# Exploration of Shared SNPs in Thaps
## trunc-unfiltered

July 26, 2017

Rambling exploration of SNP positions shared between two or more of the isolates. Code is included to document it thoroughly, (even if largely uninteresting to anyone else), and I will summarize it as I go.

# Contents

# 1 History

This was added to SVN 1/26/2014; not sure when it was started, but earliest related emails I see are from 1/21/14.

```
r413 | ruzzo | 2014-01-26 08:22:37 -0800 (Sun, 26 Jan 2014) | 2 lines

adding shared-snp analysis.
```

# 2 Preliminaries

NOTE: Some comments in code and some parts of the text, especially specific numbers and general conclusions, are based on Unqfiltered, Chr1, Medium stringency (i.e., "[[2]]" below) analysis. The broad picture does not appear to change with other choices, but details do, and the text is neither fully parameterized nor fully updated, so proceed with caution.

Load utility R code; do setup:

```r
source('../../../R/wlr.R') # load util code; path relative this folder or sibling in scripts/larrys

## Running as: ruzzo @ bicycle.cs.washington.edu; SVN Id, I miss you.  $Id: wlr.R  2017-07-21 or later $

setup.my.wd('shared-snps') # set working dir; UPDATE if this file moves, or if COPY/PASTE to new file
setup.my.knitr('figs-knitr/')
generic.setup('figs-mine/')
```

# 3 Major Analysis/Performance Parameters.

Choices here control how this file is processed, what data is analyzed, speed, etc. Set them carefully before running "make." Major choices are:

1. WHICH SNP TABLES ARE LOADED??? The logical vector `load.tb` selects the desired combination of SNP tables to load, in the order `full.unfiltered, chr1.unfiltered, full.qfiltered, chr1.qfiltered`. E.g., `load.tb=(T, F, T, F)` loads *full* tables for *both* q- and un-qfiltered data. Primary analysis is only performed on one of them, but the others are retained for comparison/debugging.

2. WHICH MAIN ANALYSIS??? If multiple tables are loaded, which is used for the main analysis? Parameter `pri` is a permutation of 1:4, corresponding to `load.tb`; the first loaded table in that order becomes the analysis focus. The default `pri=c(1,2,3,4)` looks at un-q-filtered data in preference to q-filtered, and full tables in preference to Chr1 within each group.

   (Choice of data for the "Table 1" coverage summary in section 5 is independent of this; full genome data is prefered over Chr 1 for both q- and unq-filtered reads; change `tset.picker` calls near the end of that section to modify this.)

3. CLEAR CACHE??? `clear.cache=T` forces Knitr cache removal at the start of the run; especially important if the previous parameters have changed since the last run.

4. HOW MANY BOOTSTRAP REPLICATES??? The variable `nboot` is a major performance factor; 1000 reps takes several hours. Set to 5 for debug and quick look; 100 or more for final run.

5. TRUNCATE TABLES TO Chrs ONLY??? I.e., remove mitochondrial-, plastid-, and BD- contigs.

The following code chunk sets the first four parameters based on where it's run. To prototype/debug on a laptop, faster is better—run on Chr1 with small `nboot`; when run on the linux servers, I typically do full genomes, more replicates. Just override them if these defaults don't work for you.

```r
# for Makefile, params can be command line args, else base on system; see wlr.r for details.
# load.tb order: full.un, chr1.un, full.qfil,  chr1.qfil

params <- pick.params(
  mac   = list(load.tb=c(F,T,F,F), pri=1:4, clear.cache=F, nboot=  1, trunc.tables=T), # quick on lap
 #linux = list(load.tb=c(F,F,F,T), pri=1:4, clear.cache=F, nboot=  5, trunc.tables=T), # quick qfil on server
  linux = list(load.tb=c(T,F,T,F), pri=1:4, clear.cache=T, nboot=101, trunc.tables=T)  # full on server
)

# Alternatively, edit/uncomment the following to override the above as needed
#params<-pick.params(default=list(load.tb=c(T,T,T,T),pri=1:4,clear.cache=T,nboot=1000,trunc.tables=T))
print(params)

# $load.tb
# full.unf chr1.unf  full.qf  chr1.qf
#     TRUE    FALSE     TRUE    FALSE
#
# $pri
# [1] 1 2 3 4
#
# $clear.cache
# [1] TRUE
#
# $nboot
# [1] 101
#
# $trunc.tables
# [1] TRUE
```

CLEAR CACHE??!! Some code chunks use the knitr cache, but extent of cache consistency checks unknown. If in doubt, delete "cache/" (knitr's) directory to force rebuild. T/F set in params above will/won't force removal (actually, rename):

```r
decache(params$clear.cache)

# Rename of 'cache' to 'cache86293' returned TRUE .
```

If still in doubt, also manually remove "00common/mycache/" (mine).
Load the main SNP data file(s) based on the parameters set in section 3.

```r
# short names to keep the following chunk compact
tb <- params$load.tb
tset <- list(NULL, NULL, NULL, NULL) # tset = 'table set'
```

```r
# see wlr.R for load paths
if(tb[1]){tset[[1]] <- load.snp.tables(use.chr1.tables = FALSE, data.name='full.tables.01.26.14')}

# Loading full tables from ../../../data/ungit-data/full.tables.01.26.14.rda ...Loaded.
# ../00common/mycache/snp.tables.chr1.unqfiltered.rda saved.

if(tb[2]){tset[[2]] <- load.snp.tables(use.chr1.tables = TRUE , data.name='full.tables.01.26.14')}
if(tb[3]){tset[[3]] <- load.snp.tables(use.chr1.tables = FALSE, data.name='full.tables.02.25.15')}

# Loading full tables from ../../../data/ungit-data/full.tables.02.25.15.rda ...Loaded.
# ../00common/mycache/snp.tables.chr1.qfiltered.rda saved.
# Bandaiding qfiltered tables...

if(tb[4]){tset[[4]] <- load.snp.tables(use.chr1.tables = TRUE , data.name='full.tables.02.25.15')}
```

Grrr! I should have excluded non-Chr contigs from full genome runs. Rather than change tons of code below to add mask params, I'm just going to truncate the tables, as follows. (See notes in wlr.r::make.mask for assumptions.)

```r
if(params$trunc.tables){
  for(i in 1:4){
    if(!is.null(tset[[i]])){
      first.mito <- match("mitochondria.fasta", tset[[i]][[7]]$Chr)
      if(!is.na(first.mito)){ # will be NA for Chr1 tables
        for(j in 1:7){
          # hmmm... slow; wonder whether head(tset[[i]][[j]],first.mito-1) is faster;
          # ok, simple tests suggest not: system.time(head(data.frame(1:1e7,1:1e7),5e6))
          tset[[i]][[j]] <- tset[[i]][[j]][1:(first.mito-1),]
        }
      }
    }
  }
} else {
    cat('***\n*** DID YOU *REALLY* WANT UNTRUNCATED TABLES???\n***\n')
}
```

The tersely-named `tset` list is sometimes convenient, but give them more descriptive names, too.

```r
snp.tables.full.unfiltered <- tset[[1]]; names(tset)[1] <- 'snp.tables.full.unfiltered'
snp.tables.chr1.unfiltered <- tset[[2]]; names(tset)[2] <- 'snp.tables.chr1.unfiltered'
snp.tables.full.qfiltered  <- tset[[3]]; names(tset)[3] <- 'snp.tables.full.qfiltered'
snp.tables.chr1.qfiltered  <- tset[[4]]; names(tset)[4] <- 'snp.tables.chr1.qfiltered'
```

Main analysis may just use one of the potentially 4 table sets. Pick it according to the priority specified in section 3, using the shorter name 'snp.tables' for this default choice.

```r
snp.tables <- tset.picker(priority=params$pri, table.set=tset)
```

```r
# Sanity check: unlike unqfiltered tables, bug in early code gave qfiltered ones different numbers
# of rows per strain, which breaks much code.  Verify this is no longer happening.
check.eq.nrows <- function(tables){
  if(!is.null(tables)){
    nrow.snp.tables <- unlist(lapply(tables,nrow))
    print(nrow.snp.tables)
    if(all(nrow.snp.tables == nrow.snp.tables[1])){
      cat('OK, all strains have same number of rows.\n')
    } else {
      cat('***\n*** Warning: Different strains have different numbers of rows! ***\n***\n')
    }
  }
}

dummy<-lapply(tset, check.eq.nrows)

#     1007     1012     1013     1014     1015     3367     1335
# 31301782 31301782 31301782 31301782 31301782 31301782 31301782
# OK, all strains have same number of rows.
#     1007     1012     1013     1014     1015     3367     1335
# 31301782 31301782 31301782 31301782 31301782 31301782 31301782
# OK, all strains have same number of rows.
```

Which tables have we got?:

```r
# 'which.snp.tables' return summary of which tables, either as a char string (default), e.g.
# "Chr1-qfiltered", or as vector of 2 strings, e.g. c("full","unfiltered").
cat('This analysis uses: (', paste(unlist(lapply(tset,which.snp.tables)),collapse=', '), ') SNP tables.\n')

# This analysis uses: ( trunc-unfiltered, NULL, trunc-qfiltered, NULL ) SNP tables.

cat('Main shared SNP analysis focuses on', which.snp.tables(snp.tables), '\n')

# Main shared SNP analysis focuses on trunc-unfiltered
```

A LATEX hack: I want which.snp.tables info in doc title/page headers, but it is unknown until now, so the following writes a command definition `\whichsnptables` into the .aux file, which is read during the *next* LATEX run, when `\begin{document}` is processed:

```
\makeatletter
\immediate\write\@auxout{\noexpand\gdef\noexpand\whichsnptables{trunc-unfiltered}}
\makeatother
```

Subsequent analysis was initially all directed at Chr1. In general, I have *not* updated the discussion to reflect genome-wide analysis.

```
if(exists('snp.tables.chr1.qfiltered') && exists('snp.tables.chr1.unqfiltered')){
  # If have both, where is new unequal to old?
  uneq <- snp.tables.chr1.qfiltered[[1]]$Ref[1:chr1.len] != snp.tables.chr1.unqfiltered[[1]]$Ref[1:chr1.len]
  cat('Sum uneq:', sum(uneq,na.rm=T), '\n')
  cat('Sum NA:  ', sum(is.na(uneq)),  '\n')
  print(which(is.na(uneq))[1:10])
  seecounts(which(is.na(uneq))[1:4],snp.tables=snp.tables.qfiltered,debug=F)
}
```

In brief, "`snp.tables`" will be a list of 7 data frames, one per strain, giving read counts for each nucleotide at each position, SNP calls, etc.:

```
names(snp.tables)

# [1] "1007" "1012" "1013" "1014" "1015" "3367" "1335"

str(snp.tables[[1]])

# 'data.frame': 31301782 obs. of  15 variables:
#  $ chr   : Factor w/ 66 levels "BD10_65","BD11_74",..: 39 39 39 39 39 39 39 39 39 39 ...
#  $ pos   : int  1 2 3 4 5 6 7 8 9 10 ...
#  $ snp   : int  0 0 0 0 0 0 0 0 0 0 ...
#  $ Chr   : chr  "Chr1" "Chr1" "Chr1" "Chr1" ...
#  $ Pos   : int  1 2 3 4 5 6 7 8 9 10 ...
#  $ Ref   : chr  "T" "C" "C" "A" ...
#  $ Cov   : num  1 3 4 5 7 7 10 12 13 15 ...
#  $ a     : num  0 0 1 0 0 0 0 0 1 0 ...
#  $ g     : num  0 0 0 0 0 0 0 0 0 0 ...
#  $ c     : num  0 0 0 0 0 0 0 0 0 0 ...
#  $ t     : num  0 0 0 0 0 0 0 0 0 0 ...
#  $ n     : num  0 0 0 0 0 0 0 0 0 0 ...
#  $ .match: num  1 3 3 5 7 7 10 12 12 15 ...
#  $ exon  : logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
#  $ indel : logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
```

Just for background, also load the desert tables:

```
# from svn+ssh://ceg1.ocean.washington.edu/var/svn/7_strains/trunk/code/snpNB/data
#load('../../../data/ungit-data/des.rda')
load('../../../data/des.rda')
```

What's the total length of all deserts in each strain? Big deserts (defined as "big.threshold" or longer)?

```
some.desert.stats <- function(big.threshold=0){
  desert.len <- unlist(lapply(des,function(x){sum(unlist(lapply(x,function(y){sum(y[,'Length'])})))}))
  bigdes.len <- unlist(lapply(des,function(x){sum(unlist(lapply(x,function(y){
                                         sum(y[y[,'Length']>=big.threshold,'Length'])})))}))
  rbind(desert.len, desert.pct=round( desert.len / genome.length.constants()$genome.length.trunc * 100),
        bigdes.len, bigdes.pct=round( bigdes.len / genome.length.constants()$genome.length.trunc * 100))
}
some.desert.stats(big.threshold=50000)

#              tp1007    tp1012   tp1013  tp1014    tp1015 thapsIT    tp1335
# desert.len 11146526 11332566 5801763 9464213 11251426 6780300 10883723
```

```
# desert.pct        36        36        19        30        36        22        35
# bigdes.len   3495805   3936973     55365   3627235   3727061     57119   4046934
# bigdes.pct        11        13         0        12        12         0        13
```

I.e., looking at all deserts, about 1/3 of L-clade, 1/5 of H-clade are in deserts, whereas, looking at the largest deserts ($> 50k$), only about 12% in L-clade (and none in H-clade). Note that the rough stats above include artifactual "deserts" created by gaps in the reference sequence, large genomic deletions, etc. A more careful analysis of this is found in nc-snps.rnw.

# 4   Refined SNP Calls

## 4.1   Method

It is appropriate that SNP calls should be conservative, to avoid many false positives, but, when a position is called a SNP in one isolate, we often see a significant number of reads for the same non-reference nucleotide at that position in other isolates, even if they are not called as SNPs. On the other hand, we sometimes see a position called a SNP in two or more isolates, but with *different* pairs of nucleotides, potentially suggesting technical errors. Analysis in this section attempts to refine the SNP calls by looking for issues such as these by looking at all 7 isolates jointly, at each position called a SNP in any of them.

For a given strain, the following function returns a vector of 0:4 to indicate which nonreference nucleotide has the maximum read count at the corresponding position. The values 1..4 indicate that the max count occurred at A, G, C, T, resp. (Ties are resolved arbitrarily ($a < g < c < t$), which possibly deserves further attention.) The value 0 means all nonreference counts are below threshold, based *either* on absolute count *or* as a fraction of coverage. Default only excludes 0 counts.

```r
nref.nuc.new <- function(strain=1, mask=T, thresh.count=0, thresh.rate=0.0){
        # get read count for max nonref nuc
        nref <- apply(snp.tables[[strain]][mask, c('a', 'g', 'c', 't')], 1, max)
        # where does nref count match a (g,c,t, resp) count
        as <- ifelse(nref == snp.tables[[strain]][mask,'a'],1,0)
        gs <- ifelse(nref == snp.tables[[strain]][mask,'g'],2,0)
        cs <- ifelse(nref == snp.tables[[strain]][mask,'c'],3,0)
        ts <- ifelse(nref == snp.tables[[strain]][mask,'t'],4,0)
        # most positions will show 3 zeros and one of 1:4, so max identifies max nonref count;
        # ties broken arbitrarily  (a<g<c<t)
        merge <- pmax(as,gs,cs,ts)
        # but if max nonref count is zero or below threshold, return 0
        merge[nref == 0 | nref < thresh.count] <- 0
        merge[nref/snp.tables[[strain]][mask,'Cov'] < thresh.rate] <- 0
        return(merge)
}
```

Get union and intersection of the sets of called SNPs. ("$snp" is 0/1.) Also, 5-way (L-clade) and 4-way (L-excluding Gyre).

```r
# 4-way union/intersection
u4.snps <- snp.tables[[1]]$snp
i4.snps <- snp.tables[[1]]$snp
for(i in c(2,5,7)) {
        u4.snps <- pmax(u4.snps, snp.tables[[i]]$snp)
        i4.snps <- pmin(i4.snps, snp.tables[[i]]$snp)
}
# 5-way: add gyre
u5.snps <- pmax(u4.snps, snp.tables[[4]]$snp)
i5.snps <- pmin(i4.snps, snp.tables[[4]]$snp)
# 7-way
union.snps     <- pmax(u5.snps, snp.tables[[3]]$snp, snp.tables[[6]]$snp)
intersect.snps <- pmin(i5.snps, snp.tables[[3]]$snp, snp.tables[[6]]$snp)
nu4snps <- sum(u4.snps)
nu5snps <- sum(u5.snps)
ni4snps <- sum(i4.snps)
```

```
ni5snps <- sum(i5.snps)
nusnps  <- sum(union.snps)
nisnps  <- sum(intersect.snps)
c(n4u=nu4snps, n5u=nu5snps, n7u=nusnps, n4i=ni4snps, n5i=ni5snps, n7i=nisnps)

#    n4u     n5u     n7u     n4i     n5i     n7i
# 196296 197799 474613 128683  70687  15186
```

There are nusnps=474613 positions called as SNPs in one or more strains (but only nisnps=15186 that are shared among all 7). Note that the 4-way union is only modestly larger (1.5254229 times larger) than the 4-way intersection, emphasizing the inherent similarities among these SNP sets. The corresponding 5-way numbers show that Gyre adds relatively little to the 5-way union vs the 4-way union, whereas it removes a fair bit from the 5-way intersection. However, much of that loss is simply because Gyre has fewer called SNPs: only 89184 vs 128683 in the 4-way intersection, and they are highly concordant:

```
sum(snp.tables[[4]]$snp*i4.snps)/sum(snp.tables[[4]]$snp)

# [1] 0.7925973
```

So, a likely source of the Gyre's difference in called SNPs is technical (lower read coverage, higher read error rate) rather than biological.

Inclusion of the 2 H-clade members, however, causes more dramatic changes in both union and intersection numbers. I examine all these relationships in more detail below, but first I examine what I believe to be a significant source of technical error in these comparisons—erroneous SNP calls, especially false negative calls.

It is appropriate that SNP calls should be conservative, to avoid many false positives, but, if a position is called a SNP in one strain, we often see a significant number of reads for the same non-reference nucleotide at that position in other strains, even if they are not called as SNPs. For my purposes below, these will be considered "shared SNPs," based on three different levels of permissiveness. Note that, e.g., $\geq 84\%$ of all positions have zero reads for any non-reference nucleotide, and only a small fraction have 2 or more non-reference reads:

```
nonmatch <- rbind(
  unlist(lapply(snp.tables,function(x){sum(x$Cov-x$.match == 0)})),
  unlist(lapply(snp.tables,function(x){sum(x$Cov-x$.match == 1)})),
  unlist(lapply(snp.tables,function(x){sum(x$Cov-x$.match == 2)})),
  unlist(lapply(snp.tables,function(x){sum(x$Cov-x$.match == 3)})),
  unlist(lapply(snp.tables,function(x){sum(x$Cov-x$.match >= 4)})),
  unlist(lapply(snp.tables,function(x){sum((x$Cov-x$.match)[union.snps==0] >= 4)}))
)/nrow(snp.tables[[1]])*100
rownames(nonmatch) <- c('% ==0','% ==1','% ==2','% ==3','% >=4', '% >=4, nonSNP')
nonmatch

#                      1007        1012        1013        1014        1015        3367        1335
# % ==0          92.4326481  88.6711338  84.5965383  86.25661312  90.4442629  86.3165937  84.9608722
# % ==1           6.3661967   9.4731028  12.3462396  11.94901300   7.9893279  10.9701135  11.8031619
# % ==2           0.4436073   0.9107788   1.4480773   1.14379111   0.6581766   1.1767158   1.8092900
# % ==3           0.0908830   0.1568633   0.2616369   0.19303054   0.1231208   0.2238563   0.4539230
# % >=4           0.6666649   0.7881213   1.3475079   0.45755222   0.7851119   1.3127208   0.9727529
# % >=4, nonSNP   0.1041826   0.1892065   0.3693783   0.09624053   0.1760091   0.3790263   0.3876297
```

Build a table of max non-reference nucleotides at each position in the union.snps set. The three criteria are

- [[1]]: any non-zero count at any coverage is considered significant

- [[2]]: (count $\geq 2$ and count/coverage $\geq 0.05$) is considered significant

- [[3]]: (count $\geq 4$ and count/coverage $\geq 0.10$) is considered significant

In all three cases, the nonref nucleotide must also be consistent across all strains passing that threshold; see below.

```
non.refs <- vector('list',4)
for(i in 1:4){
  non.refs[[i]] <- matrix(0, nrow=nusnps, ncol=7)
```

```
  colnames(non.refs[[i]]) <- names(snp.tables)
  rownames(non.refs[[i]]) <-
          paste(snp.tables[[1]]$chr[union.snps==1], ':', snp.tables[[1]]$pos[union.snps==1], sep='')
}
for(j in 1:7){
  non.refs[[1]][,j] <- nref.nuc.new(j, mask=union.snps==1, thresh.count=0, thresh.rate=0.00)
  non.refs[[2]][,j] <- nref.nuc.new(j, mask=union.snps==1, thresh.count=2, thresh.rate=0.05)
  non.refs[[3]][,j] <- nref.nuc.new(j, mask=union.snps==1, thresh.count=4, thresh.rate=0.10)
}
```

For comparison, I want to look at unfiltered SAMTools SNP calls. In complete opposition to the measures of consistency imposed above, I'm going to simply force this into the "non.refs" structure constructed above by imagining that any position called a SNP in any strain has its max nonref count on "A", so any given position called a SNP in any strain will automatically be declared "consistent." This will allow the tree-code, etc. given below to work in a uniform way (even though interpretation of the results is different.) Results will be jammed into a 4th component of the "non.refs" list; i.e., we have a 4th criterion:

  • [[4]]: all called SNPs at a given position are considered "consistent."

As this case was a late addition to the analysis, the commentary throughout this document has not necessarily been updated to reflect that this case is distinct from the first three.

```
for(j in 1:7){
  non.refs[[4]][,j] <- snp.tables[[j]]$snp[union.snps==1]
}
```

```
str(non.refs[[4]])

#  num [1:474613, 1:7] 0 0 0 0 0 0 0 0 1 0 ...
#  - attr(*, "dimnames")=List of 2
#   ..$ : chr [1:474613] "Chr1:333" "Chr1:417" "Chr1:435" "Chr1:438" ...
#   ..$ : chr [1:7] "1007" "1012" "1013" "1014" ...
```

"non.refs" indicates, among those positions in the union of all called SNPS having any non-reference read count above the thresholds listed above, the non-ref nucleotide having the highest read count in each strain. If, for a given position, the max of this code is the same as the min (among non-zero values), then every strain having over-threshold nonref reads in that position, in fact has most non-reference reads on the *same* nucleotide. These are defined as the "consistent" SNPs.

```
find.consistent <- function(nr){
  nr.max <- apply(nr,1,max)
  nr.min <- apply(nr,1,function(x){ifelse(max(x)==0,0,min(x[x>0]))})
  return(nr.min == nr.max)
}
consistent  <- lapply(non.refs, find.consistent)
```

## 4.2   Save them

```
# wrap this in a data structure to be cached:
Description <- [2753 chars quoted with ''']

refined.snps <-
  list(Description=Description,

       Data=list(
         based.on.which.snp.tables=which.snp.tables(),
         number.union.snps=nusnps,
         number.intersection.snps=nisnps,
         non.ref.nucleotide=non.refs,
         consistent.snps=consistent),
```

```r
        Code=list(
          get.snps = function(strain, stringency=2){
            # return nusnps x 1 Bool vector of consistent SNPs @ specified stringency & strain
            return(refined.snps$Data$consistent.snps[[stringency]] &
                   refined.snps$Data$non.ref.nucleotide[[stringency]][,strain] > 0)
          }
          ,
          get.snp.locs.char = function(strain, stringency=2){
            # return char vector of locations of consistent SNPs @ specified stringency & strain
            snps <- refined.snps$Code$get.snps(strain, stringency)
            return(names(snps)[snps])
          }
          ,
          get.snp.locs.df = function(strain, stringency=2){
            # return data frame (Chr/Pos) of locations of consistent SNPs @ specified stringency & strain
            snplist <- strsplit(refined.snps$Code$get.snp.locs.char(strain, stringency), ':', fixed=TRUE)
            # strsplit returns long list of 2-vectors, 1st=chr, 2nd=char position
            df <- data.frame(Chr=          unlist(lapply(snplist,function(x){return(x[1])})),
                             Pos=as.integer(unlist(lapply(snplist,function(x){return(x[2])}))),
                             stringsAsFactors = FALSE)
            return(df)
          }
        )
    )

# dont't clobber existing .rda, but save if absent.  (delete to re-save)
# result for trunc, unfiltered tables saved to "data" else "mycache"
if(which.snp.tables() == 'trunc-unfiltered'){
  rda.refined <- '../../../data/refined.snps-trunc-unfiltered.rda'
} else {
  rda.refined <- paste('../00common/mycache/refined.snps', which.snp.tables(), 'rda', sep='.')
}
if(file.exists(rda.refined)){
  cat('Pre-existing file', rda.refined, 'unchanged.\n')
} else {
  cat('Saving', rda.refined, '...')
  save(refined.snps, file=rda.refined, compress=TRUE)
  cat('Saved.\n')
}

# Saving ../../../data/refined.snps-trunc-unfiltered.rda ...Saved.
```

Knitr seems to be failing to format the long char string above, which says:

```r
cat(refined.snps$Description)

# Contents of this .rda file:
#
#   * Description: this text
#
#   * Data -- 5 items defining refined SNPs, at 4 different stringency levels, as defined
#     in shared-snps.rnw:
#
#     * based.on.which.snp.tables: {"Chr1","full","trunc"}-{"unfiltered","qfiltered"},
#       depending on which snp tables were used to build this data. ("trunc" = all Chrs.)
#
#     * number.union.snps: the total number of SNPs (SAMtools calls) in the union of SNPs
#       across all 7 strains.
#
#     * number.intersection.snps: similar, for the 7-way intersection.
#
#       nusnps/nisnps are easily recalculated from the data below, but their inclusion
#       may be convenient, e.g., to quickly see if the .rda represents the full genome
#       (nusnps=488848), or the chr 1 subset (nusnps=47499); (redundant with "based.on...";
#       numbers above are for unfiltered, perhaps slightly different if qfiltered)
#
#     * non.ref.nucleotide: 4 arrays, each nusnps x 7, of values 0..4 (0..1 in the 4th
```

```
#        array).  In the 1st 3 arrays, 0 means the given position in the given strain did
#        not have nonreference read counts above the corresponding filtering threshold,
#        i.e., is NOT a refined SNP in that strain, whereas 1..4 mean that it did pass
#        threshold, for A,C,G,T resp.  In the 4th array, this value is just 1/0,
#        indicating is/is not a called SNP in that strain.
#
#    * consistent.snps: 4 Bool vectors of length nusnps flagging positions whose nonref
#        nucs (wrt to the 4 filtering criteria) are deemed *consistent* across
#        all 7 strains.  For the 1st 3, this means all nonzero entries of non.ref.nuc
#        are equal, i.e., nonref read counts passing threshold are on the SAME nonref
#        nucleotide in all strains having over-threshold counts.  Just for comparison
#        and uniformity of data structures, the 4th is all TRUE, i.e., union of SNPs
#        across all strains, without any regard for thresholds or consistency.
#
#        In short, the refined SNPs according to our medium filtering criteria are
#        strains/positions where consistent.snps[[2]]==TRUE and non.ref.nucleotide[[2]]>0.
#
#        Rownames in both non.ref.nucs and consistent define location, e.g. "Chr1:333".
#
#   * Code -- simple routines to extract refined SNPs in (potentially) convenient formats:
#
#      * get.snps(strain, stringency=2)
#        returns nusnps x 1 Bool vector of consistent SNPs @ specified stringency in
#        given strain
#
#      * get.snp.locs.char(strain, stringency=2)
#        returns n x 1 char vector of locations of consistent SNPs @ specified stringency
#        in given strain, e.g. "Chr1:1234", where n == sum(get.snps(...))
#
#      * get.snp.locs.df(strain, stringency=2){
#        As above, but returns data frame (char vector Chr, int vector Pos) with the same info.
```

```
str(consistent[[1]])
```

```
#  Named logi [1:474613] TRUE FALSE TRUE FALSE TRUE TRUE ...
#  - attr(*, "names")= chr [1:474613] "Chr1:333" "Chr1:417" "Chr1:435" "Chr1:438" ...
```

```
consistent.count <- unlist(lapply(consistent, sum)) ; consistent.count

# [1] 358872 468117 470970 474613

inconsistent.count <- consistent.count[4] - consistent.count; inconsistent.count

# [1] 115741   6496   3643      0

inconsistent.percent <- inconsistent.count/consistent.count[4]*100; inconsistent.percent

# [1] 24.3863948  1.3686941  0.7675727  0.0000000
```

I.e., of the 474613 positions in which a SNP is called, 358872 are consistent by my loose definition, and 470970 are consistent by my tightest definition. The increase in concordance supports the view that the loose definition is too loose. Perhaps misleadingly, these counts include positions that are "consistent SNPs" in only one strain; more below. (*TODO* I suspect, but have not yet systematically checked, that most of the rest are positions with low coverage and/or very low read counts on the mixture of non-reference nucleotides.)

## 4.3 Examples: Consistent

Here are a few (nonrandomly selected) prototypical consistent SNPs:

```
esnps <- names(consistent[[2]][consistent[[2]]])
esnps2 <- as.integer(unlist(lapply(strsplit(esnps[c(7,11:13,92)],':',fixed=TRUE),function(x){x[2]})))
seecounts(esnps2,snp.tables=snp.tables)
```

```
#       chr  pos Ref Strain  A   G  C   T SNP  exon indel nrf rat
# 1  Chr1  567   T
# 2                  1007  0   0  2  29   0  TRUE FALSE
# 3                  1012  0   0 15  44   1  TRUE FALSE
# 4                  1013  0   0 15  97   0  TRUE FALSE
# 5                  1014  0   1  0  33   0  TRUE FALSE
# 6                  1015  0   0  9  43   1  TRUE FALSE
# 7                  3367  0   0 16  46   1  TRUE FALSE
# 8                  1335  0   0  2 116   0  TRUE FALSE
# 9  Chr1 1053   A
# 10                 1007 39   0  0   5   0  TRUE FALSE
# 11                 1012 55   0  0  12   0  TRUE FALSE
# 12                 1013 17   1  0  40   0  TRUE FALSE
# 13                 1014 25   0  0   5   0  TRUE FALSE
# 14                 1015 38   0  1  20   1  TRUE FALSE
# 15                 3367 13   0  0   9   0  TRUE FALSE
# 16                 1335 71   1  0  46   1  TRUE FALSE
# 17 Chr1 1055   G
# 18                 1007  0  41  0   2   0  TRUE FALSE
# 19                 1012  1  63  0   8   0  TRUE FALSE
# 20                 1013  1  62  0   8   0  TRUE FALSE
# 21                 1014  1  26  0   8   1  TRUE FALSE
# 22                 1015  0  44  0  14   0  TRUE FALSE
# 23                 3367  0  27  0   0   0  TRUE FALSE
# 24                 1335  0  78  0  40   1  TRUE FALSE
# 25 Chr1 1176   G
# 26                 1007  2  67  0   0   0 FALSE FALSE
# 27                 1012  1  68  0   0   0 FALSE FALSE
# 28                 1013 29  73  0   0   0 FALSE FALSE
# 29                 1014  1  52  0   0   0 FALSE FALSE
# 30                 1015  4 103  0   0   0 FALSE FALSE
# 31                 3367 11   8  0   0   1 FALSE FALSE
# 32                 1335  1 206  0   0   0 FALSE FALSE
# 33 Chr1 8670   A
# 34                 1007 19   0  0   7   0  TRUE FALSE
# 35                 1012 36   0  0  12   0  TRUE FALSE
# 36                 1013 44   0  0  12   0  TRUE FALSE
# 37                 1014 10   0  0   7   0  TRUE FALSE
# 38                 1015 24   0  0  11   1  TRUE FALSE
# 39                 3367 18   0  0   0   0  TRUE FALSE
# 40                 1335 27   0  0   6   0  TRUE FALSE
```

## 4.4   Examples: Inconsistent

Here is a brief look at some *in*-consistent positions. E.g., Chr1:2013 shows nontrivial counts on 3 alleles in Wales, as do 2319, 3286, 5002, 5433, whereas 7878 shows a different alternate allele in Italy than in Wales.

```
unc <- names(consistent[[2]][!consistent[[2]]])
unc2 <- as.integer(unlist(lapply(strsplit(unc[1:10],':',fixed=TRUE),function(x){x[2]})))
seecounts(unc2,snp.tables=snp.tables)

#       chr   pos Ref Strain  A   G   C  T SNP  exon indel nrf rat
# 1  Chr1  2013   T
# 2                  1007  4   0   0 20   0  TRUE FALSE
# 3                  1012  8   0   0 34   0  TRUE FALSE
# 4                  1013  9  12   0 16   1  TRUE FALSE
# 5                  1014  1   0   0 19   0  TRUE FALSE
# 6                  1015 13   0   0 24   1  TRUE FALSE
# 7                  3367 10   0   0 36   0  TRUE FALSE
# 8                  1335 20   0   0 68   1  TRUE FALSE
# 9  Chr1  2319   C
# 10                 1007  0  29  22  0   1  TRUE FALSE
# 11                 1012  0  54  26  0   1  TRUE FALSE
# 12                 1013 19  19  18  0   1  TRUE FALSE
# 13                 1014  0  25  19  0   1  TRUE FALSE
```

```
# 14                     1015   0  54  29   0   1   TRUE FALSE
# 15                     3367   5   0  43   0   0   TRUE FALSE
# 16                     1335   0 132  48   0   1   TRUE FALSE
# 17 Chr1  3286    T
# 18                     1007   4   0   1 17   0   TRUE FALSE
# 19                     1012   9   0   3 45   0   TRUE FALSE
# 20                     1013 39   1  38 12   1   TRUE FALSE
# 21                     1014   4   0   6 27   0   TRUE FALSE
# 22                     1015 11   0   7 37   0   TRUE FALSE
# 23                     3367   8   0  39 10   0   TRUE FALSE
# 24                     1335 15   0   4 75   0   TRUE FALSE
# 25 Chr1  5002    T
# 26                     1007   0  15   0 12   0   TRUE FALSE
# 27                     1012   1  23   0 26   1   TRUE FALSE
# 28                     1013 21  11   0 39   0   TRUE FALSE
# 29                     1014   0   8   0 12   0   TRUE FALSE
# 30                     1015   0  19   0 16   1   TRUE FALSE
# 31                     3367   0   0   0 35   0   TRUE FALSE
# 32                     1335   0  57   0 60   0   TRUE FALSE
# 33 Chr1  5433    G
# 34                     1007   0  50   0   3   0   TRUE FALSE
# 35                     1012   0  78   0   5   0   TRUE FALSE
# 36                     1013 18  47   0 14   1   TRUE FALSE
# 37                     1014   9  19   0   0   1   TRUE FALSE
# 38                     1015   7  63   0   2   0   TRUE FALSE
# 39                     3367   8  54   0   0   0   TRUE FALSE
# 40                     1335   6 109   0   4   0   TRUE FALSE
# 41 Chr1  7858    C
# 42                     1007   0   0  48   0   0   TRUE FALSE
# 43                     1012   0   1  61   0   0   TRUE FALSE
# 44                     1013   0   0 131 10   0   TRUE FALSE
# 45                     1014   0   0  34   0   0   TRUE FALSE
# 46                     1015   0   0  74   0   0   TRUE FALSE
# 47                     3367 20   0   8   0   1   TRUE FALSE
# 48                     1335   0   0 120   0   0   TRUE FALSE
# 49 Chr1  8914    A
# 50                     1007 23   0   0   2   0   TRUE FALSE
# 51                     1012 29   0  15   0   1   TRUE FALSE
# 52                     1013 25   0   6   0   0   TRUE FALSE
# 53                     1014 22   0   0   0   0   TRUE FALSE
# 54                     1015 31   0   5   2   0   TRUE FALSE
# 55                     3367   8   0   0   1   0   TRUE FALSE
# 56                     1335 68   0   7   0   0   TRUE FALSE
# 57 Chr1  8974    C
# 58                     1007   0   2   6   0   0   TRUE FALSE
# 59                     1012   0   2  17   0   0   TRUE FALSE
# 60                     1013 10  22   4   0   1   TRUE FALSE
# 61                     1014   0   1  10   0   0   TRUE FALSE
# 62                     1015   0   2  15   0   0   TRUE FALSE
# 63                     3367   2   0   3   0   0   TRUE FALSE
# 64                     1335   0  11  49   0   0   TRUE FALSE
# 65 Chr1 10099    T
# 66                     1007 17   0   0 29   0   TRUE FALSE
# 67                     1012 48   0   0 36   0   TRUE FALSE
# 68                     1013   0   2   6 68   0   TRUE FALSE
# 69                     1014 34   0   0 26   0   TRUE FALSE
# 70                     1015 41   0   0 38   0   TRUE FALSE
# 71                     3367   0   1   0 14   0   TRUE FALSE
# 72                     1335 55   0   0 68   1   TRUE FALSE
# 73 Chr1 15154    A
# 74                     1007 25   0   0   0   0 FALSE FALSE
# 75                     1012 56   0   0   1   0 FALSE FALSE
# 76                     1013 10   0  38 10   1 FALSE FALSE
# 77                     1014 26   0   0   0   0 FALSE FALSE
# 78                     1015 37   0   0   0   0 FALSE FALSE
# 79                     3367 19   0   0 13   1 FALSE FALSE
# 80                     1335 70   0   0   3   0 FALSE FALSE
```

## 4.5   Examples: Homozygous nonref

And at some *homozygous nonreference* positions (defined to be those with nonref fraction > 0.75):

```
hnr <- lapply(snp.tables, function(x){x$.match/x$Cov < 0.25})    # find them
hnr <- lapply(hnr, function(x){ifelse(is.na(x),FALSE,x)})         # remove NA
unlist(lapply(hnr,sum))                                           # count per strain

#  1007  1012  1013  1014  1015  3367  1335
#  6619  7645 62072   440  3593 72356   558
```

Hmm, in L-clade, excluding the ref isolate (1335) this tracks time-in culture to some degree; Maybe many of these are in hemizygous regions. Next two chunks lifted from nc-snps to get tables for hemi-deletion.

```
cnv.chronly <- load.cnv.tables('../../../data/cnv.txt', chrs.only=TRUE)

str(cnv.chronly)

# 'data.frame': 1956 obs. of  11 variables:
#  $ strain  : Factor w/ 7 levels "IT","tp1007",..: 3 3 3 3 3 3 3 3 3 3 ...
#  $ chr     : Factor w/ 65 levels "BD1_7","BD10_65",..: 38 38 38 38 38 38 38 38 38 38 ...
#  $ start   : int  10601 112001 215001 358901 536501 554801 673401 781801 806901 853201 ...
#  $ end     : int  13500 116500 221100 370300 538600 559300 685000 787400 811100 855600 ...
#  $ length  : int  2900 4500 6100 11400 2100 4500 11600 5600 4200 2400 ...
#  $ filtered : logi  FALSE FALSE FALSE TRUE FALSE FALSE ...
#  $ type    : Factor w/ 1 level "CNVnator": 1 1 1 1 1 1 1 1 1 1 ...
#  $ cov_ratio: num  0.63738 1.54893 1.65381 0.00204 0.68486 ...
#  $ dup_frac : num  0.41188 0.00908 0.01178 0.97997 0.0211 ...
#  $ iStart  : num  10601 112001 215001 358901 536501 ...
#  $ iEnd    : num  13500 116500 221100 370300 538600 ...

cnv.chronly[c(1:4,nrow(cnv.chronly)+c(-1,0)),]                    ## first/last few rows

#       strain   chr  start    end length filtered      type  cov_ratio   dup_frac    iStart      iEnd
# 1     tp1012  Chr1  10601  13500   2900    FALSE CNVnator 0.63738000 0.41187900     10601     13500
# 2     tp1012  Chr1 112001 116500   4500    FALSE CNVnator 1.54893000 0.00907677    112001    116500
# 3     tp1012  Chr1 215001 221100   6100    FALSE CNVnator 1.65381000 0.01178470    215001    221100
# 4     tp1012  Chr1 358901 370300  11400     TRUE CNVnator 0.00204431 0.97997300    358901    370300
# 1955 tp1335 Chr24 259901 278000  18100    FALSE CNVnator 1.41458000 0.38091100  31264334  31282433
# 1956 tp1335 Chr24 286901 289800   2900    FALSE CNVnator 1.74941000 0.74228100  31291334  31294233
```

```
get.cnv.dels <- function(cov.thresh.lo = 0.0,
                         cov.thresh.hi = 0.8,
                         cnv,
                         snp.tables = NULL,
                         DEBUG = FALSE
){
  # build list of 7 Bool vectors of genome length, with i-th == T iff
  #   * i-th pos is 'NA' in genome seq (if snp.tables are provided), or
  #   * in CNVnator call for coverage in half-open [cov.thresh.lo, hi), and
  #   * not marked 'filtered' by CNVnator
  cnv.deletions <- vector(mode='list',7)                 # make list of bool vectors
  if(is.null(snp.tables)){
    # if no tables, assume full
    t.len <- genome.length.constants()$genome.length.trunc
  } else {
    t.len <- nrow(snp.tables[[1]])
  }
  for(st in 1:7){
    if(is.null(snp.tables)){
      cnv.deletions[[st]] <- logical(t.len)              # all F
    } else  {
      cnv.deletions[[st]] <- is.na(snp.tables[[st]]$Pos[1:t.len])  # NA positions in genome
    }
  }
  strain.names <- c(paste('tp10',c('07',12:15),sep=''),'IT','tp1335')
```

```r
  names(cnv.deletions) <- strain.names
  for(i in 1:nrow(cnv)){
    if(!cnv$filtered[i] &&
       cnv$cov_ratio[i] >= cov.thresh.lo &&
       cnv$cov_ratio[i] <  cov.thresh.hi)
    {
      if(DEBUG){
        print(cnv[i,])
        print(as.character(cnv$strain[i]))
      }
      # following ASSUMES no CNVnator call crosses a chromosome bdry, & that
      # t.len ends at chr end (typically chr1 or chr24)
      if(cnv$iEnd[i] <= t.len){
        cnv.deletions[[as.character(cnv$strain[i])]]][cnv$iStart[i]:cnv$iEnd[i]] <- TRUE
      }
    }
  }
  return(cnv.deletions)
}

# sanity check:
cnv.dels.38 <- get.cnv.dels(0.3, 0.8, cnv.chronly, snp.tables = NULL)
unlist(lapply(cnv.dels.38,sum)) # does it match low.length.38 in tic ?

#  tp1007  tp1012  tp1013  tp1014  tp1015      IT  tp1335
# 1672500 1781500 1383600 1313700  988400  320900 1453000

# 1672500 1781500 1399400 1313700 988400 336500 1453000 <== low.length.38 from tic (circa page 8)
# 1672500 1781500 1399400 1313700 988400 336500 1453000 <== low.length.38 from tic (pg9, 6/28/17)
rm(cnv.dels.38)
```

Slight discrepancy in H-clade that I should hunt down, but basically OK. (hmm; maybe untrunc tbls.)

```r
# the ones we want for the current analysis:
hemi.masks <- get.cnv.dels(0.3, 0.8, cnv.chronly, snp.tables=snp.tables)

rbind(
  homnr        = unlist(lapply(hnr,sum)),
  hemi         = unlist(lapply(hemi.masks, sum)),
  homnr.unhemi = unlist(lapply(list(1,2,3,4,5,6,7), function(i){sum(hnr[[i]] & !hemi.masks[[i]])}))
)

#                  1007    1012    1013    1014    1015    3367    1335
# homnr            6619    7645   62072     440    3593   72356     558
# hemi          1834990 1940024 1527725 1472095 1134652 480817 1596965
# homnr.unhemi     4441    4220   59390     434    2683   71732     537


# based on the thought that hnr in 1335 may reflect errors in the ref seq,
# are they shared with others?
unlist(lapply(hnr, function(x){sum(x & hnr[[7]])}))                 # hnr shared with 1335

# 1007 1012 1013 1014 1015 3367 1335
#  271  306  311  243  332  330  558

# answer: around 300 in each strain, of 558 in NY, genomewide,
# so that seems like a plausibly important factor.

hnr.lclade <- hnr[[1]] | hnr[[2]] | hnr[[4]] | hnr[[5]] | hnr[[7]]  # union over L-clade
sum(hnr.lclade)                                                    # count all in L-clade

# [1] 11010


sum(hnr[[3]] | hnr[[6]])                                           # present in H-clade

# [1] 104450
```

```
sum(hnr[[3]] & hnr[[6]])                                                 # shared in H-clade

# [1] 29978

# look at a few in L-clade
w.hnr.l <- which(hnr.lclade)
seecounts(w.hnr.l[1:10],snp.tables=snp.tables)

#      chr    pos  Ref Strain     A     G     C     T SNP   exon indel nrf rat
# 1  Chr1   5397    C
# 2                     1007     0     0    24    27   1   TRUE FALSE
# 3                     1012     0     0    34    40   1   TRUE FALSE
# 4                     1013     0     0    12    42   0   TRUE FALSE
# 5                     1014     1     0    30    28   1   TRUE FALSE
# 6                     1015     0     0    33    35   1   TRUE FALSE
# 7                     3367     0     0    20    38   1   TRUE FALSE
# 8                     1335     0     0    29    98   1   TRUE FALSE
# 9  Chr1  20071    T
# 10                    1007    22     0     0    15   1  FALSE FALSE
# 11                    1012   109     0     0    41   1  FALSE FALSE
# 12                    1013    28     0     0    33   1  FALSE FALSE
# 13                    1014    76     0     0    29   1  FALSE FALSE
# 14                    1015   130     0     0    28   1  FALSE FALSE
# 15                    3367    27     0     0    28   0  FALSE FALSE
# 16                    1335    95     0     0    57   0  FALSE FALSE
# 17 Chr1  25350    G
# 18                    1007   104    31     0     0   1  FALSE FALSE
# 19                    1012   171    53     0     0   1  FALSE FALSE
# 20                    1013   209    87     1     0   0  FALSE FALSE
# 21                    1014    19    32     0     0   0  FALSE FALSE
# 22                    1015    91    44     0     0   1  FALSE FALSE
# 23                    3367   397    94     0     0   0  FALSE FALSE
# 24                    1335    80    64     0     0   0  FALSE FALSE
# 25 Chr1  26205    T
# 26                    1007    50     0     0    20   1  FALSE FALSE
# 27                    1012   104     0     0    33   1  FALSE FALSE
# 28                    1013   224     0     0    69   1  FALSE FALSE
# 29                    1014    23     0     1    16   1  FALSE FALSE
# 30                    1015    88     0     0    33   1  FALSE FALSE
# 31                    3367   143     0     0    41   1  FALSE FALSE
# 32                    1335   196     0     0    67   1  FALSE FALSE
# 33 Chr1  90942    C
# 34                    1007     0     0    15     0   0  FALSE FALSE
# 35                    1012     0     0    33     0   0  FALSE FALSE
# 36                    1013     0     0    46     0   0  FALSE FALSE
# 37                    1014     0     0    16     0   0  FALSE FALSE
# 38                    1015     0     0     7    25   1  FALSE FALSE
# 39                    3367     0     0    56     0   0  FALSE FALSE
# 40                    1335     0     0    70     0   0  FALSE FALSE
# 41 Chr1 149447    T
# 42                    1007     0     0     0     1   0  FALSE FALSE
# 43                    1012     0     1     0     0   0  FALSE FALSE
# 44                    1013     0     0     0     1   0  FALSE FALSE
# 45                    1014     0     0     0     8   0  FALSE FALSE
# 46                    1015     0     0     0     2   0  FALSE FALSE
# 47                    3367     0     0     0     1   0  FALSE FALSE
# 48                    1335     0     1     0     1   0  FALSE FALSE
# 49 Chr1 149457 <NA>
# 50                    1007  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 51                    1012  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 52                    1013  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 53                    1014  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 54                    1015  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 55                    3367  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 56                    1335  <NA>  <NA>  <NA>  <NA>   0  FALSE FALSE
# 57 Chr1 156248    A
# 58                    1007     2     0    39     0   0  FALSE FALSE
# 59                    1012     7     0    67     0   0  FALSE FALSE
```

```
# 60                        1013    5    0   53    0   0 FALSE FALSE
# 61                        1014   11    0   13    0   1 FALSE FALSE
# 62                        1015    6    0   44    0   0 FALSE FALSE
# 63                        3367    9    0   66    0   0 FALSE FALSE
# 64                        1335   62    0   31    0   1 FALSE FALSE
# 65 Chr1 176517    C
# 66                        1007    0    0    0    1   0  TRUE FALSE
# 67                        1012    0    0    2    0   0  TRUE FALSE
# 68                        1013    0    0    4    0   0  TRUE FALSE
# 69                        1014    0    0    6    0   0  TRUE FALSE
# 70                        1015    0    0    0    0   0  TRUE FALSE
# 71                        3367    0    0    4    0   0  TRUE FALSE
# 72                        1335    0    0   11    0   0  TRUE FALSE
# 73 Chr1 193761    C
# 74                        1007    0    0   20   14   1 FALSE FALSE
# 75                        1012    0    0   19   31   1 FALSE FALSE
# 76                        1013    0    1    4    6   1 FALSE FALSE
# 77                        1014    0    0    9    4   1 FALSE FALSE
# 78                        1015    0    0   28   39   1 FALSE FALSE
# 79                        3367    0    0    7   11   0 FALSE FALSE
# 80                        1335    0    0   10   43   1 FALSE FALSE

# one of those is a little weird:
xx<-snp.tables[[1]][149457,]
for (i in 2:7){xx <- rbind(xx,snp.tables[[i]][149457,])}
row.names(xx)<-names(snp.tables)
# My guess is that Chr/Pos/Ref are left as NA if coverage is zero.
xx

#        chr    pos snp  Chr    Pos  Ref Cov a g c t n .match  exon indel
# 1007 Chr1 149457   0 <NA>     NA <NA>   0 0 0 0 0 0      0 FALSE FALSE
# 1012 Chr1 149457   0 <NA>     NA <NA>   0 0 0 0 0 0      0 FALSE FALSE
# 1013 Chr1 149457   0 <NA>     NA <NA>   0 0 0 0 0 0      0 FALSE FALSE
# 1014 Chr1 149457   0 Chr1 149457    G   1 0 0 0 0 0      1 FALSE FALSE
# 1015 Chr1 149457   0 <NA>     NA <NA>   0 0 0 0 0 0      0 FALSE FALSE
# 3367 Chr1 149457   0 <NA>     NA <NA>   0 0 0 0 0 0      0 FALSE FALSE
# 1335 Chr1 149457   0 Chr1 149457    G   1 0 0 1 0 0      0 FALSE FALSE
```

# 5   Table 1 stats

Here is a brief summary of per-strain SNP counts, pairwise overlaps, and other conveniently available stats, such as those shown in Table 1 of the paper.

```
snp.counts     <- matrix(NA,7,4)
snp.pctofny    <- matrix(NA,7,4)
snp.pctofself  <- matrix(NA,7,4)
snp.inter  <- matrix(NA,7,7)
snp.union  <- matrix(NA,7,7)
rownames(snp.counts)      <- names(snp.tables)
rownames(snp.pctofny)     <- names(snp.tables)
rownames(snp.pctofself)   <- names(snp.tables)
rownames(snp.inter)   <- names(snp.tables)
colnames(snp.inter)   <- names(snp.tables)
rownames(snp.union)   <- names(snp.tables)
colnames(snp.union)   <- names(snp.tables)
for(stringency in 1:4){
  cat('\nStringency', stringency, ifelse(stringency==4,'(i.e. raw SAMTools SNP calls)',''),
      ':\n--------------\n')
  for(i in 1:7){
    f.snps.i <- refined.snps$Code$get.snps(i, stringency)
    snp.counts[i,stringency] <- sum(f.snps.i)
    for(j in i:7){
      f.snps.j <- refined.snps$Code$get.snps(j, stringency)
      snp.inter[i,j] <- sum(f.snps.i & f.snps.j)
```

```
      snp.union[i,j] <- sum(f.snps.i | f.snps.j)
    }
  }
  snp.pctofny  [,stringency] <- snp.inter[,7]/snp.counts[7,stringency]
  snp.pctofself[,stringency] <- snp.inter[,7]/snp.counts[ ,stringency]
  cat('Union Counts:\n');                     print(snp.union)
  cat('Intersect Counts:\n');                 print(snp.inter)
  cat('Intersect as percent of union:\n'); print(snp.inter/snp.union*100,digits=3)
}

#
# Stringency 1  :
# --------------
# Union Counts:
#         1007    1012    1013    1014    1015    3367    1335
# 1007 175110 185299 304446 191894 189265 297494 193810
# 1012     NA 180026 306192 195182 192148 299196 196561
# 1013     NA     NA 249044 302460 307109 316432 306430
# 1014     NA     NA     NA 168167 192630 295200 195041
# 1015     NA     NA     NA     NA 181549 300380 194151
# 3367     NA     NA     NA     NA     NA 237364 299559
# 1335     NA     NA     NA     NA     NA     NA 181546
# Intersect Counts:
#         1007    1012    1013    1014    1015    3367    1335
# 1007 175110 169837 119708 151383 167394 114980 162846
# 1012     NA 180026 122878 153011 169427 118194 165011
# 1013     NA     NA 249044 114751 123484 169976 124160
# 1014     NA     NA     NA 168167 157086 110331 154672
# 1015     NA     NA     NA     NA 181549 118533 168944
# 3367     NA     NA     NA     NA     NA 237364 119351
# 1335     NA     NA     NA     NA     NA     NA 181546
# Intersect as percent of union:
#      1007  1012  1013  1014  1015  3367  1335
# 1007  100  91.7  39.3  78.9  88.4  38.6  84.0
# 1012   NA 100.0  40.1  78.4  88.2  39.5  83.9
# 1013   NA    NA 100.0  37.9  40.2  53.7  40.5
# 1014   NA    NA    NA 100.0  81.5  37.4  79.3
# 1015   NA    NA    NA    NA 100.0  39.5  87.0
# 3367   NA    NA    NA    NA    NA 100.0  39.8
# 1335   NA    NA    NA    NA    NA    NA 100.0
#
# Stringency 2  :
# --------------
# Union Counts:
#         1007    1012    1013    1014    1015    3367    1335
# 1007 182700 189244 374013 193188 196130 364691 195659
# 1012     NA 186491 375922 195230 196927 366436 196736
# 1013     NA     NA 304293 361304 377929 407688 373319
# 1014     NA     NA     NA 149531 194388 351232 188751
# 1015     NA     NA     NA     NA 190302 368934 195492
# 3367     NA     NA     NA     NA     NA 290904 363768
# 1335     NA     NA     NA     NA     NA     NA 180187
# Intersect Counts:
#         1007    1012    1013    1014    1015    3367    1335
# 1007 182700 179947 112980 139043 176872 108913 167228
# 1012     NA 186491 114862 140792 179866 110959 169942
# 1013     NA     NA 304293  92520 116666 187509 111161
# 1014     NA     NA     NA 149531 145445  89203 140967
# 1015     NA     NA     NA     NA 190302 112272 174997
# 3367     NA     NA     NA     NA     NA 290904 107323
# 1335     NA     NA     NA     NA     NA     NA 180187
# Intersect as percent of union:
#      1007  1012  1013  1014  1015  3367  1335
# 1007  100  95.1  30.2  72.0  90.2  29.9  85.5
# 1012   NA 100.0  30.6  72.1  91.3  30.3  86.4
# 1013   NA    NA 100.0  25.6  30.9  46.0  29.8
# 1014   NA    NA    NA 100.0  74.8  25.4  74.7
# 1015   NA    NA    NA    NA 100.0  30.4  89.5
```

```
# 3367    NA    NA    NA    NA    NA 100.0  29.5
# 1335    NA    NA    NA    NA    NA    NA 100.0
#
# Stringency 3  :
# --------------
# Union Counts:
#         1007   1012   1013   1014   1015   3367   1335
# 1007 172445 184829 367417 181202 191833 358522 189174
# 1012     NA 180855 371490 187550 193436 362568 192116
# 1013     NA     NA 298841 340611 373801 406440 366421
# 1014     NA     NA     NA 106701 188560 330490 177650
# 1015     NA     NA     NA     NA 185294 365440 191520
# 3367     NA     NA     NA     NA     NA 286507 357408
# 1335     NA     NA     NA     NA     NA     NA 170198
# Intersect Counts:
#         1007   1012   1013   1014   1015   3367   1335
# 1007 172445 168471 103869  97944 165906 100430 153469
# 1012     NA 180855 108206 100006 172713 104794 158937
# 1013     NA     NA 298841  64931 110334 178908 102618
# 1014     NA     NA     NA 106701 103435  62718  99249
# 1015     NA     NA     NA     NA 185294 106361 163972
# 3367     NA     NA     NA     NA     NA 286507  99297
# 1335     NA     NA     NA     NA     NA     NA 170198
# Intersect as percent of union:
#      1007  1012  1013  1014  1015  3367  1335
# 1007  100  91.1  28.3  54.1  86.5  28.0  81.1
# 1012   NA 100.0  29.1  53.3  89.3  28.9  82.7
# 1013   NA    NA 100.0  19.1  29.5  44.0  28.0
# 1014   NA    NA    NA 100.0  54.9  19.0  55.9
# 1015   NA    NA    NA    NA 100.0  29.1  85.6
# 3367   NA    NA    NA    NA    NA 100.0  27.8
# 1335   NA    NA    NA    NA    NA    NA 100.0
#
# Stringency 4 (i.e. raw SAMTools SNP calls) :
# --------------
# Union Counts:
#         1007   1012   1013   1014   1015   3367   1335
# 1007 161103 176738 343873 171675 185741 336599 180313
# 1012     NA 166089 346766 176177 186459 339458 182312
# 1013     NA     NA 247737 302322 352586 386037 339669
# 1014     NA     NA     NA  89184 179976 295574 162912
# 1015     NA     NA     NA     NA 174701 345396 184068
# 3367     NA     NA     NA     NA     NA 240413 331982
# 1335     NA     NA     NA     NA     NA     NA 153901
# Intersect Counts:
#         1007   1012   1013  1014   1015   3367   1335
# 1007 161103 150454  64967 78612 150063  64917 134691
# 1012     NA 166089  67060 79096 154331  67044 137678
# 1013     NA     NA 247737 34599  69852 102113  61969
# 1014     NA     NA     NA 89184  83909  34023  80173
# 1015     NA     NA     NA    NA 174701  69718 144534
# 3367     NA     NA     NA    NA     NA 240413  62332
# 1335     NA     NA     NA    NA     NA     NA 153901
# Intersect as percent of union:
#      1007  1012  1013  1014  1015  3367  1335
# 1007  100  85.1  18.9  45.8  80.8  19.3  74.7
# 1012   NA 100.0  19.3  44.9  82.8  19.8  75.5
# 1013   NA    NA 100.0  11.4  19.8  26.5  18.2
# 1014   NA    NA    NA 100.0  46.6  11.5  49.2
# 1015   NA    NA    NA    NA 100.0  20.2  78.5
# 3367   NA    NA    NA    NA    NA 100.0  18.8
# 1335   NA    NA    NA    NA    NA    NA 100.0

vs.stringency <- cbind(snp.counts, matrix(NA,7,1), round(snp.counts[,1:3]/snp.counts[,4]*100,1))
colnames(vs.stringency) <- c('[[1]]', '[[2]]', '[[3]]', '[[4]]', '----', '[[1]]%', '[[2]]%', '[[3]]%')

# SNPs vs filtering stringency (raw counts and as % of [[4]]).  Medium filter
# adds 10-20% in most cases.  Big exception is Gyre, where low coverage,
```

```
# high err rate and SAMTools conservatism seemed to seriously undercall:
print(vs.stringency)

#         [[1]]  [[2]]  [[3]]  [[4]]  ---- [[1]]% [[2]]% [[3]]%
# 1007 175110 182700 172445 161103   NA  108.7  113.4  107.0
# 1012 180026 186491 180855 166089   NA  108.4  112.3  108.9
# 1013 249044 304293 298841 247737   NA  100.5  122.8  120.6
# 1014 168167 149531 106701  89184   NA  188.6  167.7  119.6
# 1015 181549 190302 185294 174701   NA  103.9  108.9  106.1
# 3367 237364 290904 286507 240413   NA   98.7  121.0  119.2
# 1335 181546 180187 170198 153901   NA  118.0  117.1  110.6

# Intersect NY as % of self (vs stringency):
print(snp.pctofself*100, digits=3)

#        [,1]  [,2]  [,3]  [,4]
# 1007  93.0  91.5  89.0  83.6
# 1012  91.7  91.1  87.9  82.9
# 1013  49.9  36.5  34.3  25.0
# 1014  92.0  94.3  93.0  89.9
# 1015  93.1  92.0  88.5  82.7
# 3367  50.3  36.9  34.7  25.9
# 1335 100.0 100.0 100.0 100.0

# Intersect NY as % of NY (vs stringency):
print(snp.pctofny*100, digits=3)

#        [,1]  [,2]  [,3]  [,4]
# 1007  89.7  92.8  90.2  87.5
# 1012  90.9  94.3  93.4  89.5
# 1013  68.4  61.7  60.3  40.3
# 1014  85.2  78.2  58.3  52.1
# 1015  93.1  97.1  96.3  93.9
# 3367  65.7  59.6  58.3  40.5
# 1335 100.0 100.0 100.0 100.0
```

Quick look at coverage. Are there any NA?:

```
nacount <- NULL
for(i in 1:4){
  if(!is.null(tset[[i]])){
    nacount <- rbind(nacount,
                     unlist(lapply(tset[[i]], function(x){sum(is.na(x$Cov))}))
                     )
    rownames(nacount)[nrow(nacount)] <- names(tset)[i]
  }
}
nacount

#                             1007 1012 1013 1014 1015 3367 1335
# snp.tables.full.unfiltered     0    0    0    0    0    0    0
# snp.tables.full.qfiltered      0    0    0    0    0    0    0
```

Seemingly no. What's average in unq- vs q-filtered:

```
snp.tables.unqfil <- tset.picker(c(1,2), table.set = tset)
snp.tables.qfil   <- tset.picker(c(3,4), table.set = tset)
cov.unqfil <- unlist(lapply(snp.tables.unqfil, function(x){mean(x$Cov)}))
cov.qfil   <- unlist(lapply(snp.tables.qfil,   function(x){mean(x$Cov,na.rm=T)}))
cov.both <- rbind(cov.unqfil,cov.qfil,cov.qfil/cov.unqfil)
i <- 1
if(!is.null(snp.tables.unqfil)){
  rownames(cov.both)[i] <- which.snp.tables(snp.tables.unqfil)
  i <- i+1
}
if(!is.null(snp.tables.qfil)){
  rownames(cov.both)[i] <- which.snp.tables(snp.tables.qfil)
```

```
  i <- i+1
}
if(i==3){
  rownames(cov.both)[i] <- 'Ratio'
}
cat('Mean Coverage:\n'); cov.both

# Mean Coverage:
#                      1007       1012       1013       1014       1015       3367       1335
# trunc-unfiltered 37.0555484 70.8060724 69.6610432 33.1009373 61.5365159 64.0284488 107.7425968
# trunc-qfiltered  28.2750286 51.3249686 45.4036337 13.7261052 48.7880005 44.8042054  81.8823765
# Ratio             0.7630444  0.7248668  0.6517794  0.4146742  0.7928301  0.6997547   0.7599815
```

## 5.1   Table 1 Data

Throw together the conveniently-available Table 1 data, *in Table 1 row order*:

```
# if coverage unavailable, build NA vector
if(!is.null(cov.unqfil)){cov.unqfilv <- cov.unqfil} else {cov.unqfilv <- rep(NA,times=7)}
if(!is.null(cov.qfil  )){cov.qfilv   <- cov.qfil  } else {cov.qfilv   <- rep(NA,times=7)}
t1data.df <- data.frame(
  id        = st.locs(1:7, id=T, loc=F, date=F),
  loc       = st.locs(1:7, id=F, loc=T, date=F),
  date      = st.locs(1:7, id=F, loc=F, date=T),
  cov.unq   = cov.unqfilv,
  cov.q     = cov.qfilv,
  SNPs.4    = snp.counts[,4],
  SNPs.2    = snp.counts[,2],
  olap.ny.4 = snp.pctofny[,4]*100,
  olap.ny.2 = snp.pctofny[,2]*100
)
t1row.order <- c(7,1,2,5,3,6,4)
print(t1data.df[t1row.order,],digits=3)

#              id              loc date cov.unq cov.q SNPs.4 SNPs.2 olap.ny.4 olap.ny.2
# 1335 CCMP1335       New York 1958   107.7  81.9 153901 180187     100.0     100.0
# 1007 CCMP1007       Virginia 1964    37.1  28.3 161103 182700      87.5      92.8
# 1012 CCMP1012    W. Australia 1965    70.8  51.3 166089 186491      89.5      94.3
# 1015 CCMP1015   Puget Sound 1985    61.5  48.8 174701 190302      93.9      97.1
# 1013 CCMP1013           Wales 1973    69.7  45.4 247737 304293      40.3      61.7
# 3367 CCMP3367           Italy 2007    64.0  44.8 240413 290904      40.5      59.6
# 1014 CCMP1014 N. Pacific Gyre 1971    33.1  13.7  89184 149531      52.1      78.2
```

# 6   Shared-SNPs P-Value

Text of the main paper quotes a "p-value" for the observed degree of SNP sharing in L-clade (and/or L-clade excluding Gyre) under a null model that these isolates were sampled from a population globally in Hardy-Weinberg equilibrium. Details of this analysis are as follows.

## 6.1   SNP Concordance

Arbitrarily pick one isolate, say, $A$, as the "template". Arbitrarily pick a heterozygous (aka "SNP") position in $A$. Let $p_1$, and $q_1 = 1 - p_1$ be the frequencies in the overall population of the two nucleotides observed at that position in $A$. (Positions having 3 or 4 nucleotide variants segregating in the population are assumed to be negligibly rare.) Under the HWE null model, a second isolate $B$ will also be heterozygous at the same position with probability $2p_1q_1 \leq 1/2$. Similarly, this position will be heterozygous in a third isolate $C$ with the same probability, *independently*, and so on for isolates $D$ and $E$. Overall, (assuming HWE) the probability that a heterozygous position in $A$ is simultaneously heterozygous in the other 4 isolates is at most $1/2^4 = 1/16$. Continuing, suppose we pick a second heterozygous position in $A$, *on a different chromosome* with allele frequencies $p_2, q_2 = 1 - p_2$, say. Again assuming HWE, this

position will be a SNP in all of $B, C, D$ and $E$ with probability $(2p_2q_2)^4 \leq 1/16$, and this is independent of the first position, since segregation on different chromosomes is unlinked. Repeat this at 24 heterozygous positions in $A$, one per chromosome. Then, the number of five-way concordant positions observed should be dominated by the number observed when sampling from a binomial distribution with parameters $n = 24$ and $p = 1/16$, i.e., we expect at most $1/16 = 6.25\%$ of positions to agree, or at most $24/16 = 1.5$ five-way concordant positions in total. In sharp contrast, choosing CCMP 1014 (North Pacific Gyre) as the template, we see many more five-way concordant positions than predicted under these assumptions:

```r
gyre.count <- sum(snp.tables[[4]]$snp)
# NOTE: what we now calle "refined" SNPs were once called "filtered" SNPs and I have NOT tried
# to update variable names and annotation in the code below to reflect the terminology change...
# 'unfil.' => unfiltered for consistency; see below.
unfil.fiveway.count    <- sum( snp.tables[[4]]$snp * i4.snps)
unfil.fiveway.percent <- unfil.fiveway.count / gyre.count * 100
unfil.p.value <- pbinom(floor(unfil.fiveway.count/gyre.count*24)-1, 24, 1/16, lower.tail = FALSE)
consistency.comparison <-
  data.frame(
    fiveway.count   = unfil.fiveway.count,
    fiveway.percent = unfil.fiveway.percent,
    p.value         = unfil.p.value
  )
consistency.comparison

#   fiveway.count fiveway.percent       p.value
# 1         70687        79.25973 4.142632e-19
```

Namely, 89184 positions are called as SNPs in CCMP1014, of which 70687 or 79.2597327% are also called as SNPs in *all four* other L-clade isolates. 79.2597327% of 24 is 19.0223358, and the probability of seeing 19 or more "Heads" in 24 flips of a biased coin with $P(\text{Heads}) \leq 1/16$, i.e., our p-value under the HWE null hypothesis, is at most: $4.1426317 \times 10^{-19}$ based on this simple binomial model. This is obviously strong evidence against the null hypothesis.

This analysis is potentially overly-simplistic in four respects, addressed below.

1. "$2pq \leq 1/2$" is conservative. Neutral theory predicts that most variant nucleotides are rare in the population, so $2pq \ll 1/2$ is to be expected. This should make our quoted p-value very conservative.

2. Effect of Erroneous SNP calls. We base our analysis on *predicted* (by SAMTOOLS) heterozygous positions, not absolute-truth, which may affect our conclusions. However,

   • False negatives in $A$ are irrelevant, since we never examine those positions. (This is the motivation for using CCMP1014 as the template; it has the lowest predicted SNP rate, likely due to a high false negative rate in that sequencing run. As noted elsewhere, it had the lowest coverage and lowest sequence quality of the 7 isolates, both of which impare SNP calling.)

   • False negatives in $BCDE$ make such positions appear *non*-concordant. For our purpose, this makes our statistic more conservative since it can only deflate a statistic that we argue is nevertheless unexpectedly large.

   • False positive calls in $A$ are conservatively treated, as well: barring simultaneous false-positive calls in all of $BCDE$, such a position will appear non-concordant, again deflating the statistic. The *false* positive rates in $B, C, D$ and $E$ are unknown, but cannot exceed SAMTOOLS *total* positive rate, which is below 1% in all 7 isolates, suggesting a simultaneous $BCDE$ false positive rate $< 10^{-8}$, which will have a negligible effect.

   • A potentially more serious issue is a true positive in $A$ aligned to false positives in $BCD$ and/or $E$. (I.e., a position that is polymorphic in the population and heterozygous in A, under the HWE null model is likely to be homozygous for one of the two alleles in one or more of $BCDE$; false positive SNP calls in all of those isolates would make the site appear concordant, i.e., provide evidence against the null model.) However, (a) my impression is that SAMTOOLS is more prone to false negative calls than to false positive calls (see Section 4), and (b) we would need a high rate of false positives to turn a truely heterozygous but non-concordant $A$ call into a false "concordant" call—I'd expect at most half (especially given point 1

above) of $BCDE$ to be heterozygous, but all would need to be falsely declared heterozygous. Such a high false positive rate on $BCDE$ seems unlikely (see previous bullet), and would likely be counterbalanced by a similarly increased rate of false positives on $A$, which, as noted, tend to deflate our statistic (previous bullet again).

- Systematic errors. If there were, say, a sequence-context-dependent bias in the DNA sequencing, mapping and/or SNP-calling that tended to suggest (or hide) a SNP at some position, we're going to systematically over- (or under-) estimate concordant SNPs across isolates. The discordance of called SNPs between the L- and H-clades and within the H-clade suggests that this is not a major problem, but it is worth noting as a possibility.

3. Discordant nucleotides at "concordant" SNP positions. A "shared" SNP at a given position might be, say, G/C in one isolate vs T/C in another, reflecting an unexpected tri-allelic position in the population or a technical sequencing error. It is inappropriate to count such a "shared" SNP position as evidence against the null hypothesis, since it isn't clear that it is truely shared. Instead, I will identify such inconsistent positions, based on the "stringency [[2]]" criteria established above, and treat each as non-concordant. I.e., a position will be considered to be a "5-way concordant SNP" if and only if it was called as a SNP by SAMTOOLS (independently) in all 5 L-clade isolates, *and* shows the same dominant non-reference nucleotide in all 5, according to criteria [[2]] above. As it turns out, this correction has a very minor effect on the resulting p-value:

```
# 'unfil.' => Ignoring "consistency";  'fil.' => Filtering for "consistency":
fil.fiveway.count   <- sum((snp.tables[[4]]$snp * i4.snps)[union.snps == 1] & consistent[[2]])
fil.fiveway.percent <- fil.fiveway.count / gyre.count * 100
fil.p.value <- pbinom(floor(fil.fiveway.count/gyre.count*24)-1, 24, 1/16, lower.tail = FALSE)
# append new stats to previous table for easy comparison
consistency.comparison <-
  rbind(consistency.comparison,
        data.frame(
           fiveway.count   = fil.fiveway.count,
           fiveway.percent = fil.fiveway.percent,
           p.value         = fil.p.value
        )
  )
rownames(consistency.comparison) <- c('unfiltered', 'consistency.filtered')
consistency.comparison


#                      fiveway.count fiveway.percent      p.value
# unfiltered                   70687        79.25973 4.142632e-19
# consistency.filtered         69915        78.39411 1.976512e-17
```

In particular, it removes 0.9% of five-way consistent positions (only 772 of 70687 positions), and still shows a highly significant p-value.

4. "$P(E[X]) \neq E[P(X)]$". I'm expressing this poorly, but finding the p-value based on the *expected* number of concordant positions is somewhat non-standard. A more typical set-up would use the *actual* value of some statistic, then calculate the probability of observing a value that extreme (or more extreme) under the null model. The fundamental problem is that we have thousands of SNPs, but I don't see an easy way to use more than 24 of them at a time, because potential genetic linkage seemingly destroys statistical independence, which is key to most simple analyses. A somewhat more formal, but still non-standard, approach is the following. Suppose we randomly sample one SNP per chromosome and count the number $X$ of them that are 5-way concordant. What I outlined above calculated the p-value based on $E[X]$, the expected value of $X$, i.e., $P(E[X])$. Alternatively, we can calculate $E[P(X)]$, the expected p-value. (They are not the same.) In effect, this averages the p-values that would be seen over many different randomly-sampled sets of 24 SNPs. This is not difficult to calculate. First, the probability that we would observe $0 \leq i \leq 24$ concordant positions in a sample of 24, given that 78.39% of positions are concordant follows this binomial distribution:

```
x.equals.i.distribution <- dbinom(0:24, 24, fil.fiveway.percent/100)
print(x.equals.i.distribution, digits=3)


#  [1] 1.07e-16 9.33e-15 3.89e-13 1.04e-11 1.97e-10 2.86e-09 3.29e-08 3.07e-07 2.37e-06 1.53e-05
```

```
# [11] 8.31e-05 3.84e-04 1.51e-03 5.05e-03 1.44e-02 3.48e-02 7.11e-02 1.21e-01 1.71e-01 1.96e-01
# [21] 1.78e-01 1.23e-01 6.09e-02 1.92e-02 2.90e-03
```

Second, the p-value corresponding to $0 \leq i \leq 24$ observed concordant positions also follows a different binomial distribution:

```r
p.val.of.x.equals.i <- c(1, pbinom(0:23, 24, 1/16, lower.tail = F))
print(p.val.of.x.equals.i, digits=3)
```

```
#  [1] 1.00e+00 7.88e-01 4.48e-01 1.87e-01 5.95e-02 1.49e-02 3.01e-03 4.99e-04 6.90e-05 8.02e-06
# [11] 7.89e-07 6.60e-08 4.72e-09 2.87e-10 1.49e-11 6.59e-13 2.46e-14 7.66e-16 1.98e-17 4.14e-19
# [21] 6.88e-21 8.70e-23 7.88e-25 4.56e-27 1.26e-29
```

Finally, the expected (or "average") p-value is just the weighted average of the latter values, weighted by the former:

```r
e.of.p.of.x <- sum(x.equals.i.distribution * p.val.of.x.equals.i)
e.of.p.of.x
```

```
# [1] 6.939136e-10
```

This is still highly significant, but weaker than the $P(E[X])$ analysis, basically because $X < E[X]$ has a fair probability of occurring, and the corresponding p-value $P(X)$ rises rapidly as $X$ declines.

Another way to look at the numbers:

```r
pvdf <- data.frame(x.density=x.equals.i.distribution,
                   x.cdf=cumsum(x.equals.i.distribution),
                   pval.of.x=p.val.of.x.equals.i)
print(pvdf, digits=4)
```

```
#    x.density      x.cdf pval.of.x
# 1  1.071e-16 1.071e-16 1.000e+00
# 2  9.325e-15 9.432e-15 7.875e-01
# 3  3.891e-13 3.985e-13 4.476e-01
# 4  1.035e-11 1.075e-11 1.869e-01
# 5  1.972e-10 2.080e-10 5.950e-02
# 6  2.862e-09 3.070e-09 1.490e-02
# 7  3.289e-08 3.596e-08 3.010e-03
# 8  3.068e-07 3.428e-07 4.994e-04
# 9  2.366e-06 2.709e-06 6.899e-05
# 10 1.526e-05 1.797e-05 8.015e-06
# 11 8.306e-05 1.010e-04 7.887e-07
# 12 3.836e-04 4.846e-04 6.603e-08
# 13 1.508e-03 1.992e-03 4.716e-09
# 14 5.050e-03 7.042e-03 2.875e-10
# 15 1.440e-02 2.144e-02 1.493e-11
# 16 3.482e-02 5.626e-02 6.590e-13
# 17 7.107e-02 1.273e-01 2.456e-14
# 18 1.213e-01 2.487e-01 7.662e-16
# 19 1.712e-01 4.199e-01 1.977e-17
# 20 1.962e-01 6.161e-01 4.143e-19
# 21 1.780e-01 7.941e-01 6.877e-21
# 22 1.230e-01 9.170e-01 8.701e-23
# 23 6.085e-02 9.779e-01 7.884e-25
# 24 1.920e-02 9.971e-01 4.556e-27
# 25 2.903e-03 1.000e+00 1.262e-29
```

E.g., row 9 in that table says that the concordance rate (78%) is so high that a sample of 24 SNPs will almost always have 9 or more five-way concordant positions (probability of fewer is only 2.709e-06), while under the null model, seeing 9 or more is very unlikely (probability at most 6.899e-05).  ***AM I OFF-BY-ONE INTERPRETING ROW 9 HERE??***

## 6.2   Notes

In earlier drafts, an analog of the above analysis was based on the concordance of *refined* SNPs. This now seems to me to be questionable, since the "refined" SNP calling makes SNPs called across L-clade non-independent. OTOH, the above analysis seems valid: SAMTOOLS was run on each isolate independently, and likewise "criterion [[2]]" is evaluated independently in each strain, and is being used here solely to remove SNP predictions, not to add them. "Systematic errors" as outlined above remain a potential problem, but again discordance with/within H-clade suggests that this is of limited concern.

For completeness, I did a similar analysis including a sample of H-clade comparisons: Gyre vs Italy, NY vs Italy, NY vs Italy+Wales, and of Italy vs Wales. As expected, none of these show a statistically significant p-value, although the $\approx 40\%$ concordance in the 2-way comparisons, while $< 1/2$ as predicted, is a bit higher than I expected based on "neutral theory implies many rare variants." (I did not bother to include "criterion[[2]] filtering" in these calculations.)

```
# 'gi.twoway' => gyre vs italy 2-way concordance;
# 'ni.twoway' => new york vs italy 2-way concordance;
# not bothering with criterion[[2]] filtering
gi.twoway.count    <- sum(snp.tables[[4]]$snp * snp.tables[[6]]$snp)
gi.twoway.percent <- gi.twoway.count / gyre.count * 100
gi.p.value <- pbinom(floor(gi.twoway.count/gyre.count*24)-1, 24, 1/2, lower.tail = FALSE)
ny.count <- sum(snp.tables[[7]]$snp)
ni.twoway.count    <- sum(snp.tables[[7]]$snp * snp.tables[[6]]$snp)
ni.twoway.percent <- ni.twoway.count / ny.count * 100
ni.p.value <- pbinom(floor(ni.twoway.count/ny.count*24)-1, 24, 1/2, lower.tail = FALSE)
niw.threeway.count    <- sum(snp.tables[[7]]$snp * snp.tables[[6]]$snp * snp.tables[[3]]$snp)
niw.threeway.percent <- niw.threeway.count / ny.count * 100
niw.p.value <- pbinom(floor(niw.threeway.count/ny.count*24)-1, 24, 1/4, lower.tail = FALSE)
it.count <- sum(snp.tables[[6]]$snp)
iw.twoway.count    <- sum(snp.tables[[6]]$snp * snp.tables[[3]]$snp)
iw.twoway.percent <- iw.twoway.count / it.count * 100
iw.p.value <- pbinom(floor(iw.twoway.count/it.count*24)-1, 24, 1/2, lower.tail = FALSE)
consistency.comparison <-
  rbind(consistency.comparison,
        data.frame(
          fiveway.count   = c(gi.twoway.count,   ni.twoway.count,   niw.threeway.count,   iw.twoway.count),
          fiveway.percent = c(gi.twoway.percent, ni.twoway.percent, niw.threeway.percent, iw.twoway.percent),
          p.value         = c(gi.p.value,        ni.p.value,        niw.p.value,          iw.p.value)
        )
  )
colnames(consistency.comparison)[1:2] <- c('552232way.count', '552232way.percent') # old col names misleading
rownames(consistency.comparison)[3:6] <- c('gyre.vs.italy', 'new.york.vs.italy',   # new rows
                                           'ny.vs.it.plus.wales', 'it.vs.wales')
consistency.comparison

#                       552232way.count 552232way.percent      p.value
# unfiltered                      70687          79.25973 4.142632e-19
# consistency.filtered            69915          78.39411 1.976512e-17
# gyre.vs.italy                   34023          38.14922 9.242052e-01
# new.york.vs.italy               62332          40.50136 9.242052e-01
# ny.vs.it.plus.wales             35796          23.25911 7.533516e-01
# it.vs.wales                    102113          42.47399 8.462719e-01
```

## 6.3   P-Value: The Bottom Line

So, what to say in the body of the paper? $E[P(X)]$ is highly significant, and conservative, but complex to explain. $P(E[X])$ is simpler to explain, but may be criticized as misleading if we aren't very careful in that explanation. I'm slightly leaning towards the last option, but want to sleep on it and draft the key sentence or two before settling.

# 7   Sharing

The following analysis looks at the sharing patterns among the consistent SNPs. I assume that shared SNPs reflect shared ancestry, and that SNPs accumulate slowly over time. Then, in outline, the story is consistent with what we

have seen in other analyses—there seem to be 3 groups: 1013 (Wales) in one, 3367 (Italy) in another, and the other 5 in a third, with some hints as to the order of divergence. A caveat is that in a sexual population, non-shared SNPs do not immediately imply non-shared ancestry; they may merely reflect Hardy-Weinberg capturing a homozygous state in one isolate vs the other. (Or read errors, etc.) Thus, if we are right that the H-isolates retain sex, then the large number of "private" SNPs in H may be at least partially due to HWE.

Analysis is broken into cases based on how many strains share a particular SNP.

## 7.1   Code

To categorize SNPs by sharing patterns, first convert the 7-way consistent sharing pattern into a 7-bit binary number, and tabulate based on that:

```
# convert (n x 7) 0-1 matrix to n vector of 0-127
tobin <- function(x){
  bin <- integer(nrow(x)) # initialized to 0
  for(i in 1:7){
    bin <- bin*2 + as.integer(x[,i]>0)
  }
  return(bin)
}

# get full set of patterns
snp.pattern.all <- lapply(non.refs,tobin)
# prune to just the consistent ones
snp.pattern <- snp.pattern.all
for(i in 1:3){
  snp.pattern[[i]][!consistent[[i]]] <- NA
}

# analogous to built-in ``table'' but simpler.  Count entries in an integer
# vector sharing values in a (smallish) range.  Result is a 2-column matrix with
# the shared values in col 1 and count of occurrences of that value in col 2.
# Out-of-range values cause subscript error.
mytable <- function(vec, therange=range(vec,na.rm=T)){
  counts <- matrix(0,nrow=therange[2]-therange[1]+1,ncol=2,dimnames=list(NULL,c('val','count')))
  counts[1:nrow(counts),1] <- therange[1]:therange[2]
  for(i in 1:length(vec)){
    if(!is.na(vec[i])){
      counts[vec[i]-therange[1]+1,2] <- counts[vec[i]-therange[1]+1,2] + 1
    }
  }
  return(counts)
}

pattern.counts <- lapply(snp.pattern, function(x){mytable(x,c(0,127))})
```

To display the results, build a data frame whose i-th row, $0 \le i \le 127$ shows one of the 128 possible sharing patterns, with counts of the numbers of consistent, shared SNPs with that pattern according to criteria c1-c3.

```
tobitvec <- function(x){
  bitvec <- integer(7)
  for(i in 7:1){
    bitvec[i] <- x %% 2
    x <- x %/% 2
  }
  return(bitvec)
}

flg <- function(x){
  return(ifelse(x==1,'X',''))
}

pat.summary <- function(listOfTbls){
  mydf <- data.frame(pat=0:127,sharedBy=NA,
```

```
                          tp1007='',tp1012='',tp1013='',tp1014='',tp1015='',tp3367='',tp1335='',
                          count1=NA,count2=NA,count3=NA,count4=NA,stringsAsFactors=F)

  for(i in 1:128){
    bvec <- tobitvec(i-1)
    mydf[i,'sharedBy']=sum(bvec)
    mydf[i,'tp1007']=flg(bvec[1])
    mydf[i,'tp1012']=flg(bvec[2])
    mydf[i,'tp1013']=flg(bvec[3])
    mydf[i,'tp1014']=flg(bvec[4])
    mydf[i,'tp1015']=flg(bvec[5])
    mydf[i,'tp3367']=flg(bvec[6])
    mydf[i,'tp1335']=flg(bvec[7])
  }

  for(i in 1:length(listOfTbls)){
    tbl <- listOfTbls[[i]]
    if(!is.null(tbl)){
      mydf[,9+i] <- tbl[,2]   ## count1/2/3/4 are columns 10/11/12/13 in mydf
      #for(j in 1:length(tbl)){
      #   k <- as.integer(rownames(tbl)[j]);
      #   mydf[k+1,9+i] <- tbl[j]   ## count1/2/3 are columns 10/11/12
      #}
    }
  }

  mydf$pat <-as.octmode(mydf$pat)   # display bit pattern in octal
  return(mydf)
}

pat.summaries <- pat.summary(pattern.counts)
```

## 7.2   Sanity Checks

Some sanity checking: table sums equal to number of consistent positions?

```
all(consistent.count == apply(pat.summaries[,10:13],2,sum))
```

```
# [1] TRUE
```

More sanity checking: visually inspect a pattern with small counts, specifically pattern 12, i.e., consistent SNPs shared by only strains 1014 and 1015 (2nd and 3 rows from bottom, binary code $12 = 2^3 + 2^2$). There are only 10 such positions on Chr1. Chr1 2524239 has pattern 12 under criteria c1 and c2 but not c3; Chr1 1088766 has in c2 only. Both look good. Neither position is a *called* SNP except in 1015. However, all but 1 nonreference read agree with the called SNP (the exception being one read in Wales). Both 1014 and 1015 have at least 2 non-reference reads, comprising at least 5% of coverage, and in both strains, those reads are on the same non-reference base, satisfying criterion c2. The other strains have higher coverage and/or lower non-reference counts, so they do not satisfy c2. Position 2524239 also satisfies c1, but not c3, since 2 reads out of 35 is below the 10% threshold. (It is pattern 4 inder c3, i.e., a SNP private to 1015.) Position 1088766 is also pattern 4 under c3 (2 reads out of 56 in 1335 is below both thresholds), and it is not consistent under c1, since the single A read in 1013 is discordant with the other non-reference reads.

```
unlist(lapply(snp.pattern,function(x){sum(x==12,na.rm=T)}))
```

```
# [1] 143 136 176 417
```

```
sp1 <- snp.pattern[[1]]==12
sp2 <- snp.pattern[[2]]==12
sp3 <- snp.pattern[[3]]==12
sp4 <- snp.pattern[[4]]==12
c(sum(sp1,na.rm=T), sum(sp2,na.rm=T), sum(sp3,na.rm=T), sum(sp4,na.rm=T))
```

```
# [1] 143 136 176 417
```

```
r1 <- rownames(non.refs[[1]])[which(sp1)]
r2 <- rownames(non.refs[[2]])[which(sp2)]
r3 <- rownames(non.refs[[3]])[which(sp3)]
r4 <- rownames(non.refs[[4]])[which(sp4)]

r2

#   [1] "Chr1:1088766"     "Chr1:2524239"     "Chr2:713075"      "Chr2:1464209"
#   [5] "Chr2:2406031"     "Chr2:2480466"     "Chr2:2480532"     "Chr2:2480838"
#   [9] "Chr2:2481998"     "Chr2:2483322"     "Chr2:2488863"     "Chr2:2489189"
#  [13] "Chr2:2490933"     "Chr2:2492886"     "Chr2:2492887"     "Chr2:2497794"
#  [17] "Chr2:2500122"     "Chr2:2503000"     "Chr2:2507585"     "Chr2:2507680"
#  [21] "Chr2:2510117"     "Chr2:2513923"     "Chr2:2515103"     "Chr2:2516669"
#  [25] "Chr2:2516751"     "Chr2:2518558"     "Chr2:2518653"     "Chr2:2518980"
#  [29] "Chr2:2519285"     "Chr2:2519288"     "Chr2:2519718"     "Chr2:2520984"
#  [33] "Chr2:2521271"     "Chr2:2522648"     "Chr2:2524223"     "Chr2:2524439"
#  [37] "Chr2:2525160"     "Chr2:2525463"     "Chr2:2527916"     "Chr2:2528472"
#  [41] "Chr2:2528769"     "Chr2:2529076"     "Chr2:2529140"     "Chr2:2529186"
#  [45] "Chr2:2529432"     "Chr2:2529684"     "Chr2:2530064"     "Chr2:2530216"
#  [49] "Chr2:2530239"     "Chr2:2530294"     "Chr2:2530768"     "Chr2:2530896"
#  [53] "Chr2:2531114"     "Chr2:2531285"     "Chr2:2531498"     "Chr2:2531567"
#  [57] "Chr2:2532173"     "Chr2:2532365"     "Chr2:2533028"     "Chr2:2533171"
#  [61] "Chr2:2533440"     "Chr2:2534441"     "Chr2:2535121"     "Chr2:2535122"
#  [65] "Chr2:2535314"     "Chr2:2535493"     "Chr2:2535503"     "Chr2:2535509"
#  [69] "Chr2:2535862"     "Chr2:2536242"     "Chr2:2537201"     "Chr2:2537864"
#  [73] "Chr2:2537917"     "Chr2:2538072"     "Chr2:2538498"     "Chr2:2539318"
#  [77] "Chr2:2543595"     "Chr2:2545615"     "Chr2:2545798"     "Chr2:2546865"
#  [81] "Chr2:2546991"     "Chr2:2547055"     "Chr2:2547086"     "Chr2:2547120"
#  [85] "Chr2:2547155"     "Chr2:2547212"     "Chr2:2547248"     "Chr2:2547318"
#  [89] "Chr2:2547554"     "Chr2:2547938"     "Chr2:2547944"     "Chr2:2548131"
#  [93] "Chr2:2549281"     "Chr2:2551574"     "Chr2:2551930"     "Chr2:2554708"
#  [97] "Chr2:2554860"     "Chr2:2555005"     "Chr2:2555203"     "Chr2:2555820"
# [101] "Chr3:496665"      "Chr4:901220"      "Chr4:983962"      "Chr4:1086210"
# [105] "Chr4:1086589"     "Chr5:7509"        "Chr5:141375"      "Chr5:1397904"
# [109] "Chr6:1034519"     "Chr7:399475"      "Chr8:556556"      "Chr10:95217"
# [113] "Chr12:422344"     "Chr12:458461"     "Chr13:963939"     "Chr14:56058"
# [117] "Chr15:417704"     "Chr16a:39914"     "Chr16a:39917"     "Chr16a:394030"
# [121] "Chr17:461465"     "Chr18:673261"     "Chr19a_19:303090" "Chr19a_19:308244"
# [125] "Chr19b_31:4468"   "Chr19b_31:138559" "Chr19c_29:64170"  "Chr19c_29:64811"
# [129] "Chr19c_29:65720"  "Chr20:15766"      "Chr20:230994"     "Chr20:486431"
# [133] "Chr22:249009"     "Chr22:380816"     "Chr23:274291"     "Chr24:114599"

c1 <- as.integer(unlist(lapply(strsplit(r1[1:min(20,length(r1))],':',fixed=TRUE),function(x){x[2]})))
c2 <- as.integer(unlist(lapply(strsplit(r2[1:min(20,length(r2))],':',fixed=TRUE),function(x){x[2]})))
c3 <- as.integer(unlist(lapply(strsplit(r3[1:min(20,length(r3))],':',fixed=TRUE),function(x){x[2]})))
c4 <- as.integer(unlist(lapply(strsplit(r4[1:min(20,length(r4))],':',fixed=TRUE),function(x){x[2]})))

c1

#  [1]  198498   914018 1317406 1481838 1501481 1878058 2145849 2388286 2524239 2718093   62676
# [12]  393166   458314  713075 1416054 2148271 2149651 2310069 2406031 2480466

c2

#  [1] 1088766 2524239  713075 1464209 2406031 2480466 2480532 2480838 2481998 2483322 2488863
# [12] 2489189 2490933 2492886 2492887 2497794 2500122 2503000 2507585 2507680

c3

#  [1]  371484 1210354 1886633 2264683 2898352  207186  903516 1264023 1276745 1464904 1464905
# [12] 2229060 2347253 2406031 2439655 2480532 2480838 2483322 2488863 2489189

c4

#  [1]  518347  691730  767408 1049906 1390437 2072951 2254059 2254789 2264683 2823796 2898352
# [12] 2998868   77394   77407  155680  761325  968120 1182096 1222176 1264023

seecounts(c2,snp.tables=snp.tables)
```

```
#       chr      pos Ref Strain   A   G   C   T SNP  exon indel nrf rat
# 1   Chr1 1088766   G
# 2                      1007   0  32   1   0   0 FALSE FALSE
# 3                      1012   0  39   1   0   0 FALSE FALSE
# 4                      1013   1  74   0   0   0 FALSE FALSE
# 5                      1014   0  26   2   0   0 FALSE FALSE
# 6                      1015   0  38   9   0   1 FALSE FALSE
# 7                      3367   0  36   1   0   0 FALSE FALSE
# 8                      1335   0  54   2   0   0 FALSE FALSE
# 9   Chr1 2524239   C
# 10                     1007   0   0  37   0   0  TRUE FALSE
# 11                     1012   0   0  47   0   0  TRUE FALSE
# 12                     1013   0   0  62   0   0  TRUE FALSE
# 13                     1014   0   0  33   2   0  TRUE FALSE
# 14                     1015   0   0  11  15   1  TRUE FALSE
# 15                     3367   0   0  41   0   0  TRUE FALSE
# 16                     1335   0   0  95   0   0  TRUE FALSE
# 17  Chr1  713075   T
# 18                     1007   0   0   0  43   0  TRUE FALSE
# 19                     1012   0   0   0 115   0  TRUE FALSE
# 20                     1013   1   0   0  89   0  TRUE FALSE
# 21                     1014   0   0   0  63   0  TRUE FALSE
# 22                     1015   0   0   0  97   0  TRUE FALSE
# 23                     3367   0   0   0  75   0  TRUE FALSE
# 24                     1335   0   0   0 149   0  TRUE FALSE
# 25  Chr1 1464209   T
# 26                     1007   0   0   0  26   0 FALSE FALSE
# 27                     1012   0   0   0  57   0 FALSE FALSE
# 28                     1013   0   0   0  36   0 FALSE FALSE
# 29                     1014   0   0   0  25   0 FALSE FALSE
# 30                     1015   0   0   0  40   0 FALSE FALSE
# 31                     3367   0   0   0  52   0 FALSE FALSE
# 32                     1335   0   0   0 104   0 FALSE FALSE
# 33  Chr1 2406031   C
# 34                     1007   0   0  29   0   0  TRUE FALSE
# 35                     1012   0   0  52   0   0  TRUE FALSE
# 36                     1013   0   0  71   0   0  TRUE FALSE
# 37                     1014   0   0  20   0   0  TRUE FALSE
# 38                     1015   0   0  51   0   0  TRUE FALSE
# 39                     3367   0   0  60   0   0  TRUE FALSE
# 40                     1335   0   0  97   0   0  TRUE FALSE
# 41  Chr1 2480466   A
# 42                     1007  33   0   0   0   0  TRUE FALSE
# 43                     1012  57   0   0   1   0  TRUE FALSE
# 44                     1013  56   0   0   0   0  TRUE FALSE
# 45                     1014  24   0   0   0   0  TRUE FALSE
# 46                     1015  62   0   0   0   0  TRUE FALSE
# 47                     3367  43   0   0   1   0  TRUE FALSE
# 48                     1335  99   0   0   0   0  TRUE FALSE
# 49  Chr1 2480532   G
# 50                     1007   0  35   0   0   0  TRUE FALSE
# 51                     1012   0  50   0   0   0  TRUE FALSE
# 52                     1013   0  70   0   0   0  TRUE FALSE
# 53                     1014   0  19   0   0   0  TRUE FALSE
# 54                     1015   0  43   0   0   0  TRUE FALSE
# 55                     3367   0  46   0   0   0  TRUE FALSE
# 56                     1335   0 113   0   0   0  TRUE FALSE
# 57  Chr1 2480838   T
# 58                     1007   0   0   0  16   0  TRUE FALSE
# 59                     1012   0   0   0  30   0  TRUE FALSE
# 60                     1013   0   0   0  46   0  TRUE FALSE
# 61                     1014   0   0   0  26   0  TRUE FALSE
# 62                     1015   0   0   0  23   0  TRUE FALSE
# 63                     3367   0   0   0  28   0  TRUE FALSE
# 64                     1335   0   0   0 100   0  TRUE FALSE
# 65  Chr1 2481998   T
# 66                     1007   0   0   0  34   0  TRUE FALSE
```

```
# 67                       1012   0   0   0  54   0  TRUE FALSE
# 68                       1013   0   0   0  90   0  TRUE FALSE
# 69                       1014   0   0   0  36   0  TRUE FALSE
# 70                       1015   0   0   0  59   0  TRUE FALSE
# 71                       3367   0   0   0  64   0  TRUE FALSE
# 72                       1335   0   0   0 100   0  TRUE FALSE
# 73   Chr1 2483322   A
# 74                       1007  35   0   0   0   0  TRUE FALSE
# 75                       1012  59   1   0   0   0  TRUE FALSE
# 76                       1013 117   0   0   0   0  TRUE FALSE
# 77                       1014  53   0   0   0   0  TRUE FALSE
# 78                       1015  78   0   0   0   0  TRUE FALSE
# 79                       3367  71   0   1   0   0  TRUE FALSE
# 80                       1335 131   0   0   0   0  TRUE FALSE
# 81   Chr1 2488863   C
# 82                       1007   0   0  29   0   0 FALSE FALSE
# 83                       1012   0   0  55   0   0 FALSE FALSE
# 84                       1013   0   0  49   0   0 FALSE FALSE
# 85                       1014   0   0  27   0   0 FALSE FALSE
# 86                       1015   0   0  48   0   0 FALSE FALSE
# 87                       3367   0   0  63   0   0 FALSE FALSE
# 88                       1335   0   0  90   0   0 FALSE FALSE
# 89   Chr1 2489189   C
# 90                       1007   0   0  37   0   0 FALSE FALSE
# 91                       1012   0   0  88   0   0 FALSE FALSE
# 92                       1013   0   0  60   0   0 FALSE FALSE
# 93                       1014   0   0  45   0   0 FALSE FALSE
# 94                       1015   0   0  68   0   0 FALSE FALSE
# 95                       3367   0   0  39   0   0 FALSE FALSE
# 96                       1335   0   0 132   0   0 FALSE FALSE
# 97   Chr1 2490933   G
# 98                       1007   0  35   0   0   0 FALSE FALSE
# 99                       1012   0  72   0   0   0 FALSE FALSE
# 100                      1013   0  60   0   0   0 FALSE FALSE
# 101                      1014   0  25   1   0   0 FALSE FALSE
# 102                      1015   0  44   0   0   0 FALSE FALSE
# 103                      3367   0  47   1   1   0 FALSE FALSE
# 104                      1335   0  71   0   0   0 FALSE FALSE
# 105 Chr1 2492886   T
# 106                      1007   0   0   0  34   0 FALSE FALSE
# 107                      1012   0   0   0  98   0 FALSE FALSE
# 108                      1013   0   1   0  60   0 FALSE FALSE
# 109                      1014   0   0   0  37   0 FALSE FALSE
# 110                      1015   0   0   0  75   0 FALSE FALSE
# 111                      3367   0   0   0  73   0 FALSE FALSE
# 112                      1335   0   0   0 125   0 FALSE FALSE
# 113 Chr1 2492887   G
# 114                      1007   0  33   0   0   0 FALSE FALSE
# 115                      1012   0  95   0   0   0 FALSE FALSE
# 116                      1013   0  59   0   0   0 FALSE FALSE
# 117                      1014   0  36   0   0   0 FALSE FALSE
# 118                      1015   0  72   0   0   0 FALSE FALSE
# 119                      3367   0  71   0   0   0 FALSE FALSE
# 120                      1335   0 125   0   0   0 FALSE FALSE
# 121 Chr1 2497794   T
# 122                      1007   0   0   0  43   0  TRUE FALSE
# 123                      1012   0   0   0  77   0  TRUE FALSE
# 124                      1013   0   0   0  76   0  TRUE FALSE
# 125                      1014   0   0   0  25   0  TRUE FALSE
# 126                      1015   0   0   0  75   0  TRUE FALSE
# 127                      3367   0   0   0  63   0  TRUE FALSE
# 128                      1335   0   0   0 130   0  TRUE FALSE
# 129 Chr1 2500122   A
# 130                      1007  33   0   0   0   0 FALSE FALSE
# 131                      1012  71   0   0   0   0 FALSE FALSE
# 132                      1013  58   0   0   0   0 FALSE FALSE
# 133                      1014  31   0   0   0   0 FALSE FALSE
```

```
# 134                            1015  48    0    0    0    0 FALSE FALSE
# 135                            3367  46    0    0    0    0 FALSE FALSE
# 136                            1335  74    0    0    0    0 FALSE FALSE
# 137 Chr1 2503000   T
# 138                            1007   0    0    0   35    0 FALSE FALSE
# 139                            1012   0    0    0   65    0 FALSE FALSE
# 140                            1013   0    0    0   83    0 FALSE FALSE
# 141                            1014   0    0    0   33    0 FALSE FALSE
# 142                            1015   0    0    0   54    0 FALSE FALSE
# 143                            3367   0    0    0   56    0 FALSE FALSE
# 144                            1335   0    0    0   65    0 FALSE FALSE
# 145 Chr1 2507585   A
# 146                            1007  40    0    0    0    0  TRUE FALSE
# 147                            1012  65    0    0    0    0  TRUE FALSE
# 148                            1013  56    0    0    0    0  TRUE FALSE
# 149                            1014  31    0    0    0    0  TRUE FALSE
# 150                            1015  47    0    0    0    0  TRUE FALSE
# 151                            3367  78    0    1    0    0  TRUE FALSE
# 152                            1335 118    0    0    0    0  TRUE FALSE
# 153 Chr1 2507680   A
# 154                            1007  30    0    0    0    0 FALSE FALSE
# 155                            1012  82    0    0    0    0 FALSE FALSE
# 156                            1013  62    0    0    0    0 FALSE FALSE
# 157                            1014  32    1    0    0    0 FALSE FALSE
# 158                            1015  70    0    0    0    0 FALSE FALSE
# 159                            3367  87    0    0    0    0 FALSE FALSE
# 160                            1335 125    1    0    0    0 FALSE FALSE
```

Position 1088766, however, in a good example of the situation that motivated this analysis—one strain has a G/C SNP and 5 of the other 6 strains have nonreference reads consistent with that SNP. Although, excluding 1015, the nonreference read counts are not high enough to justify a SNP call in any strain considered in isolation, the fact that they *consistently* agree with the 1015 SNP suggests that they are real. One alternative hypothesis is that there is some sequence-dependent bias at this locus that favors misreading a G as a C. On the other hand, one could equally well posit a shared SNP, and a locus-dependant bias that *supresses* C reads, explaining the unbalanced readout that we observe. However, it is hard to reconcile either view with the significant strain-specific patterns that we see in the shared SNPs (as seen below). I think a more likely explanation is that (a) there are some number of relatively rare SNPs present in each of the sampled populations, (b) some of these SNPs happened to be present in one or two cells of the roughly 5-10 cells that we believe constituted the founding population of the culture grown for sequencing, and (c) stochastic effects during culture growth and during sequencing may have further perturbed the apparent frequency of each variant, but the bottom line is that the above-threshold presence of consistent non-reference reads is evidence for shared SNPs at the population level (and the proportions of such reads represent estimates of the population-level frequencies of the variants, albeit a noisy estimate at any specific position).

An aside: I was curious to see whether there is any consistent pattern to positions that are called consistent SNPs in all but Italy, so I repeated the above, basically. My summary is that coverage in Italy tends to be below average in these positions, but otherwise they don't stand out. For the record:

```
abit <- snp.pattern[[2]]==125
abit[is.na(abit)]<-F
sum(abit)

# [1] 14648


rabit <- rownames(non.refs[[2]])[which(abit)]
rabits <- rabit[1:20]
cabit <- as.integer(unlist(lapply(strsplit(rabits,':',fixed=TRUE),function(x){x[2]})))
cabit

#  [1]  1244  1575  6485  7181  7220  7661  8144  8208  8518  8552  8567  8670  8685 14361 15254
# [16] 15280 16103 17587 18904 25546


seecounts(cabit,snp.tables=snp.tables)

#     chr   pos Ref Strain   A   G   C   T SNP  exon indel nrf rat
# 1  Chr1 1244   G
# 2                    1007   3  30   0   0   0  TRUE FALSE
# 3                    1012   5  54   0   0   0  TRUE FALSE
```

```
# 4                        1013  15  47   0   0   1  TRUE FALSE
# 5                        1014   3  30   0   0   0  TRUE FALSE
# 6                        1015  21  68   0   0   1  TRUE FALSE
# 7                        3367   0  10   0   0   0  TRUE FALSE
# 8                        1335 108 108   0   0   1  TRUE FALSE
# 9   Chr1  1575   G
# 10                       1007  26  11   0   0   0  TRUE FALSE
# 11                       1012  49  27   0   0   0  TRUE FALSE
# 12                       1013  19  28   0   0   0  TRUE FALSE
# 13                       1014  15  17   0   0   0  TRUE FALSE
# 14                       1015  46  42   0   0   1  TRUE FALSE
# 15                       3367   0  11   0   0   0  TRUE FALSE
# 16                       1335  37  99   0   0   0  TRUE FALSE
# 17  Chr1  6485   G
# 18                       1007  26  20   0   0   0  TRUE FALSE
# 19                       1012  33  39   0   0   0  TRUE FALSE
# 20                       1013  54  48   0   0   0  TRUE FALSE
# 21                       1014  13  10   0   0   0  TRUE FALSE
# 22                       1015  34  41   0   0   1  TRUE FALSE
# 23                       3367   0  42   0   0   0  TRUE FALSE
# 24                       1335  71  69   0   0   0  TRUE FALSE
# 25  Chr1  7181   G
# 26                       1007   0  37  31   0   0  TRUE FALSE
# 27                       1012   0  66  36   0   0  TRUE FALSE
# 28                       1013   0  30  86   0   0  TRUE FALSE
# 29                       1014   0  19   8   0   0  TRUE FALSE
# 30                       1015   0  44  33   0   1  TRUE FALSE
# 31                       3367   0  33   0   0   0  TRUE FALSE
# 32                       1335   0  94  78   0   0  TRUE FALSE
# 33  Chr1  7220   C
# 34                       1007  17   0  26   6   0  TRUE FALSE
# 35                       1012  45   0  31  14   0  TRUE FALSE
# 36                       1013 112   1  41  14   0  TRUE FALSE
# 37                       1014  16   1  16   2   0  TRUE FALSE
# 38                       1015  66   0  26   7   1  TRUE FALSE
# 39                       3367   0   0  24   0   0  TRUE FALSE
# 40                       1335  68   0  51  25   0  TRUE FALSE
# 41  Chr1  7661   T
# 42                       1007   0   0  10  14   0  TRUE FALSE
# 43                       1012   0   0   7  24   0  TRUE FALSE
# 44                       1013   0   0  32  23   1  TRUE FALSE
# 45                       1014   0   0   8  11   0  TRUE FALSE
# 46                       1015   0   0   6  41   0  TRUE FALSE
# 47                       3367   0   0   0   8   0  TRUE FALSE
# 48                       1335   0   0   8  42   0  TRUE FALSE
# 49  Chr1  8144   G
# 50                       1007  10  16   0   1   0  TRUE FALSE
# 51                       1012  19  28   0   0   1  TRUE FALSE
# 52                       1013  63  67   0   0   0  TRUE FALSE
# 53                       1014   7  12   0   0   0  TRUE FALSE
# 54                       1015  18  28   0   0   0  TRUE FALSE
# 55                       3367   0   7   0   0   0  TRUE FALSE
# 56                       1335  17  58   0   0   1  TRUE FALSE
# 57  Chr1  8208   G
# 58                       1007   0  15   0   8   1  TRUE FALSE
# 59                       1012   0  28   0  16   0  TRUE FALSE
# 60                       1013   0  24   0  63   0  TRUE FALSE
# 61                       1014   0  15   0   4   0  TRUE FALSE
# 62                       1015   0  25   0  13   1  TRUE FALSE
# 63                       3367   0   9   0   1   0  TRUE FALSE
# 64                       1335   0  49   0  21   1  TRUE FALSE
# 65  Chr1  8518   T
# 66                       1007   0   0  20  18   1 FALSE FALSE
# 67                       1012   0   0  45  30   1 FALSE FALSE
# 68                       1013   0   0  57  75   1 FALSE FALSE
# 69                       1014   0   0  10  32   0 FALSE FALSE
# 70                       1015   0   0  41  18   1 FALSE FALSE
# 71                       3367   0   0   0  11   0 FALSE FALSE
# 72                       1335   0   0 120  71   1 FALSE FALSE
# 73  Chr1  8552   G
# 74                       1007   3  13   0   0   0  TRUE FALSE
# 75                       1012  21  31   0   1   0  TRUE FALSE
# 76                       1013  33  35   0   0   1  TRUE FALSE
# 77                       1014   7  15   0   0   0  TRUE FALSE
# 78                       1015  14  22   0   0   0  TRUE FALSE
# 79                       3367   0  28   0   0   0  TRUE FALSE
# 80                       1335  27  59   0   0   0  TRUE FALSE
# 81  Chr1  8567   A
# 82                       1007  16  18   0   0   1  TRUE FALSE
# 83                       1012  34  35   0   0   1  TRUE FALSE
```

```
# 84                     1013  66  75   0    0   1  TRUE FALSE
# 85                     1014   9   4   0    0   0  TRUE FALSE
# 86                     1015  17  31   0    0   1  TRUE FALSE
# 87                     3367  29   0   0    0   0  TRUE FALSE
# 88                     1335  59  44   0    0   1  TRUE FALSE
# 89  Chr1  8670   A
# 90                     1007  19   0   0    7   0  TRUE FALSE
# 91                     1012  36   0   0   12   0  TRUE FALSE
# 92                     1013  44   0   0   12   0  TRUE FALSE
# 93                     1014  10   0   0    7   0  TRUE FALSE
# 94                     1015  24   0   0   11   1  TRUE FALSE
# 95                     3367  18   0   0    0   0  TRUE FALSE
# 96                     1335  27   0   0    6   0  TRUE FALSE
# 97  Chr1  8685   G
# 98                     1007   7  16   0    0   0  TRUE FALSE
# 99                     1012  12  37   0    0   0  TRUE FALSE
# 100                    1013  18  30   0    0   1  TRUE FALSE
# 101                    1014   5  32   0    0   0  TRUE FALSE
# 102                    1015  11  35   0    0   1  TRUE FALSE
# 103                    3367   0  12   0    0   0  TRUE FALSE
# 104                    1335   5  45   0    0   0  TRUE FALSE
# 105 Chr1 14361   A
# 106                    1007  29   7   0    0   0 FALSE FALSE
# 107                    1012  54   6   0    0   0 FALSE FALSE
# 108                    1013  28  12   0    0   1 FALSE FALSE
# 109                    1014  22   2   1    0   0 FALSE FALSE
# 110                    1015  51   9   0    0   0 FALSE FALSE
# 111                    3367  12   1   0    0   0 FALSE FALSE
# 112                    1335  64   8   0    0   0 FALSE FALSE
# 113 Chr1 15254   T
# 114                    1007  11   0   0   22   1 FALSE FALSE
# 115                    1012  28   0   0   53   1 FALSE FALSE
# 116                    1013  39   0   0   66   1 FALSE FALSE
# 117                    1014   3   0   0   14   1 FALSE FALSE
# 118                    1015  18   0   0   39   1 FALSE FALSE
# 119                    3367   0   0   0   89   0 FALSE FALSE
# 120                    1335  15   0   0   63   1 FALSE FALSE
# 121 Chr1 15280   T
# 122                    1007   0  14   0   32   1 FALSE FALSE
# 123                    1012   0  31   0   53   1 FALSE FALSE
# 124                    1013   0   6   0  102   0 FALSE FALSE
# 125                    1014   0   3   1   40   0 FALSE FALSE
# 126                    1015   0  22   1   51   1 FALSE FALSE
# 127                    3367   0   0   0   74   0 FALSE FALSE
# 128                    1335   0  26   0  109   1 FALSE FALSE
# 129 Chr1 16103   A
# 130                    1007  12   0  14    0   1 FALSE FALSE
# 131                    1012  50   0  19    0   1 FALSE FALSE
# 132                    1013  29   0  15    0   1 FALSE FALSE
# 133                    1014  28   0   2    0   0 FALSE FALSE
# 134                    1015  37   0  10    0   1 FALSE FALSE
# 135                    3367  41   0   0    0   0 FALSE FALSE
# 136                    1335  56   0  12    0   0 FALSE FALSE
# 137 Chr1 17587   A
# 138                    1007  22   2   0    0   0 FALSE FALSE
# 139                    1012  62   6   0    1   0 FALSE FALSE
# 140                    1013  22  12   0    0   1 FALSE FALSE
# 141                    1014  22   2   0    0   0 FALSE FALSE
# 142                    1015  29   3   0    0   0 FALSE FALSE
# 143                    3367  20   1   0    0   0 FALSE FALSE
# 144                    1335  82  11   0    0   0 FALSE FALSE
# 145 Chr1 18904   T
# 146                    1007   0   5   0   34   0 FALSE FALSE
# 147                    1012   0   4   0   39   0 FALSE FALSE
# 148                    1013   0   9   0   21   0 FALSE FALSE
# 149                    1014   0   3   0   21   0 FALSE FALSE
# 150                    1015   0   9   0   48   0 FALSE FALSE
# 151                    3367   0   5   0   96   0 FALSE FALSE
# 152                    1335   0  27   0   73   1 FALSE FALSE
# 153 Chr1 25546   A
# 154                    1007  31   0   0   14   1 FALSE FALSE
# 155                    1012  64   0   0   22   1 FALSE FALSE
# 156                    1013  20   0   0   50   1 FALSE FALSE
# 157                    1014  22   0   1   18   1 FALSE FALSE
# 158                    1015  64   0   0   18   1 FALSE FALSE
# 159                    3367  73   0   0    0   0 FALSE FALSE
# 160                    1335  80   0   0    5   0 FALSE FALSE
```

More sanity: there are 83 sites on Chr1 shared by zero strains in the tightest condition. (I.e., SAMTOOLS called

it a SNP, but the read counts/proportions fall below our 3rd threshold). Are they due to low coverage? Seemingly yes:

```
zp3 <- snp.pattern[[3]] == 0
zr3 <- rownames(non.refs[[3]])[which(zp3)]
zc3 <- as.integer(unlist(lapply(strsplit(zr3[1:min(100,length(zr3))],':',fixed=TRUE),function(x){x[2]})))
zc3

#   [1]    91284   127986   161271   196862   196864   199166   282391   289344   289363   314132   314661
#  [12]   438976   447253   475823   501830   501975   504462   652889   657955   692139   709443   762174
#  [23]   826899   856950   875379   913014   938651   967184  1036942  1100300  1113225  1181146  1203203
#  [34]  1210360  1212223  1224082  1270250  1270251  1348311  1431628  1473437  1516083  1526912  1628300
#  [45]  1637082  1686331  1736789  1763837  1782580  1967158  2024930  2075603  2098145  2110716  2194162
#  [56]  2242316  2258647  2261176  2325671  2376777  2432898  2441781  2498706  2550796  2554565  2581374
#  [67]  2614631  2619528  2659281  2675254  2691279  2703771  2737914  2744068  2802553  2842231  2846930
#  [78]  2906880  2931365  2948653  2957936  3014028  3016252    31184   101081   109502   195069   198570
#  [89]   208189   278516   292413   297200   320853   349833   357243   357245   403824   418951   422130
# [100]   459508
```

```
seecounts(zc3[1:5], snp.tables=snp.tables)
```

```
#      chr   pos Ref Strain   A  G   C   T SNP  exon indel nrf rat
# 1  Chr1  91284   T
# 2                    1007    0  0   0  17   0 FALSE FALSE
# 3                    1012    0  0   0  38   0 FALSE FALSE
# 4                    1013    2  0   0  13   0 FALSE FALSE
# 5                    1014    0  0   0  20   0 FALSE FALSE
# 6                    1015    0  0   0  35   0 FALSE FALSE
# 7                    3367    3  0   0  12   1 FALSE FALSE
# 8                    1335    0  0   0  47   0 FALSE FALSE
# 9  Chr1 127986   A
# 10                   1007   47  0   0   0   0  TRUE FALSE
# 11                   1012   92  0   0   0   0  TRUE FALSE
# 12                   1013   19  1   0   0   0  TRUE FALSE
# 13                   1014   73  0   0   0   0  TRUE FALSE
# 14                   1015   83  0   0   0   0  TRUE FALSE
# 15                   3367   13  3   0   0   1  TRUE FALSE
# 16                   1335  160  0   0   0   0  TRUE FALSE
# 17 Chr1 161271   A
# 18                   1007   31  0   0   0   0  TRUE FALSE
# 19                   1012   47  0   0   0   0  TRUE FALSE
# 20                   1013   18  3   0   0   0  TRUE FALSE
# 21                   1014   30  0   0   0   0  TRUE FALSE
# 22                   1015   59  0   0   0   0  TRUE FALSE
# 23                   3367    8  3   0   0   1  TRUE FALSE
# 24                   1335  102  0   0   0   0  TRUE FALSE
# 25 Chr1 196862   C
# 26                   1007    0  0  10   0   0 FALSE FALSE
# 27                   1012    0  0  22   0   0 FALSE FALSE
# 28                   1013    0  0   8   2   0 FALSE FALSE
# 29                   1014    0  0  14   0   0 FALSE FALSE
# 30                   1015    0  0  18   0   0 FALSE FALSE
# 31                   3367    1  0   4   3   1 FALSE FALSE
# 32                   1335    0  0  18   0   0 FALSE FALSE
# 33 Chr1 196864   T
# 34                   1007    0  0   0  11   0 FALSE FALSE
# 35                   1012    0  0   1  23   0 FALSE FALSE
# 36                   1013    3  0   0   8   1 FALSE FALSE
# 37                   1014    0  0   0  12   0 FALSE FALSE
# 38                   1015    1  0   1  19   0 FALSE FALSE
# 39                   3367    3  0   0   4   1 FALSE FALSE
# 40                   1335    0  0   1  19   0 FALSE FALSE
```

## 7.3   Main Analysis

Turning to the main analysis, there is a large increase in the number of consistent positions between the loose and medium stringency levels; medium and tight are similar in most respects. The likely interpretation is that the loose

criterion is including many "SNPs" induced by read errors, and that either of the tighter criteria are successfully filtering them out. In the interest of simplicity, the narrative below will focus on the shared SNPs at the medium stringency level (the "count2" column in the data frame), although the numbers for all three (sometimes all 4) are displayed. Also note that the prose and some comments in the code were based on the Chr1 analysis, and so may occasionally be off-target for the whole-genome data.

```r
# Show a subset of pat.summaries, optionally with totals of count_i in last row, and optionally
# aggregating low-count rows as ``Other''
#
#   sharedBy=c(2,4) selects SNPs shared by 2 or 4 strains,
#   subset=as.octmode('35') select those with sharing pattern a subset (optionally proper) of this
#   split=as.octmode('14') additionally restricts to patterns stradling split/subset minus split
#   c2.thresh=42 suppresses printout of rows with count2 < 42
#   restrict.to=c(0,42,127) restrict to these 3 rows
showgroup <- function(p.summ=pat.summaries, sharedBy=0:7, subset=127, split=NULL, proper.subset=F,
                      total=T, c2.thresh=0, fourteenth=F, restrict.to=NULL){
  # pick just those bit patterns that are subsets of 'subset'
  pick <- bitwAnd(0:127,bitwNot(subset))==0
  if(proper.subset){
    pick[subset+1] <- F
  }
  if(!is.null(split)){ # AND that stradle left/right subtrees
    cosplit <- bitwAnd(subset,bitwNot(split))
    pick <- pick & bitwAnd(0:127,split)!=0 & bitwAnd(0:127,cosplit)!=0
  }
  # and have desired shareBy counts
  pick <- pick & (p.summ$sharedBy %in% sharedBy)
  # and are among the set of interest
  if(!is.null(restrict.to)){
    pick <- pick & (0:127 %in% restrict.to)
  }
  # find rows with low counts
  pick.low <- pick & (p.summ$count2 < c2.thresh)
  # now show them
  show <- p.summ[pick & ! pick.low,]
  # rename columns just to narrow the printouts
  colnames(show) <- c('Pat','ShrBy','1007', '1012', '1013', '1014', '1015', '3367', '1335',
                      'count1', 'count2', 'count3','count4')
  show[,1] <- format(show[,1])   # convert octal col to char so can override in last row(2)
  nlow <- sum(pick.low)
  if(nlow > 0){
    n <- nrow(show)+1
    lows <- apply(p.summ[pick.low,10:13],2,sum)
    show[n,10:13] <- lows
    show[n,1:9] <- ''
    row.names(show)[n] <- 'Other'
    if(fourteenth){
      # do this: add 14th col just to hold this comment:
      show <- cbind(show,' ='', stringsAsFactors=F)
      show[n,14] <- paste('(', nlow, 'rows w/ c2 <', c2.thresh, ')')
    } else {
      ## or this (looks a bit funky, but fits across page without line-wrap):
      show[n,1:8] <-c('(', nlow, 'rows', 'w/', 'c2', '<', c2.thresh, ')')
    }
  }
  if(total){
    n <- nrow(show)+1
    tots <- apply(show[,10:13],2,sum)
    show[n,10:13] <- tots
    show[n,1:9] <- ''
    row.names(show)[n] <- 'Total'
    if(ncol(show)==14){show[n,14]<-''}
  }
  return(show)
}
```

First, are there any SNPs that are not "consistent SNPs?" Yes, a few in c3. As noted above, they seem to be mainly

low-coverage positions.

```
showgroup(pat.summaries,0,total=F)  # chr1 totals: 0 0 83

#   Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 1   0     0                                        0      2   1164      0
```

Next, look at completely shared SNPs, those found in all 7 strains.

```
showgroup(pat.summaries,7,total=F)  # Chr1 count1 = 8593, count2 = 7054, count3 = 4790 c4=1641

#       Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 128 177     7    X    X    X    X    X    X    X   82193  67223  46524  15186
```

I.e., of the 468117 consistent positions, 67223 or 14.4% are shared by all 7 strains.

Next look at singletons, aka private SNPs—SNPs that are called in one strain and no other strain has a significant number of non-ref reads at that position. Presumably these are variants that arose in a given population after it separated from the others.

```
showgroup(pat.summaries,1)   # chr1 totals: 9669   18865   19670   23574

#         Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 2       001     1                              X    199    620   1157   2260
# 3       002     1                         X       41774  84335  88149 105614
# 5       004     1                    X            921   2070   2578   4608
# 9       010     1               X                 208    559    714   1231
# 17      020     1          X                    47772  93481  96798 113191
# 33      040     1     X                          285    611   1031   2450
# 65      100     1  X                             121    321    542   2005
# Total                                          91280 181997 190969 231359
```

The import of shared/private SNPs changes between sexual and asexual populations. Presumably asexuals slowly gain and rarely lose private SNPs; shared ones predate separation of the lineages. In sexual lineages, however, SNPs may be rather freely "gained" or "lost," merely by recombination (converting between homo- and heterozygous in the sample we sequenced). Thus, the low private counts for the 5 L-isolates compared to the large count of het positions overall suggest that (a) they are asexual, and (b) none of them has been isolated from the others for very long (if at all). Conversely, the high counts for Italy and Wales suggest that (a) if asexual, they have been separated from each other and from the rest for a long time, but (b) if sexual, there is little surprise: we have ≈160K SNPs shared between the two (90K just in those two (below), plus 70K shared by all 7), and ≈90K additional positions that are het in one but not the other. These are close to, but not exactly equal to, the 1:2:1 ratios we would naively expect from two samples of a single HWE population. The most parsimonious explanation seems to be that the H-clade is sexual, but perhaps some het positions private to each population separates them.

Aside: counts of "consistent" SNPs minus these singletons yeilds count of shared SNPs:

```
singlets <- apply(pat.summaries[pat.summaries$sharedBy==1,10:13],2,sum)
rbind(consistent=consistent.count,singlets=singlets,shared=consistent.count-singlets)

#            count1 count2 count3 count4
# consistent 358872 468117 470970 474613
# singlets    91280 181997 190969 231359
# shared     267592 286120 280001 243254
```

The slightly higher count of shared positions in the medium case further supports this choice for subsequent analysis.

Next look at consistent SNPs shared between just a pair of isolates.

```
showgroup(pat.summaries,2) # chr 1 counts: 7641   9549   9472   6924

#         Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 4       003     2                         X    X   2266    281    432    587
# 6       005     2                    X         X    210    410    854   1407
# 7       006     2                    X    X       1282    119    261    590
```

| # | Pat | ShrBy | 1007 | 1012 | 1013 | 1014 | 1015 | 3367 | 1335 | count1 | count2 | count3 | count4 |
|---|-----|-------|------|------|------|------|------|------|------|--------|--------|--------|--------|
| # 10 | 011 | 2 | | | | X | | | X | 358 | 535 | 384 | 827 |
| # 11 | 012 | 2 | | | | X | | X | | 2060 | 158 | 59 | 93 |
| # 13 | 014 | 2 | | | | X | X | | | 143 | 136 | 176 | 417 |
| # 18 | 021 | 2 | | | X | | | | X | 2445 | 154 | 260 | 402 |
| # 19 | 022 | 2 | | | X | | | X | | 55300 | 87584 | 84944 | 58009 |
| # 21 | 024 | 2 | | | X | | X | | | 1406 | 178 | 371 | 625 |
| # 25 | 030 | 2 | | | X | X | | | | 2257 | 180 | 59 | 93 |
| # 34 | 041 | 2 | | X | | | | | X | 31 | 85 | 230 | 368 |
| # 35 | 042 | 2 | | X | | | | X | | 1429 | 117 | 224 | 394 |
| # 37 | 044 | 2 | | X | | | X | | | 72 | 215 | 895 | 1809 |
| # 41 | 050 | 2 | | X | | X | | | | 31 | 25 | 54 | 105 |
| # 49 | 060 | 2 | | X | X | | | | | 1651 | 98 | 247 | 388 |
| # 66 | 101 | 2 | X | | | | | | X | 19 | 33 | 75 | 314 |
| # 67 | 102 | 2 | X | | | | | X | | 887 | 98 | 126 | 351 |
| # 69 | 104 | 2 | X | | | | X | | | 39 | 105 | 356 | 1196 |
| # 73 | 110 | 2 | X | | | X | | | | 19 | 20 | 50 | 150 |
| # 81 | 120 | 2 | X | | X | | | | | 1009 | 117 | 116 | 309 |
| # 97 | 140 | 2 | X | X | | | | | | 591 | 1150 | 1281 | 2144 |
| # Total | | | | | | | | | | 73505 | 91798 | 91454 | 70578 |

I.e., of the 91798 paired SNPs, 87584 or 95.4% are found between Italy and Wales, with comparatively few shared between any other pairs (only).

SNPs shared among exactly 3 isolates are relatively rare. (The 5 trios containing both Italy and Wales predominate in the loose set, probably because they share many pairs that become triples with the addition of a few read errors.)

```
showgroup(pat.summaries,3) # chr 1 counts: 1438    294    671   1034
```

| # | Pat | ShrBy | 1007 | 1012 | 1013 | 1014 | 1015 | 3367 | 1335 | count1 | count2 | count3 | count4 |
|---|-----|-------|------|------|------|------|------|------|------|--------|--------|--------|--------|
| # 8 | 007 | 3 | | | | | X | X | X | 152 | 146 | 278 | 557 |
| # 12 | 013 | 3 | | | | X | | X | X | 350 | 253 | 183 | 338 |
| # 14 | 015 | 3 | | | | X | X | | X | 776 | 1050 | 757 | 1389 |
| # 15 | 016 | 3 | | | | X | X | X | | 109 | 62 | 65 | 152 |
| # 20 | 023 | 3 | | | X | | | X | X | 3322 | 481 | 838 | 533 |
| # 22 | 025 | 3 | | | X | | X | | X | 197 | 168 | 333 | 522 |
| # 23 | 026 | 3 | | | X | | X | X | | 1794 | 395 | 771 | 789 |
| # 26 | 031 | 3 | | | X | X | | | X | 361 | 255 | 178 | 361 |
| # 27 | 032 | 3 | | | X | X | | X | | 2845 | 237 | 112 | 86 |
| # 29 | 034 | 3 | | | X | X | X | | | 144 | 116 | 125 | 219 |
| # 36 | 043 | 3 | | X | | | | X | X | 175 | 70 | 104 | 133 |
| # 38 | 045 | 3 | | X | | | X | | X | 116 | 409 | 1369 | 1656 |
| # 39 | 046 | 3 | | X | | | X | X | | 103 | 124 | 367 | 604 |
| # 42 | 051 | 3 | | X | | X | | | X | 46 | 55 | 79 | 126 |
| # 43 | 052 | 3 | | X | | X | | X | | 101 | 11 | 27 | 22 |
| # 45 | 054 | 3 | | X | | X | X | | | 22 | 79 | 171 | 292 |
| # 50 | 061 | 3 | | X | X | | | | X | 169 | 67 | 98 | 115 |
| # 51 | 062 | 3 | | X | X | | | X | | 2161 | 289 | 447 | 469 |
| # 53 | 064 | 3 | | X | X | | X | | | 114 | 148 | 390 | 601 |
| # 57 | 070 | 3 | | X | X | X | | | | 123 | 21 | 18 | 24 |
| # 68 | 103 | 3 | X | | | | | X | X | 83 | 25 | 51 | 143 |
| # 70 | 105 | 3 | X | | | | X | | X | 35 | 113 | 283 | 805 |
| # 71 | 106 | 3 | X | | | | X | X | | 42 | 63 | 127 | 377 |
| # 74 | 111 | 3 | X | | | X | | | X | 16 | 9 | 19 | 139 |
| # 75 | 112 | 3 | X | | | X | | X | | 55 | 12 | 16 | 26 |
| # 77 | 114 | 3 | X | | | X | X | | | 8 | 38 | 69 | 365 |
| # 82 | 121 | 3 | X | | X | | | | X | 87 | 19 | 29 | 73 |
| # 83 | 122 | 3 | X | | X | | | X | | 1325 | 191 | 220 | 354 |
| # 85 | 124 | 3 | X | | X | | X | | | 66 | 81 | 167 | 400 |
| # 89 | 130 | 3 | X | | X | X | | | | 73 | 12 | 9 | 27 |
| # 98 | 141 | 3 | X | X | | | | | X | 65 | 109 | 208 | 519 |
| # 99 | 142 | 3 | X | X | | | | X | | 342 | 419 | 477 | 755 |
| # 101 | 144 | 3 | X | X | | | X | | | 193 | 1079 | 2430 | 4432 |
| # 105 | 150 | 3 | X | X | | X | | | | 47 | 44 | 78 | 238 |
| # 113 | 160 | 3 | X | X | X | | | | | 306 | 359 | 421 | 712 |
| # Total | | | | | | | | | | 15923 | 7009 | 11314 | 18353 |

Four-way sharing is more common, but dominated by the coastal (i.e., non-Gyre) L-clade isolates. This is likely a

reflection of the strong 5-way sharing among the L-clade, from which the Gyre commonly drops out due to the lower coverage/higher error rate in that sequencing run.

```
showgroup(pat.summaries,4) # chr 1 counts: 564   1346   2552   3479
```

| # | Pat | ShrBy | 1007 | 1012 | 1013 | 1014 | 1015 | 3367 | 1335 | count1 | count2 | count3 | count4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # 16 | 017 | 4 | | | | X | X | X | X | 344 | 361 | 250 | 564 |
| # 24 | 027 | 4 | | | X | | X | X | X | 441 | 602 | 1120 | 771 |
| # 28 | 033 | 4 | | | X | X | | X | X | 1174 | 1062 | 725 | 306 |
| # 30 | 035 | 4 | | | X | X | X | | X | 495 | 529 | 373 | 503 |
| # 31 | 036 | 4 | | | X | X | X | X | | 381 | 367 | 287 | 211 |
| # 40 | 047 | 4 | | X | | | X | X | X | 93 | 178 | 485 | 708 |
| # 44 | 053 | 4 | | X | | X | | X | X | 46 | 38 | 36 | 56 |
| # 46 | 055 | 4 | | X | | X | X | | X | 480 | 709 | 750 | 971 |
| # 47 | 056 | 4 | | X | | X | X | X | | 17 | 50 | 65 | 88 |
| # 52 | 063 | 4 | | X | X | | | X | X | 325 | 167 | 265 | 194 |
| # 54 | 065 | 4 | | X | X | | X | | X | 88 | 208 | 528 | 582 |
| # 55 | 066 | 4 | | X | X | | X | X | | 263 | 432 | 944 | 851 |
| # 58 | 071 | 4 | | X | X | X | | | X | 30 | 17 | 14 | 28 |
| # 59 | 072 | 4 | | X | X | X | | X | | 158 | 50 | 35 | 31 |
| # 61 | 074 | 4 | | X | X | X | X | | | 28 | 51 | 60 | 116 |
| # 72 | 107 | 4 | X | | | | X | X | X | 36 | 64 | 105 | 330 |
| # 76 | 113 | 4 | X | | | X | | X | X | 18 | 8 | 8 | 66 |
| # 78 | 115 | 4 | X | | | X | X | | X | 103 | 141 | 138 | 604 |
| # 79 | 116 | 4 | X | | | X | X | X | | 5 | 8 | 19 | 101 |
| # 84 | 123 | 4 | X | | X | | | X | X | 162 | 85 | 77 | 124 |
| # 86 | 125 | 4 | X | | X | | X | | X | 41 | 79 | 142 | 283 |
| # 87 | 126 | 4 | X | | X | | X | X | | 124 | 214 | 297 | 425 |
| # 90 | 131 | 4 | X | | X | X | | | X | 17 | 5 | 5 | 52 |
| # 91 | 132 | 4 | X | | X | X | | X | | 86 | 27 | 19 | 38 |
| # 93 | 134 | 4 | X | | X | X | X | | | 15 | 24 | 25 | 143 |
| # 100 | 143 | 4 | X | X | | | | X | X | 46 | 55 | 86 | 190 |
| # 102 | 145 | 4 | X | X | | | X | | X | 3312 | 9777 | 18140 | 23189 |
| # 103 | 146 | 4 | X | X | | | X | X | | 121 | 455 | 958 | 1795 |
| # 106 | 151 | 4 | X | X | | X | | | X | 33 | 77 | 71 | 220 |
| # 107 | 152 | 4 | X | X | | X | | X | | 27 | 21 | 24 | 67 |
| # 109 | 154 | 4 | X | X | | X | X | | | 896 | 1611 | 1430 | 1738 |
| # 114 | 161 | 4 | X | X | X | | | | X | 72 | 88 | 102 | 207 |
| # 115 | 162 | 4 | X | X | X | | | X | | 1552 | 1955 | 1909 | 1014 |
| # 117 | 164 | 4 | X | X | X | | X | | | 165 | 603 | 1060 | 1752 |
| # 121 | 170 | 4 | X | X | X | X | | | | 21 | 26 | 21 | 69 |
| # Total | | | | | | | | | | 11215 | 20144 | 30573 | 38387 |

Five-way sharing is much more common, and is strongly dominated by the 5 L-clade isolates.

```
showgroup(pat.summaries,5) # chr 1 counts: 3969   5047   4624   6125
```

| # | Pat | ShrBy | 1007 | 1012 | 1013 | 1014 | 1015 | 3367 | 1335 | count1 | count2 | count3 | count4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # 32 | 037 | 5 | | | X | X | X | X | X | 2087 | 1987 | 1386 | 620 |
| # 48 | 057 | 5 | | X | | X | X | X | X | 227 | 231 | 205 | 324 |
| # 56 | 067 | 5 | | X | X | | X | X | X | 422 | 685 | 1836 | 1151 |
| # 60 | 073 | 5 | | X | X | X | | X | X | 155 | 89 | 72 | 38 |
| # 62 | 075 | 5 | | X | X | X | X | X | | 200 | 185 | 233 | 328 |
| # 63 | 076 | 5 | | X | X | X | X | | X | 108 | 158 | 188 | 128 |
| # 80 | 117 | 5 | X | | | X | X | X | X | 42 | 40 | 35 | 241 |
| # 88 | 127 | 5 | X | | X | | X | X | X | 133 | 253 | 399 | 482 |
| # 92 | 133 | 5 | X | | X | X | | X | X | 56 | 31 | 16 | 52 |
| # 94 | 135 | 5 | X | | X | X | X | X | | 116 | 121 | 96 | 235 |
| # 95 | 136 | 5 | X | | X | X | X | | X | 41 | 71 | 63 | 106 |
| # 104 | 147 | 5 | X | X | | | X | X | X | 1372 | 3155 | 5536 | 10001 |
| # 108 | 153 | 5 | X | X | | X | | X | X | 35 | 32 | 30 | 96 |
| # 110 | 155 | 5 | X | X | | X | X | | X | 33045 | 38232 | 26997 | 30602 |
| # 111 | 156 | 5 | X | X | | X | X | X | | 492 | 681 | 566 | 735 |
| # 116 | 163 | 5 | X | X | X | | | X | X | 271 | 263 | 302 | 316 |
| # 118 | 165 | 5 | X | X | X | | X | | X | 1875 | 3928 | 6825 | 9715 |
| # 119 | 166 | 5 | X | X | X | | X | X | | 621 | 1958 | 3252 | 2688 |
| # 122 | 171 | 5 | X | X | X | X | | | X | 30 | 29 | 18 | 70 |

```
# 123    172     5    X    X    X    X         X          105      95      59      86
# 125    174     5    X    X    X    X    X                567     789     656     782
# Total                                                  42000   53013   48770   58796
```

Six-way sharing is also common, with the sets *ex*cluding Gyre, Italy, or Wales having the most mutually-shared SNPs.

```
showgroup(pat.summaries,6) # chr 1 counts:  4166    4741    5312   4722

#       Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 64    077    6         X    X    X    X    X    X     920     872     917     485
# 96    137    6    X         X    X    X    X    X     394     325     275     333
# 112   157    6    X    X         X    X    X    X   13257   11725    8245   12202
# 120   167    6    X    X    X         X    X    X    8614   16443   28402   15091
# 124   173    6    X    X    X    X         X    X     140      93      74     114
# 126   175    6    X    X    X    X    X         X   17128   14648   10156   12697
# 127   176    6    X    X    X    X    X    X        2303    2825    2133    1032
# Total                                             42756   46931   50202   41954
```

# 8   Trees

So, overall, the picture looks like a long shared history (67223 7-way shared positions), followed by a split of the 5 L-isolates from the 2 H-isolates, then a long shared history in the 5 (38232 quintuples), in parallel with a long shared history in H- (87584 pairs), then separate histories in Italy and Wales (>84335 "private" SNPs in each, although again if they are sexual, many of these just reflect HWE), and very limited differentiation among the 5 L-isolates.

Branch lengths of course depend on filtering criteria used (and, of course, full vs Chr1 differ by about a factor of 10), but the tree *topology* appears to be fairly stable. Various versions are drawn below, exactly to explore how robust this story is. I think we should go with "medium stringency" SNP filtering (based on un-qfiltered reads).

NOTE: Much of this analysis make less sense for q-filtered read data, since (a) the point of the SNP filtering was to try to correct for noise in the raw reads, which may (or may not; haven't looked closely, yet) be largely fixed by qfiltering (e.g., "loose" or no SNP filtering may be more appropriate, post-q-filtering, esp. if we had re-run SAMTools to call SNPs based on the q-filtered reads), and (b) tree topology *does* appear to change, in that Gyre's coverage has been so sharply reduced by qfiltering that it clearly stands aside from the others (and that's confirmed by bootstrap), but this also seems to be clearly a technical rather than a biological artifact. SO, code below will run on q-filtered data, but *is not tuned to it*. Likewise, most comments in the prose below were made to describe the un-q-filtered data, and *are misleading and in some cases flatly wrong* for qfiltered data, but it doesn't seem worthwhile to bother with a rewrite...

Trees are coded in newick format, which doesn't seem to tolerate line-breaks; print with line-wrap:.

```
# wrap a long char string across multiple lines in printout
cat.hardwrap <- function(str,width=80){
  while(nchar(str)>width){
    cat(substr(str,1,width),'\n')
    str <- substr(str,width+1,nchar(str))
  }
  cat(str,'\n')
}
```

Trees are built as follows. Code for drawing, especially, is specific to the topology of the medium tree, and placement of some of the figure elements have been hand-optimized for this case; drawings for the other variants will not be as pretty.

```
# set up for tree figs

# the newick parser in ape seems to be confused by commas and parens in
# tip names, and blanks are not allowed, so replace by *, <, >, _, resp.
newick.name <- function(name){
  name <- gsub(' ', '_', name, fixed=TRUE)
  name <- gsub(',', '*', name, fixed=TRUE)
  name <- gsub('(', '<', name, fixed=TRUE)
```

```r
  name <- gsub(')', '>', name, fixed=TRUE)
  return(name)
}
# undo above changes
newick.name.undo <- function(name){
 #name <- gsub('_', ' ', name, fixed=TRUE) # unnecessary; ape plot routine handles this one
  name <- gsub('*', ',', name, fixed=TRUE)
  name <- gsub('<', '(', name, fixed=TRUE)
  name <- gsub('>', ')', name, fixed=TRUE)
  return(name)
}

# make a newick string from tree;  see it below
# 'pre' is prefixed to ccmpid; 'nb' optionally included;
# 'alt' can be used instead of pre/ccmp/nb/where for less formal labeling
# 'newstyle'==T => new node label: [nb_]where[(pre-less-id)]
# 'newstyle'==F => old node label: [nb_][pre id]where
newickize <- function(tree,pre='CCMP',nb=TRUE,alt=F,newstyle=TRUE){
  if(is.null(tree$where)){
    # not a leaf; paste together newick from subtrees
    sub1 <- newickize(tree$sub1,pre=pre,nb=nb,alt=alt,newstyle=newstyle)
    sub2 <- newickize(tree$sub2,pre=pre,nb=nb,alt=alt,newstyle=newstyle)
    new <- paste( '(', sub1, ',', sub2, ')', sep='')
    if(!is.null(tree$length)){
      # internal node, add length
      return(paste(new, ':', tree$length, sep=''))
    } else {
      # top level; escape blanks and add trailing ';'
      return(paste(gsub(' ', '_', new), ';', sep=''))
    }
  } else {
    # a leaf; build label and branch length
    if(alt){
      # label is just alt; if alt omitted, default to where
      new <- newick.name(ifelse( is.null(tree$alt), tree$where, tree$alt ))
    } else {
      if(newstyle){
        # new node label = [nb_]where[(pre-less-id)]
        new <- ifelse( nb && !is.null(tree$nb), paste(tree$nb, '_', sep =''), '' )
        new <- newick.name(paste(new, tree$where, sep=''))
        new <- ifelse( is.null(tree$id), new, paste(new, '_(', tree$id, ')', sep='') )
        new <- newick.name(new)
      } else {
        # old style node label = [nb_][pre id]where
        new <- ifelse( nb && !is.null(tree$nb), paste(tree$nb, '_', sep =''), '' )
        new <- ifelse( is.null(tree$id), new, paste(new, pre, tree$id, '_', sep='') )
        new <- newick.name(paste(new, tree$where, sep=''))
      }
    }
    #add length to either
    new <- paste(new, ':', tree$length, sep='')
  }
  return(new)
}

# Make a tree as nested lists, **based on the chr1, count2 topology**, but using any of the counts.
#    Internal nodes have subtrees sub1/2 and length
#    Root has sub1/2, but no length
#    Leaves have where, length, optionally, id, alt, nb.  (Omit id for 'outgroup'. Use 'alt' for less formal
#     labeling in cartoon version; it defaults to 'where'. Use 'nb' to add abcde annotations for legend.)
# The single parameter v is any of the 4 count vectors contained in pat.summaries (most conveniently
# indexed in octal).  E.g., make.tree(pat.summaries[,'count2']) reproduces the count2 tree.
# (This was previously built by hand-pasting the edge lengths; tree.by.hand is retained in appendix
# for comparison, & its counts are in comments below).
#
make.tree <- function(v){
  pat.count <- function(pat, pat.counts=v){return(pat.counts[1+strtoi(pat,8)])}
  thetree <-
    list(
      sub1 = list(
        sub1 = list(
          sub1 = list(id=3367, length=pat.count('002'), where='Venice, Italy', alt='Venice'), #8813
          sub2 = list(id=1013, length=pat.count('020'), where='Wales, UK'),                    #9652
          length=pat.count('022')),                                                            #9365
        sub2 =  list(
          sub1 = list(
            sub1 = list(
              sub1 = list(id=1007, length=pat.count('100'), nb='e', where='Virginia, USA'),     #30
              sub2 = list(id=1012, length=pat.count('040'), nb='d', where='Perth, W. Australia', alt='Perth'), #61
              length=pat.count('140')),                                                         #19
```

```
          sub2 = list(
            sub1 = list(id=1015, length=pat.count('004'),nb='c', where='Washington, USA', alt='Puget Sound'), #207
            sub2 = list(id=1335, length=pat.count('001'), nb='b', where='New York, USA',   alt='NY'), #41
              length=pat.count('005')),                                                            #18
          length=pat.count('145')),                                                               #1005
        sub2 = list(id=1014, length=pat.count('010'), nb='a', where='N. Pacific Gyre'),        #61
          length=pat.count('155')),                                                               #3912
      length=pat.count('177')),                                                                   #7054
    sub2 = list(length=0, where='outgroup')
  )
  return(thetree)
}
```

Code to plot a tree given newick description. Again, code is somewhat general, but has some specializations tied to the medium-stringency, full-genome, un-qfiltered data.

```
# run following 2 lines after an R upgrade
# update.packages()
# install.packages("ape")
library(ape)
show.tree <- function(newick.str=newick.medium,
                      col.edge  ='darkblue', lwd.edge =2,
                      col.elabel='darkblue',                cex.elabel=0.8, font.elabel=3,
                      col.arrow ='red',      lwd.arrow=1.5, cex.arrow =0.9, font.arrow =4,
                      col.clade ='black',    lwd.clade=1,   cex.clade =1.0, font.clade =3,
                      col.legbox='beige',                   cex.legend=0.8,
                      col.tip   ='darkblue',                                font.tip   =4,
                      plusx=FALSE, pltdebug=FALSE, total.snps=consistent.count[2],
                      straight.arrow=FALSE){

  ####
  #
  # ADJUST NEWICK & GET LENGTHS, COORDINATES
  #
  newick.str.noout <- sub('outgroup','_',newick.str) # Hide outgroup ('_' prints as blank)
  the.tree <- read.tree(text=newick.str.noout)

  ## nasty hack: ape's newick parser seems to be confused by commas, () in tip labels, so
  ## newickize replaced them by '*<>'; before plotting, I want to convert them back, and hope
  ## this doesn't break anything else...  And if a revised version of ape changes the internal
  ## representation of a tree, this may need to be redone.
  the.tree$tip.label <- newick.name.undo(the.tree$tip.label)

  # extract branch lengths as char string of comma-separated numbers via pattern matching hack:
  # lengths always preceded by colon
  lengths.ch <- strsplit(paste(newick.str,':'),'[^0-9][^:]*:')[[1]]

  # then convert to ints, dropping empty string at front
  lengths.int <- scan(what=integer(),quiet=T,sep=',',text=lengths.ch[-1])

  # then to data frame with named rows; a..g are terminal branches; others are internal.
  # a..e match legend in plot; f/g = wales/italy.  lengths appear in postfix order of
  # newick tree, and ape draws the 1st of them at the bottom of the plot.
  lmed <- data.frame(lengths=lengths.int,
                     row.names=c('g','f','fg','e','d','de','c','b','bc','bcde','a','abcde','all','out'))

  # extract counts needed for legend:
#leg.counts <- c(      61, 41,207, 61, 30, 1005,  18, 19) #by hand, medium chr1
  leg.counts <- lmed[c('a','b','c','d','e','bcde','bc','de'),1]
  discord <- total.snps - sum(lmed$lengths)

  #tree.labels <- list( ## x,y,text; coords are all picked by eye
  #   3000, 3.62, paste(lmed['all'  ,1], 'shared by 7', sep='\n'), # 7054
  #   8900, 5.75, paste(lmed['abcde',1], 'by 5'        , sep='\n'), # 3912
  #  12000, 1.50, paste(lmed['fg'   ,1], 'shared by 2', sep='\n'), # 9365
  #  21000, 2.00, paste(lmed['f'    ,1], 'only\nin Wales'),        # 9652
  #  21000, 1.00, paste(lmed['g'    ,1], 'only\nin Italy'),        # 8813
  #  11500, 4.50, '*')
  # automating x-placement, below; retain above for comparison...
  tip <- integer(7)   # x coords of tree tips
  tip[1] <-sum(lmed[c('all','fg','g'),1])
  tip[2] <-sum(lmed[c('all','fg','f'),1])
  tip[3] <-sum(lmed[c('all','abcde','bcde','de','e'),1])
  tip[4] <-sum(lmed[c('all','abcde','bcde','de','d'),1])
  tip[5] <-sum(lmed[c('all','abcde','bcde','bc','c'),1])
  tip[6] <-sum(lmed[c('all','abcde','bcde','bc','b'),1])
  tip[7] <-sum(lmed[c('all','abcde','a'),1])

  inode <- integer(5) # x coords of (some) internal nodes
```

```
  inode[1] <- 0                                    # root
  inode[2] <- lmed['all',1]                        # lca of all
  inode[3] <- sum(lmed[c('all','fg'),1])           # lca H-clade
  inode[4] <- sum(lmed[c('all','abcde'),1])        # lca L-clade
  inode[5] <- sum(lmed[c('all','abcde','bcde'),1]) # lca L-clade, nonGyre
  tree.labels <- list( ## x,y,text; y coords partially picked by eye
    sum(inode[c(1,2)])/2, 3.62, paste(lmed['all'  ,1], 'shared by 7', sep='\n'), # 7054
    sum(inode[c(2,4)])/2, 5.75, paste(lmed['abcde',1], 'by 5'       , sep='\n'), # 3912
    sum(inode[c(2,3)])/2, 1.50, paste(lmed['fg'   ,1], 'shared by 2', sep='\n'), # 9365
    (inode[3]+tip[2])/2,  2.00, paste(lmed['f'    ,1], 'only\nin 1013'),         # 9652
    (inode[3]+tip[1])/2,  1.00, paste(lmed['g'    ,1], 'only\nin 3367'),         # 8813
    sum(inode[c(4,5)])/2, 4.35, '* ')


  tree.labels <- list( ## x,y,text; y coords partially picked by eye
    sum(inode[c(1,2)])/2, 3.62, paste(lmed['all'  ,1], 'in 7', sep='\n'), # 7054
    sum(inode[c(2,4)])/2, 5.75, paste(lmed['abcde',1], 'in 5', sep='\n'), # 3912
    sum(inode[c(2,3)])/2, 1.50, paste(lmed['fg'   ,1], 'in 2', sep='\n'), # 9365
    (inode[3]+tip[2])/2,  2.00, paste(lmed['f'    ,1], 'only\nin 1013'),  # 9652
    (inode[3]+tip[1])/2,  1.00, paste(lmed['g'    ,1], 'only\nin 3367'),  # 8813
    sum(inode[c(4,5)])/2, 4.35, '* ')

  ####
  #
  # BOGUS PLOT
  #
  # a messy bit: need string widths to set xlim; but strwidth needs x-scale so must plot first.
  # M plot completely invisible, overlay 2nd plot via par(new=F...) .
  #
  # PROVISIONALLY set x.lim here at about 30% wider than tree; fine tune it for the real plot
  # based on strwidth(tip labels) below.
  #
  provisional.tree.x.lim <- 1.3 * max(tip) # <== PROVISIONAL plot width
  plot(0,0, type='n', bty='n', xaxt='n', yaxt='n', xlab='', ylab='', xlim=c(0,provisional.tree.x.lim), ylim=c(0,7))

  tiplabel.x <- integer(7)
  for(i in 1:7){
    # see warning above about internals of the.tree; labels have '_', printed as ' '.
    tiplabel.x[i] <- tip[i]+strwidth(gsub('_',' ',the.tree$tip.label[i],fixed=T), font=font.tip)
  }


  # visually show tip coords & max x to debug placement issues
  plt.debug <- function(tree.x.lim, tip, tiplabel.x, spx=NULL, spy=NULL){
    if(pltdebug){ # F to hide/T to show debug
      cat('Tip labels:', paste(the.tree$tip.label,sep='',collapse='/'), '\n')
      axis(2) # useful only for placing labels
      for(i in 1:7){
        points(c(tip[i],tiplabel.x[i]),c(i,i)) # debug: do I have right tip coordinates?
      }
      lines(rep(tree.x.lim,2),c(0,7)) # where is right edge?
      if(!is.null(spx)){
        points(spx,spy) # show spline control points, for tweaking
      }
    }
  }

  plt.debug(provisional.tree.x.lim, tip, tiplabel.x)

  label.end.H <- max(tiplabel.x[1:2])
  label.end.L <- max(tiplabel.x[3:7])
  clade.dx <- strwidth('x') # space between clade marker line and its label
  xdel <- 3*clade.dx        # space between labeled clade tips and marker line

  tree.x.lim <- 1.03*(max(tiplabel.x)+xdel)  # <== FINAL plot width
  tree.y.lim <- 7
  if(pltdebug){cat('Plot width hacking:', provisional.tree.x.lim, tree.x.lim, tree.x.lim/1.03/max(tip), clade.dx)}

  par(new=T)  # I.e., NOT starting a new plot

  ####
  #
  # REAL PLOT
  #
  plot(the.tree,
       x.lim = c(0, tree.x.lim),
       y.lim = c(0, tree.y.lim),
       font=font.tip, label.offset=100,          # bold-italic, nudged slightly right
       tip.color=col.tip, edge.color=col.edge,
       edge.width=lwd.edge,
```

```r
      edge.lty=c(1,1,1,1, 1 ,1,1,1,1,1,1,1,1,0)     # 5th is bottleneck edge; 14th is outgroup
    )
lines(00+c(0,0),c(3.5,6),col='white',lwd=6)         # Hide vertical line to outgroup
axis(1, pos=0.25, at=seq(0,25,by=5)*10^round(log10(max(tip)/25)))

if(pltdebug){text(tip[1]+100, 1.0, 'Venice, Italy (3367)', adj=0, font=font.tip)}

####
#
# BOTTLENECK ANNOTATION
#
# spline/elipse control points (spy/y) & tweaks thereto (dx/y)
dx <- 0.01 * tree.x.lim
dy <- .04
spx <- c(7400, 7400, 9900, 10500) # by eye, chr1, for comparison
spx <- c(inode[2]+dx,inode[2]+dx,inode[4]-3*dx,inode[4]-dx)
spy <- c( 3.8,   3.9,  5.6-dy,    5.6-dy)

plt.debug(tree.x.lim, tip, tiplabel.x, spx, spy)

if(T){
  #elipse version, defined by rect thru 2 middle pts of spx/y
  spf<-function(x){
    ifelse(x <= spx[2], spy[1],
           ifelse(x >= spx[3], spy[4],
                  spy[2]+(spy[3]-spy[2])*sqrt(pmax(0,1-((x-spx[3])/(spx[3]-spx[2]))^2))))
  }
} else {
  # spline version
  spf <- splinefun(spx,spy,method='hyman')
}
serx <- seq(spx[1],spx[length(spx)],length.out=50)
sery <- spf(serx)
tailx <- spx[1]
taily <- spy[1]
headx <- spx[4]
heady <- spy[4]
textx <- (headx+tailx)/2+(headx-tailx)*(-.01)
texty <- (heady+taily)/2+(heady-taily)*(-.10)
bottle.txt <- "inbreeding\nLoH / LoS"
if(!straight.arrow){
  arrows(headx,heady,headx+tree.x.lim*1e-3,heady, length=.1,col=col.arrow,lwd=lwd.arrow)
  lines(rev(serx), rev(sery), lty=c(5,1),col=col.arrow, lwd=lwd.arrow)
  textangle <- 66
  textadj <- c(0,0)
} else {
  # Tweak positioning slightly; visualize a rectangle from 7-node to base of L-clade;
  # center text, rotated, on diagonal towards L-clade; ditto the straight arrow.
  llx <- inode[2] # the aforementioned rectangle
  urx <- inode[4]
  lly <- 3.62
  ury <- 5.75
  # rect(llx,lly,urx,ury) # show rect for debug
  textx <- (llx+urx)/2    # center text
  texty <- (lly+ury)/2
  textangle <- atan(grconvertY(ury-lly,to='dev')/grconvertX(urx-llx,to='dev'))*360/(2*pi)
  textadj <- c(0.50, 0.43) #tweak position; ".5" = center in x , ".43" raises, THEN rotate.
  alpha <- .78 # fraction along diag at which arrow begins
  beta  <- .95 # ... and ends
  arrows((1-alpha)*llx + alpha*urx,
         (1-alpha)*lly + alpha*ury,
         (1-beta)*llx  + beta*urx,
         (1-beta)*lly  + beta*ury, length=.1,col=col.arrow,lwd=lwd.arrow,angle=25)
}
if(T){
  text(textx, texty, bottle.txt, srt=textangle, font=font.arrow, cex=cex.arrow,
       col=col.arrow, adj=textadj)
} else {
  # experiment at wrapping text along curved path; unpretty, but retain for now, maybe revisit
  bottlec <- strsplit(bottle,split=NULL)[[1]]
  for(i in 1:length(bottlec)){
    text(xser[i],yser[i],bottlec[i], srt=65, font=4, cex=.7, col=col.arrow)
  }
}

####
#
# CLADE ANNOTATION
#
clade.L.x <- label.end.L + xdel
```

```r
  clade.H.x <- label.end.H + xdel
  dy <-.33
  lines(rep(clade.L.x,2),c(3-dy,7+dy),lwd=lwd.clade,col=col.clade)
  lines(rep(clade.H.x,2),c(1-dy,2+dy),lwd=lwd.clade,col=col.clade)
  text(clade.L.x+clade.dx,5.0,'L-clade',srt=90,font=font.clade,cex=cex.clade,col=col.clade)
  text(clade.H.x+clade.dx,1.5,'H-clade',srt=90,font=font.clade,cex=cex.clade,col=col.clade)

  ####
  #
  # LEGEND
  #
  # parameter plusx controls whether we try to annotate b/c (+) and d/e (x) sharing in tree; I think
  # it looks cluttered, rather than adding clarity, so I vote no, but code is here, in case.  "Logic,"
  # if any, for my symbol choice is that + overlaid on x looks like the * at the next level; this
  # analogy is more visible if we use pch 3/4/8 rather than Courier or Helvetica chars, but probably
  # should use same in both tree & legend, which will take a modicum of additional work.
  legend.text <- c('a: only in 1014  ',
                   'b: only in 1335  ',
                   'c: only in 1015  ',
                   'd: only in 1012  ',
                   'e: only in 1007  ',
                   '*: shared by bcde',
                   paste(ifelse(plusx,'+:','  '),'shared by b/c '),
                   paste(ifelse(plusx,'x:','  '),'shared by d/e ')
  )

  legend.text <- c('a: only in 1014 ',
                   'b: only in 1335 ',
                   'c: only in 1015 ',
                   'd: only in 1012 ',
                   'e: only in 1007 ',
                   '*: in bcde        ',
                   paste(ifelse(plusx,'+:','  '),'in bc          '),
                   paste(ifelse(plusx,'x:','  '),'in de          '),
                   'Discordant SNPs '
  )
  legend.text <- paste(legend.text,format(c(leg.counts,discord),width=4),sep=' - ')
  legend.text <- paste(legend.text,' ') # add a little more right margin in box
  opar <- par(family='mono',cex=cex.legend)
  legend('topright', legend=legend.text, cex=cex.legend, inset=c(0.05,0), bg=col.legbox, box.col=col.legbox)
  par(opar)
  if(plusx){
    points(tree.labels[[16]],tree.labels[[17]]+.14,pch=8,col=col.elabel)
    points(tree.labels[[16]]+200,tree.labels[[17]]+1,pch=3,col=col.elabel)
    points(tree.labels[[16]]+200,tree.labels[[17]]-1,pch=4,col=col.elabel)
  }

  ####
  #
  # EDGE LENGTHS
  #
  for(i in seq(1,length(tree.labels)-ifelse(plusx,5,2),by=3)){
    if(F){ # T for \n in edge labels; F to remove (except "by 5")
      text(tree.labels[[i]], tree.labels[[i+1]], tree.labels[[i+2]])
    } else {
      # points(tree.labels[[i]], tree.labels[[i+1]], pch=3,col='green') # for debugging
      text(tree.labels[[i]], tree.labels[[i+1]], sub('\n([^z])',' \\1', tree.labels[[i+2]]),
           pos=3, offset=.4, font=font.elabel, col=col.elabel,cex=cex.elabel)
    }
  }
}
if(FALSE){#for debug convenience
  pdf(paperfig.path, width=8,height=5,onefile=TRUE,family='Helvetica',fonts='Courier',pointsize=10)
  show.tree(newick.medium, total.snps=consistent.count[2], pltdebug=F,straight.arrow=T)
  dev.off()
}
```

```r
caption <- function(stringency,which.tables=which.snp.tables(string.val=F)){
  caption.where <- '(UNKNOWN genome subset).'
  if(which.tables[1]=='Chr1') {caption.where <- 'on Chr1.'}
  if(which.tables[1]=='full') {caption.where <- 'genome-wide.'}
  if(which.tables[1]=='trunc'){caption.where <- 'all Chrs.'}
  cap.stringency <- c(
    'loose SNP filters.',
    'medium SNP filters.',
    'strict SNP filters.',
    'unfiltered SNPs.')
```

```
   cap <- paste('Tree based on', which.tables[2], 'reads and', cap.stringency[stringency],
                ' ``Lengths\'\' are numbers of shared/private SNPs', caption.where)
   return(cap)
}
```

Trees based on all four SNP filtering criteria are shown below. Their topologies are exactly the same, although the branch lengths are different. In all four, the length of the branch labeled "*" is probably inflated by lower coverage and higher error rate in 1014, which may mask further legitimate sharing between it and the other L-isolates. The branch lengths among the other 4 are too short for their topology to be convincing without a more rigorous analysis (e.g., a bootstrap test), but detail there is irrelevant to the story.

My sense is that the "medium" version is the best for the paper, made here and shown in Fig 1. In theory, this should look exactly like Fig 3, but something is apparently different between Knitr and direct-to-pdf. (Increasing fig.width in Knitr's chunk headers from 8 (as in the pdf call below) to 9 helps somewhat, but probably still best to make the paper fig directly rather than via Knitr.)

```
###
#
# MAKE PROTOTYPE PDF FOR PAPER, *AND* SAVE DATA NEEDED TO BUILD IT
#
w.s.t. <- which.snp.tables()
if(w.s.t. == 'trunc-unfiltered'){
  rda.Description <- 'This .rda contains data to generate Fig 3; see shared.snps.rnw for details.'
  save(rda.Description, w.s.t., pat.summaries, consistent.count, file='Fig3-data.rda')
  paperfig.path <- paste('figs-mine/paperfig-medium-tree-', w.s.t., '--Fig3proto.pdf', sep='')
} else {
  paperfig.path <- paste('figs-mine/paperfig-medium-tree-', w.s.t., '.pdf', sep='')
}
pdf(paperfig.path, width=8,height=5,onefile=TRUE,family='Helvetica',fonts='Courier',pointsize=10)
newick.medium <- newickize(make.tree(pat.summaries[,'count2']))
show.tree(newick.medium, total.snps=consistent.count[2], pltdebug=F,straight.arrow=T)
dev.off()

# pdf
#   2
```

```
# fig.paths for knitr chunks below;  .h for "hand-made" trees; plain for automatic chr1/full versions
myfigpath   <- paste(getwd(), '/figs-knitr/newick-', which.snp.tables(), '-', sep='')
myfigpath.h <- paste(getwd(), '/figs-knitr/newick-', sep='')
```

Figure 2, i.e., criteria [[1]]:

```
newick.loose <- newickize(make.tree(pat.summaries[,'count1']))
show.tree(newick.loose, total.snps=consistent.count[1])
```

Figure 3, i.e. [[2]]:

```
# newick.medium <- newickize(tree.by.hand)
# simple.newick.medium <- newickize(tree.by.hand,alt=TRUE)
newick.medium <- newickize(make.tree(pat.summaries[,'count2']))
simple.newick.medium <- newickize(make.tree(pat.summaries[,'count2']),alt=TRUE)
show.tree(newick.medium, total.snps=consistent.count[2])
```

Figure 4, i.e. [[3]]:

```
newick.strict <- newickize(make.tree(pat.summaries[,'count3']))
show.tree(newick.strict, total.snps=consistent.count[3])
```

Figure 5, i.e. [[4]]:

```
newick.unfiltered <- newickize(make.tree(pat.summaries[,'count4']))
show.tree(newick.unfiltered, total.snps=consistent.count[4])
```

Some other versions of the trees are included in the appendix.
Counts for all tree edges in the medium tree:

Figure 1: Proposed fig. for paper: Tree based on unfiltered reads and medium SNP filters. "Lengths" are numbers of shared/private SNPs all Chrs.

Figure 2: Tree based on unfiltered reads and loose SNP filters. "Lengths" are numbers of shared/private SNPs all Chrs.



Figure 3: Tree based on unfiltered reads and medium SNP filters. "Lengths" are numbers of shared/private SNPs all Chrs.

Figure 4: Tree based on unfiltered reads and strict SNP filters. "Lengths" are numbers of shared/private SNPs all Chrs.



Figure 5: Tree based on unfiltered reads and unfiltered SNPs. "Lengths" are numbers of shared/private SNPs all Chrs.

```
#pat.summaries[c(128,110,102,6,97,19,9,2,5,33,65,17,3),]
tree.edges <- c(128,110,102,6,97,19,9,2,5,33,65,17,3)-1
non.edges <- setdiff(0:127, tree.edges)
sg.edges <- showgroup(restrict.to=tree.edges) ; sg.edges

#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 2      001    1                                  X     199    620   1157   2260
# 3      002    1                             X        41774  84335  88149 105614
# 5      004    1                        X            921   2070   2578   4608
# 6      005    2                        X       X     210    410    854   1407
# 9      010    1                   X               208    559    714   1231
# 17     020    1              X                  47772  93481  96798 113191
# 19     022    2              X              X   55300  87584  84944  58009
# 33     040    1         X                         285    611   1031   2450
# 65     100    1    X                              121    321    542   2005
# 97     140    2    X    X                         591   1150   1281   2144
# 102    145    4    X    X               X       X   3312   9777  18140  23189
# 110    155    5    X    X          X    X       X  33045  38232  26997  30602
# 128    177    7    X    X    X    X    X    X   X  82193  67223  46524  15186
# Total                                          265931 386373 369709 361896
```

Counts for the top 10 discordant patterns, i.e., SNPs whose sharing pattern does not match any of the bifurcations in the tree:

```
tenth <- sort(showgroup(restrict.to=non.edges)[-(length(non.edges)+1),'count2'],decreasing=T)[10]
sg.non.edges <- showgroup(restrict.to=non.edges, c2.thresh = tenth) ; sg.non.edges

#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 32     037    5              X    X    X    X   X   2087   1987   1386    620
# 104    147    5    X    X              X    X   X   1372   3155   5536  10001
# 109    154    4    X    X         X    X        896   1611   1430   1738
# 112    157    6    X    X              X    X   X  13257  11725   8245  12202
# 115    162    4    X    X    X              X   1552   1955   1909   1014
# 118    165    5    X    X    X         X       X   1875   3928   6825   9715
# 119    166    5    X    X    X         X    X    621   1958   3252   2688
# 120    167    6    X    X    X         X    X   X   8614  16443  28402  15091
# 126    175    6    X    X    X    X    X       X  17128  14648  10156  12697
# 127    176    6    X    X    X    X    X    X   2303   2825   2133   1032
# Other  (   105 rows   w/   c2   < 1611   )   43236  21509  31987  45919
# Total                                        92941  81744 101261 112717
```

And percent of discordant SNPs:

```
nsge <- nrow(sg.edges)
discordv <- consistent.count - sg.edges[nsge,c('count1','count2','count3','count4')] ; discordv

#        count1 count2 count3 count4
# Total  92941  81744 101261 112717

discordv.pct <- round(discordv/consistent.count*100,1) ; discordv.pct

#        count1 count2 count3 count4
# Total   25.9   17.5   21.5   23.7
```

In short, the sharing pattern observed at 81744 or 17.5% of the 468117 medium-stringency consistent SNPs positions observed across all 7 isolates are discordant with the medium tree. (The strict tree has slightly more.)

A majority of the discordant SNPs fall into one of three patterns: 6-way sharing excluding Gyre (likely a technical artifact since the low coverage in Gyre reduces our power to detect SNPs there), or 6-way sharing excluding one of the two H-isolates (likely a reflection of sexuality in the H-clade—SNP positions in a population in Hardy-Weinberg equilibrium are fairly likely to be homozygous for the reference allele in a given individual).

```
third.biggest <- sort(showgroup(pat.summaries,6)[-8,'count2'],decreasing=T)[3]
big.three <- showgroup(pat.summaries,6,c2.thresh = third.biggest); big.three

#        Pat ShrBy 1007 1012 1013 1014  1015 3367 1335 count1 count2 count3 count4
# 112    157    6    X    X              X    X    X   X  13257  11725   8245  12202
# 120    167    6    X    X    X              X    X   X   8614  16443  28402  15091
# 126    175    6    X    X    X    X         X       X  17128  14648  10156  12697
# Other  (    4 rows   w/   c2   < 11725   )    3757   4115   3399   1964
# Total                                          42756  46931  50202  41954

big.three.frac <- sum(big.three[1:3,'count2'])/discordv$count2; big.three.frac

# [1] 0.5237816
```

I.e., 52.4% of discordant SNPs fall into one of these three categories.

Out of curiousity: what is the ratio of full genome to Chr 1 branch lengths. Except for the shortest few, generally $\approx$10x, as expected given the length of Chr 1:

```r
# (vectors derived by editing Newick strings, and in that order)
print(
  c(Italy=86155, Wales=95697, IW=89598, Virg=330,   Aust=632,   VA=1296,
    Puget=2113, NY=658,       PNY=480,  four=10059, Gyre=568,   five=39517, all=69526) /
  c(Italy=8813,  Wales=9652,  IW=9365,  Virg=30,    Aust=61,    VA=19,
    Puget=207,   NY=41,       PNY=18,   four=1005,  Gyre=61,    five=3912,  all= 7054),
  digits=3)


# Italy Wales   IW  Virg  Aust    VA Puget    NY  PNY  four  Gyre  five   all
#  9.78  9.91  9.57 11.00 10.36 68.21 10.21 16.05 26.67 10.01  9.31 10.10  9.86


round(genome.length.constants()$genome.length.trunc / genome.length.constants()$chr1.length, digits=4)

# [1] 10.2879
```

# 9   Semi-Automated Tree-Building

Slightly formalizing the process above: Look for the bifurcation of the 7 strains that maximizes the number of shared SNPs *within* each side of the partition while minimizing the number and fraction of SNPs that are shared by subsets that include at least one strain on each side of the partition. The 2/5 split is the winner, with 6418 SNPs in conflct with that partition (16% of the 39842 SNPs not shared by all 7; Chr1 data). The runner-up places the Gyre in a group by itself (7079 = 18% in conflict).

```r
treepart <- function(p.summ=pat.summaries, root=127, verbose=T, stringency='count2'){
  root.shared <- p.summ[root+1,stringency]
  df<-NULL
  for(i in 1:floor(root/2)){
    if(bitwAnd(i,root)==i && i < root-i){
      l1 <- showgroup(p.summ,subset=i,split=NULL,proper.subset=F,total=T)
      l  <- l1[nrow(l1),stringency]
      r1 <- showgroup(p.summ,subset=root-i,split=NULL,proper.subset=F,total=T)
      r  <- r1[nrow(r1),stringency]
      c1 <- showgroup(p.summ,subset=root,split=i,proper.subset=T,total=T)
      c  <- c1[nrow(c1),stringency]
      df <- rbind(df, data.frame(pat=i,left=l,right=r,both=l+r,cross=c,all=l+r+c,ratio=c/(l+r+c),
                                 best='',stringsAsFactors=F))
    }
  }
  df$pat<-as.octmode(df$pat)
  maxl <- which.max(df$left)
  maxr <- which.max(df$right)
  maxb <- which.max(df$both)
  minc <- which.min(df$cross)
  minr <- which.min(df$ratio)
  df$best[c(maxl,maxr,maxb,minc,minr)] <- '<'
  df$best[maxl] <- paste(df$best[maxl], 'L') # max Left
  df$best[maxr] <- paste(df$best[maxr], 'R') # max Right
  df$best[maxb] <- paste(df$best[maxb], 'B') # max Both (L+R)
  df$best[minc] <- paste(df$best[minc], 'C') # min Cross
  df$best[minr] <- paste(df$best[minr], 'O') # min ratiO (Cross/(Left+Right+Cross)
  if(verbose){
    same <- all(maxl==c(maxr,maxb,minc,minr))
    cat('root:',         format(as.octmode(root),width=3),
        '; shared:',    root.shared,
        '.  max l',     format(as.octmode(df$pat[maxl]),width=3),
        ', max r',      format(as.octmode(df$pat[maxr]),width=3),
        ', max both',   format(as.octmode(df$pat[maxb]),width=3),
        ', min cross',  format(as.octmode(df$pat[minc]),width=3),
        ', min ratio',  format(as.octmode(df$pat[minr]),width=3),
        '. \nAll the same?:',same,
        '\n')
```

```
    cat('\n')
  }
  return(df)
}
```

```
treepart()
```

```
# root: 177 ; shared: 67223 .  max l 077 , max r 010 , max both 022 , min cross 022 , min ratio 022 .
# All the same?: FALSE
#    pat    left   right   both  cross     all     ratio    best
# 1   01     622 287930 288552 112344 400896 0.2802323
# 2   02   84337 177213 261550 139346 400896 0.3475864
# 3   03   85238 104349 189587 211309 400896 0.5270918
# 4   04    2072 277815 279887 121009 400896 0.3018464
# 5   05    3102 272625 275727 125169 400896 0.3122231
# 6   06   86526  99183 185709 215187 400896 0.5367651
# 7   07   87983  97026 185009 215887 400896 0.5385112
# 8   10     561 318586 319147  81749 400896 0.2039157    < R
# 9   11    1716 279366 281082 119814 400896 0.2988655
# 10  12   85054 116885 201939 198957 400896 0.4962808
# 11  13   86743 100618 187361 213535 400896 0.5326444
# 12  14    2767 273729 276496 124400 400896 0.3103049
# 13  15    5382 271127 276509 124387 400896 0.3102725
# 14  16   87441  97314 184755 216141 400896 0.5391448
# 15  17   91097  96139 187236 213660 400896 0.5329562
# 16  20   93483 163824 257307 143589 400896 0.3581702
# 17  21   94257  94798 189055 211841 400896 0.5284188
# 18  22  265402  60429 325831  75065 400896 0.1872431 < B C O
# 19  23  266938   8065 275003 125893 400896 0.3140291
# 20  24   95731  90188 185919 214977 400896 0.5362413
# 21  25   97083  87903 184986 215910 400896 0.5385686
# 22  26  268164   4255 272419 128477 400896 0.3204746
# 23  27  271026   2732 273758 127138 400896 0.3171346
# 24  30   94222 106813 201035 199861 400896 0.4985358
# 25  31   95786  91283 187069 213827 400896 0.5333727
# 26  32  266536  17109 283645 117251 400896 0.2924724
# 27  33  270177   5553 275730 125166 400896 0.3122156
# 28  34   96722  88331 185053 215843 400896 0.5384015
# 29  35  100443  87053 187496 213400 400896 0.5323076
# 30  36  269979   2931 272910 127986 400896 0.3192499
# 31  37  278873   2084 280957 119939 400896 0.2991773
# 32  40     613 281626 282239 118657 400896 0.2959795
# 33  41    1318 271381 272699 128197 400896 0.3197762
# 34  42   85065 101681 186746 214150 400896 0.5341784
# 35  43   86121  97440 183561 217335 400896 0.5421231
# 36  44    2898 271190 274088 126808 400896 0.3163115
# 37  45    4422 267334 271756 129140 400896 0.3221284
# 38  46   87593  96322 183915 216981 400896 0.5412401
# 39  47   89792  94692 184484 216412 400896 0.5398208
# 40  50    1197 272887 274084 126812 400896 0.3163214
# 41  51    2492 269354 271846 129050 400896 0.3219039
# 42  52   85818  97951 183769 217127 400896 0.5416043
# 43  53   87755  96355 184110 216786 400896 0.5407537
# 44  54    3697 267827 271524 129372 400896 0.3227071
# 45  55    7570 266129 273699 127197 400896 0.3172818
# 46  56   88673  94747 183420 217476 400896 0.5424749
# 47  57   94104  93921 188025 212871 400896 0.5309881
# 48  60   94192  92195 186387 214509 400896 0.5350739
# 49  61   95118  88106 183224 217672 400896 0.5429638
# 50  62  266517   6162 272679 128217 400896 0.3198261
# 51  63  268442   3251 271693 129203 400896 0.3222856
# 52  64   96803  87269 184072 216824 400896 0.5408485
# 53  65   98924  85505 184429 216467 400896 0.5399580
# 54  66  270198   2099 272297 128599 400896 0.3207790
# 55  67  274929    902 275831 125065 400896 0.3119637
# 56  70   94977  88805 183782 217114 400896 0.5415719
```

```
# 57  71  96765  87113 183878 217018 400896 0.5413324
# 58  72 267758   3674 271432 129464 400896 0.3229366
# 59  73 271987   2498 274485 126411 400896 0.3153212
# 60  74  97970  85715 183685 217211 400896 0.5418138
# 61  75 103426  84756 188182 212714 400896 0.5305965
# 62  76 272458    976 273434 127462 400896 0.3179428
# 63  77 285417    323 285740 115156 400896 0.2872466      < L
```

Comparing the 5/2 split to the second-place NPG/rest split (below), the former has fewer pattern instances in conflict with the split (6418 vs 7079), as well as somewhat more random distribution of the conflicting patterns (92 vs 62 rows), whereas the 1/6 split has the majority of its conflicts (3912 of 7079, or 55%) concentrated in one pattern—the 5 NE strains. Collectively, these seem to favor the 5/2 split as the correct "history."

```
showgroup(pat.summaries,split=strtoi('022'), subset=127, proper.subset=T, c2.thresh=100)
```

| # | Pat | ShrBy | 1007 | 1012 | 1013 | 1014 | 1015 | 3367 | 1335 | count1 | count2 | count3 | count4 |
|---|-----|-------|------|------|------|------|------|------|------|--------|--------|--------|--------|
| # 4 | 003 | 2 | | | | | | X | X | 2266 | 281 | 432 | 587 |
| # 7 | 006 | 2 | | | | | X | | X | 1282 | 119 | 261 | 590 |
| # 8 | 007 | 3 | | | | | X | X | X | 152 | 146 | 278 | 557 |
| # 11 | 012 | 2 | | | | X | | X | | 2060 | 158 | 59 | 93 |
| # 12 | 013 | 3 | | | | X | | X | X | 350 | 253 | 183 | 338 |
| # 16 | 017 | 4 | | | | X | X | X | X | 344 | 361 | 250 | 564 |
| # 18 | 021 | 2 | | | X | | | | X | 2445 | 154 | 260 | 402 |
| # 20 | 023 | 3 | | | X | | | X | X | 3322 | 481 | 838 | 533 |
| # 21 | 024 | 2 | | | X | | X | | | 1406 | 178 | 371 | 625 |
| # 22 | 025 | 3 | | | X | | X | | X | 197 | 168 | 333 | 522 |
| # 23 | 026 | 3 | | | X | | X | X | | 1794 | 395 | 771 | 789 |
| # 24 | 027 | 4 | | | X | | X | X | X | 441 | 602 | 1120 | 771 |
| # 25 | 030 | 2 | | | X | X | | | | 2257 | 180 | 59 | 93 |
| # 26 | 031 | 3 | | | X | X | | | X | 361 | 255 | 178 | 361 |
| # 27 | 032 | 3 | | | X | X | | X | | 2845 | 237 | 112 | 86 |
| # 28 | 033 | 4 | | | X | X | | X | X | 1174 | 1062 | 725 | 306 |
| # 29 | 034 | 3 | | | X | X | X | | | 144 | 116 | 125 | 219 |
| # 30 | 035 | 4 | | | X | X | X | | X | 495 | 529 | 373 | 503 |
| # 31 | 036 | 4 | | | X | X | X | X | | 381 | 367 | 287 | 211 |
| # 32 | 037 | 5 | | | X | X | X | X | X | 2087 | 1987 | 1386 | 620 |
| # 35 | 042 | 2 | | X | | | | | X | 1429 | 117 | 224 | 394 |
| # 39 | 046 | 3 | | X | | | X | | X | 103 | 124 | 367 | 604 |
| # 40 | 047 | 4 | | X | | | X | X | X | 93 | 178 | 485 | 708 |
| # 48 | 057 | 5 | | X | | X | X | X | X | 227 | 231 | 205 | 324 |
| # 51 | 062 | 3 | | X | X | | | X | | 2161 | 289 | 447 | 469 |
| # 52 | 063 | 4 | | X | X | | | X | X | 325 | 167 | 265 | 194 |
| # 53 | 064 | 3 | | X | X | | X | | | 114 | 148 | 390 | 601 |
| # 54 | 065 | 4 | | X | X | | X | | X | 88 | 208 | 528 | 582 |
| # 55 | 066 | 4 | | X | X | | X | X | | 263 | 432 | 944 | 851 |
| # 56 | 067 | 5 | | X | X | | X | X | X | 422 | 685 | 1836 | 1151 |
| # 62 | 075 | 5 | | X | X | X | X | | X | 200 | 185 | 233 | 328 |
| # 63 | 076 | 5 | | X | X | X | X | X | | 108 | 158 | 188 | 128 |
| # 64 | 077 | 6 | | X | X | X | X | X | X | 920 | 872 | 917 | 485 |
| # 81 | 120 | 2 | X | | X | | | | | 1009 | 117 | 116 | 309 |
| # 83 | 122 | 3 | X | | X | | | X | | 1325 | 191 | 220 | 354 |
| # 87 | 126 | 4 | X | | X | | X | X | | 124 | 214 | 297 | 425 |
| # 88 | 127 | 5 | X | | X | | X | X | X | 133 | 253 | 399 | 482 |
| # 94 | 135 | 5 | X | | X | X | X | | X | 116 | 121 | 96 | 235 |
| # 96 | 137 | 6 | X | | X | X | X | X | X | 394 | 325 | 275 | 333 |
| # 99 | 142 | 3 | X | X | | | | X | | 342 | 419 | 477 | 755 |
| # 103 | 146 | 4 | X | X | | | X | X | | 121 | 455 | 958 | 1795 |
| # 104 | 147 | 5 | X | X | | | X | X | X | 1372 | 3155 | 5536 | 10001 |
| # 111 | 156 | 5 | X | X | | X | X | X | | 492 | 681 | 566 | 735 |
| # 112 | 157 | 6 | X | X | | X | X | X | X | 13257 | 11725 | 8245 | 12202 |
| # 113 | 160 | 3 | X | X | X | | | | | 306 | 359 | 421 | 712 |
| # 115 | 162 | 4 | X | X | X | | | X | | 1552 | 1955 | 1909 | 1014 |
| # 116 | 163 | 5 | X | X | X | | | X | X | 271 | 263 | 302 | 316 |
| # 117 | 164 | 4 | X | X | X | | X | | | 165 | 603 | 1060 | 1752 |
| # 118 | 165 | 5 | X | X | X | | X | | X | 1875 | 3928 | 6825 | 9715 |
| # 119 | 166 | 5 | X | X | X | | X | X | | 621 | 1958 | 3252 | 2688 |
| # 120 | 167 | 6 | X | X | X | | X | X | X | 8614 | 16443 | 28402 | 15091 |

```
# 125    174    5    X    X    X    X    X                    567     789     656     782
# 126    175    6    X    X    X    X    X         X    17128   14648   10156   12697
# 127    176    6    X    X    X    X    X    X     2303    2825    2133    1032
# Other   (     38 rows    w/    c2    <    100    )     5050    1815    2294    5023
# Total                                              89393   75065   90025   94037
```

```
showgroup(pat.summaries,split=strtoi('010'), subset=127, proper.subset=T, c2.thresh=100)
```

```
#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 10     011    2                   X              X      358     535     384     827
# 11     012    2                   X         X          2060     158      59      93
# 12     013    3                   X         X    X      350     253     183     338
# 13     014    2                   X    X               143     136     176     417
# 14     015    3                   X    X         X      776    1050     757    1389
# 16     017    4                   X    X    X    X      344     361     250     564
# 25     030    2              X    X                    2257     180      59      93
# 26     031    3              X    X              X      361     255     178     361
# 27     032    3              X    X         X          2845     237     112      86
# 28     033    4              X    X         X    X     1174    1062     725     306
# 29     034    3              X    X    X               144     116     125     219
# 30     035    4              X    X    X         X      495     529     373     503
# 31     036    4              X    X    X    X           381     367     287     211
# 32     037    5              X    X    X    X    X     2087    1987    1386     620
# 46     055    4         X         X    X         X      480     709     750     971
# 48     057    5         X         X    X    X    X      227     231     205     324
# 62     075    5         X    X    X    X         X      200     185     233     328
# 63     076    5         X    X    X    X    X           108     158     188     128
# 64     077    6         X    X    X    X    X    X      920     872     917     485
# 78     115    4    X              X    X         X      103     141     138     604
# 94     135    5    X         X    X    X         X      116     121      96     235
# 96     137    6    X         X    X    X    X    X      394     325     275     333
# 109    154    4    X    X         X    X               896    1611    1430    1738
# 110    155    5    X    X         X    X         X    33045   38232   26997   30602
# 111    156    5    X    X         X    X    X           492     681     566     735
# 112    157    6    X    X         X    X    X    X    13257   11725    8245   12202
# 125    174    5    X    X    X    X    X                567     789     656     782
# 126    175    6    X    X    X    X    X         X    17128   14648   10156   12697
# 127    176    6    X    X    X    X    X    X     2303    2825    2133    1032
# Other   (     33 rows    w/    c2    <    100    )     1755    1270    1424    3544
# Total                                              85766   81749   59463   72767
```

Below is the full summary of shared SNPs that do *not* directly correspond to tree splits, e.g. deep coalescence, independent coincident mutations, false positives/false negatives in the shared SNP calls, loss of SNPs in hemizygous regions, etc. (Additionally, SAMTools' SNP calls exclude positions it judges to be homozygous, and I think it operates without regard to the reference sequence, so homozygous nonreference positions, while rare except in IT/Wales, often are not called SNPs by SAMTools, but are relevant for this analysis. Provided the position is called a SNP in some other isolate, the consistency filtering we've done above should recover it, but this is still worth keeping in mind when examining the data.)

First, here are SNPs that "coalesce" on the branch from the LCA of bcde, i.e., shared among some nonempty, proper subset of bcde other than bc or de. There are 8 such patterns: any of the 4 choose 3 trios plus any of the 4 pairs having exactly one of bc.

```
sg4 <- showgroup(pat.summaries, subset=strtoi('0145'), split=5, proper.subset = F)
sg4
```

```
#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 34     041    2              X                   X       31      85     230     368
# 37     044    2              X              X            72     215     895    1809
# 38     045    3              X              X    X      116     409    1369    1656
# 66     101    2    X                             X       19      33      75     314
# 69     104    2    X                        X            39     105     356    1196
# 70     105    3    X                        X    X       35     113     283     805
# 98     141    3    X    X                        X       65     109     208     519
# 101    144    3    X    X                   X          193    1079    2430    4432
# 102    145    4    X    X                   X    X     3312    9777   18140   23189
```

```
# Total                                          3882   11925   23986   34288
```

```
sg4n <- nrow(sg4)
sg4pct <- round(sg4$count2[sg4n-1]/sg4$count2[sg4n]*100,1)
sg4pct
```

```
# [1] 82
```

So, of the 11925 SNPs found only in bcde, 82% have a sharing pattern consistent with the given tree structure.

Similarly, we analyze patterns relative to the root of the L-clade (14 patterns—any nonempty proper subset of bcde together with a):

```
sg5 <- showgroup(pat.summaries,subset=strtoi('0155'), split=8, proper.subset = F)
sg5
```

```
#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 10     011    2                  X              X      358    535    384    827
# 13     014    2                  X    X              143    136    176    417
# 14     015    3                  X    X         X      776   1050    757   1389
# 41     050    2             X    X                     31     25     54    105
# 42     051    3             X    X              X       46     55     79    126
# 45     054    3             X    X    X                22     79    171    292
# 46     055    4             X    X    X         X      480    709    750    971
# 73     110    2        X         X                     19     20     50    150
# 74     111    3        X         X              X       16      9     19    139
# 77     114    3        X         X    X                 8     38     69    365
# 78     115    4        X         X    X         X      103    141    138    604
# 105    150    3        X    X    X                     47     44     78    238
# 106    151    4        X    X    X              X       33     77     71    220
# 109    154    4        X    X    X    X                896   1611   1430   1738
# 110    155    5        X    X    X    X         X    33045  38232  26997  30602
# Total                                              36023  42761  31223  38183
```

```
sg5n <- nrow(sg5)
sg5pct <- round(sg5$count2[sg5n-1]/sg5$count2[sg5n]*100,1)
```

I.e., of the 42761 SNPs found only in abcde, 89.4% have a sharing pattern consistent with the given tree structure.

Finally, how many SNPs have patterns inconsistent with the 5-2 split, i.e., include at least one strain on each side of the 5-2 split, but not shared by all 7?

```
sg7 <- showgroup(pat.summaries, subset=127, split=strtoi('022'), proper.subset=F)
sg7
```

```
#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 4      003    2                            X    X    2266    281    432    587
# 7      006    2                       X    X         1282    119    261    590
# 8      007    3                       X    X    X     152    146    278    557
# 11     012    2             X              X         2060    158     59     93
# 12     013    3             X              X    X     350    253    183    338
# 15     016    3             X    X    X              109     62     65    152
# 16     017    4             X    X    X    X          344    361    250    564
# 18     021    2        X                   X        2445    154    260    402
# 20     023    3        X                   X    X    3322    481    838    533
# 21     024    2        X         X                  1406    178    371    625
# 22     025    3        X         X              X     197    168    333    522
# 23     026    3        X         X    X             1794    395    771    789
# 24     027    4        X         X    X    X         441    602   1120    771
# 25     030    2        X    X                       2257    180     59     93
# 26     031    3        X    X              X         361    255    178    361
# 27     032    3        X    X         X             2845    237    112     86
# 28     033    4        X    X         X    X        1174   1062    725    306
# 29     034    3        X    X    X                   144    116    125    219
# 30     035    4        X    X    X         X         495    529    373    503
# 31     036    4        X    X    X    X              381    367    287    211
# 32     037    5        X    X    X    X    X        2087   1987   1386    620
```

| # | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # 35 | 042 | 2 | | X | | | | | X | | 1429 | 117 | 224 | 394 |
| # 36 | 043 | 3 | | X | | | | | X | X | 175 | 70 | 104 | 133 |
| # 39 | 046 | 3 | | X | | | | X | X | | 103 | 124 | 367 | 604 |
| # 40 | 047 | 4 | | X | | | | X | X | X | 93 | 178 | 485 | 708 |
| # 43 | 052 | 3 | | X | | | X | | X | | 101 | 11 | 27 | 22 |
| # 44 | 053 | 4 | | X | | | X | | X | X | 46 | 38 | 36 | 56 |
| # 47 | 056 | 4 | | X | | | X | X | X | | 17 | 50 | 65 | 88 |
| # 48 | 057 | 5 | | X | | | X | X | X | X | 227 | 231 | 205 | 324 |
| # 49 | 060 | 2 | | X | X | | | | | | 1651 | 98 | 247 | 388 |
| # 50 | 061 | 3 | | X | X | | | | | X | 169 | 67 | 98 | 115 |
| # 51 | 062 | 3 | | X | X | | | | X | | 2161 | 289 | 447 | 469 |
| # 52 | 063 | 4 | | X | X | | | | X | X | 325 | 167 | 265 | 194 |
| # 53 | 064 | 3 | | X | X | | | X | | | 114 | 148 | 390 | 601 |
| # 54 | 065 | 4 | | X | X | | | X | | X | 88 | 208 | 528 | 582 |
| # 55 | 066 | 4 | | X | X | | | X | X | | 263 | 432 | 944 | 851 |
| # 56 | 067 | 5 | | X | X | | | X | X | X | 422 | 685 | 1836 | 1151 |
| # 57 | 070 | 3 | | X | X | X | | | | | 123 | 21 | 18 | 24 |
| # 58 | 071 | 4 | | X | X | X | | | | X | 30 | 17 | 14 | 28 |
| # 59 | 072 | 4 | | X | X | X | | | X | | 158 | 50 | 35 | 31 |
| # 60 | 073 | 5 | | X | X | X | | | X | X | 155 | 89 | 72 | 38 |
| # 61 | 074 | 4 | | X | X | X | X | | | | 28 | 51 | 60 | 116 |
| # 62 | 075 | 5 | | X | X | X | X | | | X | 200 | 185 | 233 | 328 |
| # 63 | 076 | 5 | | X | X | X | X | X | | | 108 | 158 | 188 | 128 |
| # 64 | 077 | 6 | | X | X | X | X | X | X | | 920 | 872 | 917 | 485 |
| # 67 | 102 | 2 | X | | | | | | X | | 887 | 98 | 126 | 351 |
| # 68 | 103 | 3 | X | | | | | | X | X | 83 | 25 | 51 | 143 |
| # 71 | 106 | 3 | X | | | | | X | X | | 42 | 63 | 127 | 377 |
| # 72 | 107 | 4 | X | | | | | X | X | X | 36 | 64 | 105 | 330 |
| # 75 | 112 | 3 | X | | | X | | | X | | 55 | 12 | 16 | 26 |
| # 76 | 113 | 4 | X | | | X | | | X | X | 18 | 8 | 8 | 66 |
| # 79 | 116 | 4 | X | | | X | | X | X | | 5 | 8 | 19 | 101 |
| # 80 | 117 | 5 | X | | | X | | X | X | X | 42 | 40 | 35 | 241 |
| # 81 | 120 | 2 | X | | X | | | | | | 1009 | 117 | 116 | 309 |
| # 82 | 121 | 3 | X | | X | | | | | X | 87 | 19 | 29 | 73 |
| # 83 | 122 | 3 | X | | X | | | | X | | 1325 | 191 | 220 | 354 |
| # 84 | 123 | 4 | X | | X | | | | X | X | 162 | 85 | 77 | 124 |
| # 85 | 124 | 3 | X | | X | | | X | | | 66 | 81 | 167 | 400 |
| # 86 | 125 | 4 | X | | X | | | X | | X | 41 | 79 | 142 | 283 |
| # 87 | 126 | 4 | X | | X | | | X | X | | 124 | 214 | 297 | 425 |
| # 88 | 127 | 5 | X | | X | | | X | X | X | 133 | 253 | 399 | 482 |
| # 89 | 130 | 3 | X | | X | X | | | | | 73 | 12 | 9 | 27 |
| # 90 | 131 | 4 | X | | X | X | | | | X | 17 | 5 | 5 | 52 |
| # 91 | 132 | 4 | X | | X | X | | | X | | 86 | 27 | 19 | 38 |
| # 92 | 133 | 5 | X | | X | X | | | X | X | 56 | 31 | 16 | 52 |
| # 93 | 134 | 4 | X | | X | X | X | | | | 15 | 24 | 25 | 143 |
| # 94 | 135 | 5 | X | | X | X | X | | | X | 116 | 121 | 96 | 235 |
| # 95 | 136 | 5 | X | | X | X | X | | X | | 41 | 71 | 63 | 106 |
| # 96 | 137 | 6 | X | | X | X | X | | X | X | 394 | 325 | 275 | 333 |
| # 99 | 142 | 3 | X | X | | | | | X | | 342 | 419 | 477 | 755 |
| # 100 | 143 | 4 | X | X | | | | | X | X | 46 | 55 | 86 | 190 |
| # 103 | 146 | 4 | X | X | | | | X | X | | 121 | 455 | 958 | 1795 |
| # 104 | 147 | 5 | X | X | | | | X | X | X | 1372 | 3155 | 5536 | 10001 |
| # 107 | 152 | 4 | X | X | | | X | | X | | 27 | 21 | 24 | 67 |
| # 108 | 153 | 5 | X | X | | | X | | X | X | 35 | 32 | 30 | 96 |
| # 111 | 156 | 5 | X | X | | | X | X | X | | 492 | 681 | 566 | 735 |
| # 112 | 157 | 6 | X | X | | | X | X | X | X | 13257 | 11725 | 8245 | 12202 |
| # 113 | 160 | 3 | X | X | X | | | | | | 306 | 359 | 421 | 712 |
| # 114 | 161 | 4 | X | X | X | | | | | X | 72 | 88 | 102 | 207 |
| # 115 | 162 | 4 | X | X | X | | | | X | | 1552 | 1955 | 1909 | 1014 |
| # 116 | 163 | 5 | X | X | X | | | | X | X | 271 | 263 | 302 | 316 |
| # 117 | 164 | 4 | X | X | X | | | X | | | 165 | 603 | 1060 | 1752 |
| # 118 | 165 | 5 | X | X | X | | | X | | X | 1875 | 3928 | 6825 | 9715 |
| # 119 | 166 | 5 | X | X | X | | | X | X | | 621 | 1958 | 3252 | 2688 |
| # 120 | 167 | 6 | X | X | X | | | X | X | X | 8614 | 16443 | 28402 | 15091 |
| # 121 | 170 | 4 | X | X | X | X | | | | | 21 | 26 | 21 | 69 |
| # 122 | 171 | 5 | X | X | X | X | | | | X | 30 | 29 | 18 | 70 |
| # 123 | 172 | 5 | X | X | X | X | | | X | | 105 | 95 | 59 | 86 |

```
# 124    173    6    X    X    X    X         X    X     140      93      74     114
# 125    174    5    X    X    X    X    X                        567     789     656     782
# 126    175    6    X    X    X    X    X         X   17128   14648   10156   12697
# 127    176    6    X    X    X    X    X    X              2303    2825    2133    1032
# 128    177    7    X    X    X    X    X    X    X   82193   67223   46524   15186
# Total                                              171586  142288  136549  109223
```

```r
sg7n <- nrow(sg7)
sg7pct <- round(sg7$count2[sg7n-1]/sg7$count2[sg7n]*100,1)
sg7pct
```

```
# [1] 47.2
```

A more compact version of that table, showing only the larger counts:

```r
thresh <- signif(.02 * sg7$count2[sg7n],1)
thresh
```

```
# [1] 3000
```

```r
showgroup(pat.summaries, subset=127, split=strtoi('022'), proper.subset=F, c2.thresh = thresh)
```

```
#        Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 104    147     5    X    X              X    X    X   1372   3155   5536  10001
# 112    157     6    X    X         X    X    X    X  13257  11725   8245  12202
# 118    165     5    X    X    X         X         X   1875   3928   6825   9715
# 120    167     6    X    X    X         X    X    X   8614  16443  28402  15091
# 126    175     6    X    X    X    X    X         X  17128  14648  10156  12697
# 128    177     7    X    X    X    X    X    X    X  82193  67223  46524  15186
# Other   (      87 rows   w/   c2    < 3000    )      47147  25166  30861  34331
# Total                                              171586 142288 136549 109223
```

So, of the 142288 SNPs found both in the L- and H-clades, 47.2% have a sharing pattern consistent with the given tree structure, i.e., are found in all 7 isolates. Among the others, three patterns dominate—(i) the 6-way pattern excluding the Gyre is the largest, plausibly explained by 7-way sharing from which the Gyre drops out due to low coverage/high error rate, (ii) the 6-way excluding Italy, and (iii) ditto for Wales. Origin of the later two cases is unclear, but may partly reflect Hardy-Weinberg—some positions that are *population-level* SNPs in those isolates will be homozygous-reference in the CCMP founder cell for IT or Wales. If I take the 7-way shared SNP count (69526) as a surrogate approximating the number of population-level SNPs in either IT or Wales that are shared with the L-clade, then I might expect, based on HWE, roughly half that number to to be lost (become homozygous) in IT, and a similar number in Wales. However, the observed counts of these positions are lower by ≈20K than I might have guessed from HWE, perhaps suggesting that IT and Wales are distinct populations, each with a pool of many thousand private polymorphisms.

In aggregate:

```r
untreelike <-
  sg7$count2[sg7n]-sg7$count2[sg7n-1] +
  sg5$count2[sg5n]-sg5$count2[sg5n-1] +
  sg4$count2[sg4n]-sg4$count2[sg4n-1]
untreelike
```

```
# [1] 81742
```

```r
consistent.count[2]
```

```
# [1] 468117
```

```r
unpct <- round(untreelike/consistent.count[2]*100,1)
unpct
```

```
# [1] 17.5
```

I.e., 81742 or 17.5% of the 468117 consistent SNPs identified (by criterion 2) across all 7 isolates are discordant with the assumed tree.

Overall, based on this data, I take the following to be obvious: (a) separation of the the H-isolates from the L-isolates (and from each other??), and (b) near-identity of the L-isolates. Due to the small counts, the exact topology among the L-isolates (esp. bcde) is uncertain, but *any* topology there is consistent with the asexual/clonal/global-expansion hypothesis, so there is little point in examining this subtree more carefully. Again, we believe the (apparent) slight separation of the Gyre from the other L-isolates is largely driven by technical artifacts (lower coverage/higher error rates) in the sequencing rather than by biological effects. However, the discord between Gyre SNPs and others is the major substantive ambiguity in the offered tree. Nevertheless, in the next section we show by a bootstrap analysis that the offered placement of Gyre with respect to the other 4 L-isolates is strongly supported by the data.

## 9.1   Bootstrap

How robust is the inferred tree? Italy/Wales seem clearly related to each other but separate from the other 5. Likewise, the 4 coastal L-isolates seem to be closely related, with little data to separate them (and perhaps little sense in trying). So, the key question here is whether the top level bifurcation is 2/5 or NPG/6. Here, we do a simple bootstrap test (on c2 numbers only) to see whether the 2/5 split is consistently the most parsimonious.

```
n2 <- sum(pattern.counts[[2]][,2]); n2

# [1] 468117
```

Conceptually, we sample, with replacement, n2=468117 SNP positions from among the 468117 positions declared consisent SNPs according to criterion c2, and recalculate the statistics examined above to see whether the 2/5 split again minimizes conflicting sharing patterns. This resampling/calculation is repeated nboot times (set near front of file). Since all that matters is the sharing pattern, this procedure is expedited by actually sampling 468117 independent integers in the range 0:127 with probabilities proportional to the patttern counts given in column 2 of pattern.counts[[2]]. The sample is then tabulated in a 128 row table analogous to pattern.summaries, for analysis by showgroups/treepart, as above.

```
boot.sample <- sample(0:127,n2,replace=T,prob=pattern.counts[[2]][,2])
str(boot.sample)

#  int [1:468117] 127 111 114 2 2 100 18 16 127 127 ...

boot.count <- mytable(boot.sample,c(0,127))
boot.count[c(1:4,125:128),] # show a few rows

#       val count
# [1,]    0     0
# [2,]    1   656
# [3,]    2 84134
# [4,]    3   277
# [5,] 124   762
# [6,] 125 14747
# [7,] 126  2758
# [8,] 127 67423

boot.counts <- list(NULL,boot.count,NULL) # dummy list with just c2 summaries
cor(pattern.counts[[2]][,2],boot.counts[[2]][,2]) # just curious - how correlated are they?

# [1] 0.9999969

boot.summaries <- pat.summary(boot.counts)
showgroup(boot.summaries,c2.thresh=400) #show a few rows

#      Pat ShrBy 1007 1012 1013 1014 1015 3367 1335 count1 count2 count3 count4
# 2    001     1                              X        NA    656    NA    NA
# 3    002     1                         X             NA  84134    NA    NA
# 5    004     1                X                       NA   2148    NA    NA
# 6    005     2                X         X             NA    422    NA    NA
# 9    010     1           X                            NA    539    NA    NA
# 10   011     2           X              X             NA    515    NA    NA
# 14   015     3           X    X         X             NA   1097    NA    NA
```

```
# 16     017    4                    X    X    X    X       NA    403    NA    NA
# 17     020    1               X                           NA  93546    NA    NA
# 19     022    2               X              X            NA  87716    NA    NA
# 20     023    3               X              X    X       NA    481    NA    NA
# 24     027    4               X         X    X    X       NA    627    NA    NA
# 28     033    4               X    X         X    X       NA   1037    NA    NA
# 30     035    4               X    X    X         X       NA    503    NA    NA
# 32     037    5               X    X    X    X    X       NA   1974    NA    NA
# 33     040    1          X                                NA    594    NA    NA
# 38     045    3               X              X         X  NA    406    NA    NA
# 46     055    4               X         X    X         X  NA    703    NA    NA
# 55     066    4               X    X         X    X       NA    425    NA    NA
# 56     067    5               X    X         X    X    X  NA    674    NA    NA
# 64     077    6               X    X    X    X    X    X  NA    886    NA    NA
# 97     140    2     X    X                                NA   1130    NA    NA
# 99     142    3     X    X                   X            NA    424    NA    NA
# 101    144    3     X    X              X                 NA   1037    NA    NA
# 102    145    4     X    X              X         X       NA   9812    NA    NA
# 103    146    4     X    X              X    X            NA    456    NA    NA
# 104    147    5     X    X              X    X    X       NA   3197    NA    NA
# 109    154    4     X    X    X         X                 NA   1658    NA    NA
# 110    155    5     X    X    X         X         X       NA  38094    NA    NA
# 111    156    5     X    X    X         X    X            NA    684    NA    NA
# 112    157    6     X    X    X         X    X    X       NA  11714    NA    NA
# 115    162    4     X    X    X              X            NA   1973    NA    NA
# 117    164    4     X    X    X         X                 NA    586    NA    NA
# 118    165    5     X    X    X         X         X       NA   3981    NA    NA
# 119    166    5     X    X    X         X    X            NA   1939    NA    NA
# 120    167    6     X    X    X         X    X    X       NA  16361    NA    NA
# 125    174    5     X    X    X    X    X                 NA    762    NA    NA
# 126    175    6     X    X    X    X    X         X       NA  14747    NA    NA
# 127    176    6     X    X    X    X    X    X            NA   2758    NA    NA
# 128    177    7     X    X    X    X    X    X    X       NA  67423    NA    NA
# Other   (     88 rows   w/    c2    <    400     )        NA   9895    NA    NA
# Total                                                     NA 468117    NA    NA
```

Tree partition analysis (and how to pluck out only the best rows based on 3 smallest cross counts and "best" criteria):

```
tp <- treepart(boot.summaries,root=127) ; tp

# root: 177 ; shared: 67423 .  max l 077 , max r 010 , max both 022 , min cross 022 , min ratio 022 .
# All the same?: FALSE
#     pat    left   right    both   cross    all     ratio     best
# 1    01     656 287705  288361  112333  400694  0.2803461
# 2    02   84134 177289  261423  139271  400694  0.3475745
# 3    03   85067 104362  189429  211265  400694  0.5272477
# 4    04    2148 277617  279765  120929  400694  0.3017989
# 5    05    3226 272454  275680  125014  400694  0.3119937
# 6    06   86410  99164  185574  215120  400694  0.5368685
# 7    07   87904  96999  184903  215791  400694  0.5385431
# 8    10     539 318600  319139   81555  400694  0.2035344    < R
# 9    11    1710 279265  280975  119719  400694  0.2987791
# 10   12   84808 117015  201823  198871  400694  0.4963164
# 11   13   86509 100637  187146  213548  400694  0.5329453
# 12   14    2827 273636  276463  124231  400694  0.3100396
# 13   15    5517 271004  276521  124173  400694  0.3098948
# 14   16   87294  97317  184611  216083  400694  0.5392719
# 15   17   91056  96122  187178  213516  400694  0.5328655
# 16   20   93546 163624  257170  143524  400694  0.3581885
# 17   21   94342  94582  188924  211770  400694  0.5285080
# 18   22  265396  60374  325770   74924  400694  0.1869856 < B C O
# 19   23  266950   8063  275013  125681  400694  0.3136583
# 20   24   95866  89864  185730  214964  400694  0.5364792
# 21   25   97225  87575  184800  215894  400694  0.5388002
# 22   26  268241   4162  272403  128291  400694  0.3201720
# 23   27  271124   2635  273759  126935  400694  0.3167879
# 24   30   94278 106731  201009  199685  400694  0.4983479
# 25   31   95861  91068  186929  213765  400694  0.5334869
# 26   32  266498  17167  283665  117029  400694  0.2920658
# 27   33  270129   5536  275665  125029  400694  0.3120311
# 28   34   96847  88094  184941  215753  400694  0.5384483
```

```
# 29   35 100593   86777 187370 213324 400694 0.5323863
# 30   36 270023    2888 272911 127783 400694 0.3189042
# 31   37 278960    2012 280972 119722 400694 0.2987866
# 32   40     594 281586 282180 118514 400694 0.2957718
# 33   41    1330 271336 272666 128028 400694 0.3195156
# 34   42   84841 101737 186578 214116 400694 0.5343629
# 35   43   85927   97505 183432 217262 400694 0.5422143
# 36   44    2979 271061 274040 126654 400694 0.3160866
# 37   45    4543 267221 271764 128930 400694 0.3217667
# 38   46   87474   96343 183817 216877 400694 0.5412534
# 39   47   89710   94700 184410 216284 400694 0.5397735
# 40   50    1156 272926 274082 126612 400694 0.3159818
# 41   51    2463 269362 271825 128869 400694 0.3216145
# 42   52   85545   98025 183570 217124 400694 0.5418699
# 43   53   87487   96444 183931 216763 400694 0.5409689
# 44   54    3750 267776 271526 129168 400694 0.3223607
# 45   55    7685 266057 273742 126952 400694 0.3168303
# 46   56   88509   94786 183295 217399 400694 0.5425562
# 47   57   94032   93945 187977 212717 400694 0.5308714
# 48   60   94228   92098 186326 214368 400694 0.5349918
# 49   61   95167   87908 183075 217619 400694 0.5431052
# 50   62 266451    6240 272691 128003 400694 0.3194532
# 51   63 268395    3266 271661 129033 400694 0.3220238
# 52   64   96941   86988 183929 216765 400694 0.5409739
# 53   65   99093   85214 184307 216387 400694 0.5400305
# 54   66 270234    2044 272278 128416 400694 0.3204840
# 55   67 275014     839 275853 124841 400694 0.3115619
# 56   70   95013   88695 183708 216986 400694 0.5415255
# 57   71   96812   86957 183769 216925 400694 0.5413732
# 58   72 267656    3757 271413 129281 400694 0.3226427
# 59   73 271878    2538 274416 126278 400694 0.3151482
# 60   74   98097   85506 183603 217091 400694 0.5417875
# 61   75 103598   84516 188114 212580 400694 0.5305295
# 62   76 272437     970 273407 127287 400694 0.3176663
# 63   77 285475     288 285763 114931 400694 0.2868299      < L
```

```r
otp <- order(tp[,'cross'])[1:3]     # 3 smallest 'cross' counts
btp <- which(tp[,'best'] != '')     # 'best' by Left/Right/Both/Cross/ratio
toptp <- unique(c(otp,btp,18,8))    # above, plus 5/2, 6/1 splits
print(tp[toptp,])                   # show the winners

#     pat    left   right    both   cross     all     ratio    best
# 18   22  265396   60374 325770   74924 400694 0.1869856 < B C O
# 8    10     539 318600 319139   81555 400694 0.2035344     < R
# 1    01     656 287705 288361 112333 400694 0.2803461
# 63   77  285475     288 285763 114931 400694 0.2868299     < L
```

Now repeat the above nboot times, and summarize results:

```r
nboot <- params$nboot #  default from params set in section 2
nboot <- ((nboot+2) %/% 4) * 4 + 1   # summary is cleaner if n mod 4 == 1, so int median/quartiles
cat('***\n*** Doing', nboot, 'bootstrap replicates.\n***\n')

# ***
# *** Doing 101 bootstrap replicates.
# ***

bcor <- numeric(nboot)
b52cross <- integer(nboot)
b61cross <- integer(nboot)
brev <- logical(nboot)
for(i in 1:nboot){
  boot.sample <- sample(0:127,n2,replace=T,prob=pattern.counts[[2]][,2])
  boot.count <- mytable(boot.sample,c(0,127))
  boot.counts <- list(NULL,boot.count,NULL) # dummy list with just c2 summaries
  boot.summaries <- pat.summary(boot.counts)
  tp <- treepart(boot.summaries,root=127, verbose=F)
  bcor[i] <- cor(pattern.counts[[2]][,2],boot.counts[[2]][,2]) # just curious - how correlated are they?
  b52cross[i] <- tp[18,'cross']
  b61cross[i] <- tp[ 8,'cross']
  brev[i] <- (b52cross[i] > b61cross[i])
```

```
  if(brev[i]){
    # show the unexpected ones; probably breaks w/ cache
    otp <- order(tp[,'cross'])[1:3]
    btp <- which(tp[,'best'] != '')
    toptp <- unique(c(otp,btp,18,8))
    print(tp[toptp,])
  }
}
# summarize:
corsummary <- t(as.matrix(c(summary(bcor),sd=sd(bcor))))
row.names(corsummary) <- 'bcor'
bdelta <- b61cross-b52cross
brevp <- 100*brev   # make it percent reversed instead of logical
thesummary <- rbind(summary(b52cross),summary(b61cross),summary(c(bdelta)),summary(brevp))
row.names(thesummary) <- c('b52cross', 'b61cross', 'b61-b52', '% rev')
thesummary <- cbind(thesummary, sd=c(sd(b52cross),sd(b61cross),sd(bdelta),sd(brevp)))
```

SUMMARY: In 101 bootstrap replicates, we saw 0 samples with the 6/1 split having fewer conflicts than the 5/2 split, and the minimum separation between them was $\approx 20$ sigma, hence highly statistically significant.

```
# 'opt' hacking is trying to force knitr to show more digits of bcor in summary, as Rstudio does, but
# it still fails...  Bottom line is the correlation seems to be  > .999 in all samples, rounds to 1.0,
# as seen in this run of 1001 samples cut/paste from Rstudio:
#          Min.        1st Qu.      Median      Mean        3rd Qu. Max.     sd
# bcor      "   0.9998" "   0.9999" "   0.9999" "   0.9999" "    1" "    1" "  0.00003462"
# > max(bcor)
# [1] 0.9999915
o.opts <- options(digits=7,width=127)
format(rbind(corsummary,thesummary),scientific=F,digits=4,drop0trailing=T)

#            Min.            1st Qu.           Median            Mean              3rd Qu.           Max.
# bcor      "   0.999979714" "   0.999991809" "   0.999994085" "   0.999993422" "   0.999996278" "   0.99999
# b52cross "74410"           "74883"          "75078"          "75043.148514851" "75192"          "75542"
# b61cross "81022"           "81564"          "81761"          "81748.564356436" "81937"          "82333"
# b61-b52  " 5857"           " 6524"          " 6694"          " 6705.415841584" " 6859"          " 7499"
# % rev    "    0"           "    0"          "    0"          "    0"           "    0"          "    0"
#            sd
# bcor      "   0.000003645"
# b52cross " 231.908446855"
# b61cross " 255.187594363"
# b61-b52  " 297.27442767"
# % rev    "    0"

options(o.opts)
```

Based on this, it is reasonable to claim that we are confident that the tree topology is as shown in the earlier figures, with the exception of the exact order of the splits with the 4 NE coastal isolates.

# 10   Notes

This section is a random brain dump of limitations of the current analysis, ideas for improvements, etc. In the main, these may not be worth doing, unless we see significant holes or get pushed by reviewers, etc, but I wanted to catalog before we forget them.

**Noise:** Various sources of "noise" in the data:

1. Read errors, low read depth — perhaps fixed by medium/strict thresholding

2. Deep coalescence

3. Skew because 1335 is the reference. (Julie notes we could partially fix this by remapping based on discovered SNPs, tho that wouldn't fix gross misassembly in 1335, e.g. collapsed or misordered tandem duplicates, or segments missing in 1335 that are present in one or more other strains, etc.; much harder to fix those, let's just hope they are rare...)

4. Varying error rates and sequencing depth among the 7. E.g., plausibly the 1000 SNPs shared by 4 but not by Gyre are a result of lower read depth (we missed a SNP that is actually present) and/or higher error rates (causing a position to appear inconsistent in gyre) in the gyre data. I can't think of a way to correct for this effect. It might be possible, perhaps by simulation, to estimate the size of the effect and see whether it could explain ≈1000 SNPs.

5. Varying numbers of founder cells in the sequencing cultures. (Again, I made some attempts at modeling this, but nothing very satisfactory yet.)

6. Tri-allelic positions where stochastic fluctuation in sequence sampling promotes the rare allele to prominence. (Julie replies: "isn't this the same as more than one founder cell? If they are diploid there should only ever be two alleles, unless there were random and very rare, thus unlikely, trisomy events?" I agree, but it is a concrete example of an effect of multiple founders that might be important. Not sure this is the most important such effect...)

7. Gaps/indels - alignments are likely to be of lower quality in the vicinity of an indel, so, maybe lower coverage/more SNPs. We ignored them. Does this add any systematic bias? e.g. if one strain had more indels than another, would this confound other analyses? unclear. Julie suggested a paper titled "Barking up the wrong tree-length: yada yada yada gap penalties"; maybe relevant?

**Other Items/Potential To Dos:**

1. any spacial structure to various sub-classes?

2. after top level split, should I reanalyze halves of partition in isolation? said another way, I think the tree-building is sensible, but not sure it's optimal.

3. if we believe no sex, then I think gain of SNP should be more common than loss of SNP, since the later can only happen by (a) mutation reverting to reference, (b) second mutation matching nonreference, (c) homologous repair (look for blocks of LOH), or (d) false negative e.g. from low read depth. Does tree-building appropriately weight the gain vs loss cases? (Does it even care?)

4. should we weight coding and/or nonsynonomous SNPs more heavily? Julie says "you do not want to weight the coding or nonsynonomous/coding SNPs because for time you want the more clock-like neutral mutations." I.e., I got this backwards. Maybe should redo tree based on noncoding SNPs only.

5. We could also do an actual parsimony analysis based on 2-state model (homozygous-ref vs not), but I'm not quite sure how to handle this in a mixed sex/nosex case.

6. Might be interesting to look at sharing just within (shared?) deserts. Given tree model above and expectation that bottleneck followed split of H- from L-clades, I would expect little or no sharing of L-clade desert SNPs with H-clade; sharing between It/Wales might suggest "desert" is actually a region under strong purifying selection (e.g. a gene); sharing/non-sharing within L-clade deserts might suggest more about evo history of the 5.

# 11 Appendix: Old Trees, etc.

Tangents, old stuff of historical interest at best, etc..

## 11.1 HWE Sharing

Tangent: As a function of nonref allele freq, assuming HWE, what is probability that nonref allele will be seen in $k$ strains, $0 \le k \le 4$ (Fig 6).

```
myfigpath.h <- paste(getwd(), '/figs-knitr/', sep='')
```
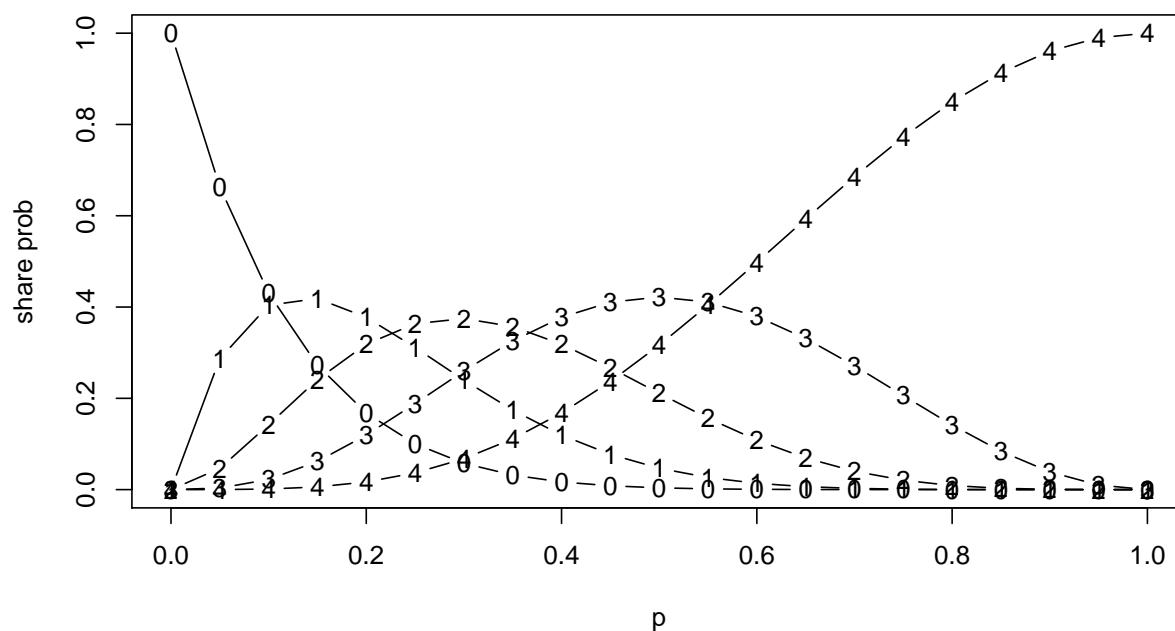
Figure 6: Sharing Probability

```
p <- (0:20)/20
q <- 1-p
r <- 2*p*q+p^2
plot(  p, 1*q^0*r^4, type='b',pch='4', ylab="share prob")
points(p, 4*q^2*r^3, type='b',pch='3')
points(p, 6*q^4*r^2, type='b',pch='2')
points(p, 4*q^6*r^1, type='b',pch='1')
points(p, 1*q^8*r^0, type='b',pch='0')
```

## 11.2   Old Tree Stuff

All based on un-q-filtered reads.

The first pass at the tree analysis was the Chr1 tree, *loose criteria* (c1); it is rendered via
`http://iubio.bio.indiana.edu/treeapp/treeprint-form.html` as Fig 7, and in newick format is:

```
newick.chr1.loose <- '(((tp3367_Italy:4551,tp1013_Wales:4954):5920,(((tp1007_Virginia:10,tp1012_Australia:29):9,(
cat.hardwrap(newick.chr1.loose)

# (((tp3367_Italy:4551,tp1013_Wales:4954):5920,(((tp1007_Virginia:10,tp1012_Austra
# lia:29):9,(tp1015_Puget_Sound:90,tp1335_NY:13):11):320,tp1014_Gyre:22):3484):859
# 3,outgroup:0);
```

Chr 1 tree based on *medium criteria* (c2) has exactly the same topology is, although the branch lengths are different. As noted earlier, the length of the branch labeled "*" is probably inflated by lower coverage and higher error rate in 1014, which may mask further legitimate sharing between it and the other L-isloates. The branch lengths among the other 4 are too short for its topology to be convincing without a more rigorous analysis (e.g., a bootstrap test).

Chr1 tree, medium criteria, in newick format:

Figure 7: Inferred Tree, based on Chr1, un-q-filtered reads, loose criteria. (Note: to visually resolve the edges among the 5, their lengths were scaled by 5x – 10x in this figure, but not in the newick description shown in the text.)

```
newick.chr1.med <- '(((tp3367_Italy:8813,tp1013_Wales:9652):9365,(((e_tp1007_Virginia:30,d_tp1012_Australia:61):1
cat.hardwrap(newick.chr1.med)

# (((tp3367_Italy:8813,tp1013_Wales:9652):9365,(((e_tp1007_Virginia:30,d_tp1012_Au
# stralia:61):19,(c_tp1015_Puget_Sound:207,b_tp1335_NY:41):18):1005,a_tp1014_Gyre:
# 61):3912):7054,outgroup:0);
```

NOTE: In early code, tree was not being recalculated; it was defined by constants in the following code chunk, hand-copied from the analysis above.

```
# tree parameters as nested lists
#   Internal nodes have subtrees sub1/2 and length
#   Root has sub1/2, but no length
#   Leaves have where, length, optionally, id, alt, nb.  (Omit id for 'outgroup'. Use 'alt' for less formal
#     labeling in cartoon version; it defaults to 'where'. Use 'nb' to add abcde annotations for legend.)
# This hand-made version is now subsumed by make.tree; retained for comparison
tree.by.hand <-
  list(
    sub1 = list(
      sub1 = list(
        sub1 = list(id=3367, length=8813, where='Venice, Italy', alt='Venice'),
        sub2 = list(id=1013, length=9652, where='Wales, UK'),
        length=9365),
      sub2 = list(
        sub1 = list(
          sub1 = list(id=1007, length=30, nb='e', where='Virginia, USA'),
          sub2 = list(id=1012, length=61, nb='d', where='Perth, W. Australia', alt='Perth'),
          length=19),
        sub2 = list(
          sub1 = list(id=1015, length=207,nb='c', where='Washington, USA', alt='Puget Sound'),
          sub2 = list(id=1335, length=41, nb='b', where='New York, USA',    alt='NY'),
          length=18),
        length=1005),
      sub2 = list(id=1014, length=61, nb='a', where='N. Pacific Gyre'),
      length=3912),
    length=7054),
  sub2 = list(length=0, where='outgroup')
  )

# historical, format example, and debug help:
oldwick.medium <- '(((CCMP3367_Italy:8813,CCMP1013_Wales:9652):9365,(((e_CCMP1007_Virginia:30,d_CCMP1012_Australia:61):19,(c_CCMP
# with simpler labeling for cartoon
simple.oldwick.medium <- '(((Italy:8813,Wales:9652):9365,(((Virginia:30,Australia:61):19,(Puget:207,NY:41):18):1005,Gyre:61):3912
cat.hardwrap(oldwick.medium)

# (((CCMP3367_Italy:8813,CCMP1013_Wales:9652):9365,(((e_CCMP1007_Virginia:30,d_CCM
# P1012_Australia:61):19,(c_CCMP1015_Puget_Sound:207,b_CCMP1335_NY:41):18):1005,a_
# CCMP1014_NPG:61):3912):7054,outgroup:0);

cat.hardwrap(simple.oldwick.medium)

# (((Italy:8813,Wales:9652):9365,(((Virginia:30,Australia:61):19,(Puget:207,NY:41)
# :18):1005,Gyre:61):3912):7054,outgroup:0);
```

Two other versions of the tree, for possible use in figs in the main paper.

Figure 8: [** as of 10/4/2015, this fig and next have stray stars on virginia, wales labels; probably due to hacking with commas in newick; not worth fixing unless we resurrect these trees for some purpose, but if so, see use of newick.name.undo in show.tree as probable fix. **]

```
tree.scale <- ifelse(which.snp.tables(string.val=F)[1]=='Chr1', 1, 10)
tree.x.lim <- 3e4 * tree.scale
the.simple.tree <- read.tree(text=simple.newick.medium)
plot(the.simple.tree, x.lim = tree.x.lim)
axis(1)
```

Figure 9:

```
plot(the.simple.tree, x.lim = tree.x.lim)
axis(1,(0:4)*7000*tree.scale,(0:4)*7000*tree.scale)
```

At some much earlier point, Tony ran the whole-genome version of the then-current code above, and manually entered tree branch lengths/legend for the resuting tree, shown in Fig 10. Code above can now automatically generate
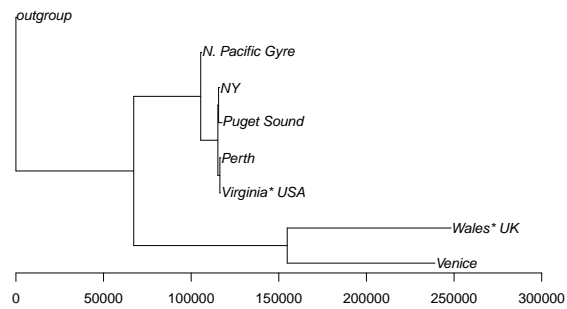
Figure 8: Tree based on unfiltered reads and medium SNP filters. "Lengths" are numbers of shared/private SNPs all Chrs. (no edge labels, nolegend)
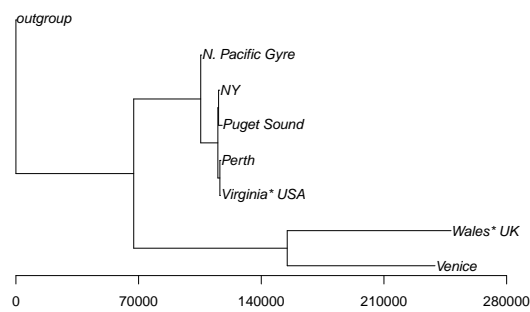


Figure 9: Tree based on unfiltered reads and medium SNP filters. "Lengths" are numbers of shared/private SNPs all Chrs. (no edge labels, no legend, short scale bar)

such a tree, but retain the following for comparison. The basic story seems clear—same topology and branch lengths scaled by about 10x, which is completely reasonable given that Chr1 is about 10% of the genome. Note that this tree is not being recalculated; it is defined by constants in the following code chunk.

```r
fullgenome.newick.medium <- '(((3367_Italy:86155,1013_Wales:95697):89598,(((e_1007_VA:330,d_1012_Australia:632):1296,(c_1015_WA:2
cat.hardwrap(fullgenome.newick.medium)

# (((3367_Italy:86155,1013_Wales:95697):89598,(((e_1007_VA:330,d_1012_Australia:63
# 2):1296,(c_1015_WA:2113,b_1335_NY:658):480):10059,a_1014_NPG:568):39517):69526,o
# utgroup:0);

legend.text <- c('a: only in 1014  ',
                 'b: only in 1335  ',
                 'c: only in 1015  ',
                 'd: only in 1012  ',
                 'e: only in 1007  ',
                 '*: shared by bcde',
                 '   shared by b/c ',
                 '   shared by d/e '
)
fullgenome.tree.x.lim <- 300000
fullgenome.counts <- c( 568, 658, 2113, 632, 330, 10059, 480, 1296 )
fullgenome.legend.text <- paste(legend.text,format(fullgenome.counts,width=5),sep=' - ')
fullgenome.tree.labels <- list( ## x,y,text
     41000,3.63,'69526\nshared by 7',
     90000,5.75,'39517\nby 5 (**)',
    115000,1.5, '89598\nshared by 2',
    210000,2.0, '95697 only\nin Wales',
    210000,1.0, '86155 only\nin Italy',
    113500,4.6, '*')
```

Figure 10:

```r
library(ape)
the.fullgenome.tree <- read.tree(text=fullgenome.newick.medium)
plot(the.fullgenome.tree, x.lim = fullgenome.tree.x.lim)
axis(1) # ; axis(2) useful only for placing labels
opar <- par(family='mono',cex=.8)
legend('topright', legend=fullgenome.legend.text)
par(opar)
for(i in seq(1,length(fullgenome.tree.labels)-2,by=3)){
  text(fullgenome.tree.labels[[i]], fullgenome.tree.labels[[i+1]], fullgenome.tree.labels[[i+2]])
}
```
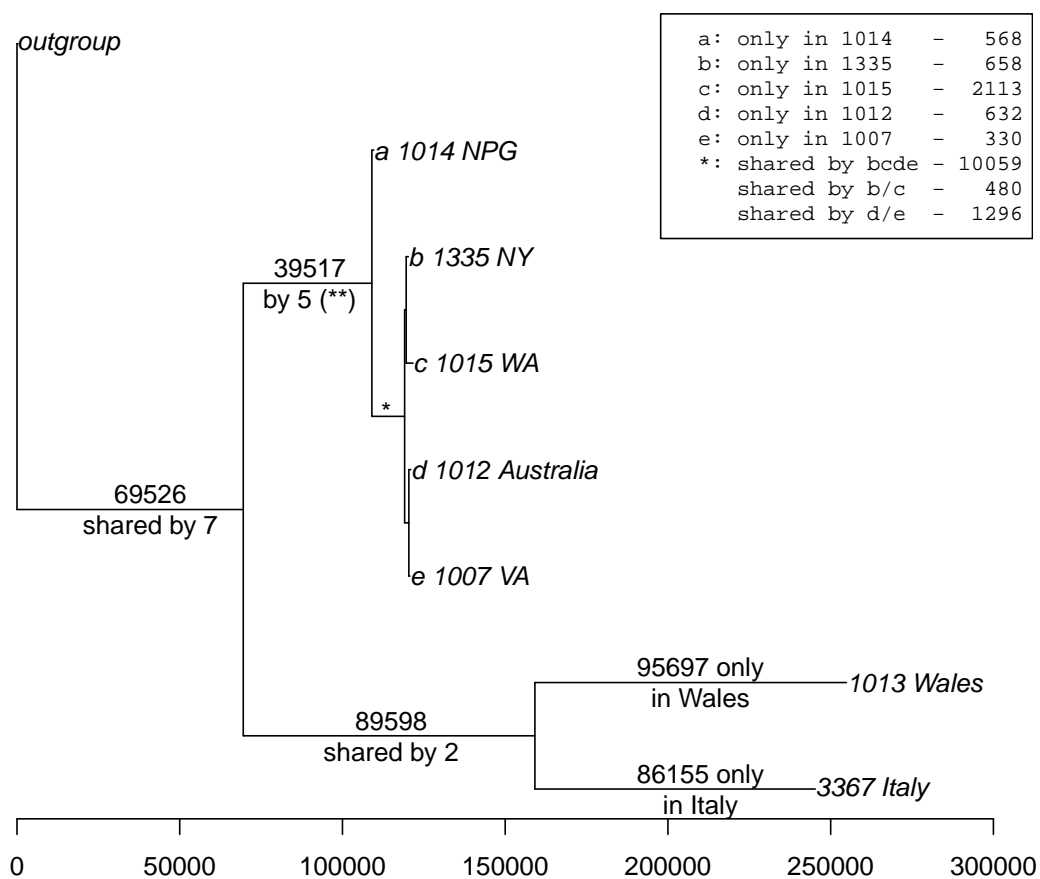
Figure 10: Tree based on unqfiltered reads and medium SNP filters. "Lengths" are numbers of shared/private SNPs genome-wide. (By-hand legacy version)