# CellPhy - The hard way (tree search, bootstrapping, mutation mapping)

Alexey Kozlov, João M Alves*, Alexandros Stamatakis & David Posada

July 2020

## 1 Introductory note

In the following supporting document, we will go through all steps required to perform a full single-cell phylogenetic analysis using **CellPhy**, including phylogenetic tree search, bootstrap support estimation and mutation mapping. For reproducibility purposes, we provide a toy single-cell dataset (from a colorectal cancer patient - CRC24) together with a step-by-step tutorial containing all information needed to generate similar results to the ones presented in **Figure 6** of the main manuscript (Kozlov *et al.* 2020).

Importantly, in order to speed up the analysis, the toy example used here is composed of 500 randomly chosen SNVs, thus corresponding to a "thinned" version of the original dataset. In any case, we followed exactly the same steps shown in this document to generate the results presented in our study.

Before starting, make sure to download and compile CellPhy, as well as our toy dataset from:

- [CellPhy's github](#)

---

## 2 Input data

CellPhy accepts two types of input:

- 1) matrix of genotypes in FASTA or **PHYLIP** format - *CRC24.ToySet.phy*

- 2) **standard VCF file** - *CRC24.ToySet.vcf*

When the input is a genotype matrix, genotypes are encoded as shown in Table S2 of Kozlov *et al.* 2020. When the input is a VCF, CellPhy can be run in two distinct modes. The first mode ("**CellPhy-EP17**") requires a VCF with at least the GT field (that stores the genotype calls), in which case CellPhy simply extracts the genotype matrix. The second mode ("**CellPhy-GL**") requires a VCF with the PL field (which stores the phred-scaled genotype likelihoods) and uses instead the likelihoods of each genotype.

Here's how they look like:

```
# PHYLIP file
$ head -n 5 CRC24.ToySet.phy
25 500
TI-TNS10    SNNCGTNCTNGAATGCNTANTGANNTYTNTNTGCG [...]
TI-TNS11    NNGCNYANTNGAGNGTGTATTGAANCNTGTNTANT [...]
TI-TNS12    NNACGTATNAGAATGNGTATNGANNCYTGNGCACG [...]
TI-TNS16    CNNNNTACTAGNATGNATATTGANTNNTGNGTANT [...]

# VCF file
$ grep -v "##" CRC24.ToySet.vcf | head -n 5
#CHROM  POS ID  REF ALT QUAL    FILTER  INFO    FORMAT  TI-TNS10    TI-TNS11 [...]
1   5107124 .   G   C   21  .   .   GT:SO:AD:BI:GQ:PL   0/1:True:6,6:0.594:21:58,37,148 ./.:.:.:.:.:. [...]
```

```
1    23800582    .    T    A    6    .    .    GT:SO:AD:BI:GQ:PL    ./.:.:.:.:.:.    ./.:.:.:.:.:.  [...]
1    25776505    .    A    G    13   .    .    GT:SO:AD:BI:GQ:PL    ./.:.:.:.:.:.    1/1:True:0,1:0.6:4:9,2,0 [...]
1    26426850    .    C    T    92   .    .    GT:SO:AD:BI:GQ:PL    0/0:True:9,0:0.6:15:5,20,284    0/0:True:9,0:0.6:15:5,20,286 [...]
```

---

# 3 Tree inference

Let us now infer a tree with default parameters using both models of CellPhy. Here, we will use a single thread and set a fixed random seed to ensure reproducibility.

## 3.1 Tree inference with "CellPhy-EP17"

This model will take as input a genotype matrix to infer a tree under the genotype error model "*GT-GTR4+FO+E*". If the input file is a PHYLIP file, CellPhy will automatically recognize the file type and proceed with the analysis. In contrast, if the input is a VCF file we need to set `--prob-msa off` which will instruct CellPhy to use the genotype field (GT) instead of the genotype likelihoods (PL).

These two runs should therefore return exactly the same result:
```
$ cellphy --msa CRC24.ToySet.phy --model GTGTR4+FO+E --seed 2 --threads 1 --prefix CRC24.PHY.EP17.Tree
$ cellphy --msa CRC24.ToySet.vcf --model GTGTR4+FO+E --prob-msa off --seed 2 --threads 1 --prefix CRC24.VCF.EP17.Tree

Analysis options:
  run mode: ML tree search
  start tree(s): random (10) + parsimony (10)

[00:00:00] Reading alignment from file: CRC24.ToySet.phy
[00:00:00] Loaded alignment with 25 taxa and 500 sites

[...]
```

In default settings, CellPhy will perform 20 tree searches using 10 random and 10 parsimony-based starting trees, and pick the best-scoring topology. While this is a reasonable choice for many practical cases, users are free to increase the number of starting trees to explore the tree space more thoroughly (e.g. for 50 tree searches, using 25 random and 25 parsimony starting trees, one would add `--tree pars{25},rand{25}` to the command line above).

Once it's finished (~ 5 minutes), CellPhy will output the ML estimates of the genotyping error and the ADO rates as well as the final tree score.
```
Optimized model parameters:

   Partition 0: noname
   Error model (ML): P17,  SEQ_ERROR: 0.000000,  ADO_RATE: 0.698852
   Rate heterogeneity: NONE
   Base frequencies (ML): 0.059424 0.051252 0.084699 0.067892 0.094797 0.246866 0.026857 0.021009 0.255205 0.091999
   Substitution rates (ML): 18.220638 18.220638 18.220638 79.667670 1000.000000 0.001000 18.220638 18.220638 18.220638 [...]

Final LogLikelihood: -9731.841911
```

If we now compare the results obtained with the two runs we can confirm that they are indeed the same:
```
$ grep "Final LogLikelihood:" CRC24.{PHY,VCF}.EP17.Tree.raxml.log
CRC24.PHY.EP17.Tree.raxml.log:Final LogLikelihood: -9731.841911
CRC24.VCF.EP17.Tree.raxml.log:Final LogLikelihood: -9731.841911
```

## 3.2 Tree inference with "CellPhy-GL"

Let's now use instead the genotype likelihood model of CellPhy. This time, our input file needs to be a VCF with an appropriate PL field - if your VCF file was generated using the **SC-Caller algorithm**, make sure to use our VCF-conversion script (*sc-caller.conversion.sh*) to add a standard PL Field to it.

```
$ cellphy --msa CRC24.ToySet.vcf --model GTGTR4+FO --seed 2 --threads 1 --prefix CRC24.VCF.GL.Tree
```

Again, once it finishes, the program will output the final tree score and write all results into different files:

```
Optimized model parameters:

   Partition 0: noname
   Rate heterogeneity: NONE
   Base frequencies (ML): 0.139958 0.190230 0.201030 0.162494 0.047445 0.102141 0.011683 0.007876 0.093882 0.043261
   Substitution rates (ML): 0.140616 0.140616 0.140616 524.497047 687.111761 74.596841 0.140616 0.140616 0.140616 0.140616 [...]


Final LogLikelihood: -30265.905611

[...]

Best ML tree saved to: CRC24.VCF.GL.Tree.raxml.bestTree
All ML trees saved to: CRC24.VCF.GL.Tree.raxml.mlTrees
Optimized model saved to: CRC24.VCF.GL.Tree.raxml.bestModel
```

---

# 4  Bootstrapping trees and tree branch support

Next we will perform a standard non-parametric bootstrap by re-sampling alignment columns and re-inferring trees for each bootstrap replicate. By default, CellPhy will conduct 50 replicate searches, but let's manually increase the number of replicates to be on the safe side.

```
$ cellphy --bootstrap --msa CRC24.ToySet.vcf --model GTGTR4+FO --seed 2 --threads 1 --bs-trees 100 --prefix CRC24.VCF.GL.Boot

[...]
[00:00:00] Data distribution: max. partitions/sites/weight per thread: 1 / 500 / 5000
[00:00:00] Starting bootstrapping analysis with 100 replicates.
[00:00:14] Bootstrap tree #1, logLikelihood: -31748.758731
[00:00:23] Bootstrap tree #2, logLikelihood: -30019.626700
[...]
[00:16:59] Bootstrap tree #99, logLikelihood: -29931.389540
[00:17:11] Bootstrap tree #100, logLikelihood: -29263.238961

Bootstrap trees saved to: CRC24.VCF.GL.Boot.raxml.bootstraps
```

Now that we have the bootstrap trees, we can map the BS support values onto the (previously inferred) best-scoring ML tree:

```
$ cellphy --support -tree CRC24.VCF.GL.Tree.raxml.bestTree \
    --bs-trees CRC24.VCF.GL.Boot.raxml.bootstraps \
    --prefix CRC24.VCF.GL.Support --threads 1

Reading reference tree from file: CRC24.VCF.GL.Tree.raxml.bestTree
Reference tree size: 25

Reading bootstrap trees from file: CRC24.VCF.GL.Boot.raxml.bootstraps
Loaded 100 trees with 25 taxa.

Best ML tree with Felsenstein bootstrap (FBP) support values saved to: CRC24.VCF.GL.Support.raxml.support
```

As stated in the log above, CellPhy will output a tree with support values (in NEWICK format) that can be visualized using any tree viewer software. Here, we will use the **ape** and **ggtree** r-packages to plot the results (Figure 1).

```
# Load packages and tree
> library(ape)
> library(ggtree)
> CRC24 <- read.tree("CRC24.VCF.GL.Support.raxml.support")

# Root tree using the Healthy sample as outgroup
> CRC24 <- root(CRC24, outgroup="Healthy", resolve.root=T)

# Plot tree and BS support (here we're ignoring branch-lengths to ease visualization)
> ggtree(CRC24, branch.length="none", ladderize=F) + geom_tiplab(size=3) + geom_nodelab(size=3, hjust=1.7)
```
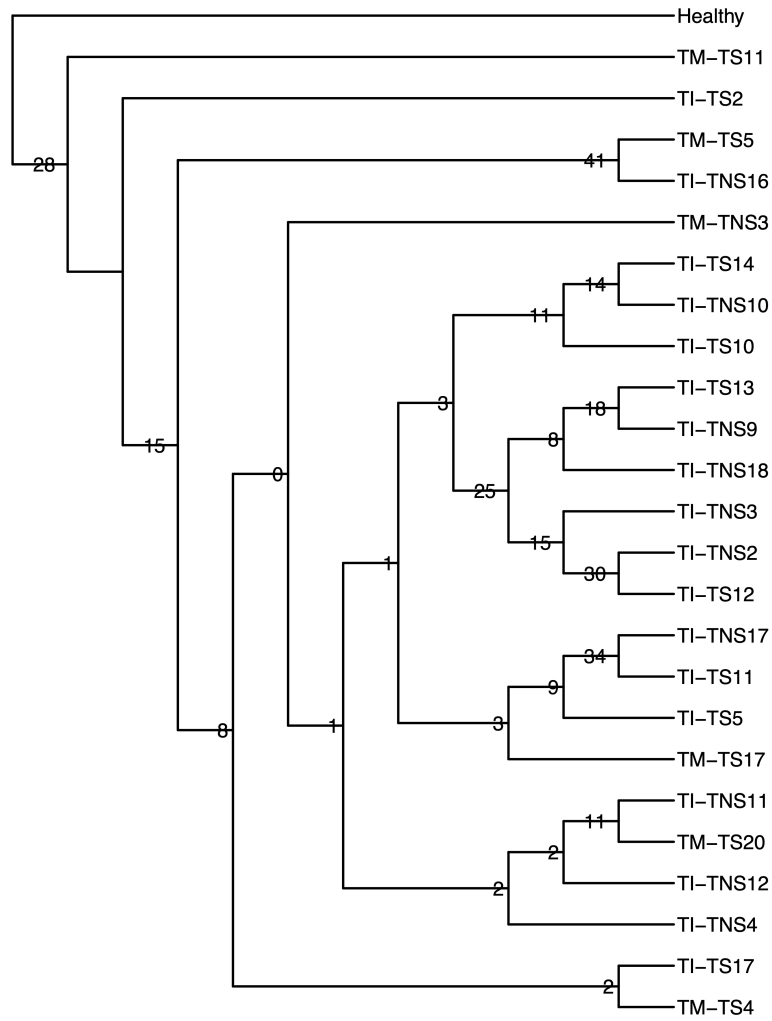
---

Figure 1: Phylogenetic tree of CRC24 dataset with bootstrap support.

# 5 Mapping mutations onto a phylogenetic tree

Cancer genomics studies are, for the most part, interested in understanding when "**driver**" mutations appeared in the malignant cell population. On this basis, we will next show how to map mutations onto a phylogenetic tree using CellPhy. Although the full VCF can be used, users are free to choose which mutations they wish to map onto the inferred phylogeny. In this tutorial, we will focus solely on a set of 15 exonic mutations from the original VCF.

Our input files will therefore consist of a "trimmed" VCF only carrying this subset of exonic mutations, together with the best tree and model estimates from our original CellPhy run (tree search).

```
$ head -n 5 CRC24.MutationsMap
#Chr    Position    GeneID
2    71042907    CLEC4F
2    142274377   LRP1B
3    33048242    GLB1
4    16764214    LDB2

$ bcftools view -T CRC24.MutationsMap CRC24.ToySet.vcf -O v -o CRC24.ToySet.Exonic.vcf

$ MODEL=CRC24.VCF.GL.Tree.raxml.bestModel
$ TREE=CRC24.VCF.GL.Tree.raxml.bestTree

$ cellphy --mutmap \
    --msa CRC24.ToySet.Exonic.vcf \
    --model $MODEL --tree $TREE --opt-branches off \
    --prefix CRC24.ToySet.Exonic.Mapped --threads 1

Branch-labeled tree saved to: CRC24.ToySet.Exonic.Mapped.raxml.mutationMapTree
Per-branch mutation list saved to: CRC24.ToySet.Exonic.Mapped.raxml.mutationMapList
```

---

# 6 Visualizing the results

Once it's done, CellPhy will output 2 distinct files:.

- (**A**) *CRC24.ToySet.Exonic.Mapped.raxml.mutationMapTree*
  → Newick tree file with indexed branches
- (**B**) *CRC24.ToySet.Exonic.Mapped.raxml.mutationMapList*
  → Text file with the number and the list of mutations per branch

We can now use the *cellphy-mutationmapping.sh* accompanying script to plot the mutations onto the phylogenetic tree. If you run this script wihtout any parameters, it will show a help message:

```
$ ./cellphy-mutationmapping.sh
Usage: ./cellphy-mutationmapping.sh raxml.mutationMapTree raxml.mutationMapList Outgroup_name Output_prefix [geneIDs]
Created by: Alexey Kozlov, Joao M Alves, Alexandros Stamatakis & David Posada - 16 June 2020
*Required files:
    -Tree
    -Mutation List
    -Outgroup name
    -Output Prefix

*Optional:
    -Gene IDs (Tab-delimited)
```

Now let's run it again but this time with the required parameters:

```
$ ./cellphy-mutationmapping.sh \
    CRC24.ToySet.Exonic.Mapped.raxml.mutationMapTree \
    CRC24.ToySet.Exonic.Mapped.raxml.mutationMapList \
    Healthy \
    CRC24.ToySet.Exonic

Generating tree plot...
Done!
```

If everything went as expected, you should have generated the following figure, in PDF format (*CRC24.ToySet.Exonic.Tree_mapped.pdf*), where the mutations are mapped onto the tree branches (**Figure 2**).



Figure 2: CRC24 phylogenetic tree with 15 exonic mutations mapped (genomic position)

If you are interested in plotting the gene names instead, you can provide a tab-delimited file (as the one we used to subset our original VCF) linking the genomic position to its gene ID:

```
$ head -n 5 CRC24.MutationsMap
#Chr    Position    GeneID
2    71042907    CLEC4F
2    142274377    LRP1B
3    33048242    GLB1
4    16764214    LDB2
```

Afterwards, we can run *cellphy-mutationmapping.sh* again, but changing the output prefix so that you don't overwrite the previous results:

```
$ ./cellphy-mutationmapping.sh \
    CRC24.ToySet.Exonic.Mapped.raxml.mutationMapTree \
    CRC24.ToySet.Exonic.Mapped.raxml.mutationMapList \
    Healthy \
    CRC24.ToySet.Exonic-GeneID \
    CRC24.MutationsMap

Converting positions to GeneID...
Done!
Generating tree plot...
Done!
```

You will notice that our tree now has the gene names displayed, instead of the genomic positions (Figure 3).



Figure 3: CRC24 phylogenetic tree with 15 exonic mutations mapped (gene names)