

PAGAN: brief tutorial (v. 0.12)

Ari Löytynoja
ari.loytynoja@gmail.com

November 26, 2010

1 Background

PAGAN is a general-purpose method for the alignment of sequence graphs. The main advantage of the graph representation of sequences is the graphs' ability to describe the uncertainty in the presence of characters at certain sequence positions. This is highly useful in phylogenetic progressive alignment and allows for capturing the different properties of insertions and deletions when alignments are iterated for multiple sequences.

PAGAN is still under development and will hopefully evolve to an easy-to-use, general-purpose method for phylogenetic sequence alignment. The graph representation has features that make it especially powerful for phylogenetic placement of sequences into existing alignments. As such a tool has direct applications in the analyses of increasingly abundant RNA-seq data and there is no satisfactory solution available, we have focused on this problem and implemented the necessary functionality first.

This document explains how to install and start using the PAGAN software.

2 Installation

PAGAN is open-source software licensed under the GPL. The C++ source code is provided at <http://www.ebi.ac.uk/~ari/pagan>. PAGAN is developed and tested on a Linux system and we can currently only provide instructions for that platform, with a special focus on the popular Ubuntu/Debian-based distributions.

Download and installation using git

The most recent version of the PAGAN source code is available from the git-repository and snapshots of this are downloadable as compressed tar-packages. The `git` software can be found at git-scm.com. On Ubuntu, it can also be installed using command:

```
sudo apt-get install git-core
```

The PAGAN source code can then be downloaded and compiled using commands:

```
git clone http://www.ebi.ac.uk/~ari/pagan
cd pagan/src
make -f Makefile.no_Qt
./pagan
```

Installation of Boost libraries

PAGAN requires one utility library from the [Boost project](#) that may not be included in standard OS installations. This library has to be installed **before** compiling the PAGAN source code. On Ubuntu, it can be installed using commands:

```
apt-cache search libboost-program-options1
```

(See which version number, ending with .0, is provided and edit the command below.)

```
sudo apt-get install libboost-program-options1.40.0
```

If your distribution does not provide Boost libraries (highly unlikely) or you are not allowed to install software on your system, you can download and install the necessary library from the Boost project repositories directly into your PAGAN source code directory. This can be achieved using following commands:

```
# make a temporary directory
mkdir ~/tmp_boost
cd ~/tmp_boost

# get the source code
wget http://sourceforge.net/projects/boost/files/boost/1.44.0/boost_1_44_0.zip/download
unzip boost_1_44_0.zip
cd boost_1_44_0

# fix the permission
chmod +x tools/jam/src/build.sh

# compile and install
sh ./bootstrap.sh --prefix=$PATH_TO_PAGAN_DIR/boost --with-libraries=program_options
./bjam install

# check that it's there and set to library path
ls $PATH_TO_PAGAN_DIR/boost/lib
export LD_LIBRARY_PATH=$PATH_TO_PAGAN_DIR/boost/lib:$LD_LIBRARY_PATH

# clean up
cd ../../
rm -r ./tmp_boost
```

Note that you need to replace `$PATH_TO_PAGAN_DIR` with the appropriate file path. You also need to copy the line `export LD_LIBRARY_PATH=...` in your `~/.bashrc` or similar such that it will automatically be set for future sessions.

3 Using PAGAN

PAGAN is a command-line program. A list of the most important program options is outputted if no arguments are provided:

```
./pagan
```

and a more complete list is given with the option `--help`:

```
./pagan --help
```

3.1 Phylogenetic multiple alignment

PAGAN is based on a progressive algorithm that aligns sequences according to a guide tree. It (currently) cannot compute a tree by itself and requires the user to provide a **rooted** binary tree relating the sequences. The leaf names in the tree and the sequence names (until the first space) in the sequence file have to match exactly. Alignment is only performed for the parts of the guide phylogeny that have sequences associated; the unnecessary branches and sequences are pruned/dropped out.

The minimal command to perform the alignment is:

```
./pagan --seqfile sequence_file --treefile tree_file
```

The resulting alignment will be written in files `outfile.fas` and `outfile.xml`. If you want to use another file name, you can specify that with option `--outfile`:

```
./pagan --seqfile sequence_file --treefile tree_file --outfile another_name
```

PAGAN will automatically add suffix `.fas` and `.xml`.

The sequence input file has to be in FASTA format and the guide tree in Newick tree format, with branch lengths as substitutions per site. The resulting alignment will be written in FASTA format and in XML-based [HSAML format](#). PAGAN supports the alignment of nucleotide and amino-acid sequences although the latter is still rather experimental and not fully optimised.

3.2 Phylogenetic placement of sequences in an existing alignment

PAGAN can reconstruct ancestral nodes for a given phylogeny based on an existing alignment. Additional sequences can then be aligned and added to this reference alignment without affecting the relative alignment of original sequences. The main application of this is phylogenetic placement of short reads, coming from NGS platforms, into existing alignments but it can, in principle, be used for any sequences, even amino-acid ones.

The minimal command to perform the alignment is:

```
./pagan --ref-seqfile ref_alignment_file --ref-treefile ref_tree_file  
--readsfile reads_file
```

The reference alignment has to be in FASTA format and the reference tree in Newick or Newick Extended (NHX) format; the reads (or sequences in general) that will be added in the reference alignment can be either in FASTA or FASTQ format. If the reference alignment consists of one sequence only, the reference guide tree is not required. Again, you can define your own output file with option `--outfile`.

Pair-end reads and 454 data

For the alignment of NGS reads, additional options `--pair-end` and `--454` are useful. The first one merges paired reads into one (separated by a spacer) before the alignment and the second models the ambiguous length of mononucleotide runs in data coming from Roche 454 platform. The pairing of pair-end reads assumes that the reads have identical names (until the first space) with the exception that the left read ends with `/1` and the right read with `/2`. In the resulting alignment, the paired sequence will have suffix `/p12`. Option `--454` only works with FASTQ-formatted data.

Overlapping pair-end reads

The length and overall quality of NGS data can be improved by using such a short fragment length that the reads starting from each end of the fragment overlap in the middle. PAGAN allows merging such reads and handling them as one sequence. Merging is done using option `--overlap-pair-end`, the reads successfully merged having suffix `/m12` in the resulting alignment.

The merging requires significant overlap between the two reads in their pairwise alignment (performed without masking). Options `--overlap-minimum` and `--overlap-identity` can be used to change the minimum length and base identity of this overlap. Shorter overlaps that show perfect base identity are also accepted; the minimum length of this can be changed with option `--overlap-identical-minimum`. The merged reads can be outputted in FASTQ format using option `--overlap-merge-file`. PAGAN can also be used just to merge the overlapping reads:

```
./pagan --overlap-merge-only --readsfile reads_file --overlap-merge-file merge_file
```

Assignment of reads to specific nodes

By default, PAGAN will add the sequences in the bottom of the alignment by aligning them against the root node. An alternative is to assign the sequences to one or more nodes, and find the node among those against which the read matches best and align it there. Sequences can be assigned to a specific node or a set of nodes using the NHX tree format with an additional tag TID. The tag identifier can be any string, in the example below we use number '001'. One can use any number and any combination of tags given that each node and each sequence has at most one tag.

As an example, let's assume

a reference alignment:

```
>A
AGCGATTG
>B
AACGATCG
>C
TGCGGTCC
```

and a read that we want to add in FASTA:

```
>D TID=001
AGCGATCG
```

or in FASTQ format:

```
@read_D_000103@15@13524@18140#0/1 TID=001
AGCGATCG
+
CCCCCCCC
```

The placement of the read in the reference alignment depends on the format of the reference tree:

- a phylogeny in plain Newick format adds the read at the root of the tree:

```
((A:0.1,B:0.1):0.05,C:0.15);
```

- a phylogeny in NHX format with one matching label adds the read to that node:

```
((A:0.1,B:0.1):0.05[&&NHX:TID=001],C:0.15);
```

- a phylogeny in NHX format with several matching labels finds the best node and adds the read to that:

```
((A:0.1,B:0.1):0.05[&&NHX:TID=001],C:0.15):0[&&NHX:TID=001];
```

For convenience, option `--test-every-node` is provided to exhaustively search through all the nodes and add the sequence at the best one. This ignores the TID tags and overrides the placement information given in the guide tree.

Typical analysis of NGS data

A typical command to perform the placement of NGS reads could be:

```
./pagan --ref-seqfile ref_alignment_file --ref-treefile ref_tree_file
        --readsfile reads_file --trim-read-ends --discard-overlapping-identical-reads
        --rank-reads-for-nodes [--pair-end] [--overlap-pair-end] [--454]
```

where the optional (and some mutually exclusive) options are in square brackets.

Depending on the output, the thresholds for masking and trimming can be adjusted using the relevant options explained below. The effect of trimming is obvious from the output; it is good to remember that bases written in lower case in the output were masked and considered as N's during the alignment. A sequence of N's matches any sequence and give rubbish alignments; if that happens, you may need to raise the masking threshold or lower the trimming thresholds.

3.3 Additional program options

Some of the options printed by `./pagan --help` relate to unfinished features and may not function properly. Only the main options are documented here.

Generic options

- Option `--silent` minimises the output (doesn't quite make it silent, though).
- Options `--ins-rate` and `--del-rate` can be used adjust the insertion and deletion rates (per base substitution). Although it would be possible to consider the two processes separately, this has not been implemented yet and the two rates are combined.
- Options `--gap-extension` and `--end-gap-extension` define the gap extension probability for regular and terminal gaps. For meaningful results, the latter should be greater (and, for pair-end data, equal to `--pair-read-gap-extension`).
- Options `--dna-kappa` and `--dna-rho` affect the DNA substitution scoring matrix; base frequencies are estimated from the data.
- Options `--scale-branches`, `--truncate-branches` and `--fixed-branches` override the branch lengths defined in the guide tree. By default, branches are truncated; this can be prohibited with `--real-branches`.
- Option `--output-ancestors` writes the parsimony-reconstructed ancestral sequences for the internal nodes of the tree. The tree indicating the nodes is written in `outfile.ancntree`.

There are many parameters related to “insertion calling”, the type and amount of phylogenetic information required to consider insertion-deletion as an insertion and thus prevent the later matching of those sites. These parameters are still experimental (although some of them are used and affect the resulting alignment) and will be described in detail later.

Phylogenetic placement

- Option `--qscore-minimum` sets the Q-score threshold for sites to be masked (replaced with N's for alignment, shown in lower case in the output); `--allow-skip-low-qscore` further allows skipping those sites (usefulness of this is uncertain). Masking is done by default (see `./pagan --help` for the default threshold) and can be disabled with option `--no-fastq` that prohibits all Q-score-based pre-processing done either by default or chosen by the user (e.g. trimming and 454-specific modelling).
- Option `--trim-read-ends` enables the trimming of FASTQ reads; `--trim-mean-qscore`, `--trim-window-width` and `--minimum-trimmed-length` define the minimum Q-score and width for the sliding window and the minimum length for the trimmed read. Trimming progresses inwards from each end until the mean Q-score for the window exceeds the score threshold; if the read is shortened below the length threshold, it is discarded. The length of removed (trimmed) fragments is indicated in the sequence name field in the format `P1ST$11:P1ET$12` where `$11` and `$12` refer to the start and end of the read. If the sequences consists of a read pair, the trimming of the right-hand side read is similarly indicated using the tag `P2`.
- Option `--rank-reads-for-nodes` ranks the sequences assigned to a node and aligns them in that order. The ranking is based on their score in the placement alignment used to decide between the alternative nodes. If sequences are assigned to one node each, one round of alignments against each node is performed to define the ranking before the final multiple alignment.
- Reads fully embedded in other reads can be discarded: `--discard-overlapping-reads` identifies reads that overlap in the placement alignment and removes ones that are fully embedded in another one; `--discard-overlapping-identical-reads` extends this by checking that the embedded read is identical (on base level) to the longer read before discarding it; `--discard-pairwise-overlapping-reads` makes all the pairwise alignments to identify embedded reads.
- Option `--reads-distance` sets the expected distance between the read and the pseudo-parent node (against which the read is aligned) and thus affects the substitution scoring used in the alignment. Having the distance very short (default), the alignment is stringent and expects high similarity.
- Reads with too few sites aligned against sites of reference sequences are discarded. (The stringency of the alignment is set using the option above.) Options `--min-reads-overlap` and `--min-reads-identity` set the required overlap and base identity for accepting the read.

The rest of the options are either not important for basic use or self-explanatory (or both).

3.4 Inference of ancestral sequences for existing alignments

The features required for phylogenetic placement of reads can be used without any reads, too. Command:

```
./pagan --ref-seqfile alignment_file --ref-treefile tree_file
```

writes the aligned sequences from `alignment` to files `outfile.fas` and `outfile.xml`, thus providing a tool to convert FASTA files to HSAML format. This maybe more interesting with option `--output-ancestors`. Command:

```
./pagan --ref-seqfile alignment_file --ref-treefile tree_file --output-ancestors
```

includes the parsimony-reconstructed internal nodes in the FASTA output, providing a efficient tool to infer gap structure for ancestral sequences based on existing alignments. The tree indicating the ancestral nodes is written in `outfile.ancree`.