>* is a rough estimate at the probability that the optimal alignment will be missed. The greater *<x>* is, the greater the acceleration, but the greater the chance of missing the optimal alignment. By default *<x>* is set as 1E-7 for both **--qdb** and **--hbanded**. **cmscore** is useful for testing for values of beta and tau that give the best trade-off of acceleration versus accuracy. To make this testing easier, multiple beta and tau values can be tested within a single **cmscore** call. The **--betas** *<x>* and **--betae** *<x>* combination and the **--taus** *<x>* and **--taue** *<x>* option combination allow the user to specify a beginning beta/tau value and an ending beta/tau value. For example, **--betas** 3 and **--betae** 5 would first align the sequences in *seqfile* with non-banded CYK, and then perform 3 additional QDB alignments, first with beta=1E-3, next with beta=1E-4 and finally with beta=1E-5. The tau options work in the same way. Currently, only values of 1E-*<x>* can be used. Summary statistics on timings and how often the optimal alignment are missed for each value of beta or tau are then printed to standard output.

Usually when comparing non-banded algorithms, the the two parse trees should be identical for any sequence, because the optimal alignment score is guaranteed. However, there can be cases of ties, where two or more different parse trees have identical scores. In such cases, it is possible for the two parse trees to differ. The parse tree selected as "optimal" from amongst the ties is arbitrary, dependent on order of evaluation in the DP traceback, and the order of evaluation for D&C vs. standard CYK is different. Thus, in its testsuite role, **cmscore** checks that the scores are within 0.01 bits of each other, but does not check that the parse trees are absolutely identical; identity can be checked for using the **--stringent** option.

Currently, the sequence file must be in FASTA format.

The sequences are treated as single stranded RNAs; that is, only the given strand of each sequence is aligned and scored, and no reverse complementing is done.

CM files are profiles of RNA consensus secondary structure. A CM file is produced by the **cmbuild** program, from a given RNA sequence alignment of known consensus structure.

Options

- h** Print brief help; includes version number and summary of all options, including expert options.
- i** Print individual timings and score comparisons for each sequence in *seqfile*. By default only summary statistics are printed.

Expert Options

- local** Turn on the local alignment algorithm, which allows the alignment to span two or more subsequences if necessary (e.g. if the structures of the query model and target sequence are only partially shared), allowing certain large insertions and deletions in the structure to be penalized differently than normal indels. The default is to globally align the query model to the target sequences.
- sub** Turn on the sub model construction and alignment procedure. For each sequence, an HMM is first used to predict the model start and end consensus columns, and a new sub CM is constructed that only models consensus columns from start to end. The sequence is then aligned to this sub CM. This "sub CM" procedure is not the same as the "sub CMs" described by Weinberg and Ruzzo.
- regress** *<f>* Save regression test information to a file *<f>*. This is part of the automated testing procedure at each release.
- stringent** Require the two parse trees to be identical; fail and return a non-zero exit code if they are not. Normally, **cmscore** only requires that the two parse trees have identical scores (within a floating point tolerance of 0.01 bits), because it is possible to have more than one parse tree with the same score.
- trees** Print the parse trees for each alignment of each sequence to standard output.
- std** Specify the first alignment algorithm as non-banded Divide and Conquer (D&C) CYK, and the second algorithm as standard CYK. This is default.
- qdb** Specify the first alignment algorithm as non-banded D&C CYK, and the second algorithm as QDB D&C CYK.
- qdbsmall** Specify the first alignment algorithm as non-banded D&C CYK, and the second algorithm as QDB standard CYK.
- qdbboth** Specify the first alignment algorithm as QDB D&C CYK and the second algorithm as QDB standard CYK.

- beta** <*x*> Set the probability tail loss beta parameter for **--qdb** to <*x*> In general higher values of beta give greater acceleration but greater chance of missing the optimal alignment. The default beta with **--qdb** is 1E-7.
- hbanded** Specify the first alignment algorithm as non-banded D&C (D&C) CYK, and the second algorithm as HMM banded standard CYK.
- beta** <*x*> Set the probability tail loss tau parameter for **--hbanded** to <*x*> In general higher values of tau give greater acceleration but greater chance of missing the optimal alignment. The default tau with **--hbanded** is 1E-7.
- hsafe** Must be used in combination with **--hbanded**. In HMM banded mode, realign any sequences with a negative alignment score using non-banded D&C CYK.
- hmmonly** Specify the first alignment algorithm as non-banded D&C (D&C) CYK, and the second algorithm as Viterbi to a CM Plan 9 HMM derived from the CM in *cmfile*. When enabled, only timing information is relevant, scoring is not.
- scoreonly** During the standard CYK algorithm stage, use the "score only" variant of the algorithm to save memory, and don't recover a parse tree. Cannot be used with **--qdb**, **--qdbsmall**, **--qdbboth**, **--hbanded**, or **--hmmonly**.
- betas** <*x*> Specify the first alignment algorithm as non-banded D&C CYK, and multiple stages of QDB CYK alignment. The first QDB alignment will use beta=1E-<*x*>, which will be the highest value of beta used. Must be used in combination with **--qdb** and **--betae**.
- betae** <*x*> Specify the first alignment algorithm as non-banded D&C CYK, and multiple stages of QDB CYK alignment. The final QDB alignment will use beta=1E-<*x*>, which will be the lowest value of beta used. Must be used in combination with **--qdb** and **--betas**.
- taus** <*x*> Specify the first alignment algorithm as non-banded D&C CYK, and multiple stages of HMM banded CYK alignment. The first HMM banded alignment will use tau=1E-<*x*>, which will be the highest value of tau used. Must be used in combination with **--qdb** and **--taue**.
- taue** <*x*> Specify the first alignment algorithm as non-banded D&C CYK, and multiple stages of HMM banded CYK alignment. The final HMM banded alignment will use tau=1E-<*x*>, which will be the lowest value of tau used. Must be used in combination with **--qdb** and **--taus**.

cmsearch - search a sequence database for RNAs homologous to a CM

Synopsis

cmsearch [*options*] *cmfile seqfile*

Description

cmsearch uses the covariance model (CM) in *cmfile* to search for homologous RNAs in *seqfile*, and outputs high-scoring alignments.

Currently, the sequence file must be in FASTA format.

CM files are profiles of RNA consensus secondary structure. A CM file is produced by the **cmbuild** program, from a given RNA sequence alignment of known consensus structure.

cmsearch output consists of alignments of all hits of score greater than zero bits sorted by decreasing score per sequence and per strand. That is, all hits for the same sequence and the same (Watson or Crick) strand are sorted, but hits across sequences or strands are not sorted. The threshold of zero bits can be changed to any positive number $\langle x \rangle$ using the **-T $\langle x \rangle$** option as described below. If the **-E $\langle x \rangle$** option is enabled E-values are calculated and all hits with an E-value less than or equal to $\langle x \rangle$ will be reported, ranked by increasing E-value. E-values are not turned on by default because they have not yet been rigorously tested to our satisfaction.

RNA homology search with CMs is slow. To speed it up, query-dependent banding (QDB) is turned on by default. QDB can be turned off with the **--noqdb** option. Briefly, QDB precalculates regions of the dynamic programming matrix that have negligible probability based on the query CM's transition probabilities. During search, these regions of the matrix are ignored to make searches faster. For more information on QDB see (Nawrocki and Eddy, PLoS Computational Biology 3(3): e56). QDB sacrifices the guarantee that the optimal alignment for any subsequence will be found, so the acceleration potentially comes at a cost to sensitivity. The beta parameter is the amount of probability mass considered negligible during band calculation, lower values of beta yield greater speedups but also a greater chance of missing the optimal alignment. The default beta is 1E-7: determined empirically as a good tradeoff between sensitivity and speed, though this value can be changed with the **--beta $\langle x \rangle$** option. The speedups for an average RNA family with the default beta of 1E-7 is about four-fold relative to a non-banded search, but in general the greater the sequence length of a family, the greater the speedup achieved by QDB.

Another option for accelerating **cmsearch** is HMM filtering. The idea is to first search the database with an HMM using HMM search algorithms that are much faster than CM search algorithms. High-scoring hits to the HMM are then searched again using the expensive CM methods. This concept was first introduced by Zasha Weinberg and Larry Ruzzo at the University of Washington. One technique Weinberg and Ruzzo put forth is to build an HMM that is a "maximum likelihood HMM" (ML HMM). This HMM is derived from the model parameters of the CM and is as similar as possible to the CM, while not modelling the interactions between base-paired consensus columns (which an HMM's regular grammar cannot do). We have implemented a version of ML HMMs within Infernal, which we call CM Plan 9 HMMs (CP9 HMMs) to distinguish them from the Plan 7 HMMs of the HMMER package. A CP9 HMM can be used to filter the database prior to search with a CM using the **--hmmfilter** option as described below. By default, when

--hmmfilter is enabled, E-values are calculated for CP9 HMM hits and those hits with E-values less than or equal to 500 survive the filter. This E-value threshold can be changed to *<x>* with the **--hmmE** *<x>* option. Also, a bit score threshold can be used by enabling the option **--hmmT** *<x>*. When this option is enabled CP9 HMM hits with a bit score greater than *<x>* survive the filter. HMM filtering can decrease sensitivity of CM searches because high-scoring CM hits may not survive the HMM filter. However, in our tests, CP9 HMM filtering often yields a very significant speed-up of about 30-40 fold or more while sacrificing a small level of sensitivity. QDB remains on when HMM filtering is turned on.

Another HMM filtering technique pioneered by Weinberg and Ruzzo is the construction and use of a "rigorous filter" HMM. All hits above a certain CM bit score threshold are guaranteed to survive the HMM filtering step. Their implementation of rigorous filters has been included with Infernal but is not available as an option to **cmsearch**. For more information see the User's Guide.

Options

- h** Print brief help; includes version number and summary of all options, including expert options.
- T** *<x>* Set the bit score cutoff for the per-sequence ranked hit list to *<x>*, where *<x>* is a positive real number. The default is 0 bits. Hits with bit scores better than (greater than) this threshold will be shown.
- E** *<x>* Set the E-value cutoff for the per-sequence/strand ranked hit list to *<x>*, where *<x>* is a positive real number. Hits with E-values better than (less than) or equal to this threshold will be shown. E-values are calculated by scoring the CM against a large number (1000 by default, but this number can be changed with the **--nsamples** option) of random synthesized sequences, and fitting an extreme value distribution (EVD) to the histogram of those scores. The sequences are randomly synthesized by first selecting a GC content from the distribution of GC contents of 100 nucleotide windows in the *seqfile* and generating a 2*W length sequence based on that GC content, where W is the window length of the model as calculated in **cm-build** and stored in the *cmfile*. This procedure for determining E-values is based on the procedure described in (Klein and Eddy, BMC Bioinformatics 4:44, 2003).

Expert Options

- informat** *<s>* Assert that the input *alifile* is in format *<s>*. Do not run Babelfish format autodetection. This increases the reliability of the program somewhat, because the Babelfish can make mistakes; particularly recommended for unattended, high-throughput runs of Infernal. *<s>* is case-insensitive. This option is a bit forward-looking; **cmsearch** currently only accepts FASTA format, but this may not be true in the future.
- toponly** Only search the top (Watson) strand of the sequences in *seqfile*. By default, both strands are searched.

- local** Turn on the local alignment algorithm, which allows the alignment to span two or more subsequences if necessary (e.g. if the structures of the query model and target sequence are only partially shared), allowing certain large insertions and deletions in the structure to be penalized differently than normal indels. The default is to globally align the query model to the target sequences. When enabled in combination with **--hmmfilter** or **--hmmonly**, the CP9 HMM local alignment algorithm is turned on, which allows alignments to start and stop at any consensus HMM node. Currently, as implemented HMM local alignment does not have a good way of modelling the CM local end behavior.
- noalign** Do not calculate and print alignments of each hit, only print locations and scores.
- window <n>** Set the scanning window width to <n>. This is the maximum length of a homologous sequence. By default, this is set within the CM file by **cmbuild** using the query-dependent band calculation (Nawrocki and Eddy, PLoS Computational Biology 3(3): e56). See the man entry for **cmbuild** for more information. This option must be enabled in combination with the **--noqdb** option, which turns off query-dependent banding (QDB). The reason is that QDB sets the window length based on an expected length distribution calculated from the transition probabilities of the CM.
- dumptrees** Dump verbose, ugly parse trees for each hit. Useful only for debugging purposes.
- nsamples <n>** When **-E** and/or **--hmmE** is/are enabled, set the number of synthetic sequences to <n>, where <n> is a positive integer. If <n> is less than about 1000, the fit to the EVD may fail. Higher numbers of <n> will give better determined EVD parameters but will take more time. The default with **-E** and/or **--hmmE** enabled is 1000: it was empirically chosen as a tradeoff between accuracy and computation time. Note that currently the same number of samples is always used to calculate both CM and CP9 HMM E-values if both the **-E** and **--hmmE** options are enabled.
- partition <n>[,<n>]** Used to calculate E-values for hits based on the GC content of the hit. Must be used in combination with the **-E** and/or **--hmmE** options. The <n> values specify partition points for percentage GC content. A separate EVD is fit to a histogram of scores from random sequences for each partition, and the relevant EVD is used to determine the E-value of each hit based on its GC content. For example, if the option **--partition 40,60** is enabled, three separate EVDs will be calculated as described for **-E** with the following exceptions. For the first EVD, GC contents for the random samples will be selected from the distribution of 100 nucleotide windows with GC content between 0 and 39 percent in *seqfile*. For the second EVD, the same is true but now from windows between 40 and 59 percent. Similarly for the third EVD, but now from windows between 60 and 100 percent. The E-values of three hits to subsequences with 35, 55, and 80 percent GC content would be calculated from the 1st, 2nd and 3rd EVDs respectively. For more detail on the partitioning procedure see (Klein and Eddy, BMC Bioinformatics 4:44, 2003).

- inside** Use a scanning Inside algorithm, instead of the default scanning CYK algorithm. With the Inside algorithm the score of a subsequence is the summed score of the all possible alignments of that subsequence to the CM, as opposed to the CYK score which is the score of the single most likely alignment to the CM. This option increases sensitivity at but slows down searches about two-fold.
- null2** Turn on the post hoc second null model that attempts to deal with the potential biased composition of hits. This option has not been rigorously tested, use at your own risk.
- learninserts** Turn off the default **cmsearch** behavior of setting all CM insert emission scores to 0 bits prior to search. This default behavior is to avoid high-scoring hits to low complexity sequence favored by high insert state emission scores. When **--learninserts** is enabled the insert emission scores are read from *cmfile* that were calculated in **cmbuild** and are unmodified prior to search.
- negsc <x>** Set the minimum CM bit score to report as <x> where <x> is any real negative number. This option is untested and very experimental, use at your own risk. It's probably only potentially useful in non-local mode, when no positive scoring real hits can be found.
- enfstart <n>** Used in combination with **--enfseq <s>** to attempt to enforce the subsequence <s> align beginning at consensus column <n>. A major limitation to this option is that the consensus columns from (<n> -1) to (<n> + length(<s>)-1) must all be modelled by MATL nodes. This may change in future versions.
- enfseq <s>** Used in combination with **--enfstart <n>** to enforce the subsequence <s> occur beginning at consensus column <n> as described above for the **--enfstart** option. By default, when this option is enabled an HMM that enforces only the subsequence <s> is used to filter the database prior to searching with the CM. This HMM filtering behavior can be turned off with the **--enfnohmm** option.
- enfnohmm** Used in combination with **--enfstart <n>** and **--enfseq <s>**. Do not filter the database first with an HMM that enforces the subsequence <s> prior to searching with the CM.
- time** Print timings for the search and histogram construction (if E-values are being used).
- rtrans** Replace the transition scores read from *cmfile* with RSEARCH transition scores as described in (Klein and Eddy, BMC Bioinformatics 4:44, 2003). Not recommended because empirically RSEARCH transition scores result in poorer performance than the default transitions.
- greedy** Resolve overlapping this with a greedy strategy instead of the default technique of using dynamic programming for optimal hit resolution. This is the technique used by RSEARCH to resolve overlapping hits. This option has not yet been implemented to work in combination with **--inside**.

- gcfile** *<f>* Print information on the GC content of *seqfile* to file *<f>*. First, the GC content of all non-overlapping 100 nucleotide windows in *seqfile* is calculated, and the the distribution of counts is normalized to frequencies for each possible GC percentage [0..100] which are printed to *<f>*.
- beta** *<x>* For query-dependent banding (QDB), which is on by default, set beta parameter to *<x>* where *<x>* is any positive real number less than 1.0. Beta is the probability mass considered negligible during band calculation. The default beta is 1E-7.
- noqdb** Turn the query-dependent banding (QDB) acceleration strategy off, it is on by default.
- qdbfile** *<f>* Read bands for QDB from file *<f>* which was output from **cmbuild** using the **--bfile** option. By default, bands are calculated within **cmsearch**. This is not a very useful option and it was only developed for testing the performance of banded search for CMs built from old versions of **cmbuild**.
- banddump** Print information on the query-dependent bands for each state to standard output.
- hmmfilter** This option is used to accelerate searches. Build a CM Plan 9 HMM from the CM in *cmfile* and use it to filter the database in *seqfile*. Only hits to the HMM with E-values less than or equal to 500 will then be searched with the CM. This E-value threshold of 500 can be changed using the **--hmmE** option or to a bit score threshold using the **--hmmT** option as explained below. The HMM hits are found using a scanning Forward HMM algorithm to determine likely end point positions (j) of hits. Then for each likely end point j, a HMM Backward scan is performed starting at j and moving backward to find the likely start point i. The region from residues i+W-1 to j-W+1 survives the filter, and is then passed to the CM to be searched.
- hmmonly** Search only with a CP9 HMM derived from *cmfile*. Do not pass HMM hits to be searched with the CM, but simply report them. By default E-values are calculated and all hits with an E-value of 50 or less are reported. Works with the **--hmmE** and **--hmmT** options the same way **--hmmfilter** does. Currently, no alignments are printed, only start and stop positions and scores of hits.
- hmmE** *<x>* Set the E-value cutoff for hits to survive the CP9 HMM filter as *<x>*, where *<x>* is a positive real number. HMM hits with E-values better than (less than) or equal to this threshold will survive the filter and be searched with the CM as described above for the **--hmmfilter** option. The default value is 50. CP9 HMM E-values are calculated in the same manner as CM E-values (see the **-E** option) except scoring is performed with the CP9 HMM instead of the CM. The number of randomly synthesized samples is 1000 by default but can be changed with the **--nsamples** option. The **--hmmE** option must be used in combination with either the **--hmmfilter** or the **--hmmonly** option.
- hmmT** *<x>* Set the bit score cutoff for hits to survive the CP9 HMM filter as *<x>* *<x>*, where *<x>* is a positive real number. There is no default value because by default E-values are calculated when **--hmmfilter** is enabled. Hits with bit scores better than

(greater than) this threshold $\langle x \rangle$ will survive the filter and be passed to the CM as described above for the **--hmmfilter** option. This option must be used in combination with either the **--hmmfilter** or the **--hmmonly** option.

--hmmnegsc $\langle x \rangle$ Set the minimum CP9 HMM bit score to survive the filter, or be returned as an HMM hit (if **--hmmonly**) as $\langle x \rangle$ where $\langle x \rangle$ is any real negative number. This option is untested and very experimental, use at your own risk. It's probably only potentially useful in non-local mode, when positive scoring CP9 HMM hits are rare. This option must be used in combination with the **--hmmT** option and either the **--hmmfilter** or the **--hmmonly** option.

References

- Brown, J. W. (1999). The ribonuclease P database. *Nucl. Acids Res.*, 27:314.
- Durbin, R., Eddy, S. R., Krogh, A., and Mitchison, G. J. (1998). *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, Cambridge UK.
- Eddy, S. R. (1998). Profile hidden Markov models. *Bioinformatics*, 14:755–763.
- Eddy, S. R. (2002). A memory-efficient dynamic programming algorithm for optimal alignment of a sequence to an RNA secondary structure. *BMC Bioinformatics*, 3:18.
- Eddy, S. R. and Durbin, R. (1994). RNA sequence analysis using covariance models. *Nucl. Acids Res.*, 22:2079–2088.
- Gerstein, M., Sonnhammer, E. L. L., and Chothia, C. (1994). Volume changes in protein evolution. *J. Mol. Biol.*, 235:1067–1078.
- Giegerich, R. (2000). Explaining and controlling ambiguity in dynamic programming. In Giancarlo, R. and Sankoff, D., editors, *Proceedings of the 11th Annual Symposium on Combinatorial Pattern Matching*, number 1848, pages 46–59, Montréal, Canada. Springer-Verlag, Berlin.
- Griffiths-Jones, S., Bateman, A., Marshall, M., Khanna, A., and Eddy, S. R. (2003). Rfam: an RNA family database. *Nucl. Acids Res.*, 31:439–441.
- Henikoff, S. and Henikoff, J. G. (1994). Position-based sequence weights. *J. Mol. Biol.*, 243:574–578.
- Karplus, K., Barrett, C., and Hughey, R. (1998). Hidden Markov models for detecting remote protein homologies. *Bioinformatics*, 14:846–856.
- Klein, R. J. and Eddy, S. R. (2003). RSEARCH: finding homologs of single structured RNA sequences. *BMC Bioinformatics*, 4:44.
- Krogh, A., Brown, M., Mian, I. S., Sjolander, K., and Haussler, D. (1994). Hidden Markov models in computational biology: Applications to protein modeling. *J. Mol. Biol.*, 235:1501–1531.
- Nawrocki, E. P. and Eddy, S. R. (2007). Query-dependent banding (QDB) for faster RNA similarity searches. *PLoS Comput. Biol.*, 3:e56.
- Weinberg, Z. and Ruzzo, W. L. (2004a). Exploiting conserved structure for faster annotation of non-coding RNAs without loss of accuracy. *Bioinformatics*, 20 Suppl. 1:I334–I341.
- Weinberg, Z. and Ruzzo, W. L. (2004b). Faster genome annotation of non-coding RNA families without loss of accuracy. *RECOMB '04*, pages 243–251.
- Weinberg, Z. and Ruzzo, W. L. (2006). Sequence-based heuristics for faster annotation of non-coding RNA families. *Bioinformatics*, 22:35–39.