

Alternative Sigma Factor Over-Expression Enables Heterologous Expression of a Type II Polyketide Biosynthetic Pathway in *Escherichia coli*

David Cole Stevens, Kyle R. Conway, Nelson Pearce, Luis Roberto Villegas-Peñaranda, Anthony Garza and Christopher N. Boddy*

File S1. Expanded experimental procedures.

Cloning, PCR amplification, and site directed mutagenesis.

A complete list of the plasmids used in this study.

PLASMID	DESCRIPTION
DCS08	<i>E. coli rpoN</i> in pCR-BLUNT
DCS11	<i>E. coli rpoN</i> in pKH22
DCS53	<i>E. coli rpoF</i> in pCR-BLUNT
DCS54	<i>E. coli rpoE</i> in pCR-BLUNT
DCS55	<i>E. coli fecI</i> in pCR-BLUNT
DCS57	<i>E. coli rpoF</i> in pKH22
DCS58	<i>E. coli rpoE</i> in pKH22
DCS59	<i>E. coli fecI</i> in pKH22
DCS61	<i>S. rimosus oxyABCD</i> cassette in pET28b
DCS62	DCS61 with <i>oxyAp</i> Δ -24 mutation
NDP4	<i>E. coli rpoS</i> in pCR-BLUNT
NDP5	<i>E. coli rpoH</i> in pCR-BLUNT
NDP6	<i>E. coli rpoS</i> in pKH22
NDP7	<i>E. coli rpoH</i> in pKH22
MRH08	<i>S. rimosus</i> oxytetracycline biosynthetic pathway in pET28b
LRVP09	pBR322 lacking <i>BLA</i>

A complete list of the primers used for PCR amplification and mutagenesis.

GENE	FORWARD	REVERSE	PROTEIN
<i>E. coli rpoN</i>	aatt catatg aacgaaccttgcaactca	aatt gaattc taaacgagtggttacg	σ^{54}
<i>E. coli rpoD</i>	aatt gctagc atggagcaaaaccgc	aatt gaattc taattcgtccaggaag	σ^{70}
<i>E. coli rpoE</i>	aatt catatg agcagcagtaacgg	aatt gaattc caacgcctgataagc	σ^E
<i>E. coli rpoF</i>	aatt catatg gtgaattcactctatac	aatt gaattc ttataactaccagtt	σ^F
<i>E. coli fecI</i>	aatt catatg tctgaccgcgccacta	aatt gaattc caacgcctgataagc	FecI
<i>E. coli rpoH</i>	gt acatatg actgacaaaatgcaaa	cac ggaattc tacgcttcaatggca	σ^H
<i>E. coli rpoS</i>	aatt catatg aagcaaggttgcaactca	aatt gaattc caaacgagtggtttacg	σ^S
<i>S. rimosus oxyABCD</i>	T7 promoter primer	cagt gaattc catagctccaggctg	oxyA,B,C,D
<i>oxyAp</i> Δ -24 mutation	cctgcgtcccct aaacggcgg gtg	gccaccgcccgttt aggggacgcagg	N/A

qPCR study.

A complete list of the primers used for qPCR analysis.

TRANSCRIPT	FORWARD	REVERSE
16s RNA	aactgcctgatggaggggg	gtctcagttccagtggtgct
<i>rpoN</i>	gtctgttcagttgtcgacg	cggtgtaaatggtgtccaa
<i>oxyB</i>	tcaccggcctcgggtgtggtc	gcgcgagccgggtcatcagg
<i>oxyF</i>	cccgcctgctccaccagcgc	cctggagtgaggacgacgaga
<i>oxyK</i>	ccctgatcaggacgccgcc	tcgcgccaccggaggcgaa

<i>oxyP</i>	ggcccatccgccgcttagg	gcaggtcggcagccccgct
<i>oxyT</i>	aggcgaaggccctgcacagc	cggacaccaggaacgcctgg

5' RACE.

A complete list of primers used for the 5' RACE experiment

Primer name	sequence
oxyA_RACE2	aggctcatcgtcatgccgca
oxyA_nested	ccgcctcgcccgtcacac
AAP	ggccacgcgtcactagtagcgggiigggiiggii
AUAP	ggccacgcgtcactagtagc

Bioinformatics

A list of organisms used in the bioinformatics analysis.

Phylum	Genome Name	Contains σ^{54}	RefSeq ID
Actinobacteria	Bifidobacterium longum NCC2705, complete genome	N	NC_004307.2
Actinobacteria	Mycobacterium tuberculosis CDC1551, complete genome	N	NC_002755.2
Actinobacteria	Mycobacterium ulcerans Agy99, complete genome	N	NC_008611.1
Actinobacteria	Nocardia farcinica IFM 10152, complete genome	N	NC_006361.1
Actinobacteria	Rubrobacter xylanophilus DSM 9941, complete genome	N	NC_008148.1
Actinobacteria	Saccharopolyspora erythraea NRRL 2338, complete genome	N	NC_009142.1
Actinobacteria	Salinispora tropica CNB-440, complete genome	N	NC_009380.1
Actinobacteria	Streptomyces avermitilis MA-4680, complete genome	N	NC_003155.4
Actinobacteria	Streptomyces coelicolor A3(2), complete genome	N	NC_003888.3
Actinobacteria	Streptomyces griseus subsp. griseus NBRC 13350, complete genome	N	NC_010572.1
Aquificae	Sulfurihydrogenibium sp. YO3AOP1, complete genome	Y	NC_010730.1
Bacteroidetes/Chlorobi	Bacteroides fragilis NCTC 9343, complete genome	Y	NC_003228.3
Bacteroidetes/Chlorobi	Candidatus Amoebophilus asiaticus 5a2, complete genome	Y	NC_010830.1
Bacteroidetes/Chlorobi	Gramella forsetii KT0803, complete genome	Y	NC_008571.1
Chlamydiae/Verrucomicrobia	Chlamydia trachomatis D/UW-3/CX, complete genome	Y	NC_000117.1
Chlamydiae/Verrucomicrobia	Chlamydophila pneumoniae AR39, complete genome	Y	NC_002179.2
Chlamydiae/Verrucomicrobia	Coralimargarita akajimensis DSM 45221 chromosome, complete genome	N	NC_014008.1
Chloroflexi	Sphaerobacter thermophilus DSM 20745 chromosome 1, complete genome		NC_013523.1
	Sphaerobacter thermophilus DSM 20745 chromosome 2, complete genome	Y	NC_013524.1
Crenarchaeota	Pyrobaculum arsenaticum DSM 13514, complete genome	N	NC_009376.1
Cyanobacteria	Gloeobacter violaceus PCC 7421, complete genome	N	NC_005125.1
Cyanobacteria	Microcystis aeruginosa NIES-843, complete genome	N	NC_010296.1
Cyanobacteria	Nostoc punctiforme PCC 73102, complete genome	N	NC_010628.1
Cyanobacteria	Synechocystis sp. PCC 6803, complete genome	N	NC_000911.1

Deinococcus-Thermus	Deinococcus deserti VCD115, complete genome	N	NC_012526.1
Euryarchaeota	Pyrococcus horikoshii OT3, complete genome	N	NC_000961.1
Fibrobacteres/Acidobacteria	Solibacter usitatus Ellin6076, complete genome	Y	NC_008536.1
Firmicutes	Bacillus cereus subsp. cytotoxis NVH 391-98, complete genome	Y	NC_009674.1
Firmicutes	Bacillus subtilis subsp. subtilis str. 168, complete genome	Y	NC_000964.3
Firmicutes	Clostridium acetobutylicum ATCC 824, complete genome	Y	NC_003030.1
Firmicutes	Clostridium kluyveri DSM 555, complete genome	Y	NC_009706.1
Firmicutes	Clostridium perfringens str. 13, complete genome	Y	NC_003366.1
Firmicutes	Lactobacillus brevis ATCC 367, complete genome	N	NC_008497.1
Firmicutes	Listeria monocytogenes str. 4b F2365, complete genome	Y	NC_002973.6
Firmicutes	Staphylococcus aureus subsp. aureus MRSA252, complete genome	N	NC_002952.2
Firmicutes	Streptococcus pyogenes M1 GAS, complete genome	N	NC_002737.1
Proteobacteria (Alpha)	Agrobacterium tumefaciens str. C58 chromosome circular, complete sequence	Y	NC_003062.2
	Agrobacterium tumefaciens str. C58 chromosome linear, complete sequence		NC_003063.2
Proteobacteria (Alpha)	Methylobacterium extorquens DM4, complete genome	Y	NC_012988.1
Proteobacteria (Alpha)	Rickettsia rickettsii str. Iowa, complete genome	N	NC_010263.1
Proteobacteria (Beta)	Neisseria meningitidis MC58, complete genome	Y	NC_003112.2
	Ralstonia solanacearum GMI1000 plasmid pGMI1000MP, complete sequence		NC_003296.1
Proteobacteria (Beta)	Ralstonia solanacearum GMI1000, complete genome		NC_003295.1
Proteobacteria (Delta)	Geobacter sulfurreducens PCA, complete genome	Y	NC_002939.4
Proteobacteria (Delta)	Myxococcus xanthus DK 1622, complete genome	Y	NC_008095.1
Proteobacteria (Delta)	Sorangium cellulosum 'So ce 56', complete genome	Y	NC_010162.1
Proteobacteria (Epsilon)	Campylobacter jejuni subsp. jejuni NCTC 11168, complete genome	Y	NC_002163.1
Proteobacteria (Epsilon)	Helicobacter pylori J99, complete genome	Y	NC_000921.1
Proteobacteria (Gamma)	Candidatus Carsonella ruddii PV, complete genome	N	NC_008512.1
Proteobacteria (Gamma)	Escherichia coli str. K-12 substr. MG1655, complete genome	Y	NC_000913.2
Proteobacteria (Gamma)	Francisella tularensis subsp. tularensis SCHU S4, complete genome	N	NC_006570.2
Proteobacteria (Gamma)	Legionella pneumophila subsp. pneumophila str. Philadelphia 1, complete genome	Y	NC_002942.5
Proteobacteria (Gamma)	Pseudomonas aeruginosa PAO1, complete genome	Y	NC_002516.2
Proteobacteria (Gamma)	Pseudomonas putida KT2440, complete genome	Y	NC_002947.3
Proteobacteria (Gamma)	Pseudomonas syringae pv. tomato str. DC3000, complete genome	Y	NC_004578.1
Proteobacteria (Gamma)	Salmonella enterica subsp. enterica serovar Typhi str. CT18, complete genome	Y	NC_003198.1
Proteobacteria (Gamma)	Shewanella oneidensis MR-1, complete genome	Y	NC_004347.1
Proteobacteria (Gamma)	Yersinia pestis CO92, complete genome	Y	NC_003143.1
Spirochaetes	Borrelia burgdorferi B31, complete genome	Y	NC_001318.1
Spirochaetes	Treponema denticola ATCC 35405, complete genome	Y	NC_002967.9
Tenericutes	Candidatus Phytoplasma mali, complete genome	N	NC_011047.1

Positional weighted matrix used to describe the σ^{54} consensus sequence

A:	12	2	0	12	139	11	55	51	46	44	38	13	4	1	9	76
C:	14	0	0	147	23	122	17	48	64	42	62	22	18	2	173	5
G:	10	184	186	6	18	10	103	69	36	35	43	15	10	181	1	17
T:	150	0	0	21	6	43	11	18	40	65	43	136	154	2	3	88

Perl Script used to analyze genomes for putative σ^{54} promoters. The script below is a modified version of the original ProimScan perl script by David J. Studholme. It has been modified to facilitate promoter prediction in genomes of >5 Mb.

```

1  #!/usr/bin/perl -w
2
3  #####
4  #
5  #  copyright      : (C) 2000-2003 by David J. Studholme      #
6  #  email          : ds2@sanger.ac.uk                        #
7  #
8  #####
9
10 #####
11 #*
12 #* This program is free software; you can redistribute it and/or modify #
13 #* it under the terms of the GNU General Public License as published by #
14 #* the Free Software Foundation; either version 2 of the License, or #
15 #* (at your option) any later version.                               #
16 #*
17 #####
18
19 ### VERSION / UPDATE LOG ###
20
21 # 7/11/02 Corrected bugs caught by Richard Pau
22 # 21/11/02 Changed to log base 2 rather than natural log
23 # 07/03 Enabled reporting of distance upstream
24 # enabled normalising of scores.
25 # 27/08/03 Switched format of matrix files.
26 # 27/08/03 Integrated Hit object into same file.
27
28 use strict;
29 use Getopt::Std;
30 use IO::File;
31 use IO::Handle;
32
33
34 $|= 1;
35
36 ### Parse command line arguments
37 ### -m specify matrix file, -g specify genome name, -c specify cut-off
38 ### -h: help, -:statistics, -n: ptt, -i unknown , -m: matrix, -g: genome -c:cutoff
39 my %option;
40 getopts( "hsnim:g:c:", \%option );
41
42 if ( $option{h} ) {
43     show_help();
44     exit;
45 }
46
47 my ($seqread);
48 my ( $seq_length, @seq_window );
49 my ( $mean_score, $total_dev, $sd, $variance );
50 my $GC_content;
51

```

```

52  ### Set threshold score above which counts as a hit
53  my $cut_off = 80;
54  if ( $option{c} ) {
55      $cut_off = $option{c};
56  }
57
58  ### Set input/output filenames
59  my $genome = "ecoli";
60  if ( $option{g} ) {
61      $genome = $option{g};
62  }
63
64  ### Remove any filename extensions
65  if ( $genome =~ m/([\w\d\.]+)\.fna/ ) {
66      $genome = $1;
67  }
68
69  ### Name of file containing sequence to be scanned
70  my $sequence_file = "$genome.fna";
71
72  ### Name of file containing details of ORFs
73  my @ptt_file;
74  unless ( $option{n} ) {
75      my $ptt_file = "$genome.ptt";
76      my $ptt_fh = new IO::File;
77      $ptt_fh->open("<$ptt_file")
78      or die "Failed to open ptt file $genome.ptt: $!\n";
79      while ( my $readline = <$ptt_fh> ) {
80          push @ptt_file, $readline;
81      }
82      $ptt_fh->close;
83  }
84
85  ### Name of matrix file
86  my $matrix_file = $option{m} or die "You must specify a matrix file\n";
87
88  ### Open error file
89
90  open OUTPUT, '>', $option{g}.".txt" or die $!;
91  STDOUT->fdopen( \*OUTPUT, 'w' ) or die $!;
92
93
94  ### Begin the scans!!
95  promscan_kl( $matrix_file, $sequence_file, \@ptt_file );
96  exit;
97
98  sub promscan_kl {
99
100      my ( $matrix_file, $sequence_file, $ptt_file ) = @_;
101
102      my %all_scores_forward;
103      my %all_scores_reverse;
104      my @hits;
105
106      my ( $sequence, $header );
107      print STDERR "Reading the sequence file into memory ... ";

```

```

108 my $seq_fh = new IO::File;
109 $seq_fh->open("<$sequence_file")
110 or die "Failed to open $sequence_file: $!\n";
111
112 while ( my $readline = <$seq_fh> ) {
113     chomp $readline;
114     $sequence .= $readline if $readline =~ m/^[A-Za-z]+$/;
115     if ( $readline =~ m/>(.*)/ ) {
116         $header = $1;
117     }
118 }
119 $seq_fh->close;
120
121 print STDERR "done.\n";
122 print STDERR "Sequence: $header\n\n";
123 print "Sequence|$header\n";
124
125 print STDERR "Splitting the sequence into a list ... ";
126 my @sequence = split //, $sequence;
127 print STDERR "done.\n";
128
129 ### Calculate the GC content of the sequence
130 my $GC_content = calculate_GC(@sequence);
131 print "G+C|$GC_content\n";
132 print STDERR "G+C = $GC_content\n";
133 my %content;
134 $content{G} = $GC_content / 2;
135 $content{C} = $GC_content / 2;
136 $content{A} = ( 1 - $GC_content ) / 2;
137 $content{T} = ( 1 - $GC_content ) / 2;
138
139 ### Read matrix file
140 my $pmatrix = read_matrix_file($matrix_file);
141 print "Matrix|".$option{m}."\n";
142 my %matrix = %$pmatrix;
143
144
145 ### Do sanity check on matrix, and calculate width of the matrix
146 die "Matrix corrupted\n"
147 unless ( @{ $matrix{A} } == @{ $matrix{C} }
148         and @{ $matrix{C} } == @{ $matrix{G} }
149         and @{ $matrix{A} } == @{ $matrix{T} } );
150 my $seq_window_length = @{ $matrix{A} };
151
152 ### Calculate maximum possible score
153 my $overall_highest = 0;
154
155 for ( my $i = 1 ; $i < $seq_window_length ; $i++ ) {
156     my $highest_score = 0;
157     my $highest_base;
158     foreach my $base ( "A", "C", "G", "T" ) {
159         if ( $matrix{$base}[$i] > $highest_score ) {
160             $highest_score = $matrix{$base}[$i];
161             $highest_base = $base;
162         }
163     }

```

```

164         $overall_highest +=
165         ( $highest_score * log( $highest_score / $content{$highest_base} ) /
166           log(2) )
167         if $highest_score;
168     }
169
170 #print STDERR "(Highest possible score before normalisation = $overall_highest)\n" ;
171 #print "(Highest possible score before normalisation = $overall_highest)\n" ;
172
173     print STDERR "Sites scoring > $cut_off\n";
174 ###     print "Sites scoring > $cut_off\n";
175
176     ### Initialise the first window of sequence to be scanned
177     for ( my $i = 0 ; $i < $seq_window_length ; $i++ ) {
178         $seq_window[$i] = $sequence[$i];
179     }
180
181     ### read sequence one base at a time
182     my $n = 0;
183
184     while ( my $seqread = uc( shift @sequence ) ) {
185         my $seq_window;
186         if ( $seqread =~ m/[ACGTN]/ ) { # sanity-check for legal bases
187
188             ### move the sequence window along one residue
189             shift(@seq_window); # removes first residue
190             push( @seq_window, $seqread ); # adds next residue
191
192             ### Assign a score to current window
193             ### This score, known as the Kullback-Liebler distance,
194             ### reflects the theoretical binding energy of the
195             ### DNA-protein interaction.
196             ### For more information see: Stormo, G.D. 2000.
197             ### DNA binding sites: representation and discovery.
198             ### Bioinformatics. 16:16-23.
199
200             ### First check the forward strand
201             my $strand = "+";
202
203             ### Calculate score. Remember, divide by log2 to get log to the base of 2
204             my $current_score = 0;
205             for ( my $i = 0 ; $i < $seq_window_length ; $i++ ) {
206                 my $base = $seq_window[$i];
207                 $current_score = $current_score + ( ( $matrix{$base}[$i] ) *
208                   log( $matrix{$base}[$i] / $content{$base} ) / log(2) )
209                 if $matrix{$base}[$i];
210             }
211
212             ### Normalise the score
213             $current_score /= $overall_highest;
214             $current_score *= 100;
215
216             if ( $option{s} ) {
217                 ### Record each score in a matrix %all_scores_forward
218                 $all_scores_forward{$n} = $current_score;
219             }

```



```

220
221     ### If this window scores greater than the cut off, then record details ...
222     if ( $current_score > $cut_off ) {
223
224         ### Convert @seq_window into a string, $seq_window
225         $seq_window = "";
226         my $q = 0;
227         while ( $seq_window[$q] ) {
228             $seq_window .= $seq_window[$q];
229             $q++;
230         }
231
232         ### Is hit in coding or non-coding region?
233         my ( $locate_match, $intergenic_size ) =
234             locate_match( $pft_file, $n, $strand, $sequence,
235                 $seq_window );
236
237         ### Round score to nearest whole number
238         $current_score = int($current_score);
239
240         ### Print details of this window
241         if ( $locate_match or $option{n} ) {
242
243             ### Create a new Hit object and set its values
244             my $hit = Hit->new();
245             $hit->position($n);
246             $hit->strand($strand);
247             $hit->score($current_score);
248             $hit->strand($strand);
249             $hit->sequence($seq_window);
250             $hit->locus($locate_match);
251
252             ### Record the existence of this Hit object in an array
253             push @hits, $hit;
254
255             printf STDERR "%8s|%1s|%s|%5d|s\n", $n, $strand,
256                 $seq_window, $current_score, $locate_match;
257         }
258     }
259
260     ### now check the Reverse strand
261     $strand = "-";
262
263     ### Calculate score. Remember divide by log2 to get log to the base of 2
264     $current_score = 0;
265     for ( my $i = 0 ; $i < $seq_window_length ; $i++ ) {
266         my $j = $seq_window_length - $i - 1;
267         my $base = revcomp( $seq_window[$j] );
268         $current_score = $current_score + ( ( $matrix{$base}[$i] ) *
269             log( $matrix{$base}[$i] / $content{$base} ) / log(2) )
270             if $matrix{$base}[$i];
271     }
272
273     ### Normalise the score
274     $current_score /= $overall_highest;
275     $current_score *= 100;

```

```

276
277     if ( $option{s} ) {
278         ### Record each score in a matrix %all_scores_forward
279         $all_scores_reverse{$n} = $current_score;
280     }
281
282     ##### print current score if greater than threshold value
283     if ( $current_score > $cut_off ) {
284
285         ### Convert @seq_window into a string, $seq_window
286         $seq_window = "";
287         my $q = 0;
288         while ( $seq_window[$q] ) {
289             $seq_window = $seq_window . $seq_window[$q];
290             $q++;
291         }
292
293         ### Convert $seq_window to Reverse Complement
294         my $revcomp = reverse($seq_window);
295         $revcomp =~ tr/ACGTacgt/TGCAtgca/;
296         $seq_window = $revcomp;
297
298         ### Is hit in coding or non-coding region?
299         my ( $locate_match, $intergenic_size ) =
300             locate_match( $pnt_file, $n, $strand, $sequence,
301                 $seq_window );
302
303         ### Round the score to nearest whole number
304         $current_score = int($current_score);
305
306         ### Print details of this window
307         ### ... but only if it is intergenic
308         if ( $locate_match or $option{n} ) {
309
310             ### Create a new Hit object and set its values
311             my $hit = Hit->new();
312             $hit->position($n);
313             $hit->strand($strand);
314             $hit->score($current_score);
315             $hit->strand($strand);
316             $hit->sequence($seq_window);
317             $hit->locus($locate_match);
318
319             ### Record the existence of this Hit object in an array
320             push @hits, $hit;
321
322             printf STDERR "%08s%01s%0s%05d%0s\n", $n, $strand,
323                 $seq_window, $current_score, $locate_match;
324         }
325     }
326
327     ### The sequence pointer proceeds by one base
328     $n++;
329 }
330 ; # end of if($seqread =~ /[ACGTN]/)

```

```

332
333     }; # end of while (read $sequence_file, $seqread, 1) loop
334
335     ### The final value of $n, the sequence-pointer,
336     ### gives the number of readable bases
337     $seq_length = $n;
338     ### print "N = $seq_length\n";
339
340     ### Create a hash of arrays, to index Hit objects by score
341     ### Cannot use a simple hash because multiple Hits may have
342     ### the same score.
343     my %hitscores;
344     foreach my $hit (@hits) {
345         my $score = $hit->score();
346         unless ( $hitscores{$score} ) {
347             $hitscores{$score} = [];
348         }
349         my $ref = $hitscores{$score};
350         my @list = @$ref;
351         push @list, $hit;
352         @$ref = @list;
353     }
354
355     ### Sort the Hit objects into descending order of scores
356     my @unordered_scores = keys(%hitscores);
357     my @ordered_scores = sort { $b <=> $a } @unordered_scores;
358     my @ordered_hits;
359     foreach my $score (@ordered_scores) {
360         my $ref = $hitscores{$score};
361         my @hits = @$ref;
362         foreach my $hit (@hits) {
363             push @ordered_hits, $hit;
364         }
365     }
366     ### Print out number of hits
367     print "Num of Hits".scalar(@ordered_hits)."n";
368
369     ### Now print out the hits ...
370     foreach my $hit (@ordered_hits) {
371         my ( $pos, $str, $seq, $score, $locus ) = (
372             $hit->position(), $hit->strand(),
373             $hit->sequence(), $hit->score(),
374             $hit->locus()
375         );
376         printf "%8s|%1s|s|%5d|s\n", $pos, $str, $seq, $score, $locus;
377     }
378
379     ### If we have opted to calculate statistics ...
380     if ( $option{s} ) {
381
382         ### Calculate mean score for all possible windows in the
383         ### query sequence
384         my $total_score = 0;
385         for ( my $i = 0 ; $i < $seq_length ; $i++ ) {
386             $total_score =

```

```

388         $total_score + $all_scores_forward{$i} + $all_scores_reverse{$i};
389     }
390
391     $mean_score = $total_score / ( $seq_length * 2 );
392     print "Mean score = $mean_score\n";
393
394     ### Calculate variance and sd of scores for all
395     ### possible windows in the query sequence
396     $total_dev = 0;
397     for ( my $i = 0 ; $i < $seq_length ; $i++ ) {
398         $total_dev =
399             $total_dev + ( ( $all_scores_forward{$i} - $mean_score )**2 ) +
400             ( ( $all_scores_reverse{$i} - $mean_score )**2 );
401     }
402
403     $variance = $total_dev / ( $seq_length * 2 );
404     print "Variance = ", $variance, "\n\n";
405     $sd = $variance**(0.5);
406 }
407 ### print time - $^T;
408
409     $seq_fh->close;
410
411 } # end of promscan_kl()
412
413 sub locate_match {
414     my $pvt_file      = shift or die "Failed to pass a PTT file";
415     my $promoter_pos  = shift or die "Failed to pass a position\n";
416     my $promoter_strand = shift or die "Failed to pass a strand\n";
417     my $sequence      = shift or die "Failed to pass a sequence\n";
418     my $seq_window    = shift
419         or die
420         "Filed to pass a seq_window ($pvt_file,$promoter_pos,$promoter_strand\n";
421
422     unless ( defined $option{n} ) {
423
424         my ( $start_pos1, $end_pos1, $strand1, $product1 ) =
425             ( 0, 0, "?", "ORIGIN" );
426         my ( $start_pos2, $end_pos2, $strand2, $product2 ) =
427             ( 0, 0, "?", "ORIGIN" );
428
429         my $promoter_description = "";
430         ### $promoter_description = "not found";
431         foreach my $read_line ( @$pvt_file ) {
432
433             # Read each window of the protein table file
434             my @fields = read_ptt_line($read_line);
435
436             unless ( $fields[0] eq "failed" ) {
437
438                 ( $start_pos1, $end_pos1, $strand1, $product1 ) =
439                     ( $start_pos2, $end_pos2, $strand2, $product2 );
440                 ( $start_pos2, $end_pos2, $strand2, $product2 ) = @fields;
441
442                 my $size_of_intergenic_region = 0;
443                 my $upstream_sequence       = "";

```

```

444
445     ### Check whether promoter falls within this window
446     if ( $promoter_position > $end_pos1
447         && $promoter_position < $start_pos2 )
448     {
449
450         ### promoter is in intergenic region...
451         my $size_of_intergenic_region_temp =
452             $start_pos2 - $end_pos1;
453
454         if ( $promoter_strand eq "+" && $strand2 eq "+" ) {
455             $size_of_intergenic_region =
456                 $size_of_intergenic_region_temp;
457             $upstream_sequence = substr(
458                 $sequence,
459                 ( $end_pos1 - 50 ),
460                 ( $size_of_intergenic_region + 49 )
461             );
462
463             my $spacer = $start_pos2 - $promoter_position -
464                 length($seq_window);
465             $promoter_description =
466                 "-N($spacer)- $product2|$size_of_intergenic_region";
467
468         }
469
470         elsif ( $promoter_strand eq "-" && $strand1 eq "-" ) {
471             $size_of_intergenic_region =
472                 $size_of_intergenic_region_temp;
473             my $seq_tmp = substr(
474                 $sequence,
475                 ( $end_pos1 - 2 ),
476                 ( $size_of_intergenic_region + 52 )
477             )
478             or die
479             "Failed to get seq_tmp! end_pos1=$end_pos1, size_of_intergenic_region=$size_of_intergenic_region\n";
480
481             my $spacer = $promoter_position - $end_pos1;
482             $promoter_description =
483                 "-N($spacer)- $product1|$size_of_intergenic_region";
484
485         }
486
487         elsif ( $promoter_strand eq "-" && $strand1 eq "+" ) {
488             $size_of_intergenic_region =
489                 $size_of_intergenic_region_temp;
490             my $seq_tmp = substr(
491                 $sequence,
492                 ( $end_pos1 - 2 ),
493                 ( $size_of_intergenic_region + 52 )
494             )
495             or die
496             "Failed to get seq_tmp! end_pos1=$end_pos1, size_of_intergenic_region=$size_of_intergenic_region\n";
497
498             my $spacer = $promoter_position - $end_pos1;
499             $promoter_description =

```

```

500 "-N($spacer)- $product1 (WRONG ORIENTATION)|$size_of_intergenic_region";
501
502         }
503
504         elsif ( $promoter_strand eq "+" && $strand1 eq "-" ) {
505             $size_of_intergenic_region =
506             $size_of_intergenic_region_temp;
507             my $seq_tmp = substr(
508                 $sequence,
509                 ( $end_pos1 - 2 ),
510                 ( $size_of_intergenic_region + 52 )
511             )
512             or die
513             "Failed to get seq_tmp! end_pos1=$end_pos1, size_of_intergenic_region=$size_of_intergenic_region\n";
514
515             my $spacer = $promoter_position - $end_pos1;
516             $promoter_description =
517             "-N($spacer)- $product1 (WRONG ORIENTATION)|$size_of_intergenic_region";
518
519         }
520     }
521
522     elsif ( ( $promoter_position >= $start_pos1 )
523             and ( $promoter_position <= $end_pos1 )
524             and ( $option{i} ) )
525     {
526         $promoter_description = "internal to $product1";
527     }
528 }
529
530 }
531 my $size_of_intergenic_region = $start_pos2 - $end_pos1;
532 return ( $promoter_description, $size_of_intergenic_region );
533 }
534 } # End of sub locate_match
535
536 sub read_ptt_line {
537
538     # Expects a scalar value containing one line of a .ptt file
539     # Extracts the 'location', strand, and 'product'
540     # from that line and returns them, in that order
541
542     my $read_line = $_[0];
543     my $i;
544     my $product;
545     my @fields;
546
547     if ( $read_line =~ /(\d+)\.\.(\d+)\s+([+-])/ ) {
548
549         # Extract 'location' and 'strand'
550         @fields = split( /\s+/, $read_line );
551         $i = 5;
552         $product = "";
553         while ( $fields[$i] ) {
554

```

```

556             # Extract 'product'
557             $product = $product . " " . $fields[$i];
558             $i++;
559         }
560         return ( $1, $2, $3, $product );
561     }
562     ; # End of 'if( $read_line =~ ....' block
563
564     ### print ("variable read_line:$read_line\n");
565     return ("failed");
566
567 }; # end of sub read_ptt_line;
568
569 sub calculate_GC {
570     my (@sequence) = @_;
571     my ( $n, $GC, $AT ) = ( 0, 0, 0 ); # initialise counters
572     my $seqread;
573
574     foreach my $seqread (@sequence) {
575         $GC++ if $seqread =~ /[GCgc]/;
576         $AT++ if $seqread =~ /[ATat]/;
577     }
578     $n = $GC + $AT;
579     my $GC_content = $GC / $n;
580     return $GC_content;
581 }
582
583 sub revcomp {
584     my $string = shift @_ or print STDERR "Called revcomp on null string!";
585     my $seq = "";
586     my @array = split( //, $string );
587     while ( my $char = pop @array ) {
588         $char =~ tr/ACGT/TGCA/;
589         $seq = $seq . $char;
590     }
591
592     #print STDERR "Revcomp of $string -> $seq\n" ;
593     return $seq;
594 }
595
596 sub read_matrix_file {
597
598     my $matrix_file = shift;
599     my %matrix;
600
601     my $fh = new IO::File;
602     $fh->open("<$matrix_file")
603     and print STDERR "Opened matrix file: $matrix_file\n"
604     or die "Could not open Configuration file $matrix_file: $!\n";
605
606     ### print "\nScoring matrix:\n\n";
607
608     while ( my $readline = <$fh> ) {
609         chomp($readline);
610
611         if ( $readline =~ m/([ACGT]):([\s\d]+)/ ) {

```

```

612         ###    print "$readline\n";
613         print STDERR "$readline\n";
614
615         my $base = $1;
616         my @values = split /\s+/, $2;
617
618         $matrix{$base} = [];
619         foreach my $value (@values) {
620             push @{$matrix{$base}}, $value;
621         }
622     }
623 }
624
625 ###    print "\n\n";
626     print STDERR "\n\n";
627
628     $fh->close;
629     print STDERR "Closed matrix file\n";
630
631     return \%matrix;
632 }
633
634 sub show_help {
635     print STDERR <<EOF;
636
637     =====
638     =====
639     $0: scan a DNA sequence against a motif specified as a scoring matrix
640     =====
641     =====
642
643     Usage: $0 <options>
644
645     Available options
646
647     -h      : Show this help.
648     -l      : Use linear scoring formula ignoring G+C content,
649               rather than the Kullback-Leibler distance.
650     -s      : Calculate statistics (this requires a lot of memory!)
651     -i      : Show matches that are internal to ORFs as well as intergenic ones.
652     -n      : Do not use a .ptt file (ie no ORF information)
653     -g <file> : Specify the input file (genome)
654               You can specify the .ppt file or the .fna file
655               or just the base name:
656               eg. $0 -g NC_000962.ppt
657                   $0 -g NC_000962.fna
658                   $0 -g NC_000962
659     -m <file> : Specify the matrix file.
660     -c <number> : Specify the threshold (cut-off) score (default is 1000).
661               You may wish to experiment with various threshold values
662               to get the optimum balance between sensitivity and accuracy.
663
664     =====
665     =====
666
667

```


668 Output is printed to the output file as well as to STDOUT (ie the screen).
 669 For each hit, ie a close match to the consensus sequence that is being sought,
 670 the output is arranged into six columns:
 671
 672 Column 1 : The position of the hit in the input sequence.
 673 Column 2 : The strand on which the hit is found (+/-).
 674 Column 3 : Sequence of the hit.
 675 Column 4 : A score for the hit. The higher the score, the better the match!
 676 Column 5 : The position of the hit relative to adjacent/downstream Open Reading Frames
 677 Column 6 : The size of the intergenic region in which the match is found.

678
 679 For more information see the PromScan website <http://www.promscan.uklinux.net>
 680

```

681 =====
682 =====
683 EOF
684 }
685
686 ### Hit.pm class of objects representing a hit,
687 ### i.e. a good match to the sought sequence
688 ### defined by the frequency matrix
689
690 package Hit;
691
692 sub new {
693     my $self = {};
694     bless $self, "Hit";
695     return $self;
696 }
697
698 sub position {
699     my $self = shift @_ ;
700     if ( my $specified = shift @_ ) {
701         $$self{POSITION} = $specified;
702     }
703     return $$self{POSITION};
704 }
705
706 sub strand {
707     my $self = shift @_ ;
708     if ( my $specified = shift @_ ) {
709         $$self{STRAND} = $specified;
710     }
711     return $$self{STRAND};
712 }
713
714 sub score {
715     my $self = shift @_ ;
716     if ( my $specified = shift @_ ) {
717         $$self{SCORE} = $specified;
718     }
719     return $$self{SCORE};
720 }
721
722 sub sequence {
723     my $self = shift @_ ;
  
```

```
724     if ( my $specified = shift @_ ) {
725         $$self{SEQ} = $specified;
726     }
727     return $$self{SEQ};
728 }
729
730 sub locus {
731     my $self = shift @_ ;
732     if ( my $specified = shift @_ ) {
733         $$self{LOCUS} = $specified;
734     }
735     return $$self{LOCUS};
736 }
737
738 sub intergenic_size {
739     my $self = shift @_ ;
740     if ( my $specified = shift @_ ) {
741         $$self{SIZE} = $specified;
742     }
743     return $$self{SIZE};
744 }
745
746 return 1; # end of package
```