LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

Interim Report on SNP analysis and forensic microarray probe design for South American hemorrhagic fever viruses, tick-borne encephalitis virus, henipaviruses, Old World Arenaviruses, filoviruses, Crimean-Congo hemorrhagic fever viruses, Rift Valley fever viruses and Japanese encephalitis viruses

C. Jaing, S. Gardner

June 12, 2012

**Disclaimer**

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

# Interim Report on SNP analysis and forensic microarray probe design for South American hemorrhagic fever viruses, tick-borne encephalitis viruses, henipaviruses, Old World Arenaviruses, filoviruses, Crimean-Congo hemorrhagic fever viruses, Rift Valley fever viruses and Japanese encephalitis viruses

**Project Title: Forensic TaqMan and Microarray Assays for Viral Genotyping**

**Contributors**

Shea Gardner and Crystal Jaing

Lawrence Livermore National Laboratory (LLNL), Livermore, CA

**Principal Investigator and Correspondent**

Crystal Jaing

925-424-6574, jaing2@llnl.gov

LLNL-TR-560677

**June 5, 2012**

# Introduction

The goal of this project is to develop forensic genotyping assays for select agent viruses, enhancing the current capabilities for the viral bioforensics and law enforcement community. We used a multipronged approach combining bioinformatics analysis, PCR-enriched samples, microarrays and TaqMan assays to develop high resolution and cost effective genotyping methods for strain level forensic discrimination of viruses. We have leveraged substantial experience and efficiency gained through year 1 on software development, SNP discovery, TaqMan signature design and phylogenetic signature mapping to scale up the development of forensics signatures in year 2. In this report, we have summarized the whole genome wide SNP analysis and microarray probe design for forensics characterization of South American hemorrhagic fever viruses, tick-borne encephalitis viruses and henipaviruses, Old World Arenaviruses, filoviruses, Crimean-Congo hemorrhagic fever virus, Rift Valley fever virus and Japanese encephalitis virus.

# Methods

SNP analyses were performed on all available full genomes or segments using kSNP. kSNP is a bioinformatics tool for sequence comparison and can scale to hundreds of bacterial or viral genomes, and can be used for finished and/or draft genomes available as unassembled contigs [1]. The method is fast to compute, finding SNPs and building a SNP phylogeny in seconds to hours depending on input. The approach can handle as input hundreds of megabases to gigabases of sequence in a single run. k=13 identified more SNP alleles than k=25. Fewer SNPs were found with the larger k because a longer length of conserved sequence surrounding the SNP is required. With these divergent viruses, shorter k means that SNPs in closer proximity to one another can be found, thus reducing the stringency for conservation surrounding a SNP. A value of k=13 for viruses should provide better resolution of unsequenced novel isolates than k=25, so all results reported below are for k=13, unless otherwise indicated. The results are summarized in Table 1.

Phylogenetic trees were created using:
1) Multiple sequence alignments (MSA), or full genome multiple sequence alignment with MUSCLE (http://www.drive5.com/muscle/; [2]) followed by maximum likelihood with 1000 bootstrap replicates using RAxML (http://icwww.epfl.ch/~stamatak/index-Dateien/software/RAxML-Manual.7.0.4.pdf ; [3])
2) SNP Hamming distance, calculated as the number of SNPs that differ between each pair of genomes and trees built using the neighbor method in PHYLIP
3) SNP RAxML, analyzing just the SNP alleles concatentated into a string for each genome, to make a SNP alignment from multiple genomes that can be treated as though it were an MSA, then using RAxML to create a maximum likelihood tree with 1000 bootstrap replicates from this SNP alignment.

Branch specific SNP allele counts are plotted on nodes, strain/genome specific allele counts are given in brackets after strain name. For each target set, the tree (full genome MSA, SNP distance, or SNP RAxML) that results in the most SNPs mapping to the tree is shown (that is, the tree that yields the fewest homoplastic SNPs).

Microarray probes were designed for every SNP. Probe design strategy maximized sensitivity and specificity based on extensive prior lab testing on a NimbleGen microarray platform, where we demonstrated 100% SNP allele call rates and 99.5% accuracy (in prep, and unpublished

reports for DHS). We determined that maximum sensitivity and SNP discrimination accuracy result if the SNP base is at the $13^{th}$ position from the 5' end of the probe (the end farthest from the array), probes are between 32 and 40 bases long, and length varies so as to equalize hybridization free energy ($\Delta G$) to the extent possible within the allowable length range. Probes shorter than 32 bases have high false negative rates, and longer probes are inefficient at discriminating single base mismatches. We found that $\Delta G$ is a better predictor of hybridization than $T_m$. Probe candidates with hybridization free energy below $\Delta G=-43$ kcal/mol were shortened until either their $\Delta G$ exceeded -43 kcal/mol or they reached the minimum 32 bases. Probes were designed around the SNP on both the plus and minus strands, for all 4 possible SNP alleles, and all surrounding sequence variants. We design probes for both the plus and minus strand; these are not the reverse complements of one another because the SNP does not lie at the center of the probe. There are probes for each of the 4 variants on each strand, so at least 8 probes per SNP locus. In addition, any sequence variation outside of the k-mer SNP context of conserved bases is captured in multiple alternative probes for that allele, so there may be more than 8 probes per SNP locus, although for a given hybridization, only the probe variant with the best signal is used for assessing the SNP allele at the $13^{th}$ position. Finally, probes are trimmed from the 3' end to remove any N's or other degenerate bases, and omitted altogether if doing so results in a probe less than 32 bases. If a probe is a subsequence of any other, only the shorter of the two is kept. If necessary to fit on the desired array format, probes can be omitted for alleles not represented in the target sequences, e.g. for biallelic SNPs half the probes can be pruned. Both full unpruned sets of probes and those reduced to the observed allele variants are provided.

## Results

The viruses analyzed have between 5 and 144 full length finished and draft genomes or segments available for analysis (Table 1). We found between 182 and 10,276 SNP loci for each organism. The detailed phylogenetic trees are shown after Table 1. Branch-specific SNP allele counts are plotted on nodes, strain/genome specific allele counts are given in brackets after strain name. For each target set, the tree that results in the most SNPs mapping to the tree is shown (that is, the tree that yields the fewest homoplastic SNPs).

Table 1. SNP analysis results summary.

| Organism | # target sequences | # SNP loci | # probes | # probes for observed alleles only | Strain resolution |
|---|---|---|---|---|---|
| Ebola | 22 | 3083 | 56,113 | 16,163 | All strains can be uniquely resolved |
| Marburg | 31 | 2922 | 65,413 | 18,951 | All strains can be uniquely resolved |
| CCHF L | 31 | 7216 | 233,005 | 62,837 | All strains can be uniquely resolved |
| CCHF M | 49 | 4720 | 144,741 | 38,729 | All strains can be uniquely resolved |
| CCHF S | 56 | 1073 | 40,437 | 11,145 | All strains can be uniquely resolved |
| RVF L | 62 | 1190 | 39,213 | 12,181 | All strains can be uniquely resolved |
| RVF M | 69 | 897 | 33,025 | 9,888 | All strains can be uniquely resolved |
| RVF S | 89 | 436 | 17,993 | 5,381 | Only 2 strains cannot be uniquely resolved on the basis of SNPs:<br><br>Rift_Valley_fever_200803166_segmentS_gi330422537<br>Rift_Valley_fever_Kenya_9800523_segmentS_gi87622508 |
| JEV | 144 | 6759 | 273,205 | 78,468 | All strains can be uniquely resolved |
| OW Arena L | 45 | 7556 | 173,157 | 44,031 | All strains can be uniquely resolved |
| OW Arena S | 54 | 4657 | 120,873 | 30,698 | There are 2 pairs of strains that cannot be uniquely resolved on the basis of SNPs:<br><br>Lassa_AV_GPC_and_nucleoprotein_NP_genes_gi46373061<br>Lassa_AV_gi354681510<br><br>Lassa_recombinant_Josiah_segmentS_gi32390330<br>Lassa_segmentS_gi23343509 |
| NW Arena S | 100 | 5410 | 139,505 | 37,249 | The following 2 pairs of strains cannot be uniquely resolved on the basis of SNP:<br><br>Machupo_Carvallo_segmentS_gi22901290<br>Machupo_Carvallo_segmentS_gi48095765<br><br>Pichinde_gi332639<br>Pichinde_gi55733698 |
| NW Arena L | 42 | 4389 | 80,973 | 20,800 | All strains can be uniquely resolved |
| Junin L | 12 | 182 | 2397 | 812 | All strains can be uniquely resolved |
| Junin S | 26 | 660 | 20,341 | 6,100 | All strains can be uniquely resolved |
| Machupo L | 5 | 373 | 6173 | 1719 | All strains can be uniquely resolved |
| Machupo S | 13 | 614 | 15,309 | 4,336 | The following 4 strains cannot be uniquely resolved on the basis of SNPs:<br><br>Machupo_Carvallo_segmentS_gi22901290<br>Machupo_Carvallo_segmentS_gi45826501<br>Machupo_Carvallo_segmentS_gi48095765<br>Machupo_segmentS_gi34365532 |
| Nipah | 9 | 684 | 10,573 | 3,090 | All strains can be uniquely resolved |
| Hendra | 10 | 437 | 5,841 | 2,280 | All strains can be uniquely resolved |
| TBEV | 67 | 10,276 | 288,269 | 78,225 | All strains can be uniquely resolved |

Abbreviations: CCHF=Crimean-Congo hemorrhagic fever, RVF=Rift Valley fever, JEV=Japanese encephalitis virus, NW Arena=New World Arenavirus, OW Arena=Old World Arenavirus, TBEV=tick-borne encephalitis virus.

Ebola

Number_SNPs: 3083
Number_Homoplastic_SNPs from MSA tree (shown below): 97
Number_Homoplastic_SNPs from SNP Hamming distance tree: 123
Number_Homoplastic_SNPs from SNP RAxML tree: 306
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
56,113 probes for all alleles
16,163 probes for observed alleles only

Marburg

Number_SNPs: 2922
Number_Homoplastic_SNPs from MSA tree: 1168
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 200
Number_Homoplastic_SNPs from SNP RAxML tree: 214
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
65,413 probes for all alleles
18,951 probes for observed alleles only

CCHF L segment

Number_SNPs: 7216
Number_Homoplastic_SNPs from MSA tree: 1422
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 1223
Number_Homoplastic_SNPs from SNP RAxML tree: 1270
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
233,005 probes for all alleles
62,837 probes for observed alleles only

CCHF M segment

Number_SNPs: 4720
Number_Homoplastic_SNPs from MSA tree: 680
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 585
Number_Homoplastic_SNPs from SNP RAxML tree: 791
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
144,741 probes for all alleles
38,729 probes for observed alleles only

CCHF S segment

Number_SNPs: 1073
Number_Homoplastic_SNPs from MSA tree: 294
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 266
Number_Homoplastic_SNPs from SNP RAxML tree: 326
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
40,437 probes for all alleles
11,145 probes for observed alleles only

RVF L segment

Number_SNPs: 1190
Number_Homoplastic_SNPs from MSA tree (shown below): 271
Number_Homoplastic_SNPs from SNP Hamming distance tree: 276
Number_Homoplastic_SNPs from SNP RAxML tree: 271
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
39,213 probes for all alleles
12,181 probes for observed alleles only

RVF M segment

Number_SNPs: 897
Number_Homoplastic_SNPs from MSA tree: 192
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 189
Number_Homoplastic_SNPs from SNP RAxML tree: 191
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
33,025 probes for all alleles
9,888 probes for observed alleles only

RVF S segment

Number_SNPs: 436
Number_Homoplastic_SNPs from MSA tree: 107
Number_Homoplastic_SNPs from SNP Hamming distance tree: 110
Number_Homoplastic_SNPs from SNP RAxML tree (shown below): 106
Only 2 strains cannot be uniquely resolved on the basis of SNPs:
Rift_Valley_fever_200803166_segmentS_gi330422537
Rift_Valley_fever_Kenya_9800523_segmentS_gi87622508
Number of SNP microarray probes:
17,993 probes for all alleles
5,381 probes for observed alleles only

JEV

Number_SNPs: 6759
Number_Homoplastic_SNPs from MSA tree: 2050
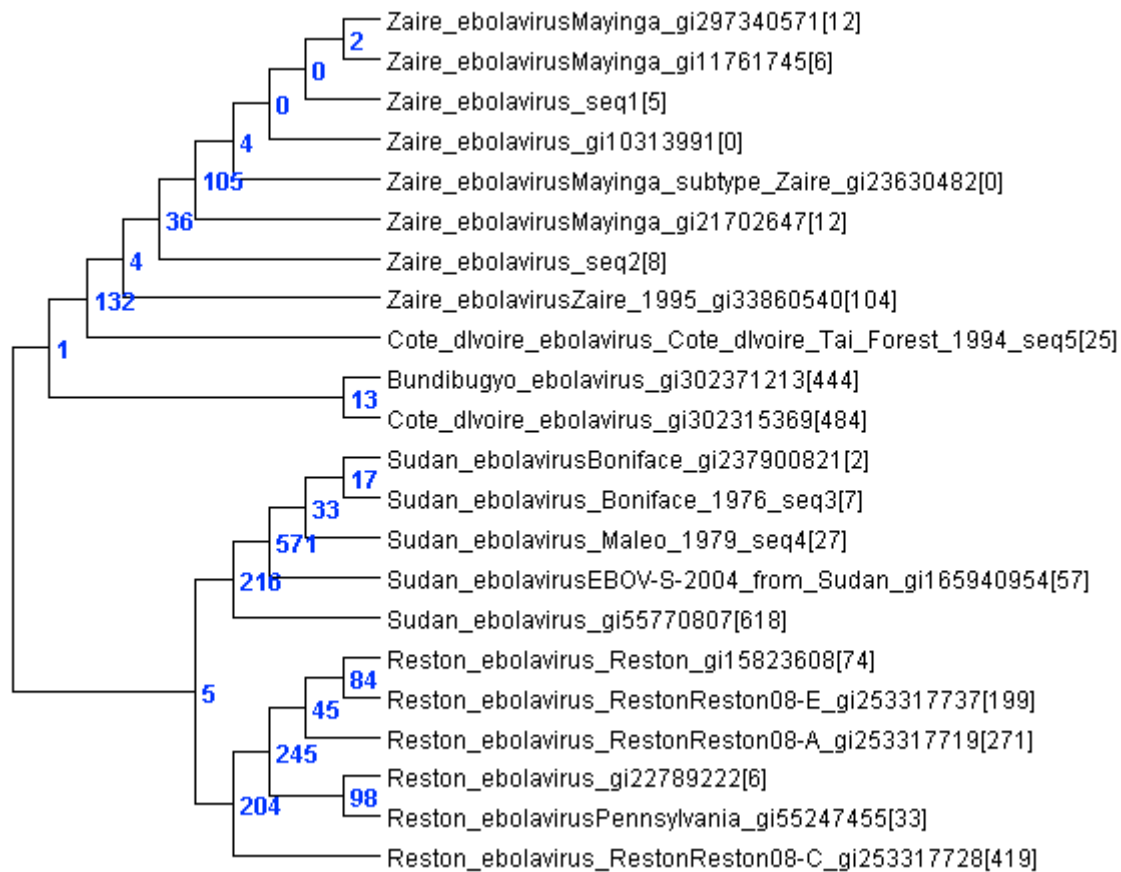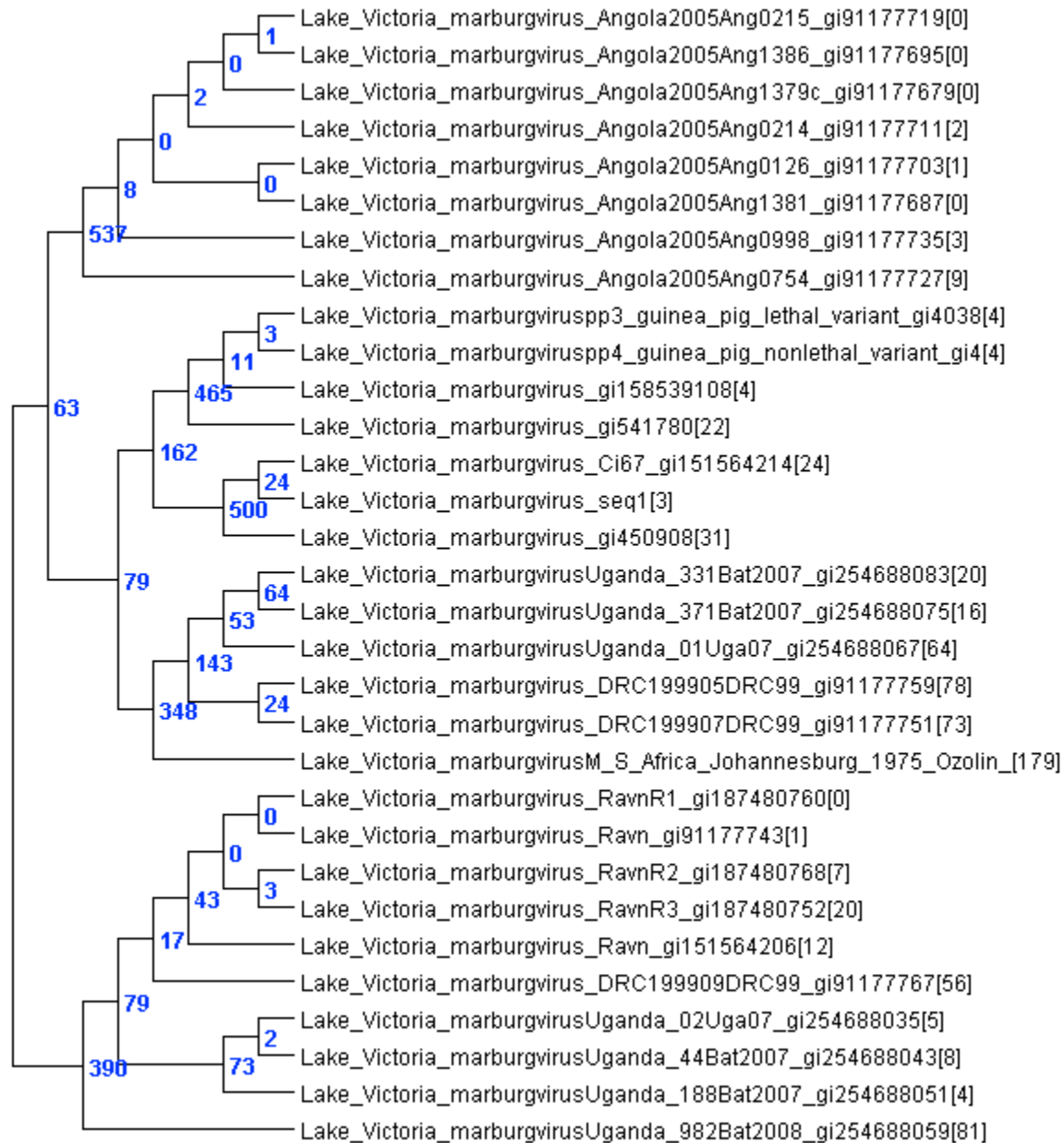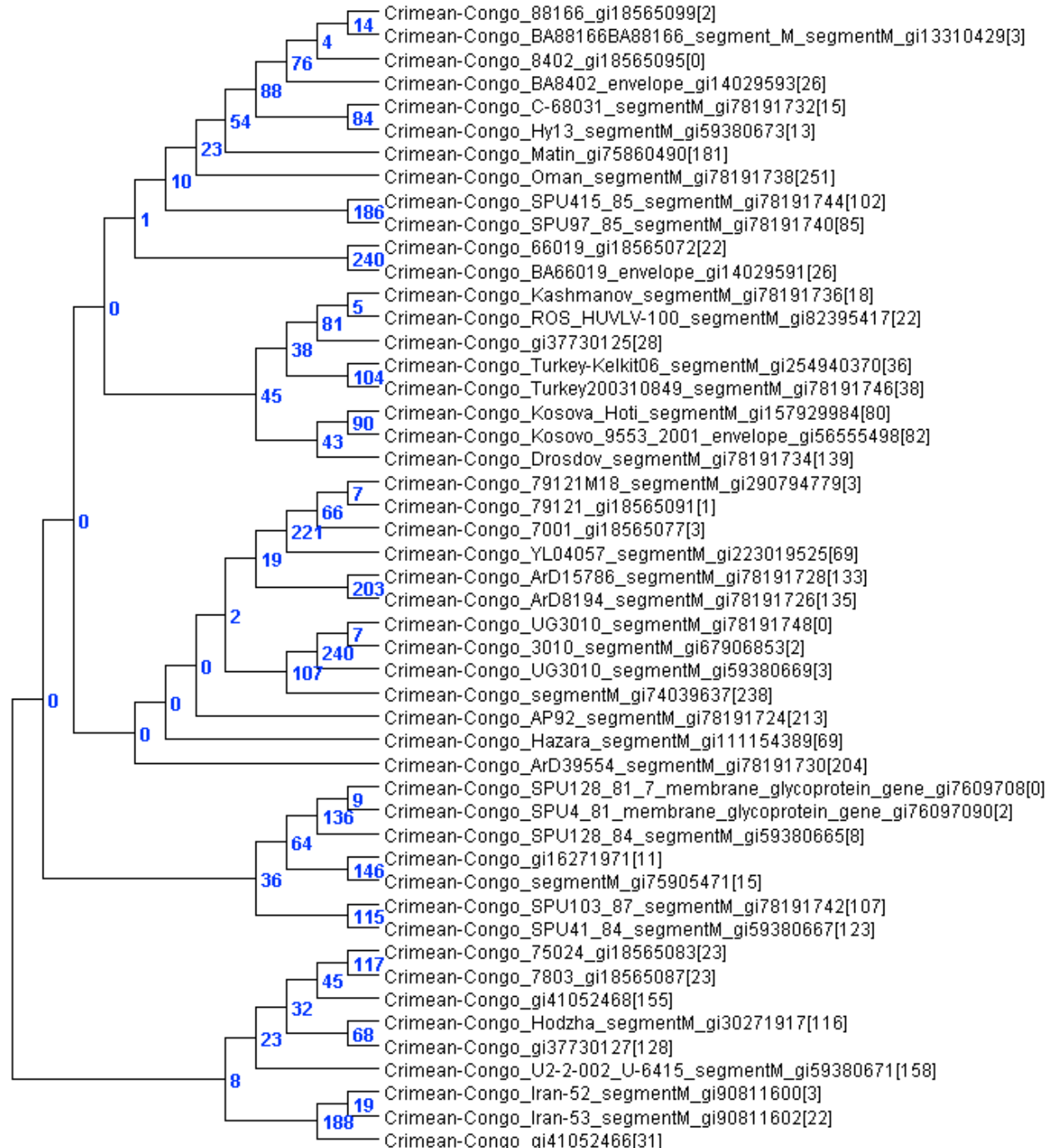Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 2018
Number_Homoplastic_SNPs from SNP RAxML tree: 2139
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
273,205 probes for all alleles
78,468 probes for observed alleles only

Japanese_encephalitis_DL04-29_gi347361102[3]
Japanese_encephalitis_SH0410_gi347301240[0]
Japanese_encephalitis_GZ042_gi347301242[0]
Japanese_encephalitis_JH0418_gi347301238[3]
Japanese_encephalitis_FJ0394_gi347301244[9]
Japanese_encephalitis_Fj02-29_gi347361104[3]
Japanese_encephalitis_HLJ02-134_gi347361110[7]
Japanese_encephalitis_FJ0339_gi347301246[2]
Japanese_encephalitis_CBH_gi347301248[1]
Japanese_encephalitis_YN98A151_gi347301250[11]
Japanese_encephalitis_LFM_gi347301254[32]
Japanese_encephalitis_ZSZ_gi347301252[3]
Japanese_encephalitis_SH3_gi347301256[9]
Japanese_encephalitis_SH045_gi347301260[4]
Japanese_encephalitis_Fj0276_gi347301262[6]
Japanese_encephalitis_HW_gi57790540[13]
Japanese_encephalitis_WHe_gi119067955[14]
Japanese_encephalitis_DL0445_gi347301236[0]
Japanese_encephalitis_HYZ_gi347301234[34]
Japanese_encephalitis_p3_gi1488030[15]
Japanese_encephalitis_CHT392_complete_genome_gi20563126[3]
Japanese_encephalitis_T1P1_complete_genome_gi20563128[7]
Japanese_encephalitis_T1P1-S1_gi34495370[2]
Japanese_encephalitis_T1P1-L4_gi34495372[1]
Japanese_encephalitis_YL_polyprotein_mRNA_gi19698517[31]
Japanese_encephalitis_gi56744200[20]
Japanese_encephalitis_gi56744202[3]
Japanese_encephalitis_gi56744204[5]
Japanese_encephalitis_gi2318126[6]
Japanese_encephalitis_gi2318128[6]
Japanese_encephalitis_gi4416166[11]
Japanese_encephalitis_GSS_gi347361108[2]
Japanese_encephalitis_NJ_2008_gi296802975[15]
Japanese_encephalitis_JaOH0566_Japan_1966_human_gi45934772[12]
Japanese_encephalitis_ML17_live_vaccine_Japan_1981_human_gi459[1
Japanese_encephalitis_YLG_gi347361118[10]
Japanese_encephalitis_HVI_gi4323271[23]
Japanese_encephalitis_TC_gi4323273[23]
Japanese_encephalitis_TL_gi4323275[40]
Japanese_encephalitis_SAT4-14-2_gi12964700[12]
Japanese_encephalitis_SA_A_genomic_RNA_gi221958[4]
Japanese_encephalitis_SAT4-12-1-7_gi15788965[18]
Japanese_encephalitis_gi687314[18]
Japanese_encephalitis_gi331331[4]
Japanese_encephalitis_SA_V_genomic_RNA_gi221960[4]
Japanese_encephalitis_gi537634[4]
Japanese_encephalitis_CH13_gi347301268[18]
Japanese_encephalitis_04940-4_gi156754270[7]
Japanese_encephalitis_057434_gi156754268[0]
Japanese_encephalitis_GP78_gi3342805[24]
Japanese_encephalitis_KPP82-39-214CT_gi307826673[73]
Japanese_encephalitis_IGIB-NIV-2009-01_gi347597909[17]
Japanese_encephalitis_014178_gi156754266[9]
Japanese_encephalitis_YN_gi347301270[17]
Japanese_encephalitis_JaTAn1_90_gi292680472[16]
Japanese_encephalitis_JaTAn2_91_gi292680474[26]
Japanese_encephalitis_K87P39_gi46519718[4]
Japanese_encephalitis_gi242346715[4]
Japanese_encephalitis_DH107_gi347301274[76]
Japanese_encephalitis_47_gi347361096[12]
Japanese_encephalitis_Ha3_gi347301272[10]
Japanese_encephalitis_CC27-L1_gi34495378[3]
Japanese_encephalitis_CC27-S8_gi34495384[2]
Japanese_encephalitis_CC27-L3_gi34495380[1]
Japanese_encephalitis_CC27-S6_gi34495382[1]
Japanese_encephalitis_CH2195LA_gi6970067[5]
Japanese_encephalitis_CH2195SA_gi6970069[1]
Japanese_encephalitis_CJN-L1_gi34495376[1]
Japanese_encephalitis_CJN-S1_gi34495374[3]
Japanese_encephalitis_JaTAn1_75_gi292680470[43]
Japanese_encephalitis_SH0601_gi146289960[21]
Japanese_encephalitis_gi9626460[30]
Japanese_encephalitis_LYZ_gi347301266[5]
Japanese_encephalitis_ZMT_gi347361124[46]
Japanese_encephalitis_CZC_gi347301258[33]
Japanese_encephalitis_B58_gi224223449[0]
Japanese_encephalitis_GB30_gi224223451[1]
Japanese_encephalitis_HB49_gi347361126[8]
Japanese_encephalitis_HB97_gi347361128[0]
Japanese_encephalitis_Nakayama_gi148009209[32]
Japanese_encephalitis_Beijing-1_gi1066797[14]
Japanese_encephalitis_Ling_gi18653908[16]
Japanese_encephalitis_TLA_gi347301264[19]
Japanese_encephalitis_Vellore_P20778_gi3406734[157]
Japanese_encephalitis_Muar_gi323145096[500]
Japanese_encephalitis_XZ0934_gi340807378[525]
Japanese_encephalitis_JKT6468_polyprotein_gene_gi30959382[521]
Japanese_encephalitis_HN0411_gi347301190[6]
Japanese_encephalitis_SC04-17_gi270272092[4]
Japanese_encephalitis_GX0523_44_gi347301192[25]
Japanese_encephalitis_GSBY0861_gi347301194[0]
Japanese_encephalitis_YN0967_gi347361094[12]
Japanese_encephalitis_HN0621_gi347301188[32]
Japanese_encephalitis_HN06129_gi347361112[12]
Japanese_encephalitis_YN82BN8219_gi347301196[10]
Japanese_encephalitis_131V_gi300089373[31]
Japanese_encephalitis_YN0623_gi347301200[20]
Japanese_encephalitis_YN0911_gi347361092[30]
Japanese_encephalitis_GX0519_gi347301198[24]
Japanese_encephalitis_HN0626_gi347301202[30]
Japanese_encephalitis_SC0415_gi347301204[36]
Japanese_encephalitis_GSBY0801_gi347361106[0]
Japanese_encephalitis_GSBY0810_gi347301208[5]
Japanese_encephalitis_HEN0701_gi218963162[19]
Japanese_encephalitis_SC0412_gi347301206[41]
Japanese_encephalitis_HN0421_gi347301210[10]
Japanese_encephalitis_JEV_sw_Mie_40_2004_gi81687247[18]
Japanese_encephalitis_GSBY0816_gi347301212[27]
Japanese_encephalitis_GZ56_gi313105111[25]
Japanese_encephalitis_GSBY0804_gi347301216[5]
Japanese_encephalitis_GSBY0827_gi347301218[5]
Japanese_encephalitis_GS07TS11_gi347301214[52]
Japanese_encephalitis_90VN70_gi302321717[21]
Japanese_encephalitis_SH53_gi347301228[27]
Japanese_encephalitis_YN79Bao83_gi347301230[32]
Japanese_encephalitis_BL06-54_gi347361100[21]
Japanese_encephalitis_YN05155_gi347301232[31]
Japanese_encephalitis_BL06-50_gi347361098[47]
Japanese_encephalitis_LN02-102_gi347361114[9]
Japanese_encephalitis_SH80_gi347301224[32]
Japanese_encephalitis_SH03103_gi347301222[5]
Japanese_encephalitis_SH03105_gi347301220[0]
Japanese_encephalitis_gi12082323[8]
Japanese_encephalitis_YN05124_gi347361120[15]
Japanese_encephalitis_JEV_sw_Mie_41_2002_gi81687252[23]
Japanese_encephalitis_SHH7M-07_gi16880524[25]
Japanese_encephalitis_XJP613_gi189086642[31]
Japanese_encephalitis_JEV_eq_Tottori_2003_gi329564799[17]
Japanese_encephalitis_JX6T_gi290759657[25]
Japanese_encephalitis_SD0810_gi347361130[17]
Japanese_encephalitis_SX09S-01_gi324985238[26]
Japanese_encephalitis_LN0716_gi347301226[22]
Japanese_encephalitis_XJ69_gi195970353[19]
Japanese_encephalitis_KV1899_gi321187331[55]
Japanese_encephalitis_YN83-Meng83-54_gi347361122[1]
Japanese_encephalitis_XZ0938_gi337263366[118]
Japanese_encephalitis_K94P05_gi5231232[95]
Japanese_encephalitis_1070_82_Subin_gi307826665[2]
Japanese_encephalitis_B-0860_82_gi307826663[12]
Japanese_encephalitis_3KPUCV569_gi307826667[45]
Japanese_encephalitis_4790-85_gi307826671[67]
Japanese_encephalitis_B-1381-85_gi307826669[54]
Japanese_encephalitis_M28_gi347361116[45]
Japanese_encephalitis_gi8163805[461]

14

Old World Arenaviruses L segment

Number_SNPs: 7556
Number_Homoplastic_SNPs from MSA tree: 698
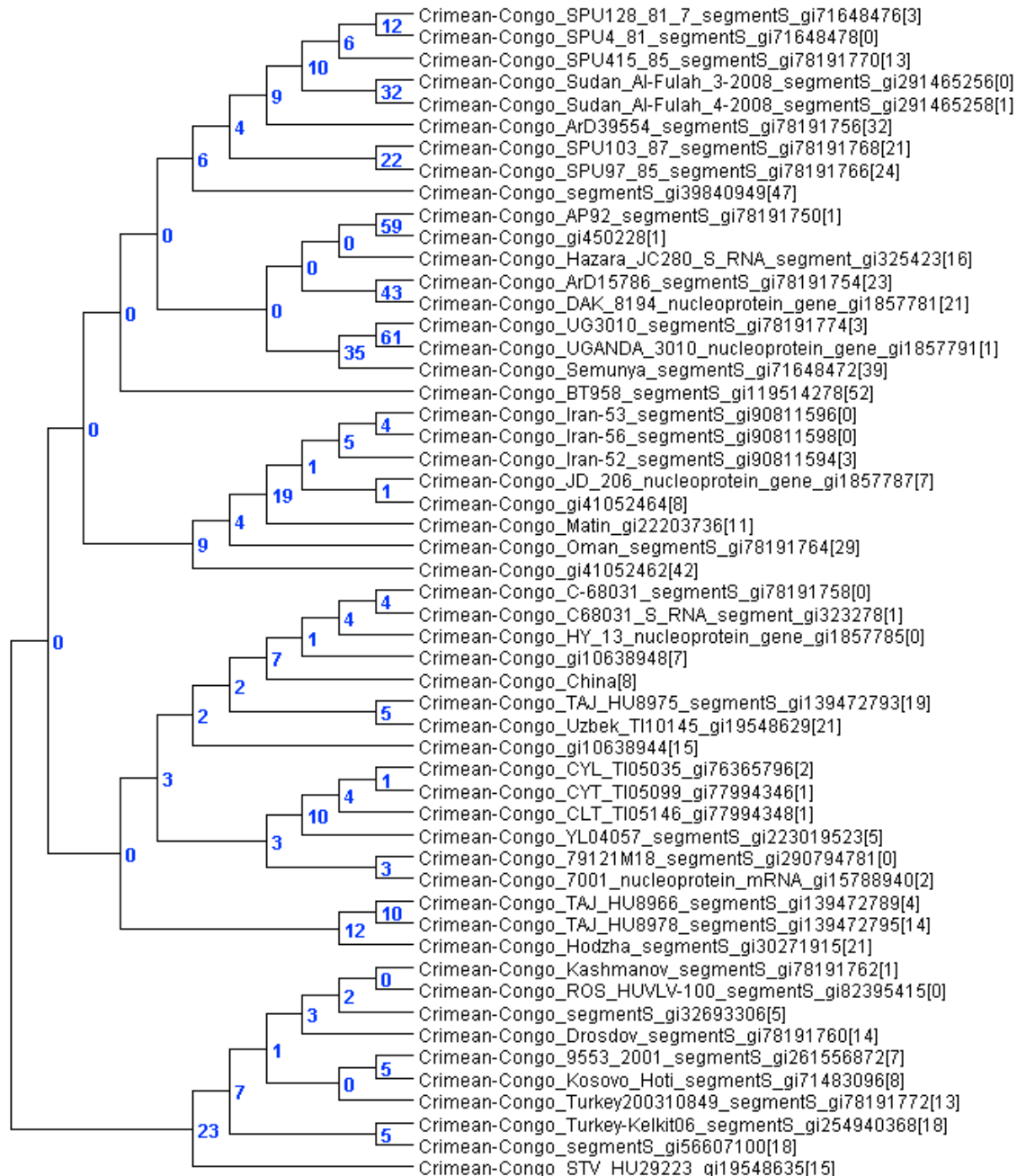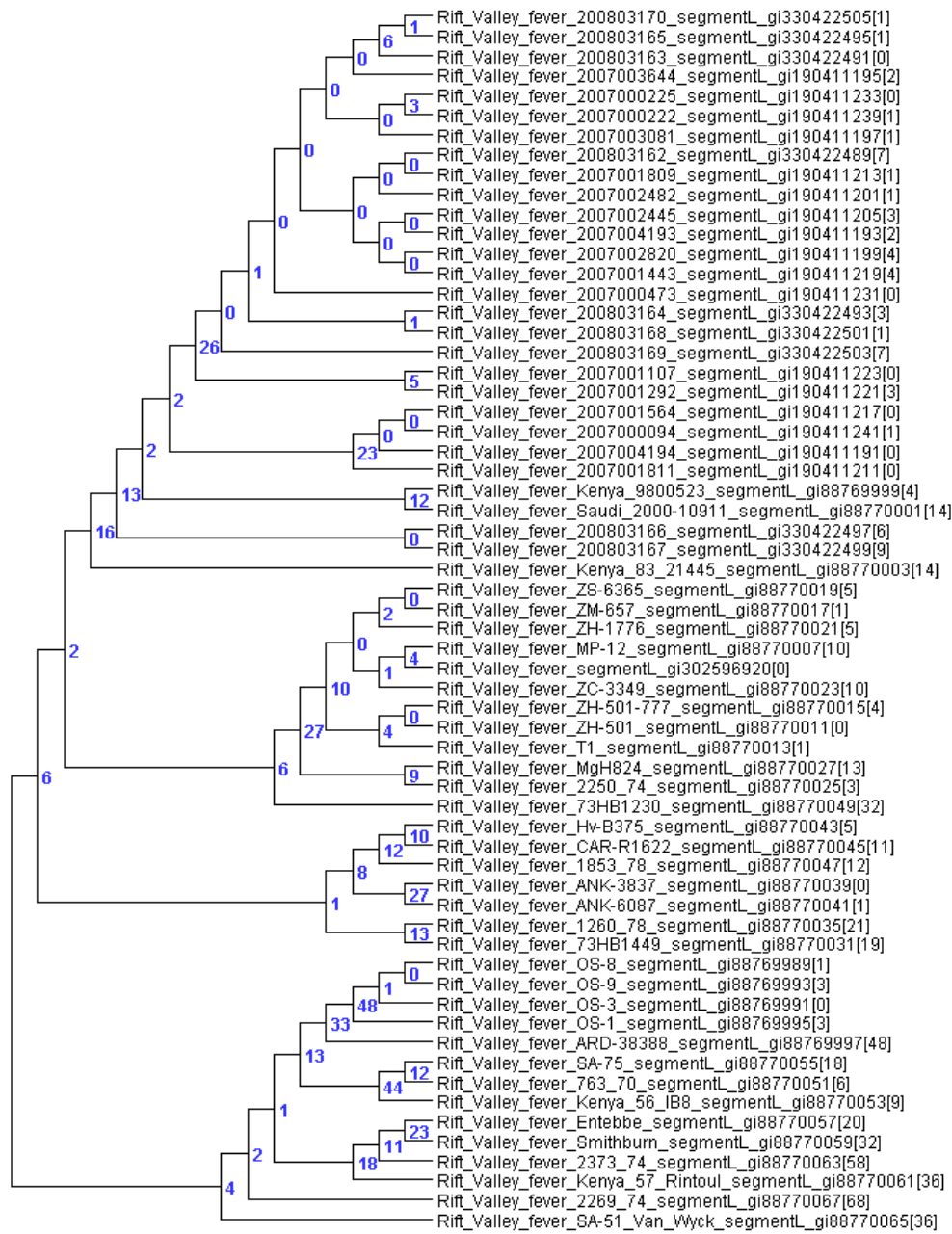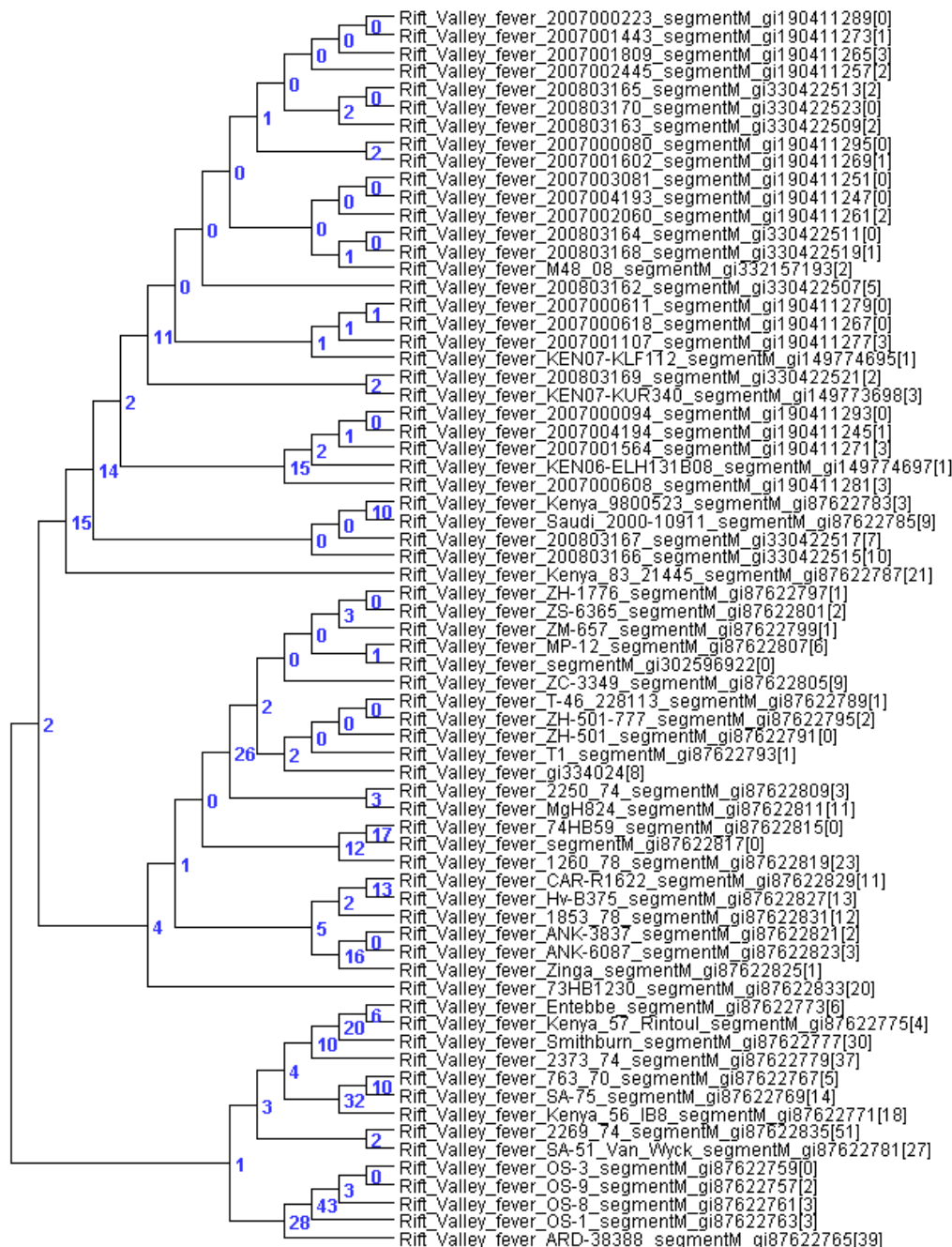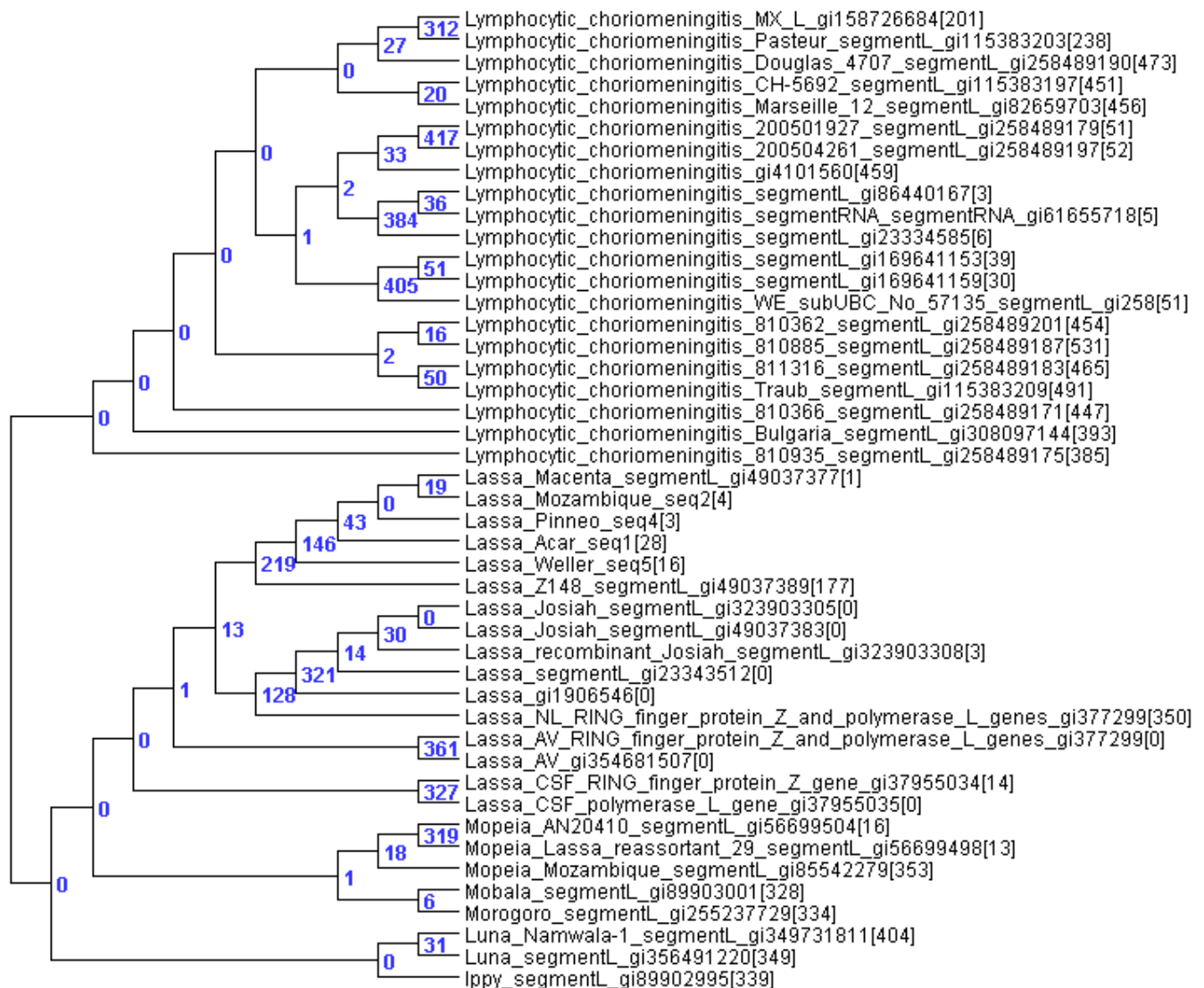Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 417
Number_Homoplastic_SNPs from SNP RAxML tree: 606
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
173,157 probes for all alleles
44,031 probes for observed alleles only

Old World Arenaviruses S segment

Number_SNPs: 4657
Number_Homoplastic_SNPs from MSA tree: 602
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 484
Number_Homoplastic_SNPs from SNP RAxML tree: 576
There are 2 pairs of strains that cannot be uniquely resolved on the basis of SNPs:

1   Lassa_AV_GPC_and_nucleoprotein_NP_genes_gi46373061
1   Lassa_AV_gi354681510
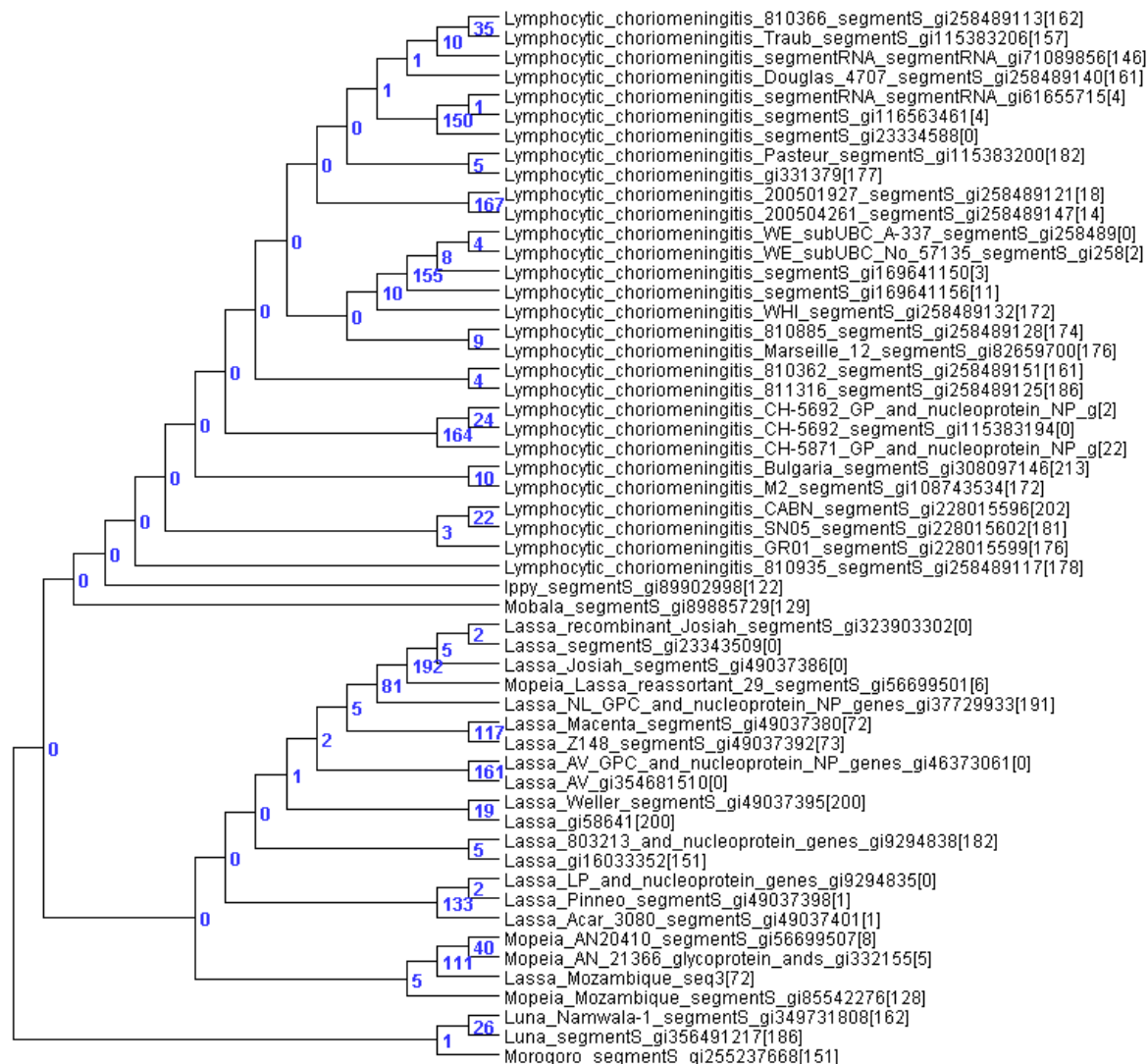
2   Lassa_recombinant_Josiah_segmentS_gi323903302
2   Lassa_segmentS_gi23343509

Number of SNP microarray probes:
120,873 probes for all alleles
30,698 probes for observed alleles only

Junin L segment

Number_SNPs: 182
Number_Homoplastic_SNPs from MSA tree: 7
Number_Homoplastic_SNPs from SNP Hamming distance tree: 5
Number_Homoplastic_SNPs from SNP RAxML tree (shown below): 4
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
2397 probes for all alleles
812 probes for observed alleles only



```
        ┌─ Junin_Candid_1_segmentL_gi227957900[2]
      0 │
    ┌─0─┤  Junin_Candid-1_segmentL_gi52222818[8]
    │ 3 └─ Junin_Candid1-wt_segmentL_gi328679013[0]
  ┌─3─┤
  │   └──── Junin_Candid1-rec_segmentL_gi328679019[2]
 ─3─┤
  │ └───── Junin_XJ44_segmentL_gi227958005[5]
  1 │
    └───── Junin_XJ39_segmentL_gi342360491[0]
        └─ Junin_XJ34_segmentL_gi342360488[1]
 0
          ┌─ Junin_MC2_segmentL_gi164519644[1]
      145 │
    ┌─4───┤ Junin_Rumero_segmentL_gi48095756[0]
 0  │     └─ Junin_rXJ13_segmentL_gi226374508[0]
    │
    ├──── Junin_XJ13_segmentL_gi226374502[0]
    └──── Junin_segmentL_gi34365526[0]
```

Junin S segment

Number_SNPs: 660
Number_Homoplastic_SNPs from MSA tree: 162
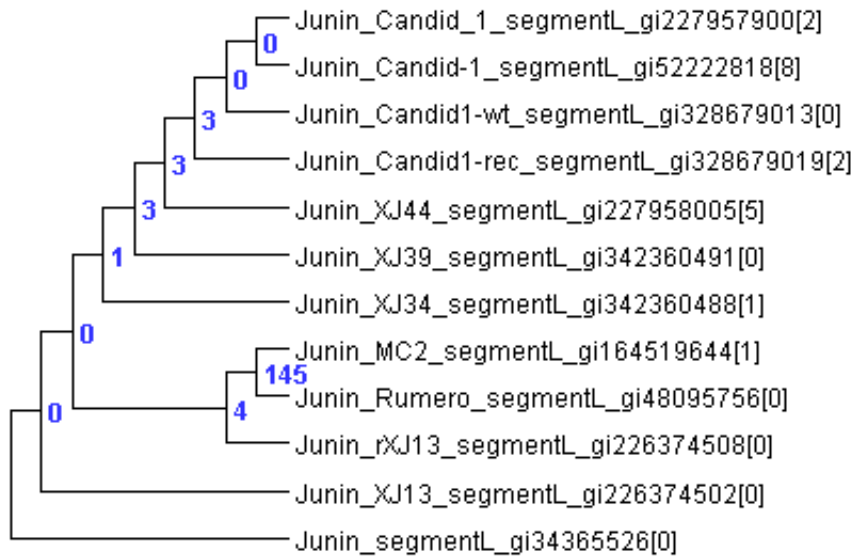Number_Homoplastic_SNPs from SNP Hamming distance tree: 161
Number_Homoplastic_SNPs from SNP RAxML tree (shown below): 160
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
20,341 probes for all alleles
6,100 probes for observed alleles only

Machupo L segment

Number_SNPs: 373
Number_Homoplastic_SNPs from MSA tree (shown below): 0
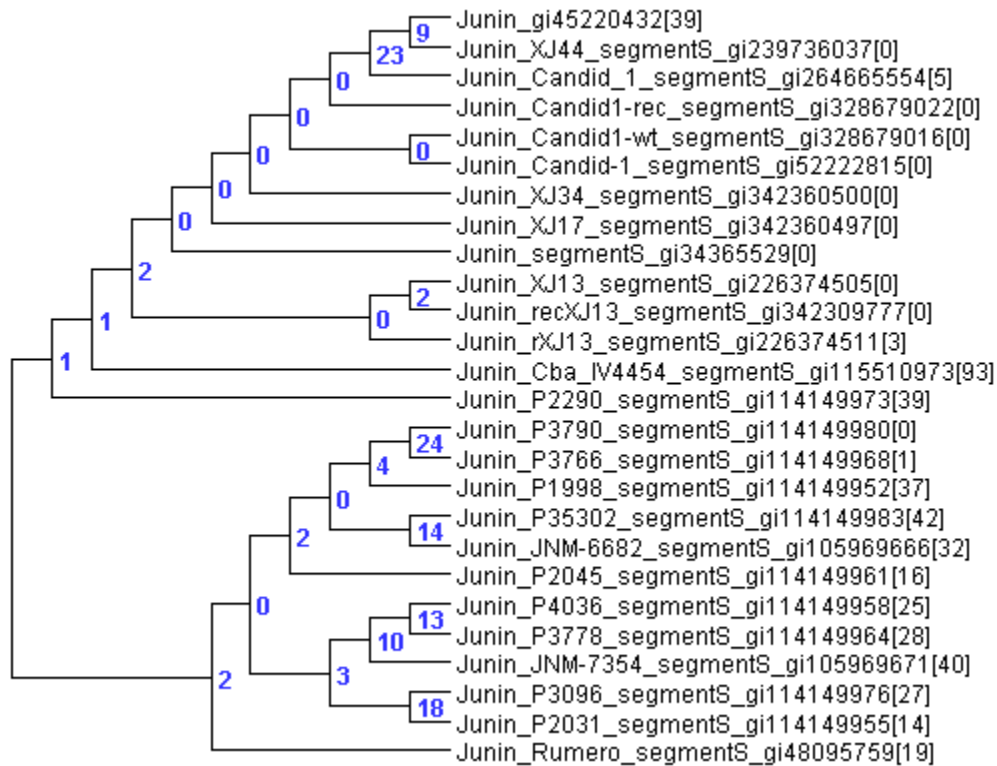Number_Homoplastic_SNPs from SNP Hamming distance tree: 0
Number_Homoplastic_SNPs from SNP RAxML tree: 0
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
6173 probes for all alleles
1719 probes for observed alleles only

Machupo S segment

Number_SNPs: 614
Number_Homoplastic_SNPs from MSA tree (shown below): 54
Number_Homoplastic_SNPs from SNP Hamming distance tree: 54
Number_Homoplastic_SNPs from SNP RAxML tree: 56
The following 4 strains cannot be uniquely resolved on the basis of SNPs:
1      Machupo_Carvallo_segmentS_gi22901290
1      Machupo_Carvallo_segmentS_gi45826501
1      Machupo_Carvallo_segmentS_gi48095765
1      Machupo_segmentS_gi34365532


Number of SNP microarray probes:
15,309 probes for all alleles
4,336 probes for observed alleles only

New World Arenaviruses L segment (including Machupo, Junin, and others)

Number_SNPs: 4389
Number_Homoplastic_SNPs from MSA tree: 79
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 50
Number_Homoplastic_SNPs from SNP RAxML tree: 233
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
80,973 probes for all alleles
20,800 probes for observed alleles only

New World Arenaviruses S segment (including Machupo, Junin, and others)

Number_SNPs: 5410
Number_Homoplastic_SNPs from MSA tree: 654
Number_Homoplastic_SNPs from SNP Hamming distance tree (shown below): 521
Number_Homoplastic_SNPs from SNP RAxML tree: 724
The following 2 pairs of strains cannot be uniquely resolved on the basis of SNP:
1       Machupo_Carvallo_segmentS_gi22901290
1       Machupo_Carvallo_segmentS_gi48095765
2       Pichinde_gi332639
2       Pichinde_gi55733698

Number of SNP microarray probes:
139,505 probes for all alleles
37,249 probes for observed alleles only

Nipah

Number_SNPs: 684
Number_Homoplastic_SNPs from MSA tree (shown below): 0
Number_Homoplastic_SNPs from SNP Hamming distance tree: 0
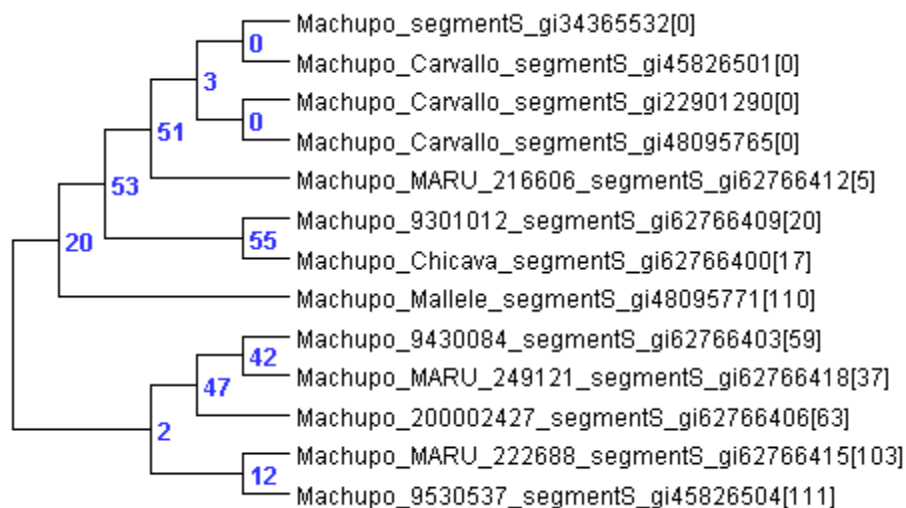Number_Homoplastic_SNPs from SNP RAxML tree: 0
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
10,573 probes for all alleles
3,090 probes for observed alleles only

Hendra

Number_SNPs: 437
Number_Homoplastic_SNPs from MSA tree (shown below): 37
Number_Homoplastic_SNPs from SNP Hamming distance tree: 38
Number_Homoplastic_SNPs from SNP RAxML tree: 37
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
5,841 probes for all alleles
2,280 probes for observed alleles only

Tick-borne encephalitis virus complex

Number_SNPs: 10,276
Number_Homoplastic_SNPs from MSA tree (shown below): 1260
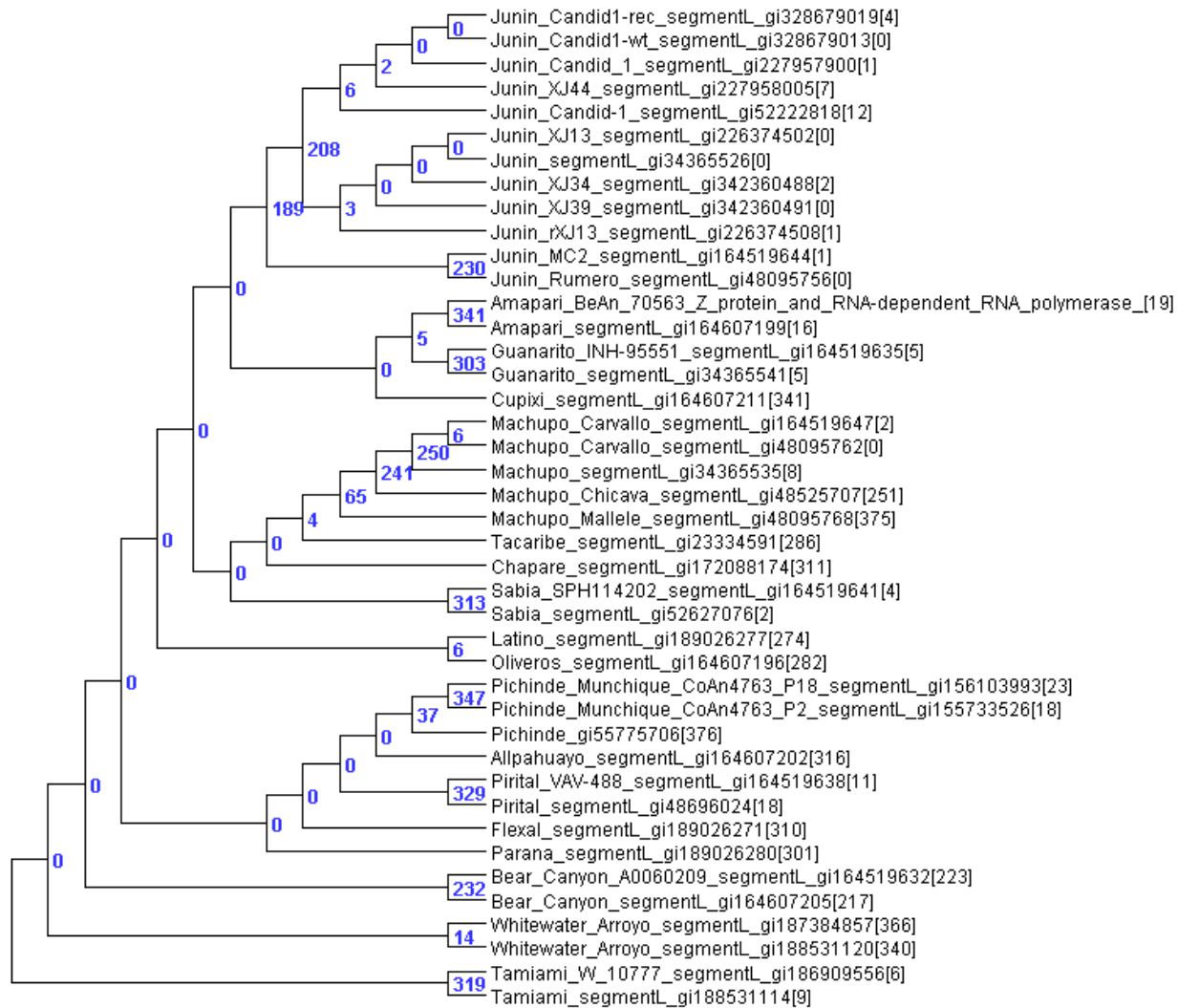Number_Homoplastic_SNPs from SNP Hamming distance tree: 1272
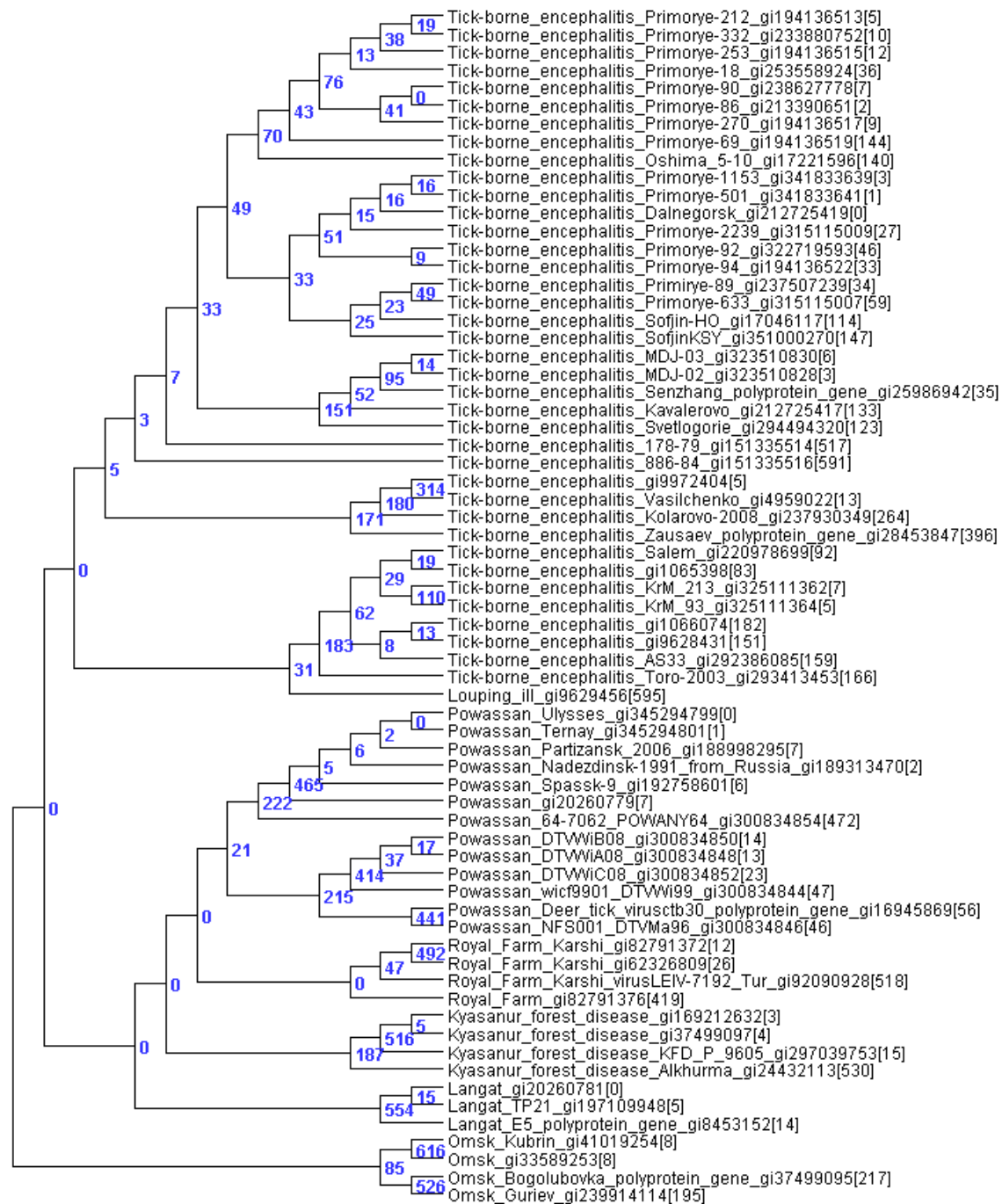Number_Homoplastic_SNPs from SNP RAxML tree: 1759
All strains can be uniquely resolved on the basis of SNPs.
Number of SNP microarray probes:
288,269 probes for all alleles
78,225 probes for observed alleles only

Tick-borne_encephalitis_Primorye-212_gi194136513[5]
Tick-borne_encephalitis_Primorye-332_gi233880752[10]
Tick-borne_encephalitis_Primorye-253_gi194136515[12]
Tick-borne_encephalitis_Primorye-18_gi253558924[36]
Tick-borne_encephalitis_Primorye-90_gi238627778[7]
Tick-borne_encephalitis_Primorye-86_gi213390651[2]
Tick-borne_encephalitis_Primorye-270_gi194136517[9]
Tick-borne_encephalitis_Primorye-69_gi194136519[144]
Tick-borne_encephalitis_Oshima_5-10_gi17221596[140]
Tick-borne_encephalitis_Primorye-1153_gi341833639[3]
Tick-borne_encephalitis_Primorye-501_gi341833641[1]
Tick-borne_encephalitis_Dalnegorsk_gi212725419[0]
Tick-borne_encephalitis_Primorye-2239_gi315115009[27]
Tick-borne_encephalitis_Primorye-92_gi322719593[46]
Tick-borne_encephalitis_Primorye-94_gi194136522[33]
Tick-borne_encephalitis_Primirye-89_gi237507239[34]
Tick-borne_encephalitis_Primorye-633_gi315115007[59]
Tick-borne_encephalitis_Sofjin-HO_gi17046117[114]
Tick-borne_encephalitis_SofjinKSY_gi351000270[147]
Tick-borne_encephalitis_MDJ-03_gi323510830[6]
Tick-borne_encephalitis_MDJ-02_gi323510828[3]
Tick-borne_encephalitis_Senzhang_polyprotein_gene_gi25986942[35]
Tick-borne_encephalitis_Kavalerovo_gi212725417[133]
Tick-borne_encephalitis_Svetlogorie_gi294494320[123]
Tick-borne_encephalitis_178-79_gi151335514[517]
Tick-borne_encephalitis_886-84_gi151335516[591]
Tick-borne_encephalitis_gi9972404[5]
Tick-borne_encephalitis_Vasilchenko_gi4959022[13]
Tick-borne_encephalitis_Kolarovo-2008_gi237930349[264]
Tick-borne_encephalitis_Zausaev_polyprotein_gene_gi28453847[396]
Tick-borne_encephalitis_Salem_gi220978699[92]
Tick-borne_encephalitis_gi1065398[83]
Tick-borne_encephalitis_KrM_213_gi325111362[7]
Tick-borne_encephalitis_KrM_93_gi325111364[5]
Tick-borne_encephalitis_gi1066074[182]
Tick-borne_encephalitis_gi9628431[151]
Tick-borne_encephalitis_AS33_gi292386085[159]
Tick-borne_encephalitis_Toro-2003_gi293413453[166]
Louping_ill_gi9629456[595]
Powassan_Ulysses_gi345294799[0]
Powassan_Ternay_gi345294801[1]
Powassan_Partizansk_2006_gi188998295[7]
Powassan_Nadezdinsk-1991_from_Russia_gi189313470[2]
Powassan_Spassk-9_gi192758601[6]
Powassan_gi20260779[7]
Powassan_64-7062_POWANY64_gi300834854[472]
Powassan_DTVWiB08_gi300834850[14]
Powassan_DTVWiA08_gi300834848[13]
Powassan_DTVWiC08_gi300834852[23]
Powassan_wicf9901_DTVWi99_gi300834844[47]
Powassan_Deer_tick_virusctb30_polyprotein_gene_gi16945869[56]
Powassan_NFS001_DTVMa96_gi300834846[46]
Royal_Farm_Karshi_gi82791372[12]
Royal_Farm_Karshi_gi62326809[26]
Royal_Farm_Karshi_virusLEIV-7192_Tur_gi92090928[518]
Royal_Farm_gi82791376[419]
Kyasanur_forest_disease_gi169212632[3]
Kyasanur_forest_disease_gi37499097[4]
Kyasanur_forest_disease_KFD_P_9605_gi297039753[15]
Kyasanur_forest_disease_Alkhurma_gi24432113[530]
Langat_gi20260781[0]
Langat_TP21_gi197109948[5]
Langat_E5_polyprotein_gene_gi8453152[14]
Omsk_Kubrin_gi41019254[8]
Omsk_gi33589253[8]
Omsk_Bogolubovka_polyprotein_gene_gi37499095[217]
Omsk_Guriev_gi239914114[195]

# References

1.    Gardner S, Slezak T: **Scalable SNP Analyses of 100+ Bacterial or Viral Genomes**. *J Forensic Res* 2010, **1**:107, doi:110.4172/2157-7145.1000107.
2.    Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy and high throughput**. *Nucleic Acids Res* 2004, **32**:1792-1797.
3.    Stamatakis A: **RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models**. *Bioinformatics* 2006, **22**:2688-2690.