

---

# **bio***bits* Documentation

**Release 1.2.0**

**Tyghe Vallard, Michael Panciera**

January 07, 2016



---

Contents

---

<b>1</b>	<b>TODO</b>	<b>3</b>
1.1	Installation . . . . .	3
1.2	Scripts . . . . .	3
1.3	AMOS . . . . .	20
1.4	CHANGELOG . . . . .	22
1.5	TODO . . . . .	24
<b>2</b>	<b>Indices and tables</b>	<b>25</b>



Various bioinformatics scripts

All documentation is hosted at <http://bio-bits.readthedocs.org/en/latest>



---

## TODO

---

- Include existing scripts

Contents:

## 1.1 Installation

It is recommended to install into a virtualenv. If you know what you are doing and don't want to install into virtualenv, then you can skip right to step 3

### 1. Setup Virtualenv

It is assumed you have virtualenv already installed. If not see <https://virtualenv.pypa.io/en/latest/installation.html>

```
virtualenv env
```

### 2. Activate virtualenv

```
. env/bin/activate
```

### 3. Install dependencies

```
pip install -r requirements.txt
```

For python 2.6 you will need to also install some additional packages

```
pip install -r requirements-py26.txt
```

### 4. Install bio\_bits

```
python setup.py install
```

## 1.2 Scripts

### 1.2.1 rename\_fasta

Many times you find you have a fasta file where the identifiers are all wrong and you want to rename them all via some mapping file.

Take the example where you have the following fasta file(example.fasta):

```
>id1  
ATGC  
>id2  
ATGC  
>id3  
ATGC
```

You want to rename each identifier(id1, id2, id3) based on a mapping you have. In a file called renamelist.csv you would have the following:

```
#From, To  
id1, samplename1  
id2, samplename2  
id3, samplename3
```

Then to rename your fasta without replacing the original file you have two options:

1. Rename without replacing original file

```
rename_fasta renamelist.csv example.fasta > renamedfasta.fasta
```

2. Rename replacing original file's contents

```
reanme_fasta renamelist.csv example.fasta --inplace
```

## Rename Mapping File Syntax

The file you specify as the rename map file is a simple comma separated text file.

The following rules apply to the format:

- The first entry is the identifier to find in the supplied fasta file.
- The second entry is what to replace the found identifier with
- Any line beginning with a pound sign(#) will be ignored by the renamer

### Missing identifiers that are in fasta but not rename file

In the case where your fasta file contains an identifier that is not in the rename map file you supply, an error will be displayed in the console telling you as such:

```
idwhatever is not in provided mapping
```

## 1.2.2 beast\_checkpoint

beast\_checkpoint is a fork of <https://gist.github.com/trvrb/5277297> that has been rewritten in python and slightly improved as the ruby script seemed to have a few errors.

It accepts any previously run or terminated beast run and will generate an xml file that essentially starts from the last generated tree/log state.

Since beast is random in nature, there does not appear to be a way to restart the run exactly from the same state that it left off.

## Example

We will use the benchmark2.xml file that comes with Beast 1.8. This file is located in:

```
BEASTv1.8.0/examples/Benchmarks/benchmark2.xml
```

First you need to fix the benchmark2.xml because each taxa has a trailing space and that is annoying

```
$> sed 's/ \"/\"/' benchmark2.xml > beast.xml
```

Now run beast for about half of the iterations and hit CTRL-C to kill it. This benchmark is set to run 1,000,000 iterations so around 500,000 you can kill it. Notice we are using a predefined seed

```
$> seed=1234567890
$> mkdir run1
$> cp beast.xml run1/beast.xml
$> beast -seed $seed -beagle_SSE beast.xml
```

Now we will want to re-run beast from that last state. We can use beast\_checkpoint to do so by supplying the original xml and the produced trees and log files. We will put the new xml into a new directory since the .trees and .log files would create an error or possibly be overwritten.

*NOTE* If your fileLog and treeFileLog do not have the same logEvery then when beast exits you may end up with more/less tree states than log states. For now you will have to manually edit the files and ensure that the last tree state matches the last log state.

---

## Todo

Could be possible to get beast\_checkpoint to check for that scenario and use the last tree state that matches the last log state

```
$> mkdir run2
$> beast_checkpoint beast.xml *.trees *.log > run2/beast.xml
```

Now you can simply just re-run beast on the new xml using the same seed

```
$> cd run2
$> beast -seed $seed -beagle_SSE beast.xml
```

## Tracer

If you name your runs sequentially as we did in the example(aka, run1, run2,...) then you can easily load all log files into tracer via the command line as follows

```
tracer run*/*.log
```

## LogCombiner

After you have run all your beast checkpointed xml files you will probably want to combine them with logcombiner which comes with beast

### 1.2.3 beast\_wrapper

Beast wrapper is intended as a helper script to run beast. At this point it just runs beast with the same arguments you would normally give to beast from the command line and just adds a estimated time left column to the console output

## Example

```
$> beast_wrapper -beagle_SSE my_beast.xml
...
state   Posterior      Prior      Likelihood      rootHeight      my_beast.ulcl.mean      location
0      -86527.5880    -6850.8316    -79676.7564    57.6772      1.16103E-3      4.86012
20000   -29044.3753    -1123.5287    -27920.8466    288.102      3.02471E-4      0.11891
40000   -25517.9525    -979.5343     -24538.4182    211.705      1.35118E-4      0.25060
60000   -24212.1250    -1040.4103    -23171.7147    188.454      1.05572E-4      0.18908
80000   -24097.9354    -1019.8099    -23078.1256    182.242      1.53593E-4      0.12857
100000  -24121.5382    -1105.6545    -23015.8837    178.060      1.26907E-4      0.10367
120000  -23930.6897    -1105.7390    -22824.9507    187.411      1.01885E-4      0.34214
140000  -23869.4856    -1087.1915    -22782.2942    178.535      8.76375E-5      0.26128
```

## 1.2.4 group\_references

group\_references splits an alignment file by reference into separate FASTQ files. group\_references takes a SAM or BAM file as input, and can optionally be given an output directory where the FASTQ files will be saved. If no output directory name is provided, the files will be saved in the new folder group\_references\_out.

```
$> group_references contigs.bam
$> group_references contigs.bam --outdir split_fastqs
```

## 1.2.5 degen

Find genes where a sequence has degenerate bases.

### How-to

**Usage:** degen.py <fasta> <options>

Options:

- gb-id=<accession\_id> Accession id for reference
- gb-file=<gbfile> Local Genbank file for reference
- tab-file=<tabfile> TSV/CSV file for reference with fields name,start,end

### Example:

```
degen sequence.fasta --gb-id 12398.91
degen sequence.fasta --gb-file tests/testinput/sequence.gb
degen sequence.fasta --tab-file tests/testinput/degen.tab
degen sequence.fasta --tab-file tests/testinput/degen.csv
```

### Output:

Gene name, degenerate position, degenerate base:

anchored capsid protein	85	R
anchored capsid protein	88	Y
membrane glycoprotein precursor	509	R
nonstructural protein NS5	8513	Y
nonstructural protein NS5	8514	Y
nonstructural protein NS5	8515	Y
anchored capsid protein	85	R
anchored capsid protein	88	Y
membrane glycoprotein precursor	509	R
nonstructural protein NS5	8513	Y
nonstructural protein NS5	8514	Y
nonstructural protein NS5	8515	Y

## Gene/Tab File

degen.tab could look like:

genename	start	stop
foo	1	2
bar	9	33

The headers do not matter, but the start field must always come before the stop field, so the below example would also be valid:

start	GENENAME	stop
1	foo	2
9	bar	33

or optionally without headers:

1	foo	2
9	bar	33

alternatively, with commas in place of tabs:

name,start,stop
foo,1,2
bar,9,33

You can also specify a coding region(CDS) in your file as well:

name,start,stop
CDS,3,33
foo,1,2
bar,9,33

## Genbank File

As downloaded from NCBI's entrez database. Use this option if you don't have internet access.

An example

LOCUS	KJ189367	10452 bp ss-RNA	linear	VRL 10-FEB-2014
DEFINITION	Dengue virus 1 isolate DENV-1/PR/BID-V8188/2010, complete genome.			
ACCESSION	KJ189367			
VERSION	KJ189367.1	GI:582052497		
DBLINK	BioProject: PRJNA31235			

KEYWORDS	.
SOURCE	Dengue virus 1
ORGANISM	Dengue virus 1
	Viruses; ssRNA viruses; ssRNA positive-strand viruses, no DNA stage; Flaviviridae; Flavivirus; Dengue virus group.
REFERENCE	1 (bases 1 to 10452)
AUTHORS	Zody, M.C., Newman, R.M., Henn, M., Munoz-Jordan, J., McElroy, K.L., Santiago, G., Poon, T.W., Charlebois, P., Weiner, B., Yang, X., Piper, M.E., Fitzgerald, M., McCowan, C., Young, S., Gargiula, S., Levin, J., Malboeuf, C., Qu, J., Ireland, A., Chapman, S.B., Murphy, C., Wortman, J., Nusbaum, C. and Birren, B.
CONSRTM	Genome Resources in Dengue Consortium; The Broad Institute Genomics Platform; The Broad Institute Genome Sequencing Center for Infectious Disease; Centers for Disease Control and Prevention Division of Vector Borne Infectious Diseases; CDC Dengue Branch Puerto Rico
TITLE	Direct Submission
JOURNAL	Submitted (22-JAN-2014) Broad Institute of MIT & Harvard, 7 Cambridge Center, Cambridge, MA 02142, USA
COMMENT	##Assembly-Data-START## Assembly Method :: Vicuna v. 1 Sequencing Technology :: Illumina ##Assembly-Data-END##
FEATURES	Location/Qualifiers
source	1..10452 /organism="Dengue virus 1" /mol_type="genomic RNA" /isolate="DENV-1/PR/BID-V8188/2010" /isolation_source="cell supernatant" /host="Homo sapiens" /db_xref="taxon:11053" /country="Puerto Rico" /collection_date="2010" /note="cell passage history: C6/36 1; cohort population: Dengue Surveillance; type: 1"
5' UTR	1..83 /note="indels in UTR have not been validated"
CDS	84..10262 /codon_start=1 /product="polyprotein" /protein_id="AHI43750.1" /db_xref="GI:582052498" /translation="MNNQRKKTGRPSFNLKRARNRVSQQLAKRFSKGLLSGQGPM KLVMAFIAFLRLFLAIAPPAGILARWSSFKNGAIKVLRGFKEISSMLNIMNRRKRSV TMLLMLLPTALAFHLTRGGEPMIVSKQERGKSLLFTSAGVNCTLIAMDLGELCE DTMTYKCPRITEAEPDDVDCWCNATDTWVTYGTCSQTGEHRREKRSVALAPHVGLGLE TRTETWMSSEGAWKQIQRVETWALRHPGFTVIAFFLAHAIGTSITQKGIIFILLMLVT PSMAMRCVGIGNRDFVEGLSGATWVDVVLHGSCVTMAKNKPILDIELLKTEVTNPA VLRKLCLIEAKISNTTDSRCPTQGEATLVEEQDANFVCRRTFVDRGWGNGCGLFGKGS LLTCAFKCVTKLEGKIVQYENLKYSVIVTVHTGDQHQVGNETTEHGTIATITPQAPT SEIQLTDYGAULTLDCSPRTGLDFNEMVLLTMEKEKSWLVHKQWFLLPLPWTSGASTSQ ETWNQRDLLVTFKTAHAKKQEVVVLGSQEGAMHTALTGATEIQTSGTTIFAGHLKCR LKMDKLTGKMSYVMCTGSFKLEKEVAETQHGTVLVQVKYEGTDAPCKIPSTQDEKG VTQNGLRITANPIVTDKEPVNIETEPPGESYIVVGAGEKALKLSWFKRGSSIGKMF EATARGARRMAILGDTAWDFSIGGVFTSVGKLVHQIFGTAYGVLFSGVSWTMKIGIG ILLTWLGLNSRSTSLSMTCIVVGMVTLYLGVMQADSGCVINWKGRELKCGSGIFVTN EVHTWTEQYKFQADSPKRLSAAGKAWEEGVCGIRSATLENIMWKQISNELNHILLE

```

NDMKFTVVVGDANGILAQGKKMIRPQPMEHKYSWKGAKIIGADIQNTTFIIDGPD
TPECPDGQRawnIWEVEDYGFVFTTNIWKLRLDSYTQMCDHRLMSAIKDSKAVHAD
MGYWIeseKNETWKLARASFIEVKTCTWPKSHTLWSNGVLESEMIIPKIYGGPISQHN
YRPGYFTQTAGPWHLGKLELFDLCEGTTVVVDEHCNRPSSLRTTVTGKIIHEWCC
RSCTLPLRFRGEDGCWYGMEIRPVKEKEENLVRSMVSAGSGEVDSFSLGILCVSIMI
EEVMRSRWSRKMLMTGTLAVFLIMQQLTWNDLIRLCIMVGANASDRMGMGTTYLAL
MATFKMRPMFAVGLLFRRLTSREVLLTIGLSLVASVELPNSLEELGDGLAMGIMMLK
LLTEFQPHQLWTLLSLLTFVKTTLSLDYAWKTTAMALSIVSLFPLCLSTTSQTTWLP
VLLGSFGCKPLTMFLITENKIWGRKSWPINEGIMAIGIVSILLSSLLKNDVPLAGPLI
AGGMLIACYVISGSSADLSLEKAAEVSWEAEHSGASHSILVEVQDDGTMKIKDEER
DDTLTILLKATLLAVSGVYPMSPATLFWVYFWQKKQRSQGVLDTPSPPEVERAVLD
NGIYRILQRGLLGRSQVGVGVFQDGVFTMWVTRGAFLMYQGKRLEPSWASVKKDLI
SYGGGWRFQGSWNTGEEVQVIATEPGKNPKNVQTPGTFTKPEGEVGAIALDFKPGTS
GSPIVNREGKIVGLYGNVTTSGTYVSAIAQAKASQEGPLPEIEDEVFKRNLTIMD
LHPGSGKTRRYLPAIVREAIRKRLRTLILAPTRVVA SEMAEALKMPIRYQTAVKSE
HTGREIVDLMCHATFTMRLSPRVPNYNMIIMDEAHFTDPASIAARGYISTRVGMGE
AAAIFMTATPPGSVEAFPQSNAVIQDEERDIPERSWNSGYDWITDFPGKTVWFVPSIK
SGNDIANCLRKNGKRVIQLSRKTFDTEYQKTKNNWDYVVTDISEMGANFRADRVID
PRRCLKPVILKDGPERVILAGPMPVTAASAAQRGRIGRNQNKEGDQYVYMGQPLNNND
EDHAHWTAEKMLLDNINTPEGIIPALFEPEREKSAAIDGEYRLGEARKTFVELMRRG
DLPVWLSSYKVASEGFQYSDRRWCDFGERNNQVLEENMDVEIWTKEGERKKLRPRWLDA
RTYSDPLALREFKEAAGRGSVSGDLILEIGKLPQHLLRAQNALDNLVMLNHSEQGG
KAYRHAMEELPDTIETLMLLALIAVLTGGVTLFFLSGKGLGKTSIGLLCVTASSALLW
MASVEPHWIAASIILEFFLMVLLIPEPDRQRTPQDNQLAYVVIIGLLFMILTVAANEMG
LLETTKKDLGIGYVAENHQATMLDVDLHPASAWTLYAVATTVITPMMRHTIENTTA
NISLTAIANQAAILMGLDKGWPISKMDIGVPLLALGCYSQVNPLTLTAAVLMLVAHYA
IIGPGLQAKATREAQKRTAACGIMKNPTVDGIVAIIDLPVVYDAKFEKQLGQIMLLILC
TSQILLMRTTWALCESITLATGPI.TTLWEGSPGKFWNTTIAVSMANIFRGSYLAGAGL
AFSLMKSLLGGGRRGTGAQGETLGEKWKRQLNQLSKSEFNTYKRSGIMEVDRSEAKEGL
KRGETTKHAVSRGTAKLRFVERNLVKPEGKVIDLGCGRGGWSYYCAGLKKVTEVKGY
TKGGPGHEEPIPMTAYGNLVLKHSGKDVFMPPEKCDTLLCDIGESSPNPTIEEGRT
LRVLKMVEPWLRGNQFCIKILNPYMPSVVETLERMQRKHGGMLVRNPLSRNSTHEMYW
VSCGTGNIVSAVNMTSRMLLNRFTMAHRKPTYERDVLDLGAGTRHVAVEPEVANLDIIG
QRIENIKNEHKSTWHYDEDNPYKTWAYHGSYEVKPGSGASSMVNVVRLLTPWDVIP
MVTQIAMTDTTPFGQQRVFKEKVDTRTPRAKRGTTQIMEVTAKWLWGFLSRNKKPRIC
TREEFTRKVRNSNAAIGAVFVDENQNSAKEAVEDERFWDLVHRRELHKQGKCATCVY
NMMGKREKKLGEGFKAKGSRAIWYMWL GARFLEFEALGFMNEDHWFSRENLSGVGEGL
LHKLGYIILRDISKIPGGNMYADDTAGWDTRVTEDDLQNEAKITDIMEPEHALLATSI
FKLTYQNKKVVRVQRPNAKNGTVMDVISRRDQRGSGQVGTYGLNTFTNMEVQLIRQMESE
GIFLPSELETPNLAERALDWLEKHGAERLKRMAISGDDCVVKPIDDRFATALTALNDM
GKVRKDIPQWEP SKGWNDWQQVPCSHHFQQLIMKDREIVVPCRNQDELVGRARVSQ
GAGWSLRETA CLGSKSYAQMQLMFHRRDLRLAANAICSAVPDVWPTSRTTWSIH AH
HQWMTTEDMLSVWNRVWIDENPWMENKTHVSSWEV PYLGKREDQWCGSLIGLTARAT
WATNIQVAINQVRLIGNENYLDYMTSMKRFKNESDSEGALW"
84..425
mat_peptide /product="anchored capsid protein"
426..923
mat_peptide /product="membrane glycoprotein precursor"
924..2408
mat_peptide /product="envelope protein"
2409..3464
mat_peptide /product="nonstructural protein NS1"
3465..4118
mat_peptide /product="nonstructural protein NS2A"
4119..4508
mat_peptide /product="nonstructural protein NS2B"
4509..6365
mat_peptide /product="nonstructural protein NS3"

```

```
mat_peptide      6366..6746
                  /product="nonstructural protein NS4A"
mat_peptide      6747..6815
                  /product="2K peptide"
mat_peptide      6816..7562
                  /product="nonstructural protein NS4B"
mat_peptide      7563..10259
                  /product="nonstructural protein NS5"
3'UTR           10263..10452
                  /note="indels in UTR have not been validated"

ORIGIN
       1 catctggacc gacaagaaca gttcgaatc ggaagcttgc ttaacgtagt tctaacagtt
       61 ttttattaga gagcagatct ctgatgaaca accaacggaa aaagacgggt cgaccgttt
      121 tcaatatgtc gaaacgcgcg agaaaccgcg tgtcaactgg ttcacagttt gcgaagagat
      181 tctcaaaaagg attgcttca ggccaaggac ccatgaaatt ggtatggct ttcatagcat
      241 ttctaaaggatt tctagccata ccccccaacag caggaatttt ggcttagatgg agctcattca
      301 agaagaatgg agcaattaaa gtgttacggg gtttcaaaaa agagatctca agcatgttga
      361 acataatgaa caggaggaaa agatccgtga ccatgctcct catgctgctg cccacagccc
      421 tggcgtttca tttgaccaca cgagggggag agccacacat gatagttgt aagcaggaaa
      481 gaggaaagtc actcttggtt aagacctctg cggggcgtcaa tatgtgcacc ctcatggcga
      541 tggacttggg agagttatgt gaggacacaa tgacctacaa atgccccgg atcaactgagg
      601 cggaaaccaga tgacgttgc tgctggtgca atgccacaga cacatgggtg acctatggg
      661 cgtgttctca aaccggcgaa caccgacgag agaaacgttc cgtggactg gccccacacg
      721 tgggacttggg tctagaaaca agaaccgaaa catggatgtc ctctgaaggc gcctggaaac
      781 aaatacaaaag agtggaaact tgggcttga gacacccagg attcacgggt atagcctttt
      841 ttttagcaca tgctatagga acatccatca ctcagaaagg gatcatttc atcttgctga
      901 tgctggtgac accatcaatg gccatgcgt gctggggaaat aggcaacaga gacttcgttgc
      961 aaggactgtc aggagcaacg tgggtggacg tggtaactgg gacacggc tgctgtcacca
     1021 ccatggcaaa aaataaacca acattggaca ttgaactctt gaagacggag gtcacgaacc
     1081 ctggcgctt gcgcaaaactg tgcatgttca gtaaaatatc aaacaccacc accgattcaa
     1141 gatgtccaaac acaaggagag gcccacactgg tggaaagaaca agacgcgaac tttgtgtgtc
     1201 gccgaacgtt tggacacaa ggtggggta atggctggg actattcgg aaggaaagtc
     1261 tattgacgtg tgccaaggatc aagtgtgtga caaaaactaga agggaaagata gttcaatatg
     1321 aaaacctaaa atattcgtg atagtcaactg tccacactgg ggaccagcac caggtgggaa
     1381 acgagaccac agaacatgg acaattgca ccataacacc tcaagctccc acgtcgaaaa
     1441 tacagctgac cgactacggc gcccctcacac tggactgtc acctagaaca gggctggact
     1501 ttaatgagat ggtgttattt gcaatgaaag aaaaatcatg gcttgtccac aaacaatgtt
     1561 ttctagactt gccactgcca tggacttcgg gggcttcaac atcccaagag accttggaaaca
     1621 gacaagattt gctggtcaca ttcaagacag ctcatgaaa gaaacaggaa gtagtcgtat
     1681 tgggatcaca ggaaggagca atgcatactg cggtactgg ggcacagaa atccagacgt
     1741 caggaacgac aacaatcttgc caggacacc tgaaatgcag actaaaaatg gataaaactga
     1801 ccttaaaggg gatgtcatat gtatgttca caggcttatt taagcttagag aaggaagtgg
     1861 ctgagaccac gcatgaaact gttctatgtc aggtcaaata tgaagggaca gacgcgcac
     1921 gcaagatccc ctttcgacc caagatgaga aaggagtgc ccagaatggg agattgataa
     1981 cagccaatcc catagttact gacaaagaaa aaccgtcaa cattgagaca gaaccaccc
     2041 ttgggtgagag ctacatcgatgt gtagggcag gcaaaaaagc ttgaaacta agctggttca
     2101 agagagggaa cagcataggg aaaatgttcg aagcaaccgc ccgaggagca cgaaggatgg
     2161 ctatcctggg agacaccgc tggacttcg gttctatagg aggagtgtt acatctgtgg
     2221 gaaaatttgtt acaccagatt ttggaaaccg catatgggt tctgttttagc ggtgtttttt
     2281 ggaccatgaa aataggaata gggattctgc tgacatgggtt gggattaaat tcaaggagca
     2341 cgtcaacttc gatgacgtgc attgttagtt gcatggtcac actgtaccta ggagtcatgg
     2401 ttcaagcgga ttggggatgt gtatgtcaact ggaagggcag agaacttaaa tgccgaagtg
     2461 gcattttgtt cactaatgaa gtccacactt ggacagagca atacaatcc caggctgact
     2521 ccccaaaaag actgtcagca gcccattggaa aggctgggaa ggaggcggtg tggaaatcc
     2581 gatcagccac gcgttttgc aacatcatgt ggaagcagat atcaaataa ttgaaccaca
     2641 ttttacttgc gaaatgacatg aaatttcacag tgggtgttagg agatgccaac ggaattttgg
     2701 cccaaaggaaa aaaaatgtt gggccacaac ccatgaaaca caaataactca tggaaaagct
     2761 ggggaaaagc taaaatcata ggacgacaca tacaaaatac cacccattt atcgacggcc
```

2821 gagacaccccc agaatgtcct gatggccaaa gagcatggaa catttggaa gttgaggact  
2881 atgggtttgg agtttcacg acaaacatat ggctgaaatt gcgtgactcc tacacccaaa  
2941 tggtgacca ccggctaata tcagctgcca tcaaggacag caaggcagtc catgctgaca  
3001 tgggtactg gatagaaagt gaaaagaacg aaacctggaa gttggcgaga gcctccttca  
3061 tagaagtcaa aacatgcacc tggccgaaat ctcacactt atggagcaat ggagtttgg  
3121 aaagtgaaat gataatccca aagatatacg gaggaccaat atctcagcac aactacagac  
3181 cagggtattt cacacaaaca gcagggccat ggcacctagg taagttggaa ctggattttg  
3241 acttgtgtga aggacccaca gtttgggtgg atgaacattt tgaaatcga ggtccatctc  
3301 tcagaaccac aacagtccaca ggaaagataa tccatgaatg gtgtgcaga tcctgcacgc  
3361 tacccccctt acgtttcaga ggagaagacg ggtgtggta tggcatggaa atcagaccag  
3421 tgaaggagaa ggaggagaat cttagtttaggt caatggctc tgcagggtca ggagaagtgg  
3481 acagttttc attaggaata ctatgcgtat caataatgtat tgaagaagtg atgagatcca  
3541 gatggagtag aaagatgctg atgactggaa cactggctgt cttccctt cttataatgg  
3601 gacaactgac atggaatgtat ctgatttaggt tatgcatcat ggtcgagct aacgcttcag  
3661 acaggatggg gatggaaaca acgtacctag cttgtatggc tacttcaaa atgagaccaa  
3721 tggcgctgt agggctatta ttccgcagac taacatccag agaagttctt ctcctaacga  
3781 ttggattaag cctggggca tccgtggagc taccaaattt cttggaggag cttagggatg  
3841 gacttgcataat gggtatcatg atgttaaaat tggactgtatg atttcagcca caccagttat  
3901 ggaccacccattt atgtctctg acatttgcataa aaacaactt ctcattggat tatgcatgg  
3961 aaacaacggc tatggcactg tctatcgat ctctcttcc tttatgcctg tctacgaccc  
4021 cccaaaaaaac aacatggctt ccgggtctgt taggatctt tggatgcaaa ccattaacca  
4081 tggccttat aacagaaaaat aaaatctggg gaaggaaaaag ttggccctc aatgaaggaa  
4141 ttatggctat tggaaatagtc agcattctac taagctcaat cctaaaaat gatgtgccgt  
4201 tggccgggccc attaataagct ggagggcatgc taatagcatg ttatgtcata tcccgtagct  
4261 cagccgattt atcattggag aaagcggctg aagtatctt ggaacaagaa gcagaacact  
4321 cccgtgcctc acacagcata ttagtagagg tccaaatgatg tggactatg aaaataaaag  
4381 atgaagagag ggtatgcacca ctaccatac tccttaaagc aacttgcgt gcaatctcag  
4441 gagttgtaccc aatgtcaata ccagcaactc tttttgtgt gttttttgg cagaaaaaga  
4501 aacagagatc aggagtgtt a tggacacac ccagccctc ggaagtggaa agagcgttc  
4561 ttgataatgg catctataga atcttgcataa gaggattgtt gggcagggtcc caagtaggag  
4621 tggagttt ccaagacggc gtgttccaca caatgtggca cggttaccagg ggagctgtcc  
4681 ttatgtacca agggaaagaga ctgaaaccaa gctggccag tggaaaaag gacttgatct  
4741 catatggagg aggttggagg ttccaaggat catggacac gggagaagaa gtgcaggtaa  
4801 tagctgtga accaggaaaa aaccccaaaa atgtacagac aacggccggc acctttaaga  
4861 ctctgtaccc cgaagttgga gccatagctc tagatttcaa accccgcaca tctggatctc  
4921 ccatctgtaccc cagagaggaa aaaatgtgg gtcgtatgg aatggagtg gtgacaacaa  
4981 gtggaaaccta cgtcgtgccc attgcccac agttaacatc acagaaggg cctctaccag  
5041 agatttgagga cgaggatattt aagaaaaagaa acttaacaat aatggacctg caccaggat  
5101 cagggaaaaac aagaagatat ctccagcc tagtccgtga ggcataaaaa aggaaactgc  
5161 gtacgttaat cctggctccc acaagagtt tgcctctga aatggcagag gcaactcaagg  
5221 gaatgccaat aagatatacg acaacagcag tgaagagtga acacacagga agggagatag  
5281 ttgacccat gtgccacgt acttttacca tgcgtcttt atccccagtg agagttccca  
5341 attacaacat gatcattatg gatgaaggac attttaccga tccagctagc atagcgccca  
5401 gagggtacat ctaaaccga gtgggtatgg gtgaagcagc tgcgtatctt atgacagcc  
5461 ctccccccagg atcgggtggag gccttccac agagcaatgc agttatccaa gatgagggaaa  
5521 gagacattcc tgagagatca tggaaactcg gtcacgt gatcaactgac tttccaggta  
5581 aaacagtctg gttgttcca agcattaaat cagggaaatga cattgccaac tggtaagaa  
5641 agaacggaaa acgggtaaatc caattgagca gaaaaaccc ttgacactgag taccagaaaa  
5701 caaaaaacaa tgactggac tatgttgcataa caacagacat ttctgaaatg gggcaaaatt  
5761 tcggggccga cagggtataa gacccaaaggc ggtgctgaa accggtataa ctaaaagatg  
5821 gtcctggacgc tggttccatc gccggaccga tgccagtgac tgcggccag tgcgtccaga  
5881 ggagaggaag aatttggaaagg aacccaaaaca aggaaggtga tcagttatgtt tataatggac  
5941 accctttaaa taatgtatg gatcacgctc atggacaga agcaaaaatg ctccttgaca  
6001 atataaacac accagaaggg atcatccag cccttttga gccagagaga gaaaagagtg  
6061 cagcaataga cggggagttac agactgcggg gagaagcaag gaaaacgttc gtggagctca  
6121 tgagaagagg agatctacca gtttgcataat cttacaaatg agcctcagaa gttttccagt  
6181 actccgacag aaggtggtgc tttgtatgggg aaaggaacaa ccaggttgc gaggagaaca  
6241 tggacgttggaa gatctggaca aaggaaggag aaaggaagaa attgcgaccc ctgtgggttgg

```
6301 acgccagaac atactctgat ccattggccc tgcgcgagtt taaagagttc gcagcaggaa
6361 gaagaagtgt ctcaggtgac ctgatattgg aaataggaa acttccacaa catttgacgt
6421 taagagccca gaatgctctg gacaaccttgg tcatgttgca caattccgaa caaggaggaa
6481 aagcctacag acatgccatg gaggaaactac cagacaccat agaaacattg atgctactag
6541 ctttgatacg tgtgttgact ggtggagtga cgctgttctt cctatcagga aaaggcctag
6601 ggaaaacatc cattggcttgc ctctgtgtga cggcctcaag cgcaactgtta tggatggcca
6661 gtgtggagcc ccattggata gcccctcca tcatactaga gttcttttg atggtgctgc
6721 tcattccaga gccagacaga cagcgcactc cacaggacaa ccagctagca tatgtggtga
6781 taggtttgtt attcatgata ctgacagtgg cagccaatga gatgggatta ttggaaacca
6841 caaagaaaaga cctggggatt ggctatgttag cccggaaaaa ccaccaacat gccacaatgc
6901 tggacgtaga cctacaccca gcttcagcct ggaccctcta tgcaagtagcc acaacagtca
6961 tcactcccat gatgagacac acaattgaaa atacaacggc aaacattcc ctgaccgcca
7021 ttgcaaatac ggcagctata ttgatgggac ttgacaaggg atggccaata tcgaagatgg
7081 acataggagt tccacttctc gccttaggggt gctattccca ggtgaaccca ttgacactga
7141 cagccgcgtt gttgatgtta gtggctcatt atgccataat tggaccagga ctgcaagcaa
7201 aggccactag agaagccaa aaaaggacag cagccggaaat aataaaaat ccaaccgtag
7261 acgggattgt tgcaatagac ttggatcctg tggtttatga tgcaaaattt gaaaaacaac
7321 taggcaaat aatgttactg atactttgtt catcacagat cctttgatg cggaccacat
7381 gggccttgc tgaatccatc acactggctt ctggaccctt gaccactctc tggagggat
7441 ctccagggaa attctggaaat accacaatag cagtgccat ggcaaattt ttcaggggaa
7501 gttatctagc aggagcagggt ctggcttctt cattgtatgaa atcttttagga ggaggttagga
7561 gaggcacggg agctcaaggg gaaacactgg gagagaaatg gaaaagacag ttgaaccaac
7621 tgagcaagtc agaattcaac acctacaaaaa ggagtggat tatggaggtg gacagatccg
7681 aagccaaaga gggactgaaa agaggagaaa caaccaaaca tgcaagtgca agaggaacag
7741 cccaaacttagt gttgttgc tggaggaaacc tcgtgaaacc agaaggaaaa gtcataagacc
7801 tcgggttgc aagaggtggc tggcatatt attgtgtgg gctgaagaaa gttactgaag
7861 tgaagggata cacaaggaa ggacctggac atgaggaacc tatcccaatg gcgacctatg
7921 gatggaaact agtaaaacta cactctggaa aggatgtatt tttatgcca cctgagaaat
7981 gtgacactct tctgtgtat attgtgttagt cctctccgaa tccaaactata gaagaaggaa
8041 gaacgttacg tggctaaaaa atggtggaaac catggctcag agggaaaccaa ttctgcataa
8101 aaatcctaaa tccttacatg ccaagtgtgg tagaaactct ggagcgaatg caaagaaaac
8161 atggagggat gctagtgcga aacccactct caagaaattt taccatgaa atgtattgg
8221 tttcatgtgg aacagggaaac attgtgtcgg cagtgaacat gacatccaga atgttactga
8281 accgattcac aatggctcac aggaagccaa catatgaaag agacgtggac ttaggcgcgt
8341 gaacaagaca tggcgtggc gaaccagagg tagccacact agatatcatt ggcagagga
8401 tagaaaaat aaaaatgaa cacaagtcaa catggcatta tgatgaggac aatccatata
8461 aaacatggc ctatcatgaa tcatatgagg tcaagccatc aggatcagcc tcatctatgg
8521 tgaatggagt ggtgagattt ctcacgaaac catggatgt catccccatg gtcacacaaa
8581 tagctatgac tgataccaca cccttggac aacagagagt gttaaagag aaagttgaca
8641 cgcgcacacc aagagaaaaa cgaggcacaac cacagattat ggaggtgaca gccaagtgt
8701 tatgggttt cctttccaga aacaaaaaaac ccagaatctg cacaagagag gagttcacaa
8761 gaaaggttag gtcaaacgcg gcaataggag cagtgttgc tggatgaaaac caatggaaact
8821 cagccaaaga agcagtggaa gacgaaaggt tttggatct tggcacaga gagagggagc
8881 ttcataaaaca gggaaaatgt ggcacgtgtc tctacaacat gatgggaaag agagagaaaa
8941 aattaggaga gtttggaaag gcaaaaggaa gtcgtcaat atgtacatg tggctggag
9001 cacgtttctt gggatgtcgaa gcccgggtt ttatgtatga agatcactgg ttttagtagag
9061 agaattcaact cagtggtgtc gaaggagaag gactgcacaa acttggatac atactcagag
9121 acatatcaaa gattccgggg gggaaatatgt atgcagatga tacagccgaa tggacacaaa
9181 gagtaacaga gggatgaccc cagaatgagg ctaaaatcac tgacatcatg gggctgaac
9241 atgctctatt ggctacgtca atttttaagc tgacttatac aaacaagggt gtgagggtgc
9301 aaagaccacg aaaaatgga accgtgtatgg atgttatatc cagacgtgtc cagagagggaa
9361 gtggacagggt cggaaactt ggtttaaata ctgttccaaat tggatggc caactaataa
9421 gacaaatgga gtctgaggga atcttttac ccagcgaatt gggaaacccccc aacctagctg
9481 agaggcgtt tggatgttca gaaaacatg gcccggaaag gctgaaacga atggcaatca
9541 gcccggatgtt gttgttgc tggatgttca gggaaatgttccaaat tggatggc caactaataa
9601 tgaatgacat gggaaaatgta agggaaagaca taccgcgtt gggaaacccccc aacctagctg
9661 atgattggca gcaagtgcct ttttggatgttca accatccca ccaactgatc atgaaggatg
9721 ggagggaaat agtggtgccca tgccgcaacc aagatgaact tggatggcagg gctagagtat
```

```

9781 cacaaggcgc cgatggagc ctgagagaaa ctgttgcct aggcaagtca tatgcacaaa
9841 tgtggcagct gatgtacttc cacaggagag acctgagact agccgctaac gctatcttt
9901 cagccgtccc agttgattgg gtcccaacca gccgcacaac ctgttcaatc catgcccacc
9961 accaatggat gacaacagaa gacatgttat cagtgtggaa tagggtttgg atagacgaaa
10021 acccatggat ggagaacaaa actcatgtat ccagtggga agaagttcca taccttagaa
10081 aaagggaaga tcaatggtgt ggatccctga taggcttgac a诶cgagggcc acctgggcca
10141 ccaacataca agtagccata aaccaagtga gaaggctcat cggaatgag aattattnag
10201 attacatgac atcaatgaag agattcaaga atgagagtga ttccgaagga gcactctggt
10261 aagtcaacac actcatgaaa taaaggaaaa tagaagatca aacaaagtaa gaagtcaagc
10321 cagattaagc catagcacgg aaagagctat gctgcctgtg agccccgtcc aaggacgtaa
10381 aatgaagtcg ggcgaaagc cacggattga gcaagccgtg ctgcctgtgg ctccatcg
10441 gggatgttagc tc
//
```

## 1.2.6 parallel\_blast

Parallel blast is a wrapper script around the blast commands as well as diamond. It utilizes GNU Parallel to run the commands in parallel by splitting up the input fasta files and distributes them across multiple subprocesses. If it detects that it is running inside of a PBS or SGE job it will run the job on multiple hosts that may be allocated to the job.

parallel\_blast requires that you have gnu parallel installed and in your environments PATH as well as diamond and/or blastn/blastx/blastp.

- diamond
- blast
- GNU parallel

### Usage

You can get all the arguments that can be supplied via the following

```
$> parallel_blast --help
```

### Examples

For the examples below assume you have an input fasta in the current directory called `input.fasta`

#### Running blastn

```
$> parallel_blast input.fasta output.blast --ninst 4 --db /path/to/nt \
--blast_exe blastn --task megablast --blast_options "--evalue 0.01"
[cmd] /path/to/parallel -u --pipe --block 10 --recstart > --sshlogin 4/: /path/to/blastn -task megablast
```

Notice how we had to quote the additional `--blast_options`

**Running diamond** Diamond v0.7.9 is the version that was tested with parallel\_blast. As diamond is still in development the options may change in future versions and parallel\_blast may not run them correctly. Please submit a new issue if you find any issues.

```
$> parallel_blast input.fasta out.blast --ninst 4 --db /path/to/diamondnr \
--blast_exe diamond --task blastx --blast_options "--tmpdir dtmp"
[cmd] /path/to/parallel -u --pipe --block 10 --recstart > --cat --sshlogin 1/: /path/to/diamond blast
```

Notice how even though we specified `--ninst 4` that `--sshlogin 1/` was used and `--threads 4` was set instead.

**Note** In recent versions of diamond, diamond outputs a daa binary file instead of a tab separated file. `parallel_blast` automatically converts the diamond output from daa to tab format for you but leaves the daa file behind(Same name as the output file you specify, but with the extension .daa)

**Command that is run** You will notice in the examples above that when you run `parallel_blast` that it outputs the command that it is running in case you want to copy/paste it and run it yourself sometime.

You might notice that the command does not include all the quoted arguments such as the `--recstart` argument which should be `--recstart ">"` as well as the `--outfmt` which should be quoted as `--outfmt "6 ..."`. If you intend on rerunning the command you will have to add the quotes manually.

### Running inside of a PBS or SGE Job

`parallel_blast` is able to detect if it is running inside of a PBS or SGE job by looking to see if `PBS_NODEFILE` or `PE_HOSTFILE` is set in the environment's variables.

If it finds either of them it will run the job by supplying `--sshlogin` for each host it finds in the file.

`PBS_NODEFILE` and `PE_HOSTFILE` have different syntax so `parallel_blast` first builds a CPU,NODENAME list from them.

#### PBS\_NODEFILE

This file is parsed and counts how many of each unique host is listed such that the following `PBS_NODEFILE`:

```
node1.localhost
node2.localhost
node2.localhost
node3.localhost
node3.localhost
node3.localhost
```

would run 1 instance on node1.localhost, 2 instances on node2.localhost and 3 instances on node3.localhost

#### PE\_HOSTFILE

This file is almost in the exact syntax that `parallel_blast` uses so it is almost a 1-to-1 mapping.

#### Diamond and multiple hosts

Since diamond utilizes threads much more efficiently than blast, for each unique host in a job only 1 instance is launched but the `-p` option is set to the number of CPUS for each host listed in the `PE_HOSTFILE` or `PBS_NODEFILE`

### 1.2.7 degen\_regions

Finds all degenerate bases in a given fasta input file that may contain multiple sequences and reports their position as well as the annotated gene name that contains them.

The fasta file must be previously aligned to the query sequence. That is, if you are using a genbank annotation file or having the script download it for you, you should have aligned all your input sequences to that sequence.

The annotation is retrieved via supplied genbank accession, genbank file path or gene tab/csv file.

## Usage

You can view the usage of degen\_regions via:

```
degen_regions --help
```

## Using Genbank Files

If you already have downloaded the genbank annotation file(typically the extension is .gb) you can use the *--gb-file* argument

The following will use the test input fasta file as well as the test input genbank file to find all degenerate bases and will put the output in a tab separated file called output.tsv

```
degen_regions -i tests/Den4_MAAPS_TestData16.fasta -o output.tsv --gb-file tests/testinput/sequence.g
```

## Fetching Genbank Files Automatically

If you want the script to automatically fetch the Genbank annotation file from the internet you can use the *--gb-id* option and specify an accession number.

```
degen_regions -i tests/Den4_MAAPS_TestData16.fasta -o output.tsv --gb-id KJ189367
```

## Using tab/csv file of gene annotation info

If you have a tab/csv file of gene annotations you can supply that using the *--tab-file* argument

You can read more about the format of the tab/csv annotation file in the [degen](#) docs

```
degen_regions -i tests/Den4_MAAPS_TestData16.fasta -o output.tsv --gb-file tests/testinput/sequence.g
```

## Manually specify CDS

You can use the *--cds* argument to set the coding region. This argument should be comma separated such as start,stop. Specifying this argument will override any other cds found in the tab file, genbank file or fetched genbank file.

The following would mark all locations as NON-CODING as you are specifying that only position 1 is coding

```
degen_regions -i tests/Den4_MAAPS_TestData16.fasta -o output.tsv --gb-file tests/testinput/sequence.g
```

## Output

The output is a simple tab separated file

seq_id	nt Position	aa position	nt composition	aa composition
721	991	331	WCA	S/T
721	1307	436	AYA	I/T
721	1826	609	AYA	I/T
721	1865	622	GRA	E/G
721	7766	2589	ARA	K/R
2055_Den4/AY618992_1/Thailand/2001/Den4_1	1927	643	RAC	D/N
2055_Den4/AY618992_1/Thailand/2001/Den4_1	2833	945	YCG	P/S
2055_Den4/AY618992_1/Thailand/2001/Den4_1	3565	1189	YAT	H/Y
2055_Den4/AY618992_1/Thailand/2001/Den4_1	6271	2091	RAA	E/K
2055_Den4/AY618992_1/Thailand/2001/Den4_1	8656	2886	YAT	H/Y
2055_Den4/AY618992_1/Thailand/2001/Den4_1	8998	3000	YAG	*/*Q
2055_Den4/AY618992_1/Thailand/2001/Den4_1	9811	3271	YCC	P/S
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10542	3515	AGN	NON-CODING
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10543	3515	NNN	NON-CODING
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10541	3514	NNN	NON-CODING
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10539	3514	NNN	NON-CODING
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10546	3516	NNN	NON-CODING
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10544	3515	NNN	NON-CODING
2055_Den4/AY618992_1/Thailand/2001/Den4_1	10542	3515	NNN	NON-CODING
1942_Den4/AY618992_1/Thailand/2001/Den4_1	4540	1514	RTA	I/V
1942_Den4/AY618992_1/Thailand/2001/Den4_1	10177	3393	MCA	P/T
1942_Den4/AY618992_1/Thailand/2001/Den4_1	10546	3516	NNN	NON-CODING
1942_Den4/AY618992_1/Thailand/2001/Den4_1	10544	3515	NNN	NON-CODING
1942_Den4/AY618992_1/Thailand/2001/Den4_1	10542	3515	NNN	NON-CODING
1942_Den4/AY618992_1/Thailand/2001/Den4_1	10542	3515	NNN	NON-CODING
1875_Den4/AY618992_1/Thailand/2001/Den4_1	1514	505	AYG	M/T
1875_Den4/AY618992_1/Thailand/2001/Den4_1	3056	1019	ARA	K/R
1875_Den4/AY618992_1/Thailand/2001/Den4_1	3058	1020	KCA	A/S
1875_Den4/AY618992_1/Thailand/2001/Den4_1	3073	1025	WTT	F/I
1875_Den4/AY618992_1/Thailand/2001/Den4_1	3491	1164	AYC	I/T
1875_Den4/AY618992_1/Thailand/2001/Den4_1	3895	1299	RTG	M/V
1875_Den4/AY618992_1/Thailand/2001/Den4_1	7445	2482	GYA	A/V
948_Den4/AY618992_1/Thailand/2001/Den4_1	2819	940	ARC	N/S
871_Den4/AY618992_1/Thailand/2001/Den4_1	2947	983	RCC	A/T
871_Den4/AY618992_1/Thailand/2001/Den4_1	3058	1020	KCA	A/S
871_Den4/AY618992_1/Thailand/2001/Den4_1	3073	1025	WTT	F/I
871_Den4/AY618992_1/Thailand/2001/Den4_1	3116	1039	GYG	A/V
871_Den4/AY618992_1/Thailand/2001/Den4_1	3181	1061	RTW	I/V
871_Den4/AY618992_1/Thailand/2001/Den4_1	3179	1060	RTW	I/V
871_Den4/AY618992_1/Thailand/2001/Den4_1	3338	1113	ART	N/S
871_Den4/AY618992_1/Thailand/2001/Den4_1	3362	1121	ARA	K/R
871_Den4/AY618992_1/Thailand/2001/Den4_1	3373	1125	WCR	S/T
871_Den4/AY618992_1/Thailand/2001/Den4_1	3371	1124	WCR	S/T
871_Den4/AY618992_1/Thailand/2001/Den4_1	4314	1439	ATV	I/M
871_Den4/AY618992_1/Thailand/2001/Den4_1	7045	2349	WCC	S/T
871_Den4/AY618992_1/Thailand/2001/Den4_1	10536	3513	GAW	NON-CODING
871_Den4/AY618992_1/Thailand/2001/Den4_1	10537	3513	YCA	NON-CODING
947_Den4/AY618992_1/Thailand/2001/Den4_1	2971	991	YTY	F/L
947_Den4/AY618992_1/Thailand/2001/Den4_1	2969	990	YTY	F/L
947_Den4/AY618992_1/Thailand/2001/Den4_1	6763	2255	YTT	F/L
1793_Den4/AY618992_1/Thailand/2001/Den4_1	223	75	MAG	K/Q
1793_Den4/AY618992_1/Thailand/2001/Den4_1	556	186	RCC	A/T
1793_Den4/AY618992_1/Thailand/2001/Den4_1	586	196	RGT	G/S
1793_Den4/AY618992_1/Thailand/2001/Den4_1	613	205	YCA	P/S
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2875	959	YCG	P/S
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2943	982	AAN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2944	982	NNG	GAPFOUND

1793_Den4/AY618992_1/Thailand/2001/Den4_1	2942	981	NNG	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2976	993	ATN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2977	993	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2975	992	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2973	992	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2980	994	NTG	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2987	996	ANN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2986	996	ANN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2989	997	NGT	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2996	999	TNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2995	999	TNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3001	1001	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2999	1000	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	2997	1000	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3004	1002	NCC	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3073	1025	NTT	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3086	1029	ARC	N/S
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3095	1032	CNG	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3116	1039	GNG	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3144	1049	GAN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3159	1054	GAN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3160	1054	NNC	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3158	1053	NNC	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3206	1069	GNC	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3235	1079	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3233	1078	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3231	1078	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3238	1080	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3236	1079	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3234	1079	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3241	1081	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3239	1080	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3237	1080	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3244	1082	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3242	1081	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3240	1081	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3247	1083	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3245	1082	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3243	1082	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3250	1084	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3248	1083	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3246	1083	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3253	1085	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3251	1084	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3249	1084	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3256	1086	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3254	1085	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3252	1085	NNN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3316	1106	NGG	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3337	1113	NAT	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3341	1114	GNA	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3408	1137	ATN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3412	1138	NTG	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3493	1165	MCC	P/T
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3509	1170	ANT	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	3837	1280	TTN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	6185	2062	ARG	K/R
1793_Den4/AY618992_1/Thailand/2001/Den4_1	6187	2063	RAR	E/K

1793_Den4/AY618992_1/Thailand/2001/Den4_1	6185	2062	RAR	E/K
1793_Den4/AY618992_1/Thailand/2001/Den4_1	6614	2205	TYT	F/S
1793_Den4/AY618992_1/Thailand/2001/Den4_1	6650	2217	ARA	K/R
1793_Den4/AY618992_1/Thailand/2001/Den4_1	8630	2877	ART	N/S
1793_Den4/AY618992_1/Thailand/2001/Den4_1	8844	2949	AAN	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	9938	3313	AYT	I/T
1793_Den4/AY618992_1/Thailand/2001/Den4_1	9941	3314	GRC	D/G
1793_Den4/AY618992_1/Thailand/2001/Den4_1	10015	3339	RTT	I/V
1793_Den4/AY618992_1/Thailand/2001/Den4_1	10087	3363	NGR	GAPFOUND
1793_Den4/AY618992_1/Thailand/2001/Den4_1	10085	3362	NGR	GAPFOUND
1901_Den4/AY618992_1/Thailand/2001/Den4_1	15	6	AAN	NON-CODING
1901_Den4/AY618992_1/Thailand/2001/Den4_1	111	38	TTN	GAPFOUND
1901_Den4/AY618992_1/Thailand/2001/Den4_1	2279	760	GYT	A/V
1901_Den4/AY618992_1/Thailand/2001/Den4_1	8798	2933	ARA	K/R
1901_Den4/AY618992_1/Thailand/2001/Den4_1	10195	3399	RAG	E/K
1901_Den4/AY618992_1/Thailand/2001/Den4_1	10366	3456	RGG	NON-CODING
1934_Den4/AY618992_1/Thailand/2001/Den4_1	15	6	AAN	NON-CODING
1934_Den4/AY618992_1/Thailand/2001/Den4_1	111	38	TTN	GAPFOUND
1934_Den4/AY618992_1/Thailand/2001/Den4_1	998	333	GMT	A/D
1934_Den4/AY618992_1/Thailand/2001/Den4_1	4515	1506	TTM	F/L
1934_Den4/AY618992_1/Thailand/2001/Den4_1	8798	2933	ARA	K/R

## 1.2.8 plot\_muts

### Usage

You can view the usage of degen\_regions via:

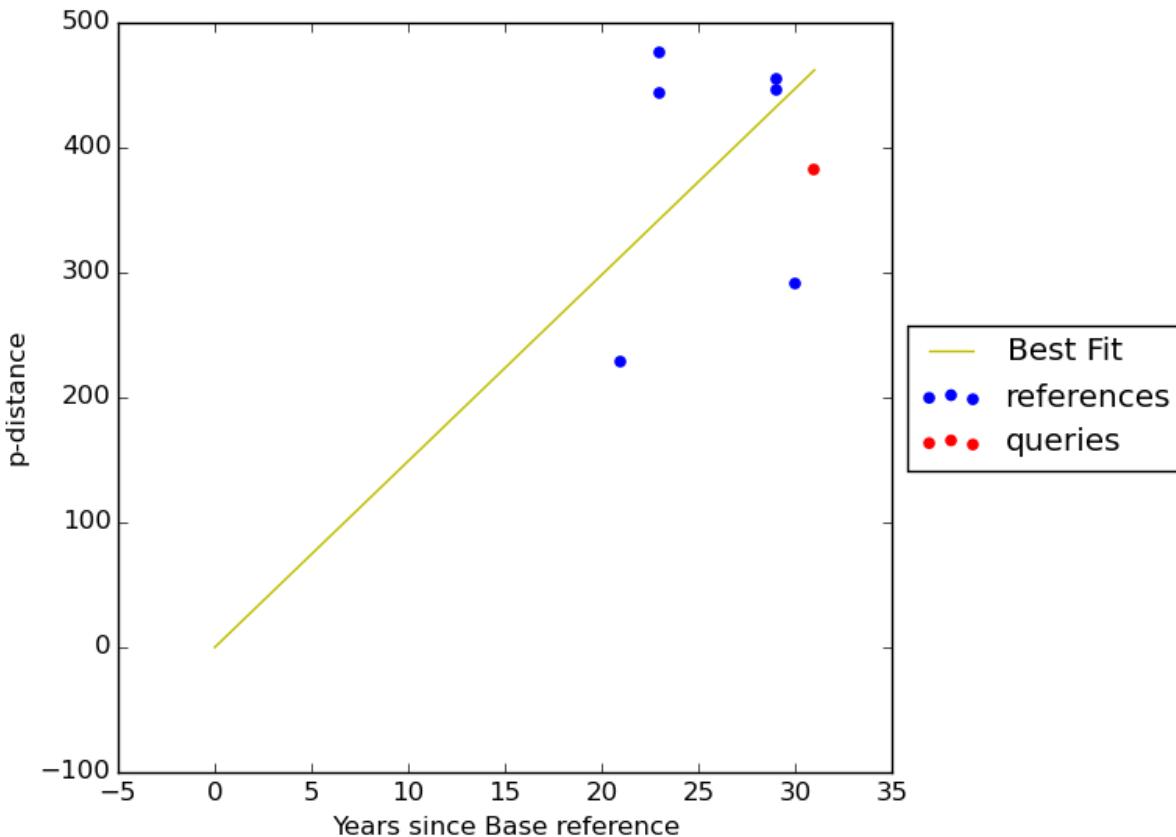
```
plot_muts --help
```

### Example

```
plot_muts --refs tests/testinput/refs.fas --query tests/testinput/query.fas --out plot.png
```

The `--out` option is optional. If it is not provided, the plot will pop up on the user's screen automatically. If this does not work, try saving the image using `--out` instead.

## Example Output



## Input File Requirements

The input must be fasta format. Both the query and ref files can have any number of sequences.

The year should be the last part of the ID, preceded by a quadruple underscore. e.g.:

```
>some|info|blah_blah____2001_09_2010  
>some____1995  
>some____09/09/2012
```

If the ID uses '/' rather than underscore, plot\_muts currently accepts the year as the *fourth* field. e.g.:

```
>some/info/blah/1995  
>some/info/blah/1995/more/info
```

### 1.2.9 fasta

fasta is a very simple script to help mangle fasta files. Currently it only supports the ability to convert fasta files that have sequences that span multiple new lines into single lines.

Later on, it may be expanded more to include even more useful fasta features.

## Usage

```
fasta --help
```

## Examples

The following examples all use the test fasta file found under `tests/testinput/col.fasta`

```
>sequence1 some description !@#$%^&* ()_+-=[ ]{}.,></?';:"  
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA  
TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTT  
GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG  
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC  
>sequence2  
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA  
TTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTT  
GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG  
CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
```

**Read fasta input from standard input** The following could output the fasta sequences as one line to your terminal(stdout) but reading from the pipe. This is useful if you want to use it in a pipeline.

```
$> cat tests/testinput/col.fasta | fasta -
```

**Read fasta input from file** The following could output the fasta sequences as one line to your terminal(stdout) as well, but reading straight from the file.

```
$> fasta tests/testinput/col.fasta
```

**Simple shell pipeline using fasta** The following is a simple shell pipeline to count how many A's there are in the sequence lines. There should be 160 since `col.fasta` is 80 characters per line and only the first line of each sequence has A and there are 2 sequences.

```
$> fasta tests/testinput/col.fasta | grep -v '>' | grep -Eo '[Aa]' | wc -l  
160
```

## 1.3 AMOS

AMOS is a file format that is similar to any assembly file format such as ACE or SAM. It contains information about each read that is used to assemble each contig.

The format is broken into different message blocks. For the Ray assembler, it produces an AMOS file that is broken into 3 types of message blocks

- RED

```
{RED  
iid:\d+  
eid:\d+  
seq:  
[ATGC]+  
.
```

```
qlt:  
[A-Z] +  
}
```

**iid** Integer identifier

**eid** Same as iid?

**seq** Sequence data

**qlt** Should be quality, but is only a series of D's from Ray assembler

- TLE

```
{TLE  
src:\d+  
off:\d+  
clr:\d+, \d+  
}
```

**src** RED iid that was used

**off** One would think offset, but unsure what it actually means

**clr** Not sure what this is either

- CTG

```
{CTG  
iid:\d+  
eid:\w+  
com:  
. * $  
. .  
seq:  
[ATGC] +  
. .  
qlt:  
[A-Z] +  
. .  
{TLE  
...  
}  
}
```

**iid** integer id of contig

**eid** contig name

**com** Communication software that generated this contig

**seq** Contig sequence data

**qlt** Supposed to be contig quality data, but for Ray it only produces D's

**TLE** 0 or more TLE blocks that represent RED sequences that compose the contig

### 1.3.1 Parsing

bio\_bits contains an interface to parse a given file handle that has been opened on an AMOS file.

To read in the AMOS file you simply do the following

```
from bio_bits import amos
a = None
with open('AMOS.afg') as fh:
    a = amos.AMOS(fh)
```

## CTG

To get information about the contigs(CTG) you can access the `.ctgs` attribute. The contigs are indexed based on their iid so to get the sequence of contig iid 1 you would do the following:

```
ctg = a.ctgs[1]
seq = ctg.seq
```

To retrieve all the reads(RED) that belong to a specific contig:

```
reads = []
for tle in ctg.tlelist:
    reads.append(a.reds[tle.src])
```

## RED

To get information about the reads(RED) you can access the `.reds` attribute. The reds are indexed based on their iid so to get the sequence of red iid 1 you would do the following:

```
red = a.reds[1]
seq = red.seq
```

If you want to convert a RED entry into anything you can use the `.format` method. The `.format` method allows you to utilize any of the properties of a RED object such as `.iid`, `.eid`, `.seq`, `.qlt`. You can see in the examples below how to do this.

### 1.3.2 Examples

Here is an example of how to convert all RED blocks into a single fastq file

```
from bio_bits import amos

# Fastq format string
fastq_fmt = '@{iid}\n{seq}\n+\n{qlt}'

with open('amos.fastq', 'w') as fh_out:
    with open('AMOS.afg') as fh_in:
        for iid, red in amos.AMOS(fh_in).reds.items():
            fq = red.format(fastq_fmt)
            fh_out.write(fq + '\n')
```

## 1.4 CHANGELOG

### 1.4.1 Version 1.3.0

- Added fasta script that removes newlines from fasta sequences

## 1.4.2 Version 1.2.1

- Fixed some python3 and python2.6 incompatability issues
- Fixed some old bio\_pieces references
- Added some simple tests for plot\_muts

## 1.4.3 Version 1.2.0

- Renamed project to bio\_bits to fix naming issue with other project
- GPL License added
- degen\_regions script added
- parallel\_blast added
- plot\_muts script added

## 1.4.4 Version 1.1.0

- Renamed parse\_contigs to group\_references to better name functionality
- group\_references now supports bam files

## 1.4.5 Version 1.0.0

- Version bump. Starting here we will employ semantic versioning
- Added version script to get version from project

## 1.4.6 Version 0.1.0

- Started project over to setup for Continuous Integration testing
- Added rename\_fasta that can rename fasta sequence identifiers based on a input rename file
- Added travis, coveralls, readthedocs
- Added amos file parser that is specific to Ray assembler amos format
- Added format functionality for amos classes such that it is easy to convert to different formats
- Added amos2fastq to pull sequences out of AMOS files organized by their contigs.
- Added vcfcat.py, a commandline app for filtering and comparing vcf files.
- Completed documentation for vcfcat
- Added beast\_checkpoint script and documentation
- Added beast\_wrapper script that prints estimated time column in beast output
- Added beast\_est\_time script that allows you to easily get estimated time left from already running beast run

## 1.5 TODO

---

### Todo

Could be possible to get beast\_checkpoint to check for that scenario and use the last tree state that matches the last log state

---

(The original entry is located in /home/docs/checkouts/readthedocs.org/user\_builds/bio-bits/checkouts/dev/docs/scripts/beast\_checkpoint.rst, line 52.)

## **Indices and tables**

---

- genindex
- modindex
- search