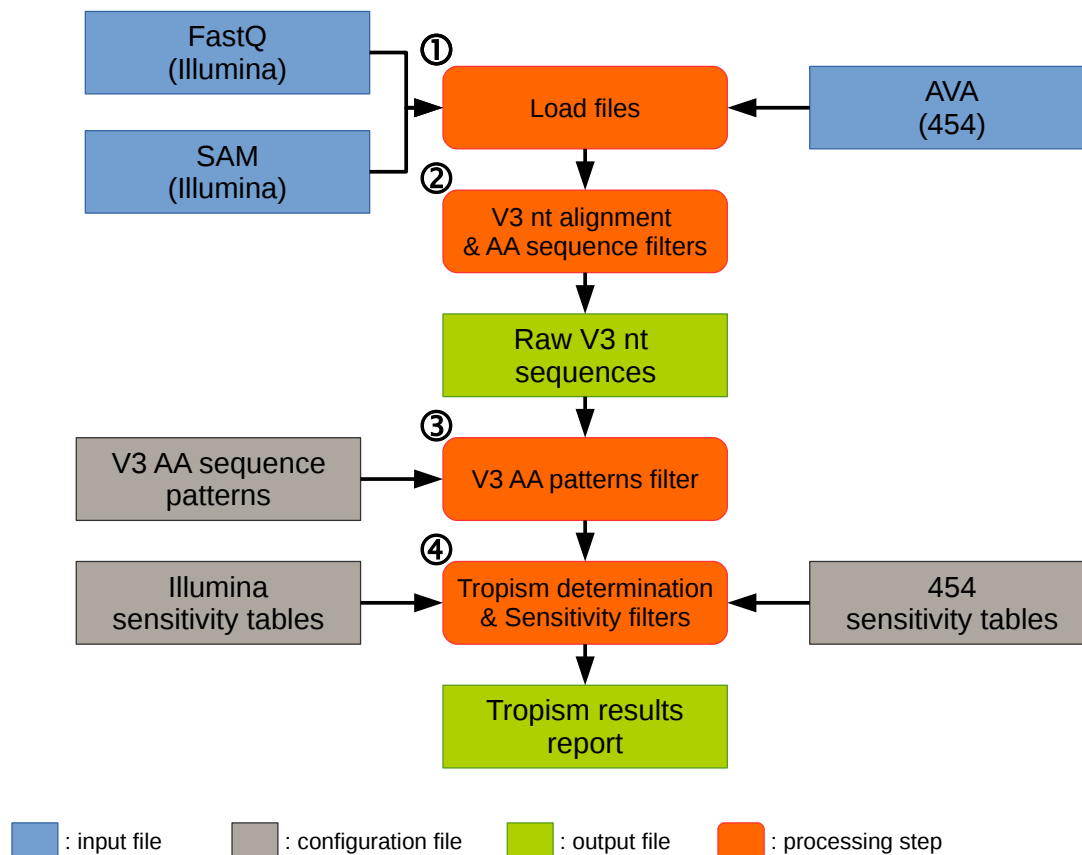# Purpose

The PyroVir software is an effective tool to interpret HIV Tropism next generation sequencing (NGS) information and guide the adaptation of anti-HIV treatment. The purpose of the software is to identify the tropism of an individual patient's HIV strains and shows whether the patient is infected with virus that enters cells using the CCR5 co-receptor, the CXCR4 co-receptor, or both.

# Software pipeline

## PyroVir pipeline



Step 1 :
The software is designed to process 454 and Illumina data files. For 454 technology the input files are AVA consensus reads output files (must contain galign_con_ prefix) and for Illumina technology FastQ or SAM files. If the subtype is known for the samples, this information can be added for a better tropism determination.

Step 2 :
The input sequences are aligned versus the BaL V3 sequence and trimmed at the boundaries of V3. The trimmed sequences are momentarily translated to amino-acid sequences and filtered. Those with no Cysteine as first amino-acid required for the disulfide bound maintaining the V3 loop, with a STOP codon, outside the size limits (between 33 and 37 AA by default) or containing ambiguous characters (N) in the nucleotide sequence are discarded.
The nucleotide remaining sequences are written in a fasta file. The sequences headers contain the count information :

```
>CON_2 Count=874
TGTACAAGACCCAATAACAATACAAGAAAAGGTATAGCTATAGGACCAGGGAGAACATTTTATACAAGAGAAAAAATAATAGGAAATATAAGACAAGCACATTGT
>CON_5 Count=758
TGTACCAGACCCAATAACAATACAAGAAAAGTATAACTATAGGACCAGGGAGAGCATTTTATACAAGAGGAGAAATAATAGGAAATATAAGACTAGCACATTGT
>CON_3 Count=539
TGTACAAGACCCAATAACAATACAAGGAAAGTATACCCATAGGACCAGGGAGAGCATTTTATACAAGAGGAGAAATAATAGGAAATATAAGACTAGCACATTGT
>CON_16 Count=270
TGTACAAGACCCAATAACAATACAAGAAAAAGTATACATATAGGACCAGGGGGAGCATTTTATACAACAGGAGGAATAATAGGAGATATAAGACAAGCACATTGT
>CON_8 Count=212
TGTACAAGACCCAATAACAATACAAAAAAAGGTATAGCTATAGGACCAGGGAGAACATTTTATACAAGAGAAAAAATAATAGGAAATATAAGACAAGCACATTGT
>CON_1 Count=202
TGTACAAGACCCAATAACAATACAAGAAAAGTATACCCATAGGACCAGGGAGAGCATTTTATACAAGAGGAGAAATAATAGGAAATATAAGACTAGCACATTGT
>CON_91 Count=188
TGTACAAGACCCAATAACAATACAAGCAAAGTATACCTATAGGACCAAGGAGAGCGTTTTATGCAACAGGAAGAATAATAGGAGATATAAGACAAGCATATTGT
...
```
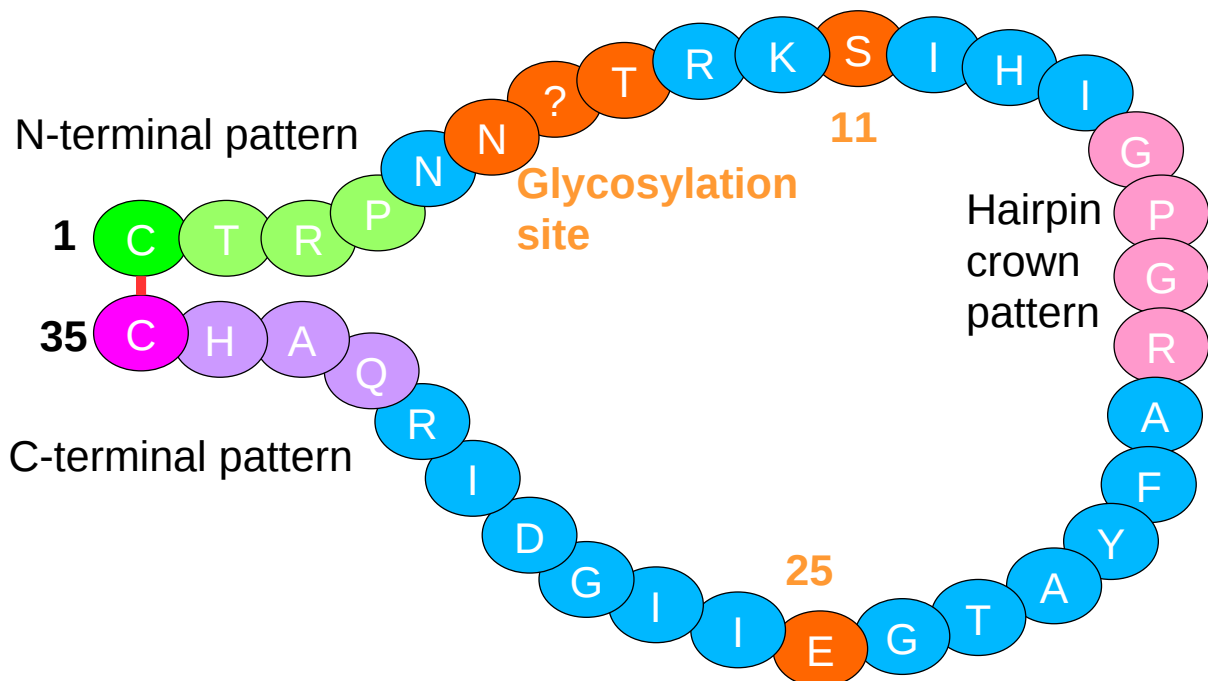
Step 3 :

With the new fasta file, each sequence is translated and filtered on the bases of three V3 specific conserved patterns. The N-terminal pattern (CTRP), the hairpin crown pattern (GPGR) and the C-terminal pattern (QAHC). An identity of 0.5 and 0.75 was allowed for substitution and insertion/deletion of an amino-acid at any of the three signature patterns. Sequences which do not match with those patterns constraints are discarded.

Step 4 :

For each remaining amino-acid sequence, the tropism is determined based on the combined rule adapted to the subtype of the sample if this field was entered in step 1, if not, the B subtype combined rule is applied.
The combined rule is based on V3 specific positions amino-acids composition (as 11, 25, glycosylation sites…) and the net electrostatic charge of the sequence as shown on the V3 loop key positions figure.

# V3 loop : Key positions



**Net electrostatic charge = (∑R + ∑K) - (∑D + ∑E)**

A stastistic filter is applied to each sequence to discard artifactual sequences resulting from PCR or

sequencing errors. Cut-off threshold are based on the results from 20 clones sequencing (Illumina and 454 technologies) which generate two matrices for each technology based on Poisson distributions. The first one is a V3 Global sensitivity with one column for the frequency values of the whole V3 sequence and the rows are the coverage. The second one is a Position-specific sensitivity matrix with all V3 positions as columns and coverage as rows.

The choosen Global and Position-specific tables for the studied sample are specific to the sequencing technology.

For R5 tropic sequences, the software uses the Global sensitivity table to discard the sequences which frequency is below the table cut-off depending on the coverage of the sequence. For X4 tropic sequences, the Position-specific table is used. The V3 position of the sequence that switch the tropism from R5 to X4 is the key position that is looked for in the column table and the sequence is discared if its frequency is below the cut-off depending on the coverage of the sequence.

Finally, a report is produced indicating the tropism X4 ratio and detailing the quasi-species with a pie-chart and a sequences table for each sample. The criteria determining the X4 tropism are presented in the sequences table. Also, the number of discarded sequences for each filtering steps are notified in the filter table.