
Scythe Documentation

Release 0.1a1

Janina Mass

June 14, 2014

CONTENTS

Scythe – Selection of Conserved Transcripts by Homology Evaluation

This documentation is under construction!

Contents:

INTRODUCTION

Scythe picks best matching transcripts for one-to-one orthologous genes from two or more species.

The goal is to provide the best, i.e. most homologous set of sequences, for a subsequent multiple sequence alignment in order to minimize sources for misalignment in an automated fashion. Scythe will perform pairwise global alignments using the Needleman-Wunsch algorithm [NE1970] as implemented in *needleall* as part of the EMBOSS package [EMB2000].

Please see the [Tutorial](#) on how to use Scythe.

Important! *needleall* appears to be broken in (the ubuntu package for) EMBOSS 6.6.0.0; this is working with *needleall* from EMBOSS 6.4.0.0.

1.1 Algorithms and Application

Scythe performs pairwise global alignments [NEED1970] to measure the similarity between transcripts. Transcripts are only compared between species.

There are three main strategies implemented in Scythe to deal with the sequences after scoring: A single-linkage approach either starting with reference transcripts as seed (*sl_ref*) and adding best matching sequences from non-reference species or starting out with the best matching transcript pair for a gene (*sl_glob*). A *reference species* is defined as a species that has only one transcript (*reference transcript*) for a gene. Every gene is processed individually, a *reference species* is only local to a gene and is automatically derived from the input data. Alternatively, the maximum-sum (*mx_sum*) approach calculates the score for all transcript pairings between the species and return a maximum-scoring set. Please note that this approach might not be feasible for large data sets.

1.1.1 mx_sum

mx_sum (maximum sum) returns an optimal solution for the problem of score maximization. There are, however, scenarios where this would not represent the desired result: In cases of single-transcript outliers, their pairwise score to all other sequences is always taken into account and may favor sequences that are less dissimilar to the outlier.

1.1.2 sl_ref

sl_ref (single linkage reference) finds similar gene models given a reference species. This single-linkage approach starts with a reference transcript as seed and adds best matching sequences from non-reference species. A *reference species* is defined as a species that has only one transcript (*reference transcript*) for a gene. Every gene is processed individually, a *reference species* is only local to a gene and is derived from the input data.

1.1.3 `sl_glob`

`sl_glob` (single linkage global) is similar to the *sl_ref* approach, but starts out with the best matching pair. This should provide a good starting point in absence of a reference gene model and may be able circumvent outliers.

TUTORIAL

2.1 Installation

2.1.1 Requirements

- Python 3

python modules

- configparser
- mysql-connector-python

```
pip install --allow-external mysql-connector-python mysql-connector-python
pip install configparser
```

In case you want to work using a virtual environment switch to that environment first (see *Install Scythe in a virtual environment*).

external software

- needleall <http://emboss.sourceforge.net/>

Install emboss on debian-based systems:

```
sudo apt-get install emboss
```

Important: needleall from (the ubuntu package for) EMBOSS 6.6.0.0 causes problems. This program was tested and works with needleall from EMBOSS 6.4.0.0.

2.1.2 Install Scythe in a virtual environment

Install virtualenv via pip:

```
pip install virtualenv
```

Create a virtual environment that uses Python 3 and activate:

```
virtualenv -p /usr/bin/python3 python3env
source python3env/bin/activate
```

Finally, install the scythe package:

```
(python3env)you@host:~$ pip install scythe
```

2.1.3 Without virtual environment

Via pip:

```
pip install scythe
```

2.2 Using Scythe

There are three ways to use Scythe: GUI, command line parameters, or configuration file.

The GUI allows you to directly obtain data from ENSEMBL, have it converted and processed by Scythe.

If you are using local files, make sure you have directories prepared that contain the *loc* and fasta files, and that you have a *grp* file ready. (See [Converters](#) on how to convert other formats to *loc* or *grp*.)

Read more about configuration files under [Configuration Files](#).

2.2.1 Graphical User Interface

Call the GUI via

```
scythe-gui.py
```

or with configuration file *conf.py*

```
scythe-gui.py conf.scy
```

Please note that Scythe will write intermediate files to the current working directory. You might want to create a new directory for each run beforehand. Data that was already downloaded from ENSEMBL will not be downloaded again. If you want to reuse this data but try different parameters, just change the output directory (via the GUI) and start Scythe from the same working directory.

2.2.2 Command Line Interface

```
Traceback (most recent call last):
  File "../scythe/scythe.py", line 11, in <module>
    import configparser
ImportError: No module named configparser
```

Please note that you cannot automatically download sequences from ENSEMBL with *scythe.py* (see [Downloading from ENSEMBL without the GUI](#)).

TROUBLESHOOTING

3.1 Error Message: *Are all gene models in your fasta files?*

Try checking the format of your fasta headers and the headers in your *grp* and *loc* files. Often, fasta headers contain more than just IDs and can be truncated. If all your fasta files follow the same scheme, calling Scythe with *-d DELIMITER -a NUM* will help.

3.1.1 Example:

fasta:

```
[...]
>ENSMLUP000000021790 pep:known_by_projection scaffold:Myoluc2.0:GL430705:85913:99944:-1 gene:ENSMLUG0
ITDGCSLFGNRSLSSESRPVS GPRTRPQRLAERLTLLGGAHSAVLPSSNRMEPPLGTQ
QGAMQPLVADDFEACLLDKVRWTRGAQRVSQMVEEVQKVIYHLTEISNQDIRFQAIPYS
LMY[...]
```

loc:

```
[...]
ENSMLUG000000013451      ENSMLUP000000021790      ENSMLUP000000012239
[...]
```

That means your delimiter (*-d*) should be " " (space character) and the id index (*-a*) should be 0 (first part of the header is relevant.) Counting starts at 0.

```
scythe.py -i DIR -g GRP -d " " -a 0
```

3.1.2 Example 2

fasta header:

```
> pep:known_by_projection|ENSMLUP000000021790
```

loc:

```
ENSMLUG000000013451      ENSMLUP000000021790      ENSMLUP000000012239
```

Call with:

```
scythe.py -i DIR -g GRP -d "|" -a 1
```


FORMAT

The formats have a very simple, human readable structure allowing users to manually modify or create them. For locus-gene model identifier relations: *loc format* For grouping of orthologous genes *grp format*

4.1 loc format

Each row of a *loc* file represents a set of gene models (GM) for a gene locus (L). Rows start with a the gene locus ID followed by its gene model identifiers. The GM in the second column in the file is treated as the reference form for its gene locus, i.e. in a case of a tie, this gene model is preferred. The columns are tab-separated.

File *Spec0.loc*:

```
Sp0L0      Sp0L0GM0      Sp0L0GM1      . . . Sp0L0GMjLF
Sp0L1      Sp1L0GM0      Sp0L1GM1      . . . Sp0L1GMkLF
.
.
.
Sp0Lm      Sp1LmGM0      Sp0LmGM1      . . . Sp0LmGMlLF
```

4.2 grp format

Each row of a *grp* file represents an orthologous group. Rows start with a unique group ID followed by orthologous gene identifiers for the respective species. Rows don't need to have the same number of columns. The gene loci order within a row is irrelevant. The program will keep the identifiers for its output. The columns are tab-separated.

File *orthogroups.grp* (Sp: species, L: orthologous gene loci):

```
0   SpaLw   SpbLw   SpcLw   . . . SpzLw
1   SpaLx   SpbLx   SpcLx   . . .
.
.
.
N   SpaLz   SpbLz   SpcLz   . . .
```


CONVERTERS

Scythe uses simple, human readable formats to store gene-transcript and ortholog information. See also [Format](#).

5.1 loc

Scripts to convert the following formats to *loc* format are included:

- gff3
- tab-separated (eg ENSEMBL BioMart)

run

```
scythe_loc_gff.py -f GFF
```

```
#----- loc output format -----#  
  
LOCUS0      TRANSCRIPT0_0      TRANSCRIPT0_1      ...      TRANSCRIPT0_n  
LOCUS1      TRANSCRIPT1_0      ...      TRANSCRIPT1_m  
.  
.  
.  
LOCUSk      TRANSCRIPTk_0      ...TRANSCRIPTk_l  
  
#-----#
```

```
#----- gff version 3 input format -----#  
example (tab is shown as \t):  
  
##gff-version 3  
L1i\texample\tgene\t1000\t9000\t.\t+\t.\tID=L1;Name=L1;Note=example  
L1.a\texample\tmRNA\t1000\t9000\t.\t+\t.\tID=L1.a;Parent=L1;Name=L1.a
```

Note that this script only relies on the "ID" and "Parent" tags, "Name" will be ignored. If "longest" is specified (eg phytozome) and == "1", the transcript will be placed on the first position for its gene.

```
#-----#
```

```
#####  
#  scythe_loc_gff.py  -f FILE.gff3  #
```

```
#####
```

```
-f, --file=gff3_FILE
-o, --output=FILE          output file [default: gff3_FILE.loc]
-h, --help                 prints this
-H, --HELP                 show help on format
```

or

```
scythe_loc_tsv.py -f FILE.tsv
```

```
#----- loc output format -----#

LOCUS0      TRANSCRIPT0_0      TRANSCRIPT0_1      ...      TRANSCRIPT0_n
LOCUS1      TRANSCRIPT1_0      ...      TRANSCRIPT1_m
.
.
.
LOCUSk      TRANSCRIPTk_0      ...TRANSCRIPTk_l

#-----#
```

example ensembl query:

```
http://www.ensembl.org/biomart/
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE Query>
<Query virtualSchemaName = "default" formatter = "TSV"
header = "0" uniqueRows = "0" count = "" datasetConfigVersion = "0.6" >
<Dataset name = "hsapiens_gene_ensembl" interface = "default" >
<Filter name = "biotype" value = "protein_coding"/>
<Attribute name = "ensembl_gene_id" />
<Attribute name = "ensembl_transcript_id" />
<Attribute name = "ensembl_peptide_id" />
<Attribute name = "cds_length" />
</Dataset>
</Query>
```

```
#####
```

```
# scythe_loc_tsv.py #
```

```
#####
```

```
-f, --file=ENSEMBLBioMart.tsv
```

```
format: 1st column: gene id, 2nd column:transcript id,
3rd column: peptide id, 4th column: cds length;
gene ids can occur multiple times
```

```
-c, --custom=COLx,COLy,COLz,... COLi in ["gene","transcript", "protein", "length"]
Use this if your file is different from the described biomart c
"cds_length" is optional but recommended, at least one of
["transcript", "protein"] need to be included
```

```
-o, --output=FILE          output file [default: ENSEMBLEBioMart.tsv.loc]
-h, --help                 prints this
-H, --HELP                 show help on format
```

```
#-----#
```

See also: [loc format](#).

5.2 grp

Please note that the *grp* converters need a concatenated *loc* file in addition to the orthology information. Scripts to convert the following formats to *grp* are included:

- orthomcl
- proteinortho
- tab-separated (eg ENSEMBL BioMart)

run

```
scythe_grp_orthomcl.py
```

```
scythe_grp_proteinortho.py
```

```
scythe_grp_tsv.py
```

See also *grp format*.

DOWNLOADING FROM ENSEMBL WITHOUT THE GUI

To download sequences (pep and cds fasta files) from ENSEMBL without the graphical user interface, use *scythe_ensembl_fasta.py* for fasta files and *scythe_ensembl_ortho_mysql.py* to download pairwise orthology information.

6.1 *scythe_ensembl_fasta*

```
usage: scythe_ensembl_fasta.py -s species1,species2 -r INT
```

options:

```
-s, --species=STR    comma-separated list of species (eg 'homo_sapiens,gorilla_gorilla')
-r, --release=NUM    ENSEMBL version (eg '75')
-d, --dir DIR        output directory [default ./]
-h, --help           prints this
```

6.2 *scythe_ensembl_ortho_mysql*

```
usage: scythe_ensembl_ortho_mysql.py -s species_1,species_2 -r INT
```

options:

```
-s, --species=STR    comma-separated list of species (eg 'homo_sapiens,gorilla_gorilla')
-r, --release=NUM    ensembl version (eg '75')
-h, --help           prints this
```

6.3 Manual merge of tab-separated files to one *.grp* file

If you have pairwise (two-species) files ready and want to merge them into a multi-species *.grp* file you can do so via the *scythe_ensembl2grp* and *scythe_mergeSubsets* scripts.

6.4 *scythe_ensembl2grp*

```
usage: scythe_ensembl2grp.py -f FILE1,FILE2 -o OUT.grp
```

```
-f, --files=STR      list of ensembl tsv files (eg sA.tsv,sB.tsv,sC.tsv)
-o, --output=FILE    output file
-h, --help           prints this
```

6.5 scythe_mergeSubsets

```
usage:  scythe_mergeSubsets.py -g groups.grp -o new.grp
```

options:

```
-g, --grp=FILE.grp
```

```
-o, --output=OUTFILE.grp      output file [default: FILE.allspec.grp]
```

```
[-r, --rename                  discard old orthogroup ids and start numbering from 0]
```

```
-h, --help                    prints this
```

```
[-n, --numspec=N             min number of species ]
```

```
-----
```

```
.grp format: GroupID          geneIDiSp1          geneIDjSp2          ...geneIDkSpn
```

CONFIGURATION FILES

The configuration files allow you save/restore your settings.

7.1 Example file

```
[Mode]
use_ensembl=yes
use_local_files=no

[Paths]
fasta_directory=/home/$USER/mydir/fa/
loc_directory=/home/$USER/mydir/loc/
grp_file=/home/$USER/mydir/mydir/test.shared_by_all.grp
output_directory=/home/$USER/mydir/ref_out/

[Cleanup]
clean_up_directories=yes

[Run_options]
num_CPU=1

[Penalties]
gap_open_cost=10
gap_extend_cost=0.5
substitution_matrix=EBLOSUM62

[Algorithm]
use_sl_glob=unset
use_sl_ref=yes
use_mx_sum=unset

[Fasta_header]
fasta_header_delimiter=" "
fasta_header_part=0
```

7.2 [Fasta_header]

In case the ids in your *loc* file are truncated or your fasta headers carry information in addition to the sequence identifier, you can define where to split the fasta header.

Example:

```
>other:xyz_sequenceID_length:123_species:x  
MYQVLQAYDWKYLHDNHDCNFQVGGADQLGNI...
```

and:

```
[Fasta_header]  
fasta_header_delimiter = "_"  
fasta_header_part = 1
```

Would result in taking only *sequenceID* into account.

7.3 [Run_options]

```
[Run_options]  
num_CPU=1
```

Will set the maximum number of CPUs to use to one.

7.4 [Algorithm]

See *Algorithms and Application*

BIBLIOGRAPHY

- [NEE1970] Needleman, Saul B.; and Wunsch, Christian D. (1970). "A general method applicable to the search for similarities in the amino acid sequence of two proteins". *Journal of Molecular Biology* 48 (3): 443–53. doi:10.1016/0022-2836(70)90057-4. PMID 5420325.
- [EMB2000] EMBOSS: The European Molecular Biology Open Software Suite (2000) Rice,P. Longden,I. and Bleasby, A.Trends in Genetics 16, (6) pp276–277