

Phylogenetics Between and Within: Seeing Transmissions and Dual Infections in HIV Deep Sequence Data

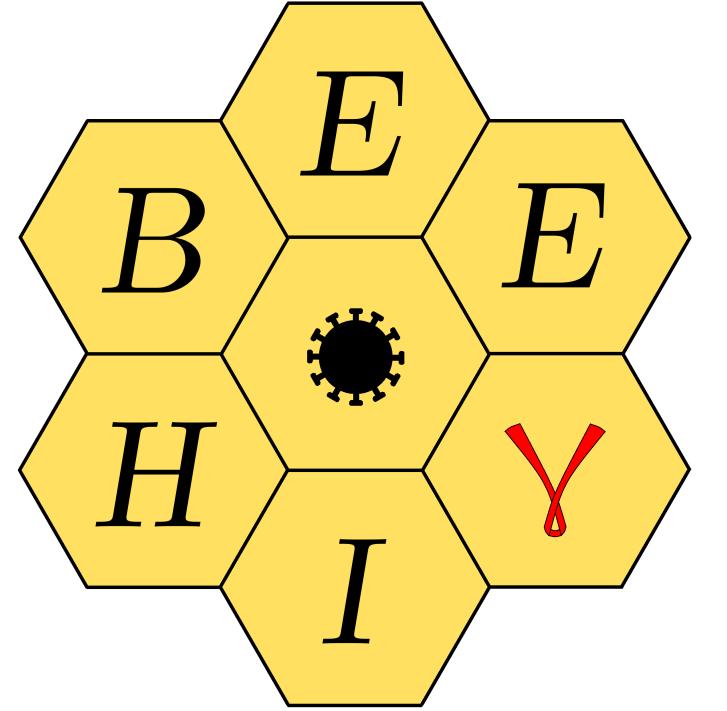
Chris Wymant^{1,2}, Matthew Hall^{1,2}, François Blanquart¹, Christophe Fraser^{1,2},
and the BEEHIVE Collaboration

¹ Department of Infectious Disease Epidemiology, Imperial College London, UK

² (from 1st July) Oxford Big Data Institute, Nuffield Department of Medicine, University of Oxford, UK

Mathematical and Computation Evolutionary Biology, June 2016: Forecasting Evolution

Notes on the talks:
ChrisWymant on Twitter



The BEEHIVE Study: Bridging the Epidemiology and Evolution of HIV in Europe



- Whole viral genome data from >4000 patients.
- Seroreconverters: 1 year window between - and + tests, or clinical evidence
- ≥ 1 viral load measurement 6-24 months after + test, pre-ART
- Wealth of clinical data (including later response to ART)

Objectives:

- 1°: the viral-molecular basis of virulence
- 2°: molecular epidemiology
- 3°: dual infections and minority variants

Imperial College London

François Blanquart
Christophe Fraser
Matthew Hall
Frank de Wolf
Chris Wymant

Amsterdam Medical Centre

Margreet Bakker
Ben Berkhout
Marion Cornelissen
Peter Reiss

Wellcome Trust Sanger Institute

Martin Hunt
Astrid Gall
Paul Kellam
Swee Hoe Ong

HIV Monitoring Foundation

Daniela Bezemer
Mariska Hillebregt
Ard van Sighem
Sima Zaheri

Karolinska University Hospital

Jan Albert

Anders Sönnernborg

Oslo University Hospital

Anne-Marte Bakken Kran
Anne Margarita Dyrhol Riise

Antwerp Institute of Tropical Medicine

Katrien Fransen

Guido Vanham

John Hopkins University

M. Kate Grabowski

Robert Koch-Institute, Berlin

Barbara Gunsenheimer-Bartmeyer

Claudia Kücherer

University Hospital Zürich

Huldrych Günthard
Roger Kouyos

Laboratory of Immunoregulation**NIAID NIH, Baltimore**

Oliver Laeyendecker

Helsinki University Hospital

Kirsi Liitsola

Matti Ristola

Université Paris Sud

Laurence Meyer

University College London

Dan Frampton

Kholoud Porter

Erasmus Medical Centre, Rotterdam

David van de Vijver

École polytechnique fédérale de Lausanne

Jacques Fellay
Istvan Bartha

Analysis Advisory Group

Sebastian Bonhoeffer

Gabriel Leventhal

Samuel Alizon

Andrew Rambaut

Oliver Pybus

Gil McVean

. . .and thanks to

Tiziano Gallo Cassarino

Nick Croucher

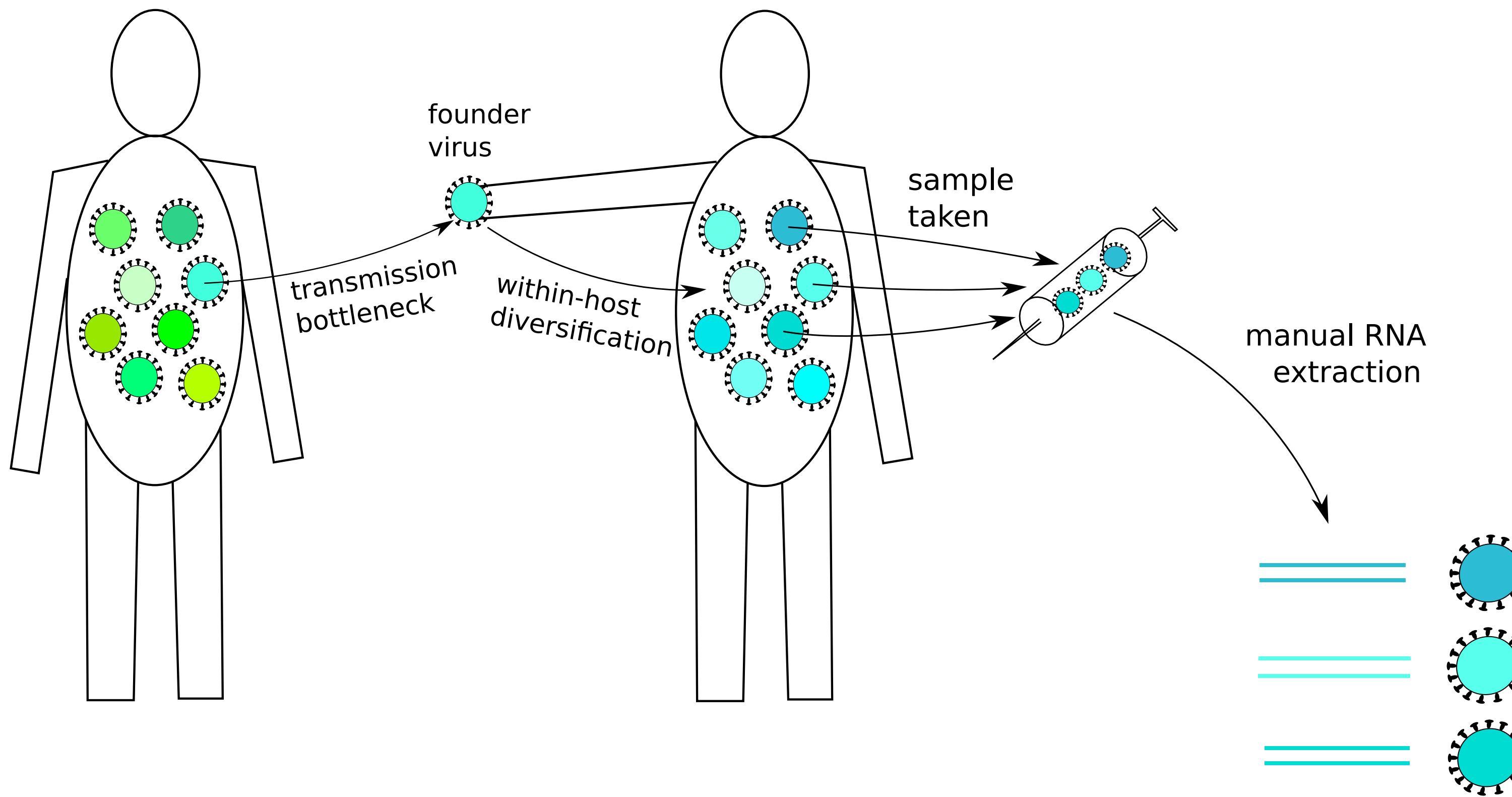
Katrina Lythgoe

Olli Ratmann



First: reconstructing the viral genotype

(Next: putting it in context.)

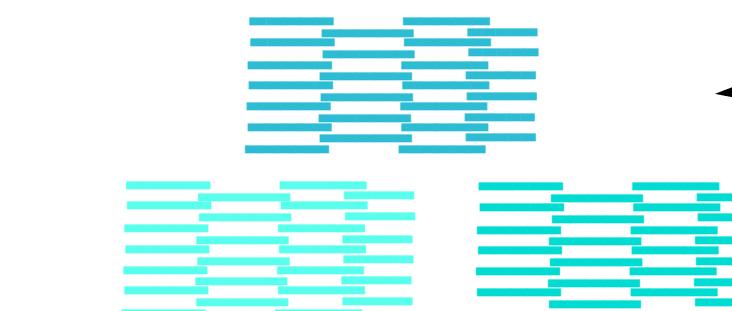


paired-end reads
(150, 250 or 300bp)



RT PCR: 4 amplicons, universal primers
(Gall et al, *Journal of Clinical Microbiology* 2012)

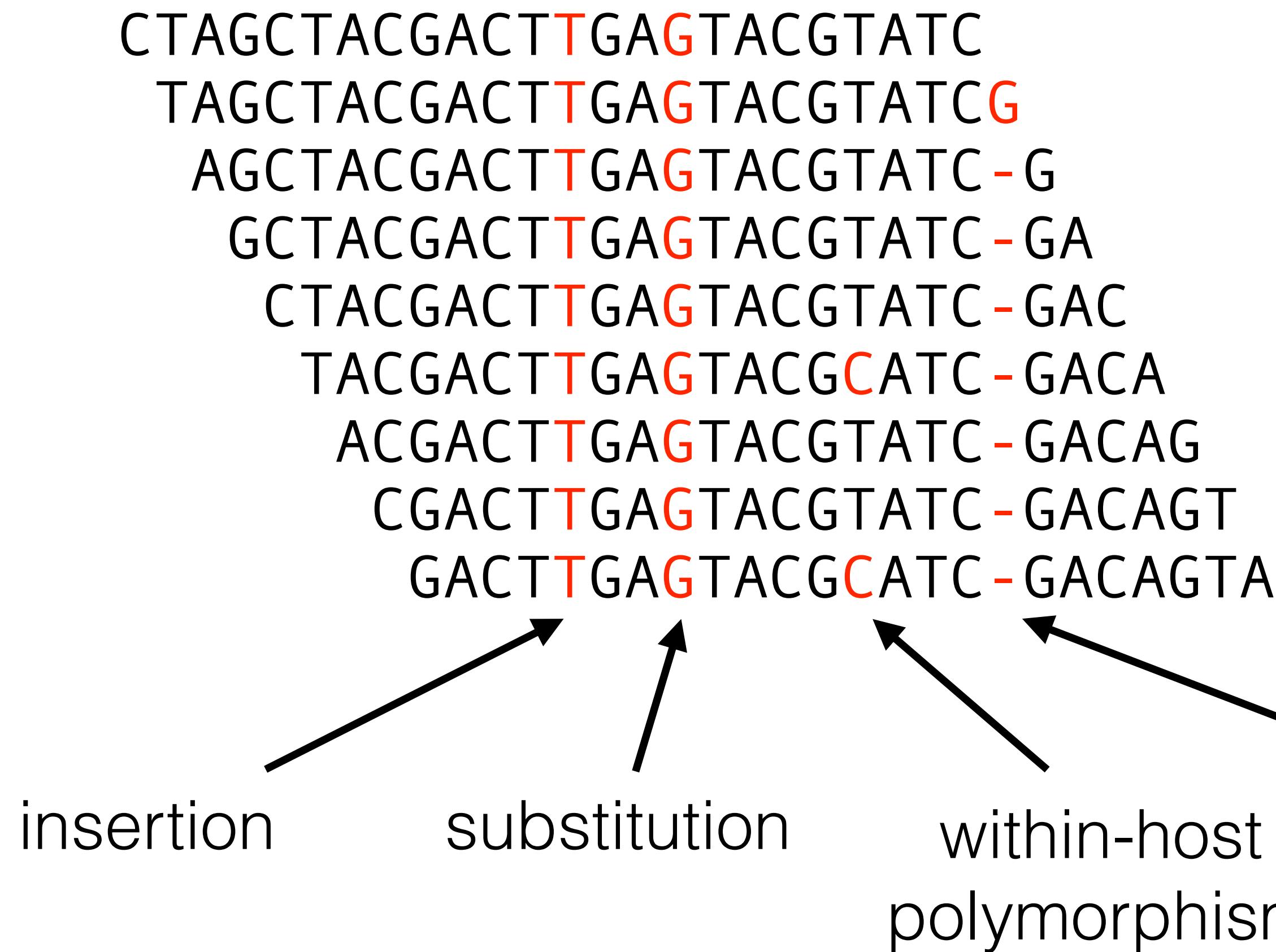
Miseq or
Hiseq



?

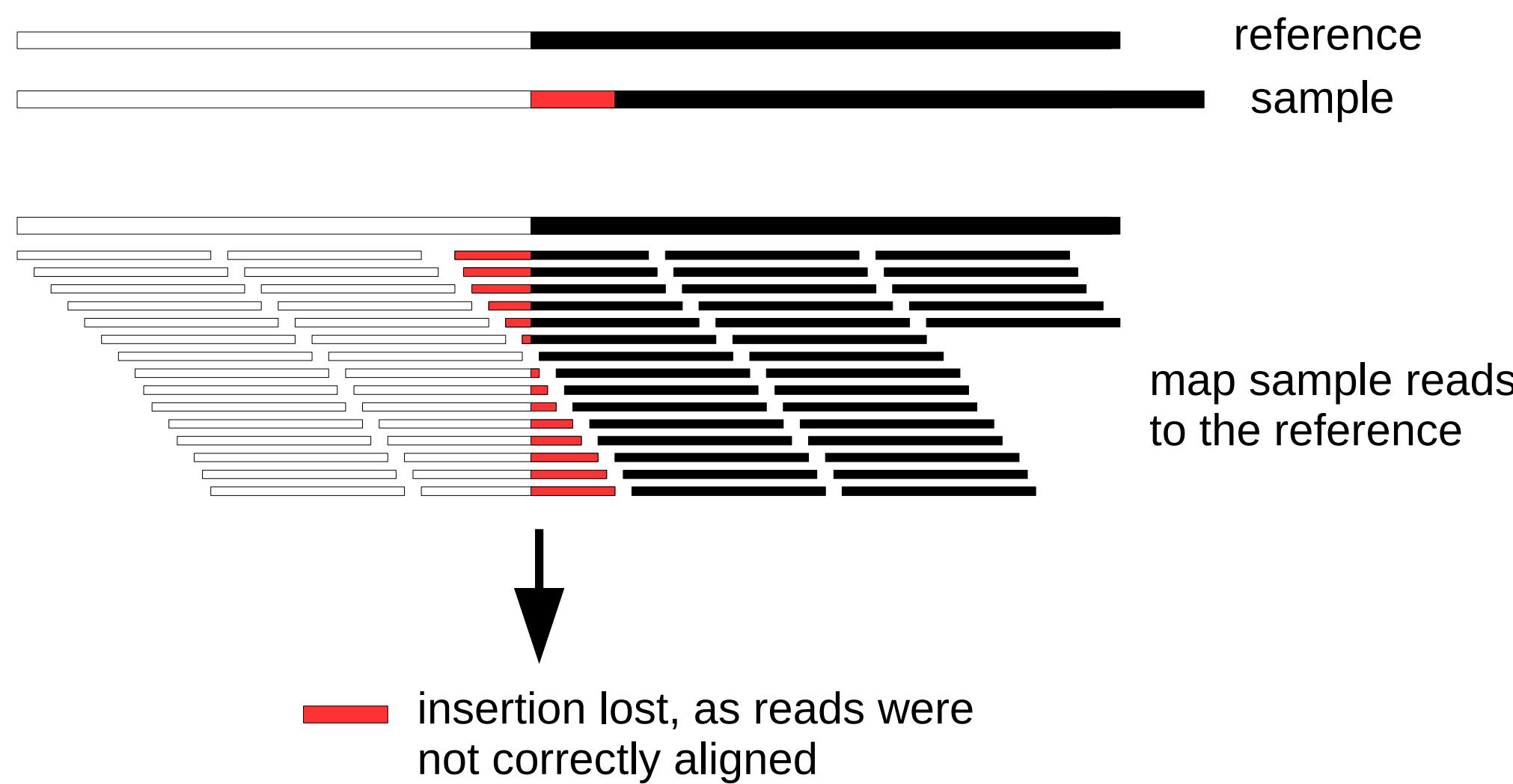
Map Reads to a Known Reference...

...CTAGCTACGACT-GAATACGTATCTGACAGTAT...



- Call the most common base amongst the reads at each position
- Retains within-host diversity information (base frequencies and linkage)
- Number of reads gives some measure of confidence

...but Mapping Induces Bias



The more different a read is from its reference, the more likely it is to be aligned incorrectly or not at all.



Systematic bias to make new sequences resemble known examples. An extremely high mutation rate makes this an acute problem for HIV.

Real data example: what the mapper should have done:

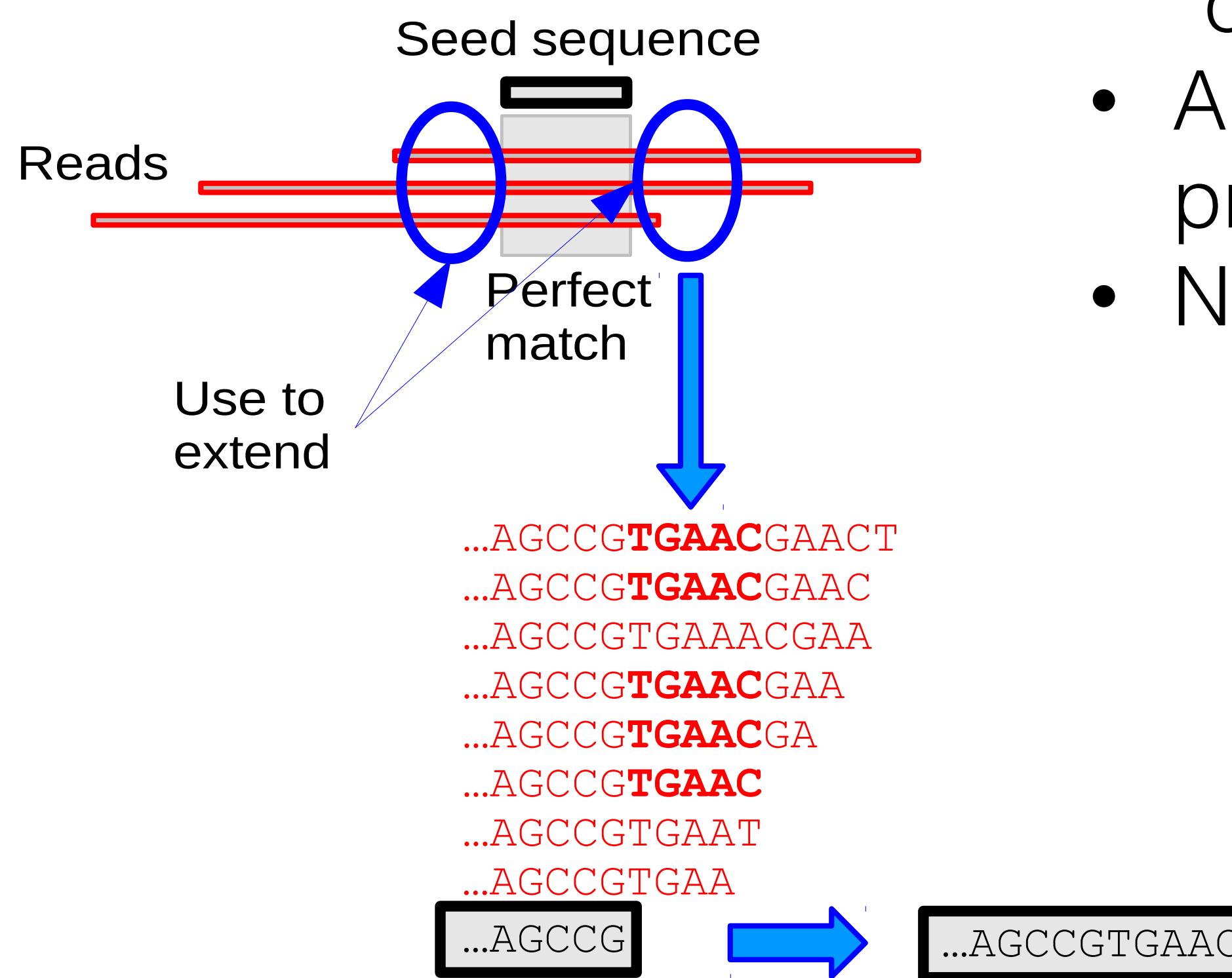
```
reference AAGTGTAGTGTGGAAGTTGACAGCAGCCTAGCACTCACAGGGCCGAGAGAACATCCGGAGTTTACAAAGACTGCTGACATC-----GAGTTTCCTACAAGGGACTTCCGCTGGGGACTTCCAGGGAAAGGCGTGGCCTGGCGGG  
read 1   AAGTGTAAATGTGGAAGTTGACAGCCGCCTAGCATTCCATCACGTAGCCGAGAGCTGCATCCGGAGTACTACAAAGACTGCTGACATCCTACAAAGACTGCTGACATCGAGCTTCTGCAAAGGGACTTCCGCTGGGGACTTCCAGG  
read 2   GACAGCCGCCTAGCATTCCATCACGTAGCCGAGAGCTGCATCCGGAGTACTACAAAGACTGCTGACATCCTACAAAGACTGCTGACATCGAGCTTCTGCAAAGGGACTTCCGCTGGGGACTTCCAGGGAGGCGTGGCCTGGCGGG
```

What it did:

```
reference AAGTGTAGTGTGGAAGTTGACAGCAGCCTAGCACTCACAGGGCCGAGAGAACATCCGGAGTTTACAAAGACTGCTGACATCAGTTCTACAAGGGACTTCCGCTGGGGACTTCCAGGGAAAGGCGTGGCCTGGCGGG  
read 1   AAGTGTAAATGTGGAAGTTGACAGCCGCCTAGCATTCCATCACGTAGCCGAGAGCTGCATCCGGAGTACTACAAAGACTGCTGACATCCTACAAAGACTGCTGACATCGAGCTTCTGCAAAGGGACTTCCGCTGGGGACTTCCAGG  
read 2   GACAGCCGCCTAGCATTCCATCACGTAGCCGAGAGCTGCATCCGGAGTACTACAAAGACTGCTGACATCCTACAAAGACTGCTGACATCGAGCTTCTGCAAAGGGACTTCCGCTGGGGACTTCCAGGGAGGCGTGGCCTGGCGGG
```

reads clipped

De Novo Assembly

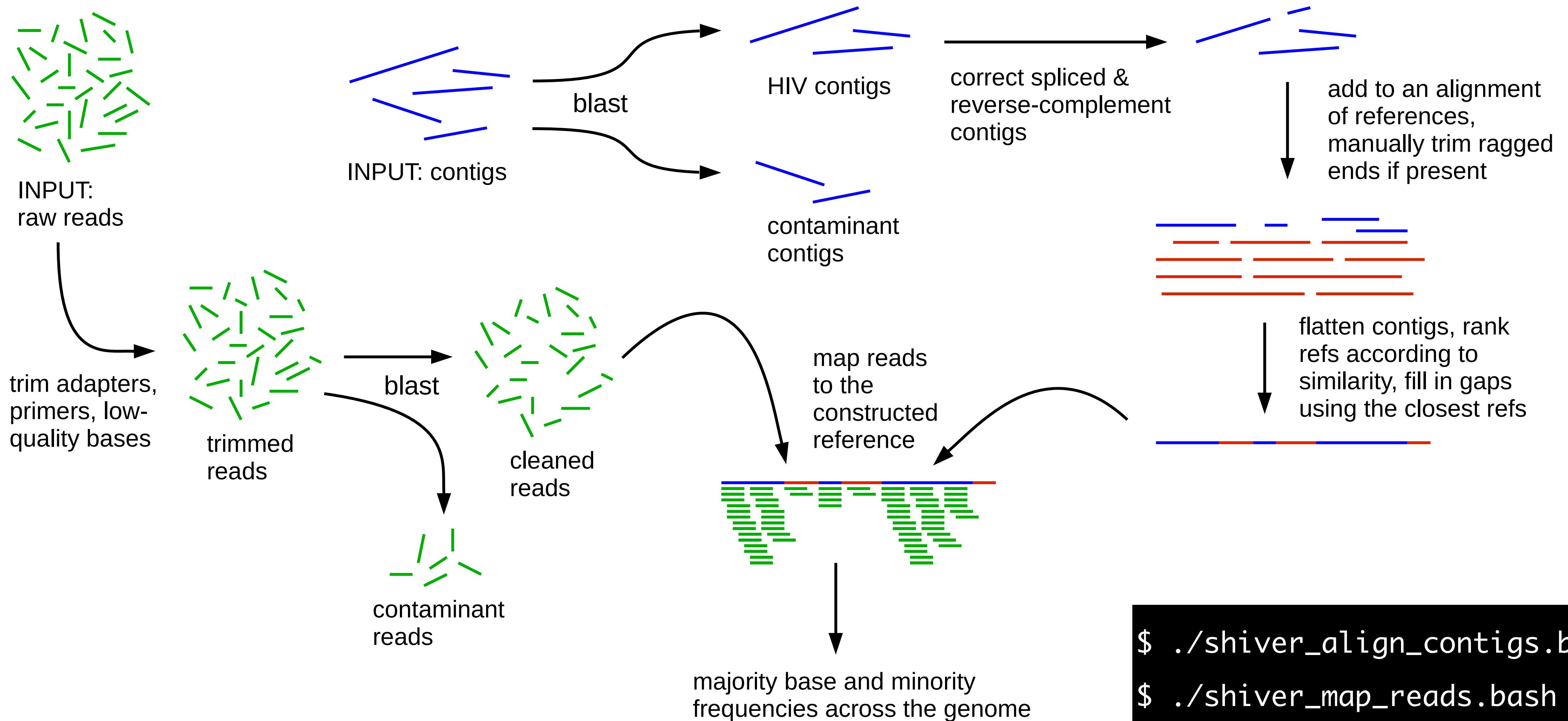


- e.g. IVA (*Hunt et al., Bioinformatics 2015*): on HIV, “outperforms all other virus de novo assemblers”.
- Align the reads to themselves, iteratively extending, producing a number of summary sequences.
- No need for a reference: no bias.

but...

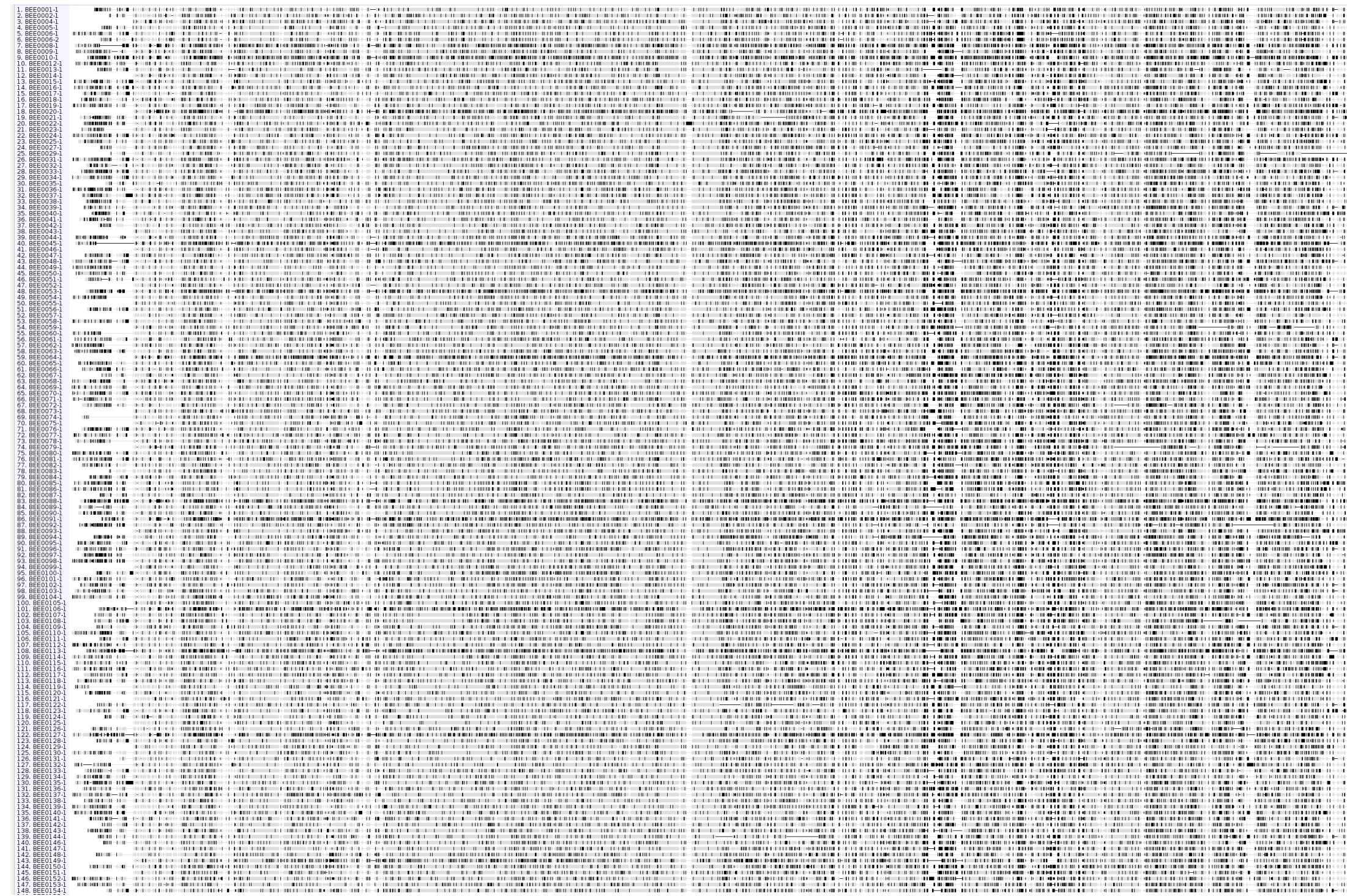
- Diversity & low coverage → no contig
- Diversity & high coverage → two or more overlapping, different contigs. Which one is ‘correct’?
- No idea of base frequencies or linkage
- Occasional assembly errors: splicing separate parts of the genome, reverse complements.

shiver - Sequences from HIV Easily Reconstructed

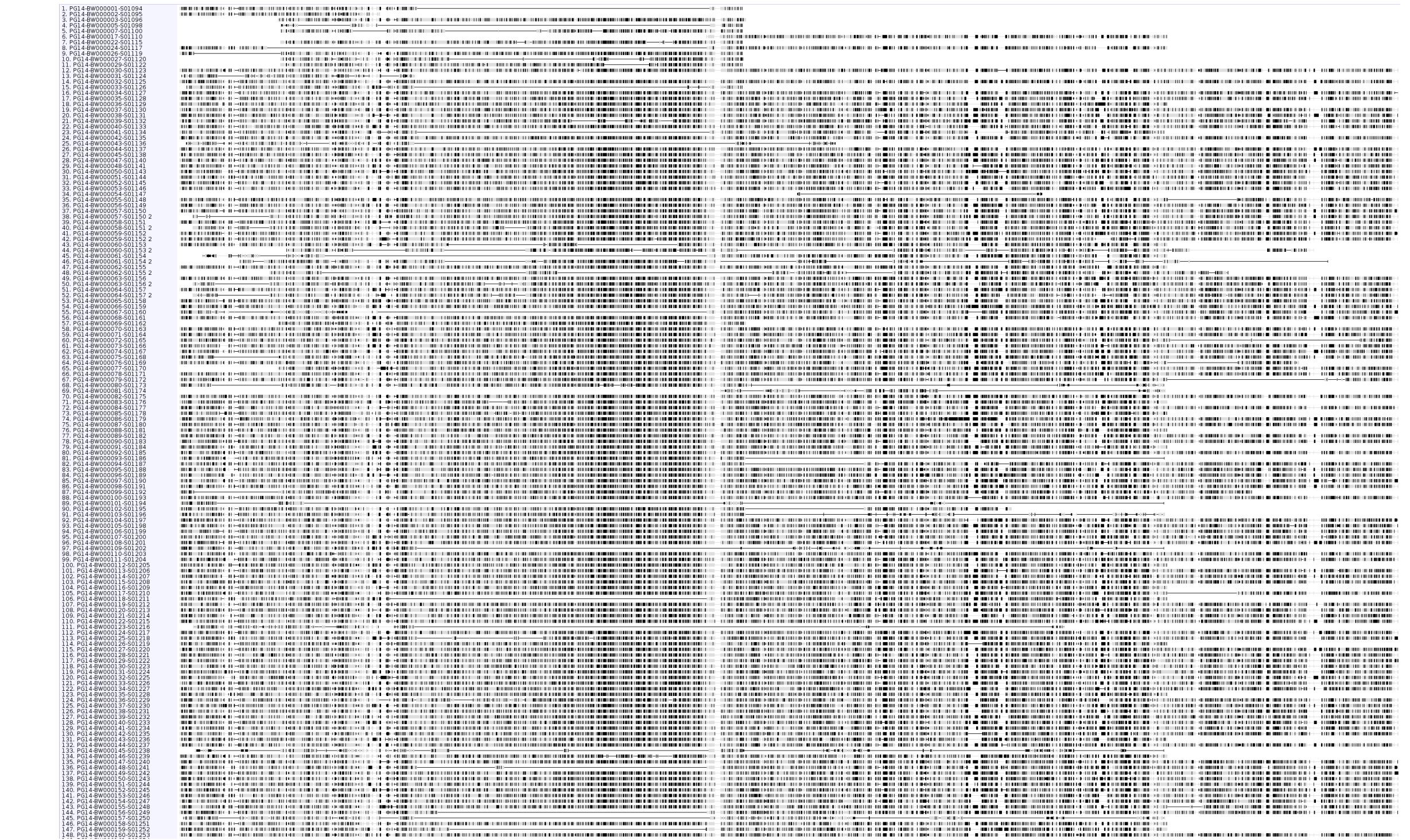


Used to reconstruct whole* genomes for

BEEHIVE: 1700+ samples



PANGEA: 4300+ samples

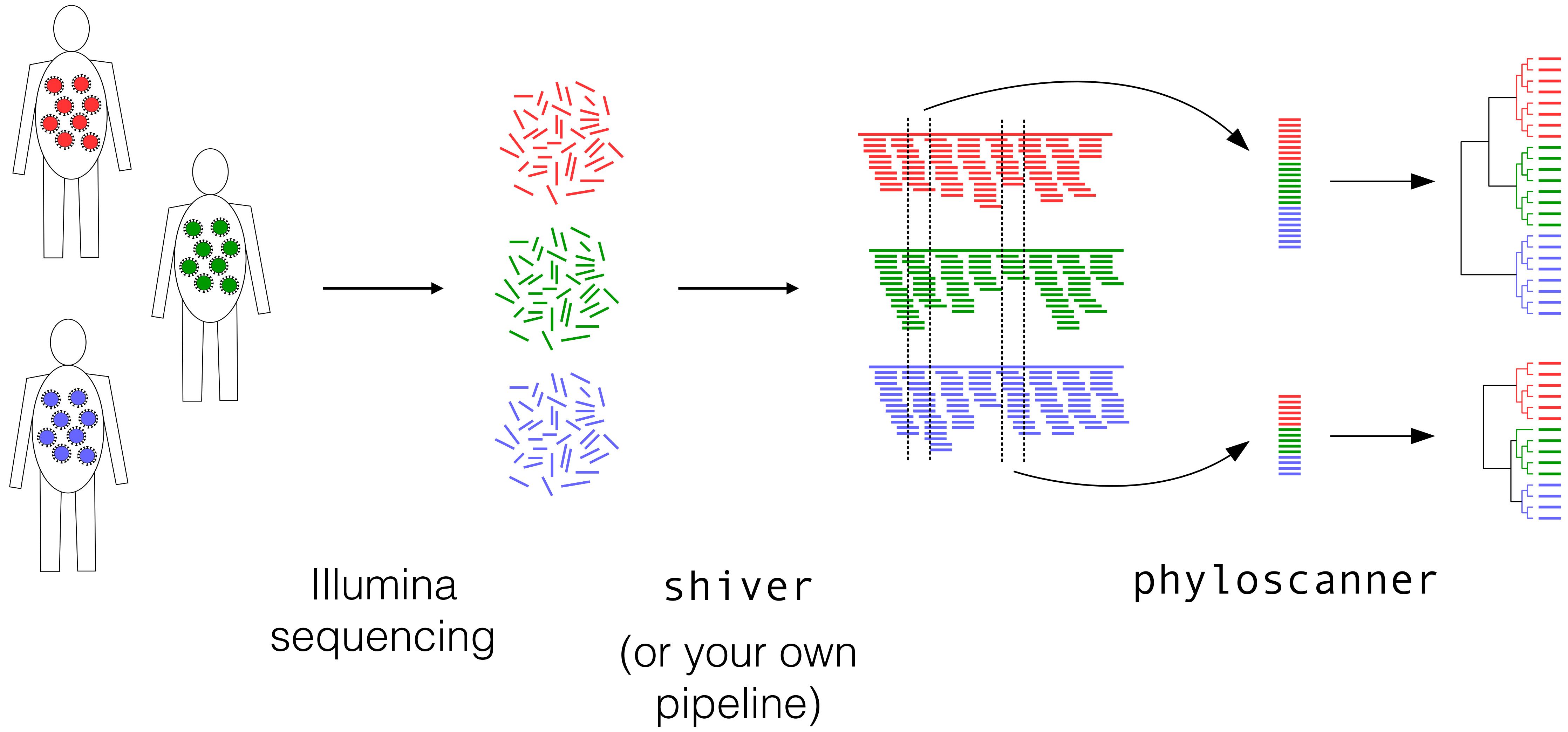


*or not, depending on the sample

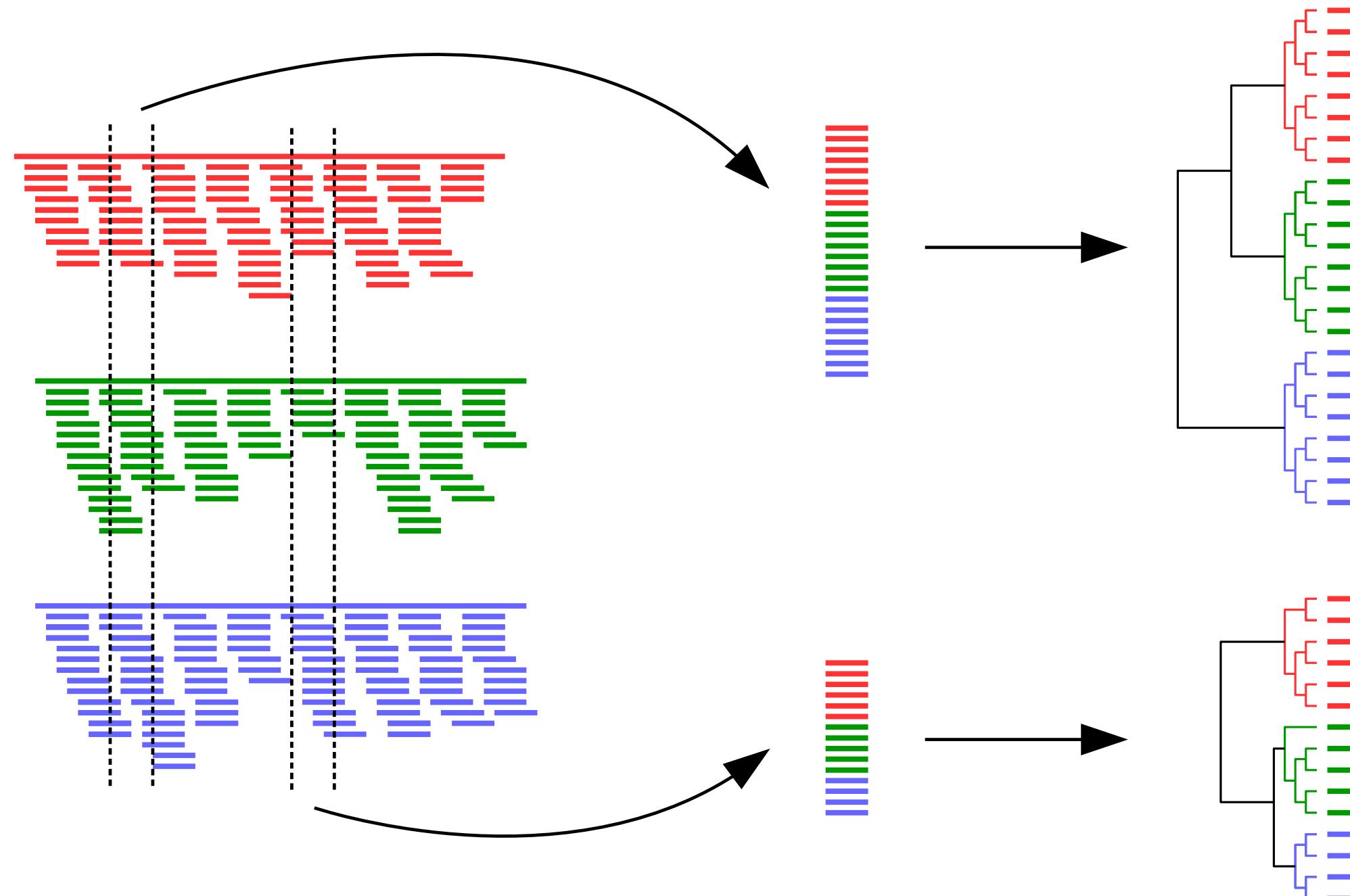
A consensus sequence does not show

- **transmission** (or at least not well; important for epidemiology),
- **dual infections** (important clinically, for bioinformatics, for evolution),
- **contamination** (important to distinguish from dual infections),
- **recombination** between strains within-host.

Put within-host diversity in context: examine within- and between-host evolution at once, across the genome.



Easy



```
$ ./phyloscanner.py MyBamFiles.txt MyRefFiles.txt --windows 1,300,200,500,...
```

red.bam
green.bam
blue.bam

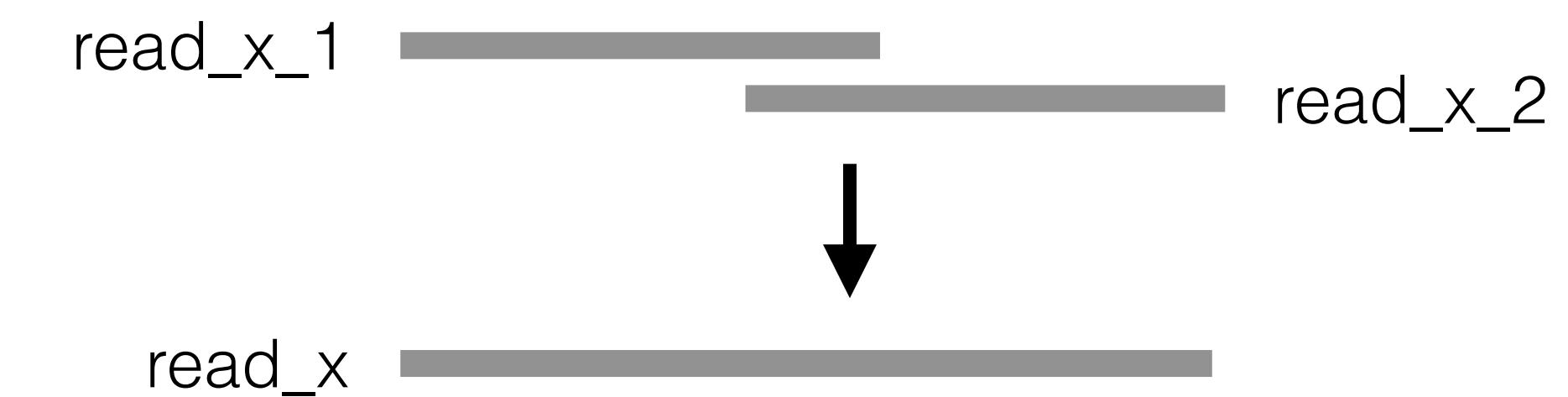
redRef.fasta
greenRef.fasta
blueRef.fasta

Genomic coordinates
of the windows to be
analysed; here,
1-300, 200-500, ...

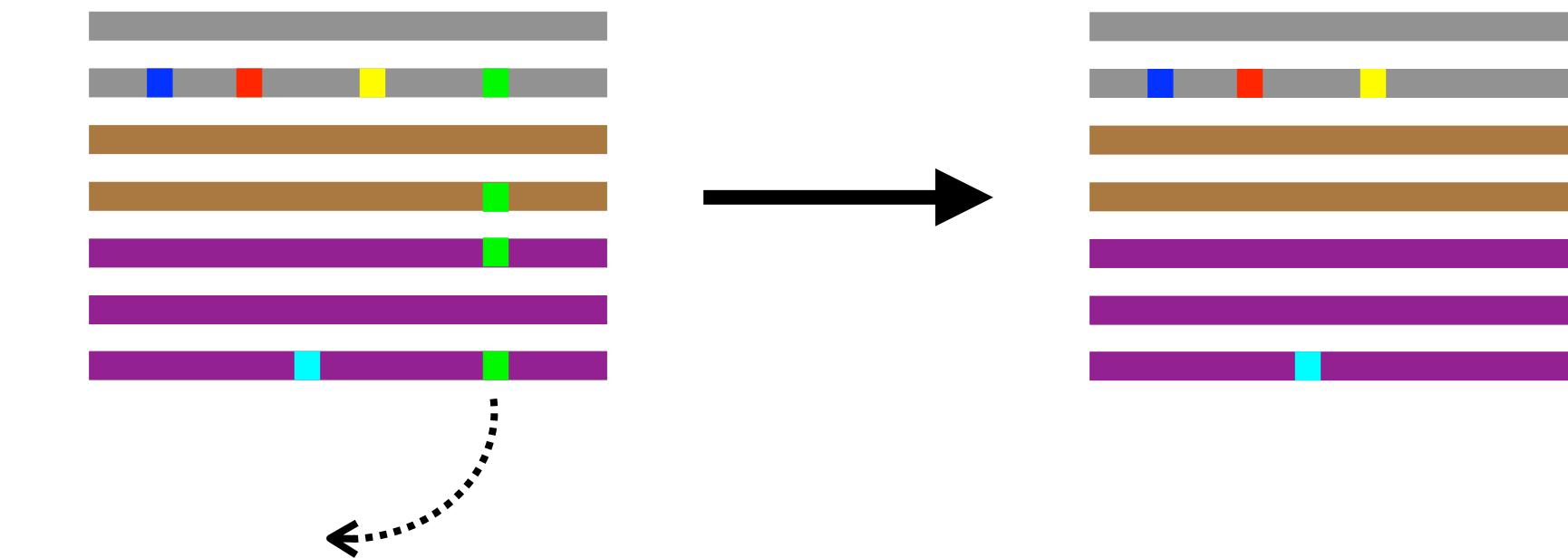
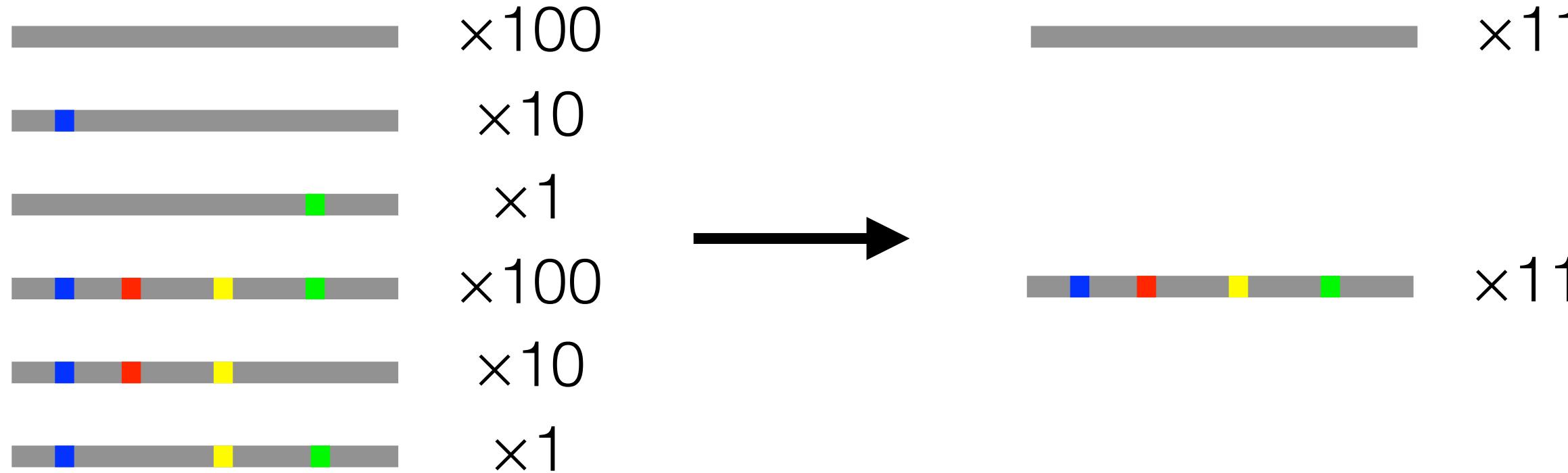
Feature-packed 1



Merge overlapping paired
reads into longer reads.



Similarity- and frequency-based
read merging, for speed.

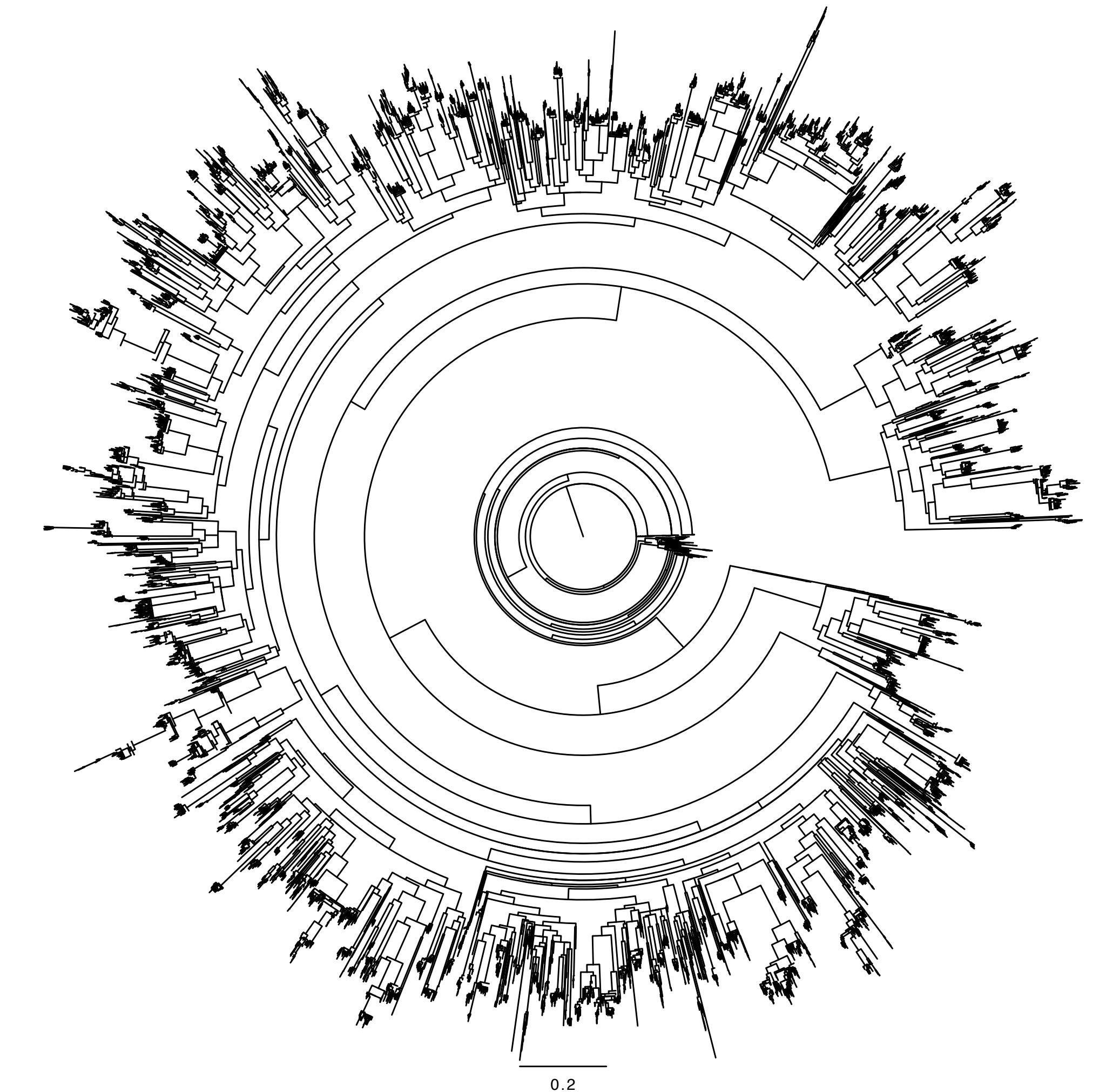


Feature-packed 2

- Include known references with the reads.
- Trim and/or discard low-quality reads.
- Minimum read count.
- Check every read for duplication in any other patient to find contamination. Remove duplicates depending on counts.
- Bootstraps.

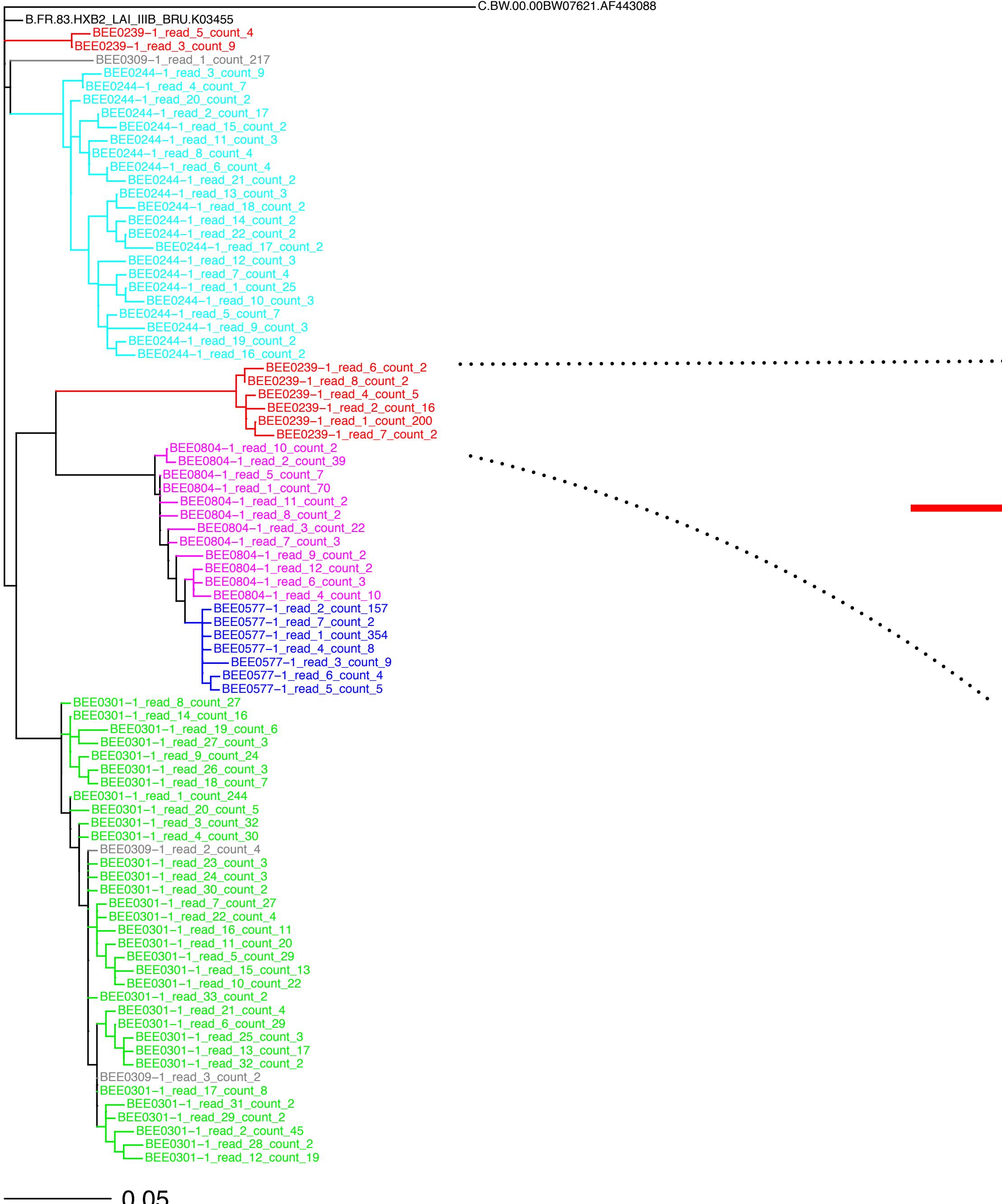
Trivially parallelisable: split the whole genome into windows, run each window as a separate job on your cluster.

<1 day for $O(1000)$ patients with coverage $O(10,000)$ (RAxML the bottleneck).



375bp at the end of the gp120 gene. 700 patients, 8574 distinct reads.
(merging threshold = 1bp, minimum read count = 2)

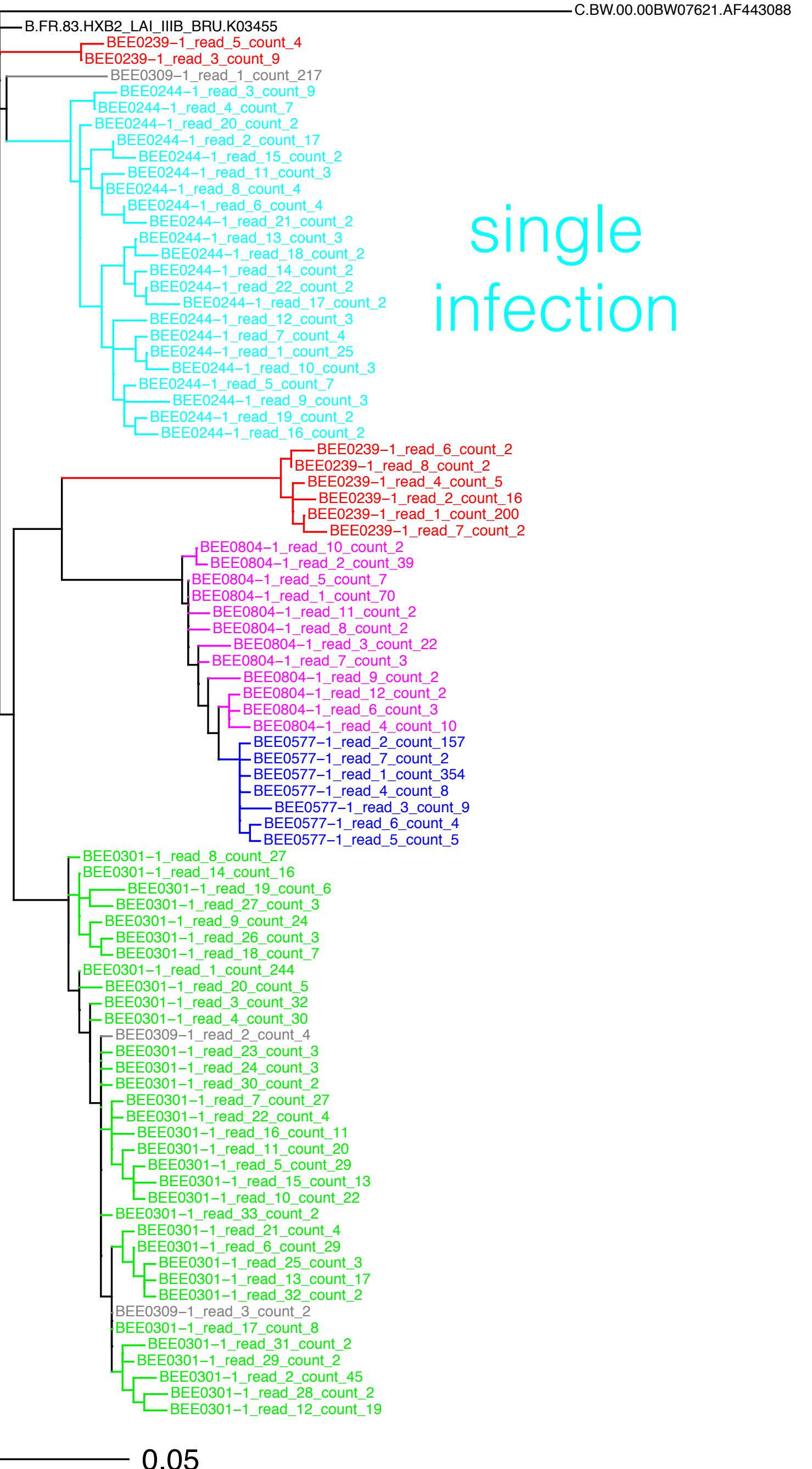
HXB2 1800-2150 (gag)



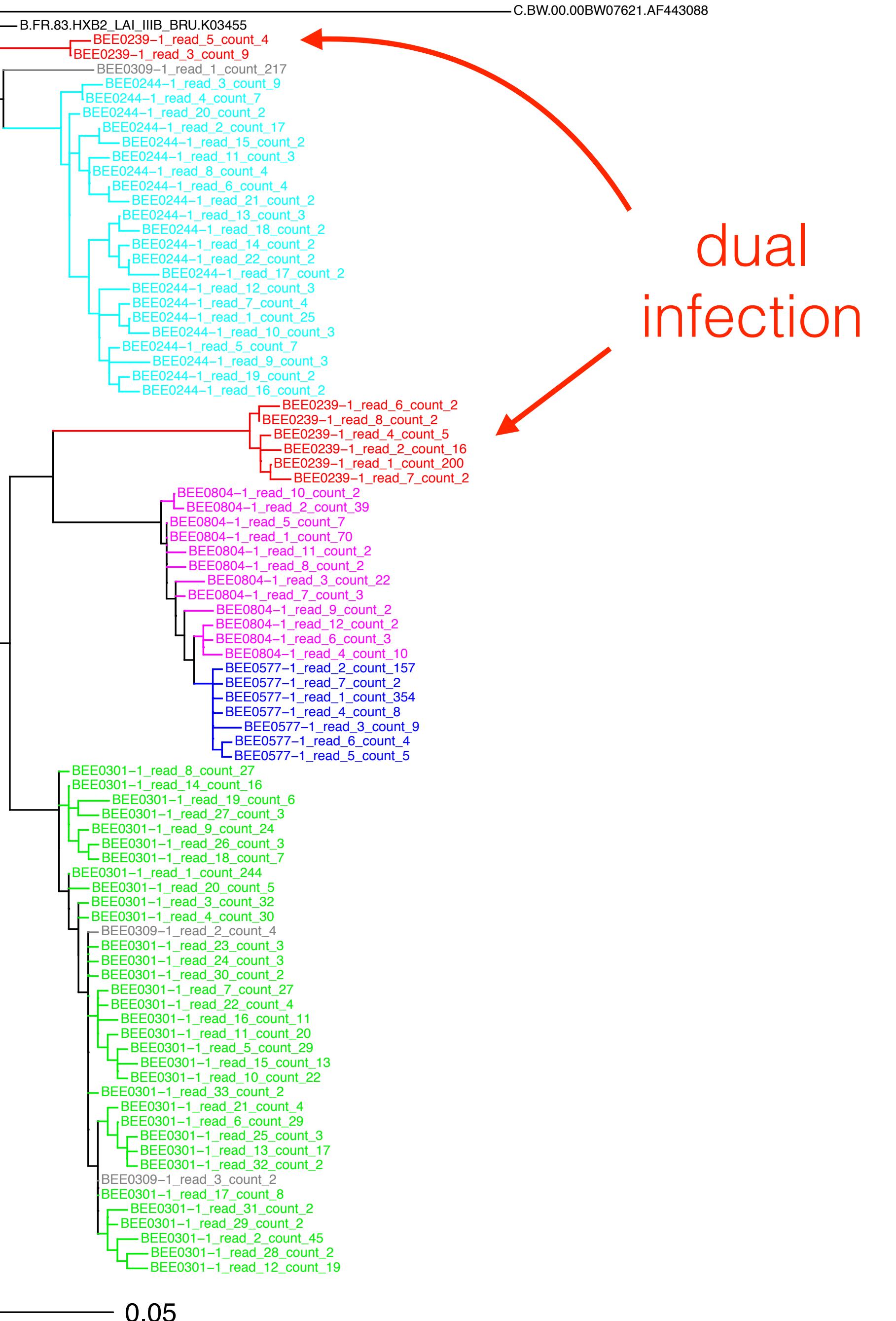
BEE0239-1_read_6_count_2
BEE0239-1_read_8_count_2
BEE0239-1_read_4_count_5
BEE0239-1_read_2_count_16
BEE0239-1_read_1_count_200
BEE0239-1_read_7_count_2

HXB2 1800-2150 (gag)

single
infection

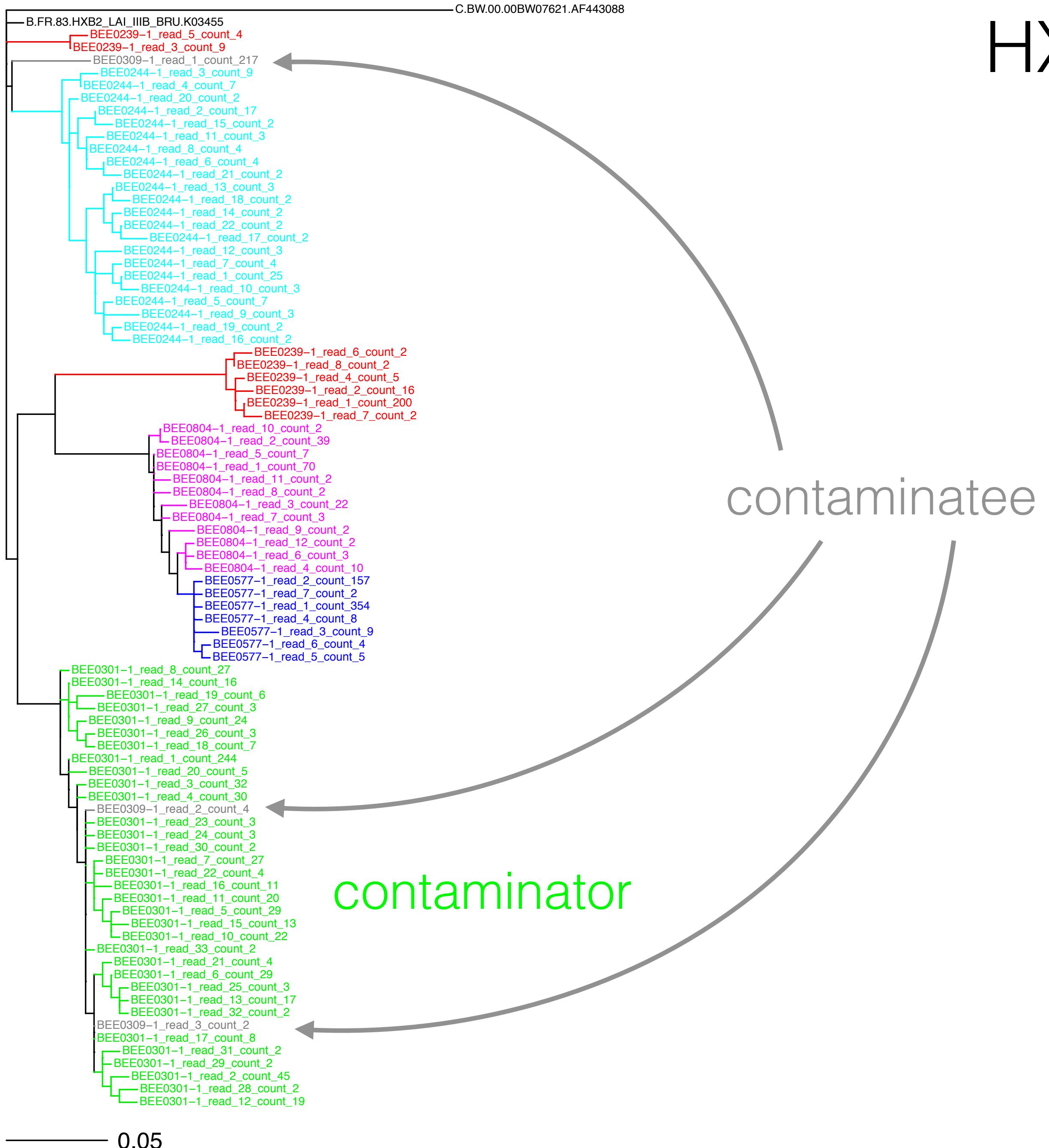


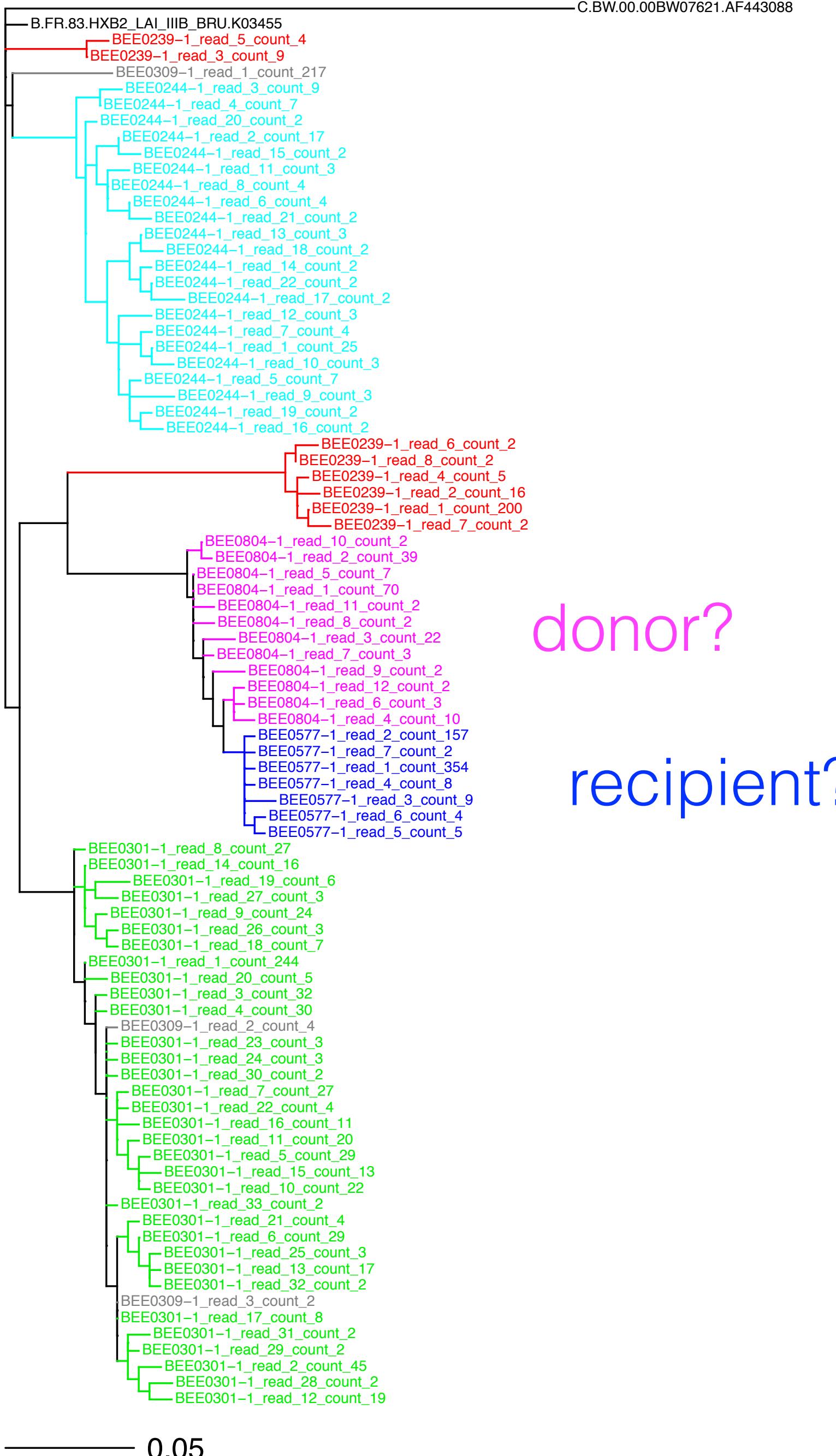
HXB2 1800-2150 (gag)



dual
infection

HXB2 1800-2150 (gag)





HXB2 1800-2150 (gag)

donor?

recipient?

Left tree: 1bp
merging,
minimum read
count = 2.

Re-run with just
the transmission
pair, no merging
and no minimum
read count for a
clearer picture

B.FR.83.HXB2_LAI_IIB.BRU.K03455 C.BW.00.BWB07621.AF443088

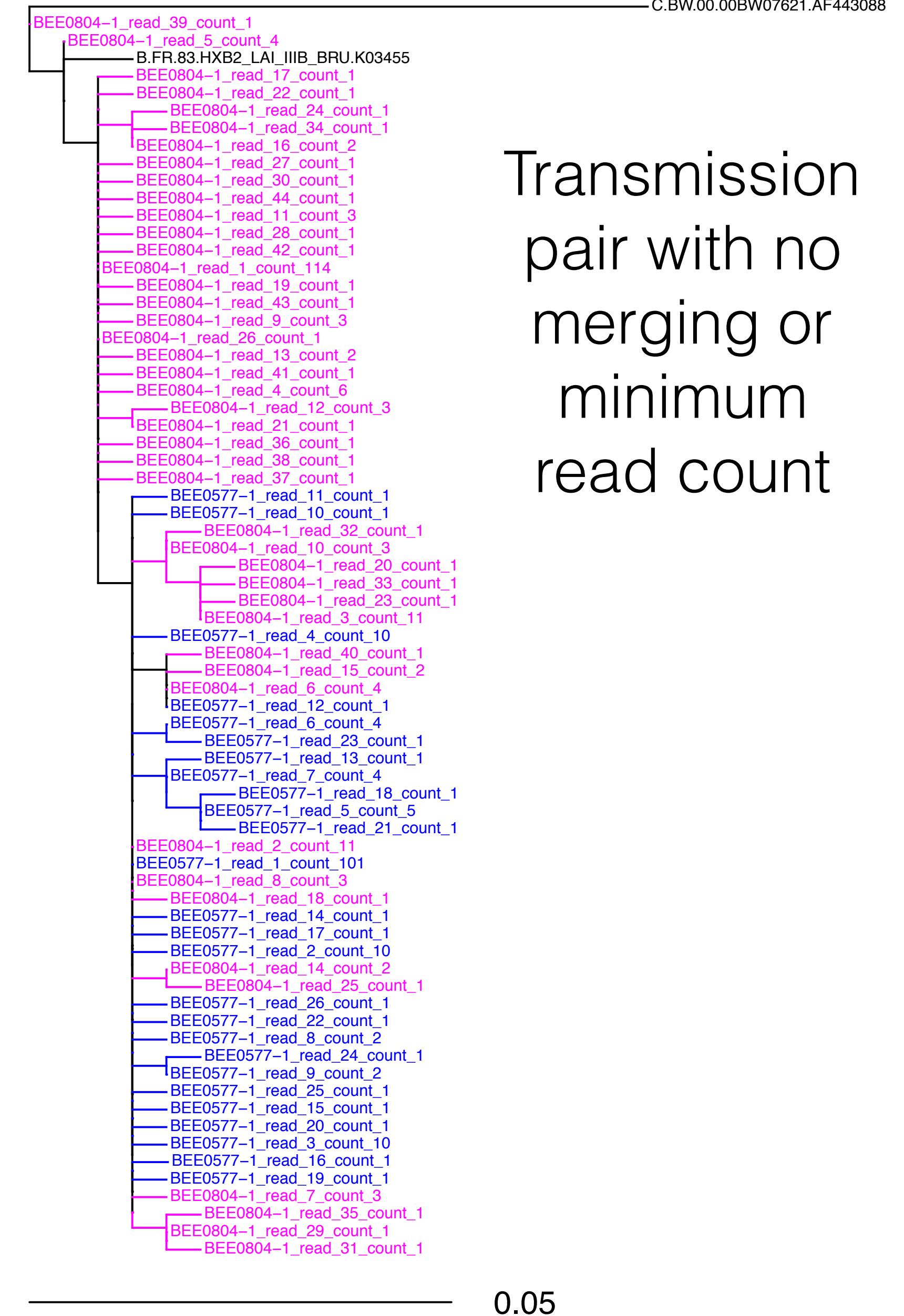
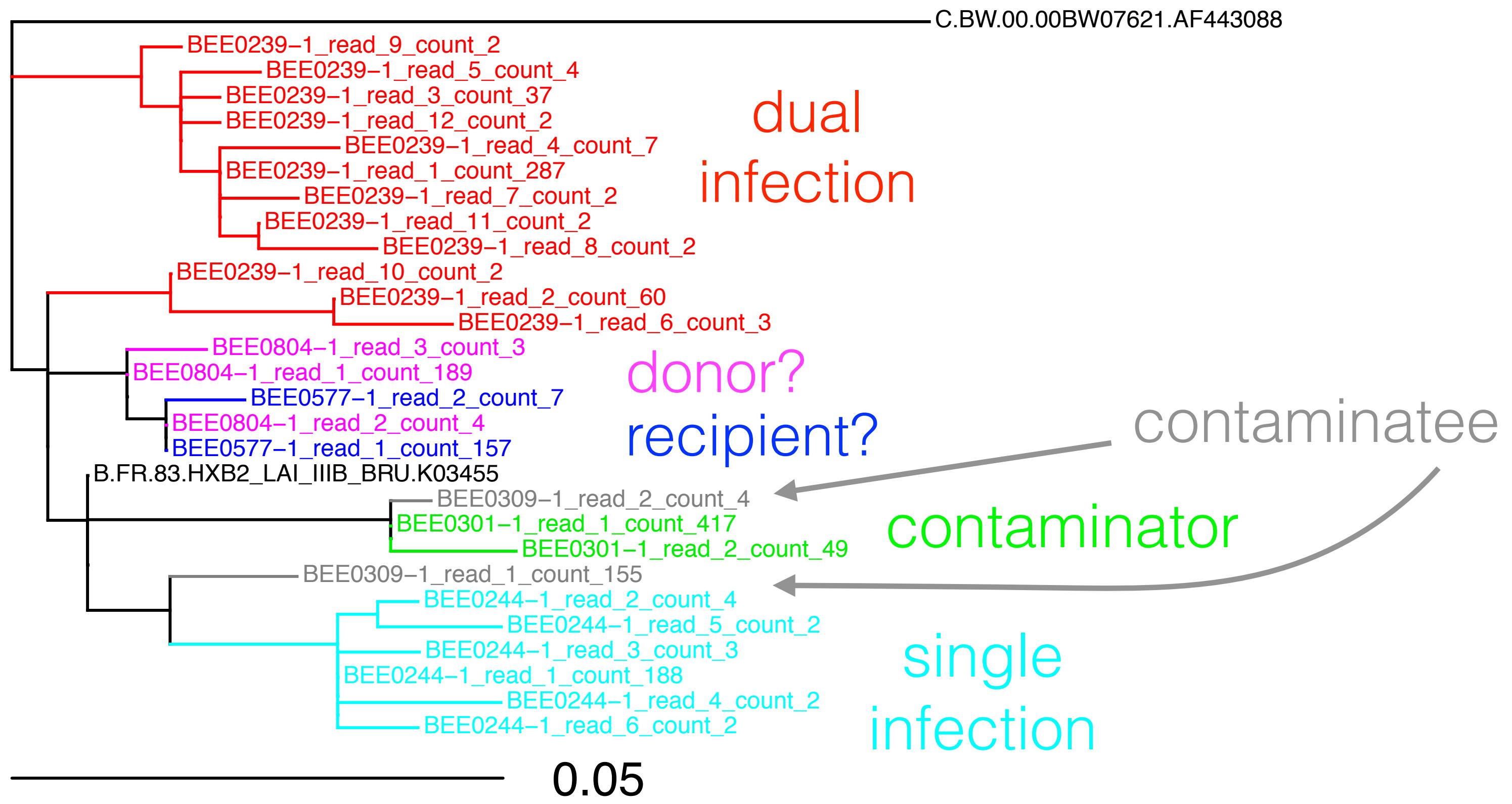
```

BEE0804-1, read_30, count_1
BEE0804-1, read_14, count_2
BEE0804-1, read_96, count_1
BEE0804-1, read_71, count_1
BEE0804-1, read_69, count_1
    BEE0804-1, read_19, count_2
    BEE0804-1, read_10, count_1
    BEE0804-1, read_56, count_1
    BEE0804-1, read_17, count_2
    BEE0804-1, read_22, count_1
    BEE0804-1, read_8, count_5
    BEE0804-1, read_34, count_1
    BEE0804-1, read_48, count_1
    BEE0804-1, read_28, count_1
    BEE0804-1, read_78, count_1
    BEE0804-1, read_90, count_1
    BEE0804-1, read_34, count_1
    BEE0804-1, read_30, count_2
    BEE0804-1, read_107, count_1
    BEE0804-1, read_91, count_1
    BEE0804-1, read_26, count_1
    BEE0804-1, read_50, count_1
    BEE0804-1, read_103, count_1
    BEE0804-1, read_1, count_40
    BEE0804-1, read_85, count_1
    BEE0804-1, read_84, count_1
    BEE0804-1, read_77, count_1
    BEE0804-1, read_63, count_1
    BEE0804-1, read_114, count_1
    BEE0804-1, read_97, count_1
    BEE0804-1, read_45, count_1
    BEE0804-1, read_37, count_1
    BEE0804-1, read_18, count_2
    BEE0804-1, read_10, count_1
    BEE0804-1, read_52, count_1
    BEE0804-1, read_12, count_3
    BEE0804-1, read_92, count_1
    BEE0804-1, read_16, count_2
    BEE0804-1, read_11, count_1
    BEE0804-1, read_19, count_2
    BEE0804-1, read_65, count_1
    BEE0804-1, read_23, count_1
    BEE0804-1, read_102, count_1
    BEE0804-1, read_25, count_1
    BEE0804-1, read_10, count_1
    BEE0804-1, read_81, count_1
    BEE0804-1, read_110, count_1
    BEE0804-1, read_67, count_1
    BEE0804-1, read_113, count_1
    BEE0804-1, read_59, count_1
    BEE0804-1, read_40, count_1
    BEE0804-1, read_44, count_1
    BEE0804-1, read_11, count_3
    BEE0804-1, read_30, count_1
    BEE0804-1, read_98, count_1
    BEE0804-1, read_53, count_1
    BEE0804-1, read_2, count_17
    BEE0804-1, read_109, count_1
    BEE0804-1, read_51, count_1
    BEE0804-1, read_35, count_1
    BEE0804-1, read_27, count_1
    BEE0804-1, read_20, count_2
    BEE0804-1, read_4, count_12
    BEE0804-1, read_89, count_1
    BEE0804-1, read_5, count_5
    BEE0804-1, read_49, count_1
    BEE0804-1, read_76, count_1
    BEE0804-1, read_43, count_1
    BEE0804-1, read_42, count_1
    BEE0804-1, read_86, count_1
    BEE0804-1, read_3, count_14
    BEE0804-1, read_58, count_1
    BEE0804-1, read_54, count_1
    BEE0804-1, read_80, count_1
    BEE0804-1, read_88, count_1
    BEE0804-1, read_91, count_1
    BEE0804-1, read_94, count_1
    BEE0804-1, read_79, count_1
    BEE0804-1, read_66, count_1
        BEE0804-1, read_13, count_2
        BEE0804-1, read_75, count_1
        BEE0804-1, read_112, count_1
        BEE0804-1, read_33, count_1
        BEE0804-1, read_6, count_6
        BEE0804-1, read_60, count_1
        BEE0804-1, read_64, count_1
        BEE0804-1, read_68, count_1
        BEE0804-1, read_55, count_1
        BEE0804-1, read_111, count_1
            BEE0804-1, read_10, count_3
            BEE0804-1, read_105, count_1
            BEE0804-1, read_9, count_3
            BEE0577-1, read_82, count_1
            BEE0577-1, read_13, count_2
            BEE0577-1, read_75, count_1
            BEE0577-1, read_112, count_1
            BEE0577-1, read_33, count_1
            BEE0577-1, read_6, count_6
            BEE0577-1, read_104, count_1
            BEE0577-1, read_98, count_1
            BEE0577-1, read_1, count_1
            BEE0577-1, read_17, count_3
            BEE0577-1, read_112, count_1
            BEE0577-1, read_44, count_1
            BEE0577-1, read_71, count_1
            BEE0577-1, read_10, count_1
            BEE0577-1, read_89, count_1
            BEE0577-1, read_61, count_1
            BEE0577-1, read_32, count_1
            BEE0577-1, read_49, count_1
            BEE0577-1, read_63, count_1
            BEE0577-1, read_60, count_1
            BEE0577-1, read_28, count_1
            BEE0577-1, read_11, count_6
            BEE0577-1, read_114, count_1
            BEE0577-1, read_111, count_1
            BEE0577-1, read_25, count_1
            BEE0577-1, read_70, count_1
            BEE0577-1, read_51, count_1
            BEE0577-1, read_72, count_1
            BEE0577-1, read_86, count_1
            BEE0577-1, read_107, count_1
            BEE0577-1, read_2, count_99
            BEE0577-1, read_86, count_1
            BEE0577-1, read_83, count_1
            BEE0577-1, read_64, count_1
            BEE0577-1, read_34, count_1
            BEE0577-1, read_60, count_1
            BEE0577-1, read_59, count_1
            BEE0577-1, read_10, count_1
            BEE0577-1, read_30, count_1
            BEE0577-1, read_20, count_2
            BEE0577-1, read_81, count_1
            BEE0577-1, read_100, count_1
            BEE0577-1, read_113, count_1
            BEE0577-1, read_3, count_1
            BEE0577-1, read_92, count_1
            BEE0577-1, read_8, count_9
            BEE0577-1, read_14, count_5
            BEE0577-1, read_21, count_1
            BEE0577-1, read_57, count_1
            BEE0577-1, read_12, count_6
            BEE0577-1, read_87, count_1
            BEE0577-1, read_87, count_1
                BEE0577-1, read_110, count_1
            BEE0577-1, read_6, count_22
            BEE0577-1, read_73, count_1
            BEE0577-1, read_57, count_1
            BEE0577-1, read_9, count_7
            BEE0577-1, read_43, count_1
            BEE0577-1, read_16, count_4
            BEE0577-1, read_69, count_1
            BEE0577-1, read_58, count_1
            BEE0577-1, read_58, count_1
            BEE0577-1, read_105, count_1
            BEE0577-1, read_94, count_1
            BEE0577-1, read_4, count_62
            BEE0577-1, read_103, count_1
            BEE0577-1, read_39, count_1
            BEE0577-1, read_29, count_1
            BEE0577-1, read_46, count_1
            BEE0577-1, read_19, count_3
            BEE0577-1, read_45, count_1
            BEE0577-1, read_45, count_1
            BEE0577-1, read_27, count_1
            BEE0577-1, read_91, count_1
            BEE0577-1, read_33, count_1
            BEE0577-1, read_10, count_7
            BEE0577-1, read_15, count_1
            BEE0577-1, read_41, count_1
            BEE0577-1, read_37, count_1
            BEE0577-1, read_90, count_1
            BEE0577-1, read_23, count_1
            BEE0577-1, read_13, count_5
            BEE0577-1, read_15, count_1
            BEE0577-1, read_108, count_1
            BEE0577-1, read_77, count_1
            BEE0577-1, read_48, count_1
            BEE0577-1, read_101, count_1
            BEE0577-1, read_35, count_1
            BEE0577-1, read_56, count_1
            BEE0577-1, read_1, count_106
            BEE0577-1, read_75, count_1
            BEE0577-1, read_55, count_1
            BEE0577-1, read_76, count_1
            BEE0577-1, read_84, count_1
            BEE0577-1, read_18, count_3
            BEE0577-1, read_65, count_1
            BEE0577-1, read_36, count_1
            BEE0577-1, read_40, count_1
            BEE0577-1, read_26, count_1
            BEE0577-1, read_39, count_1
            BEE0577-1, read_31, count_1
            BEE0577-1, read_22, count_1
            BEE0577-1, read_66, count_1

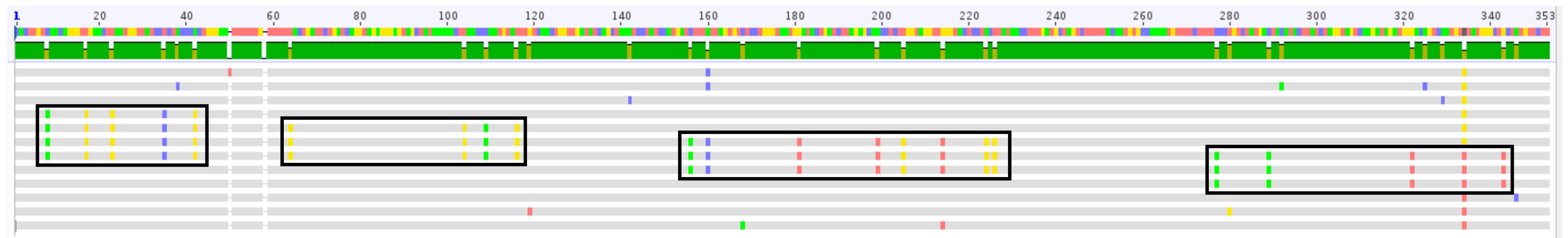
```



HXB2 2800-3150 (pol)

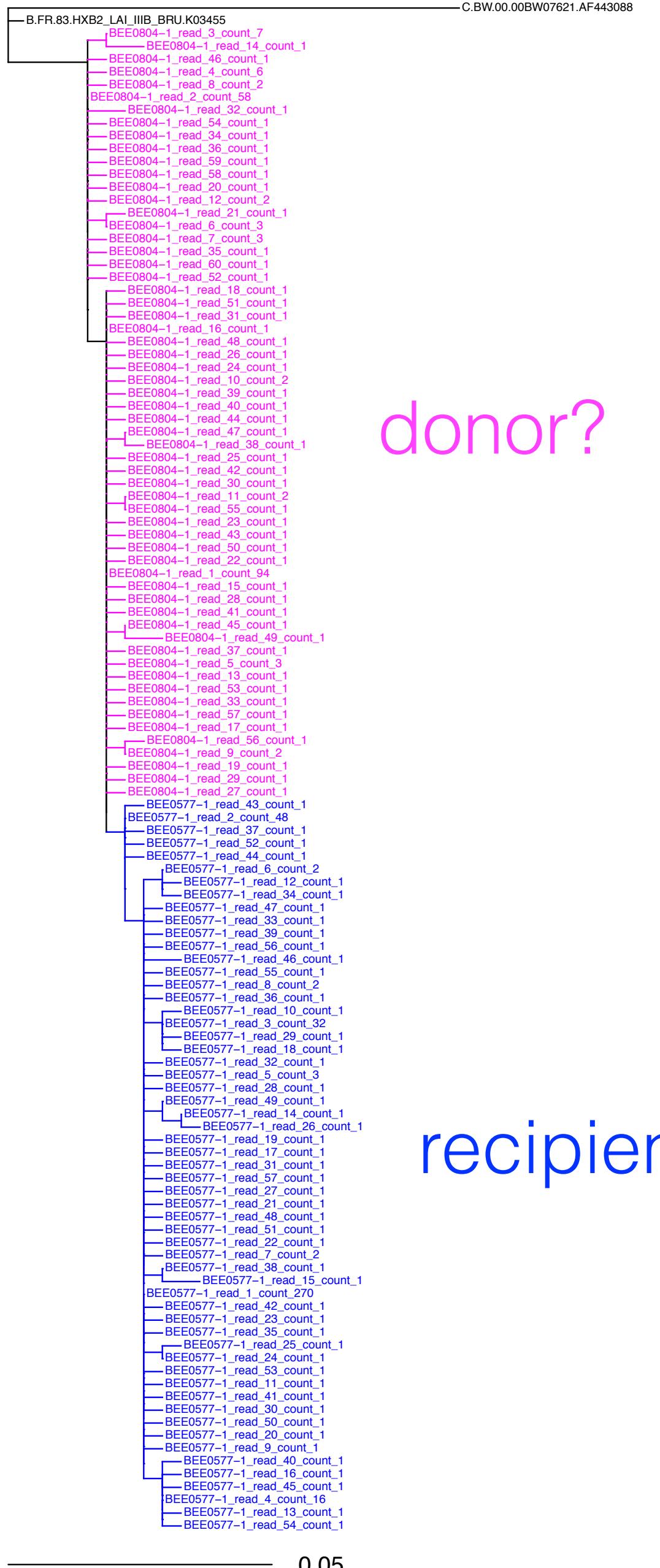
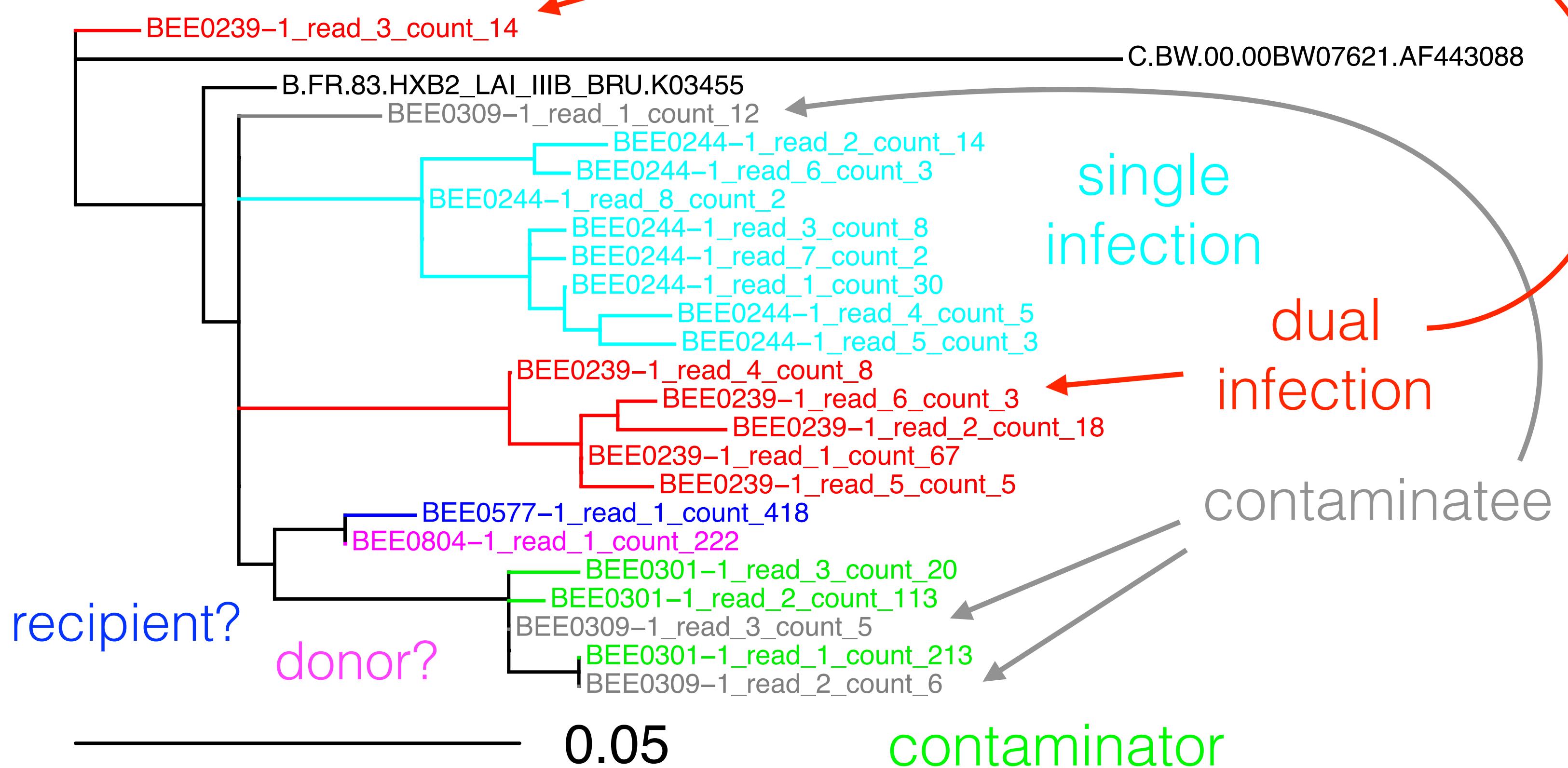


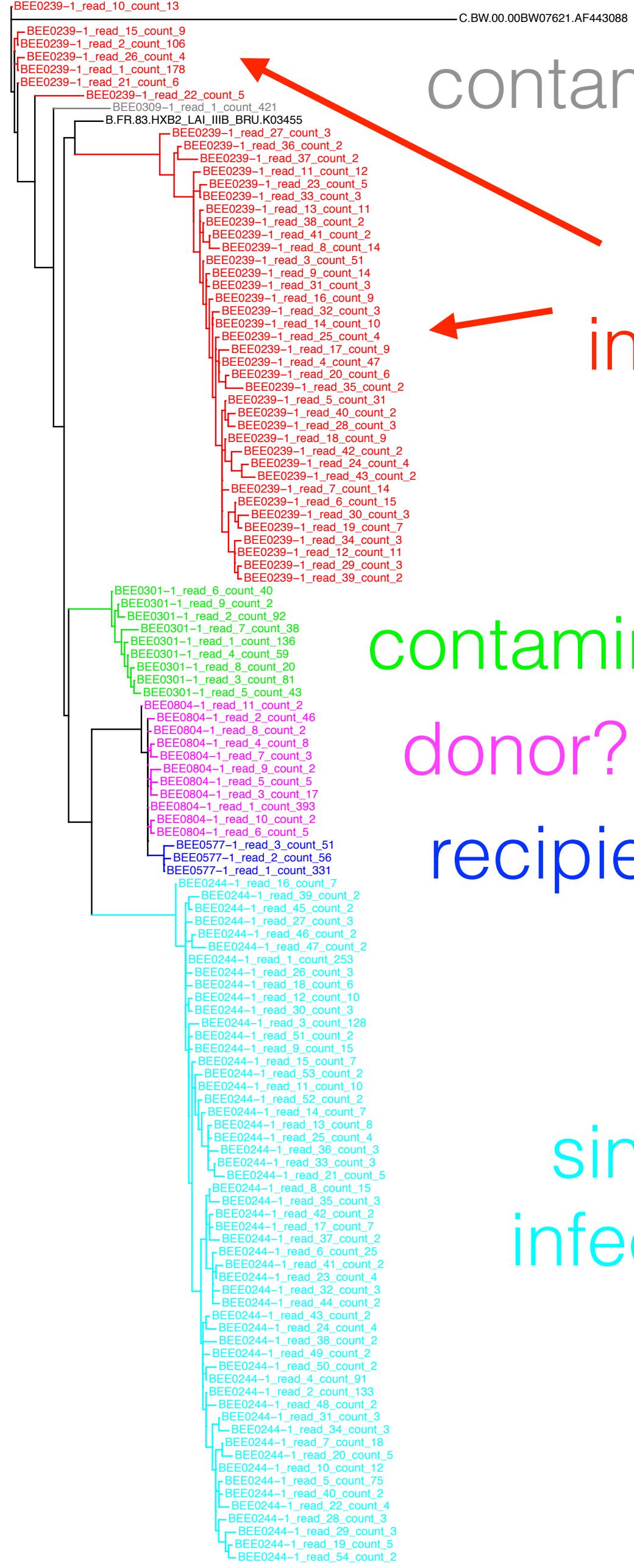
Transmission pair with no merging or minimum read count



dual infection, point of recombination

HXB2 4800-5150 (pol & vif)





contaminantee (not here)

dual
infection

contaminator (not here)

donor?

recipient?

single
infection

HXB2 6050- 6400 (vpu)

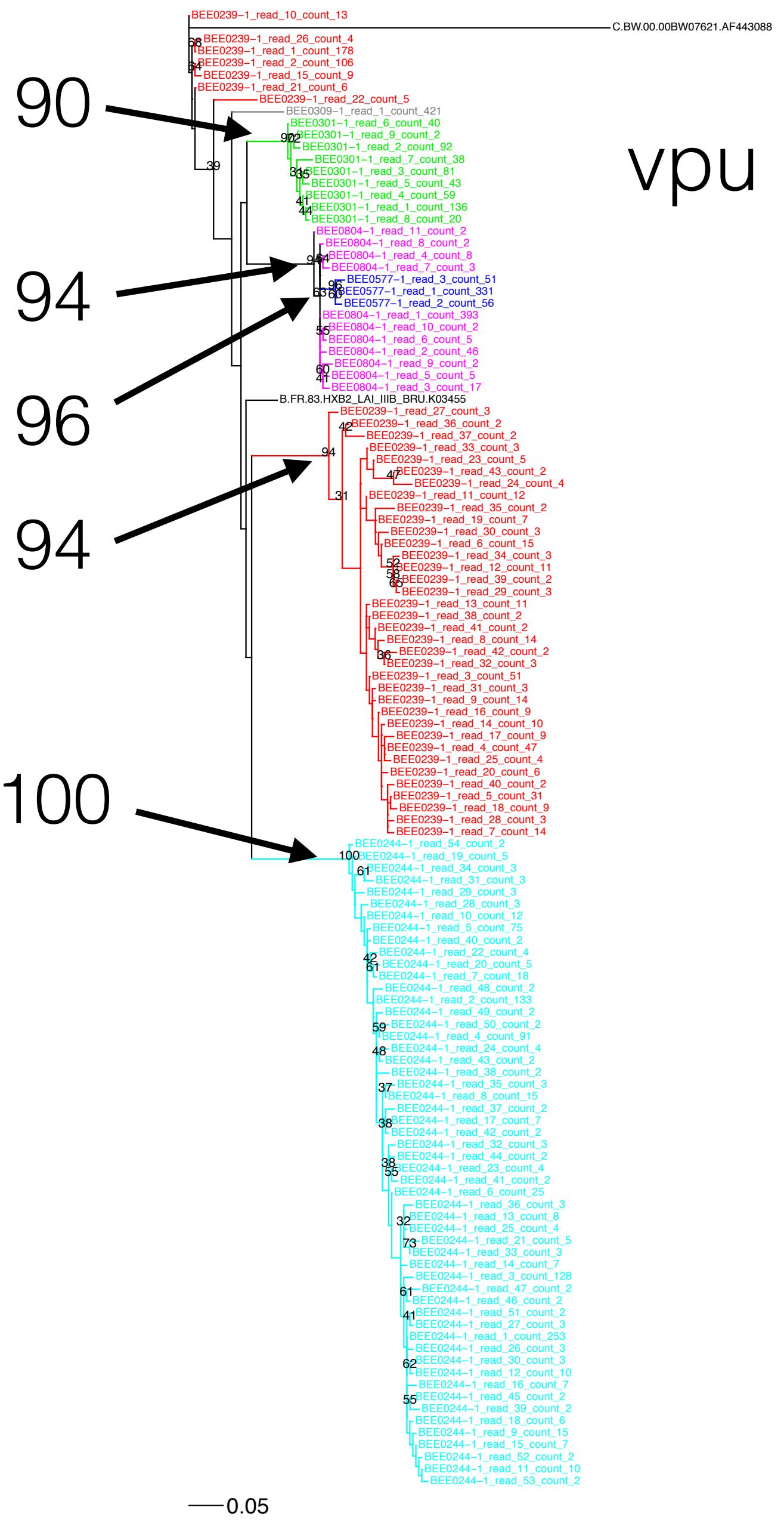
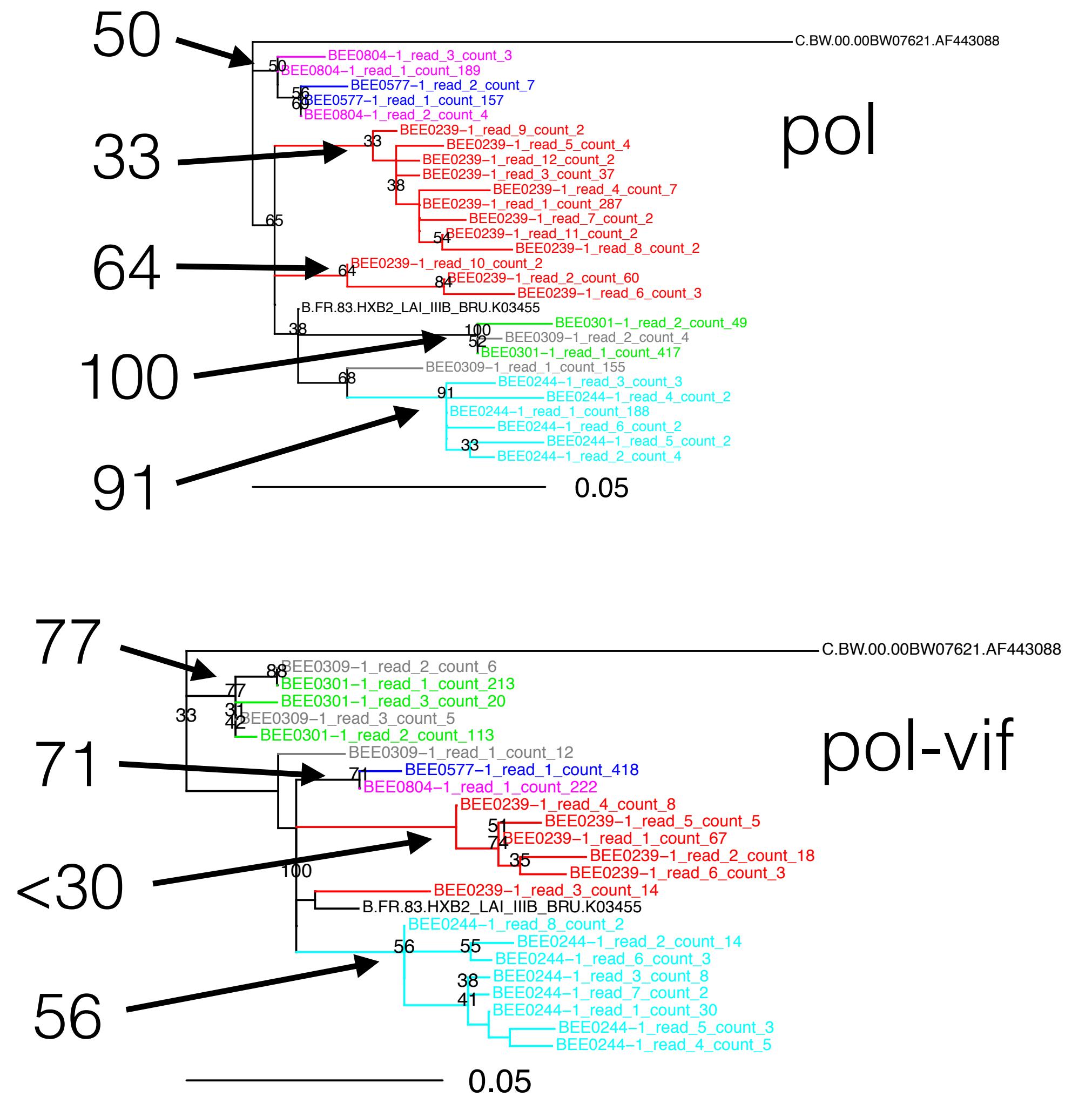
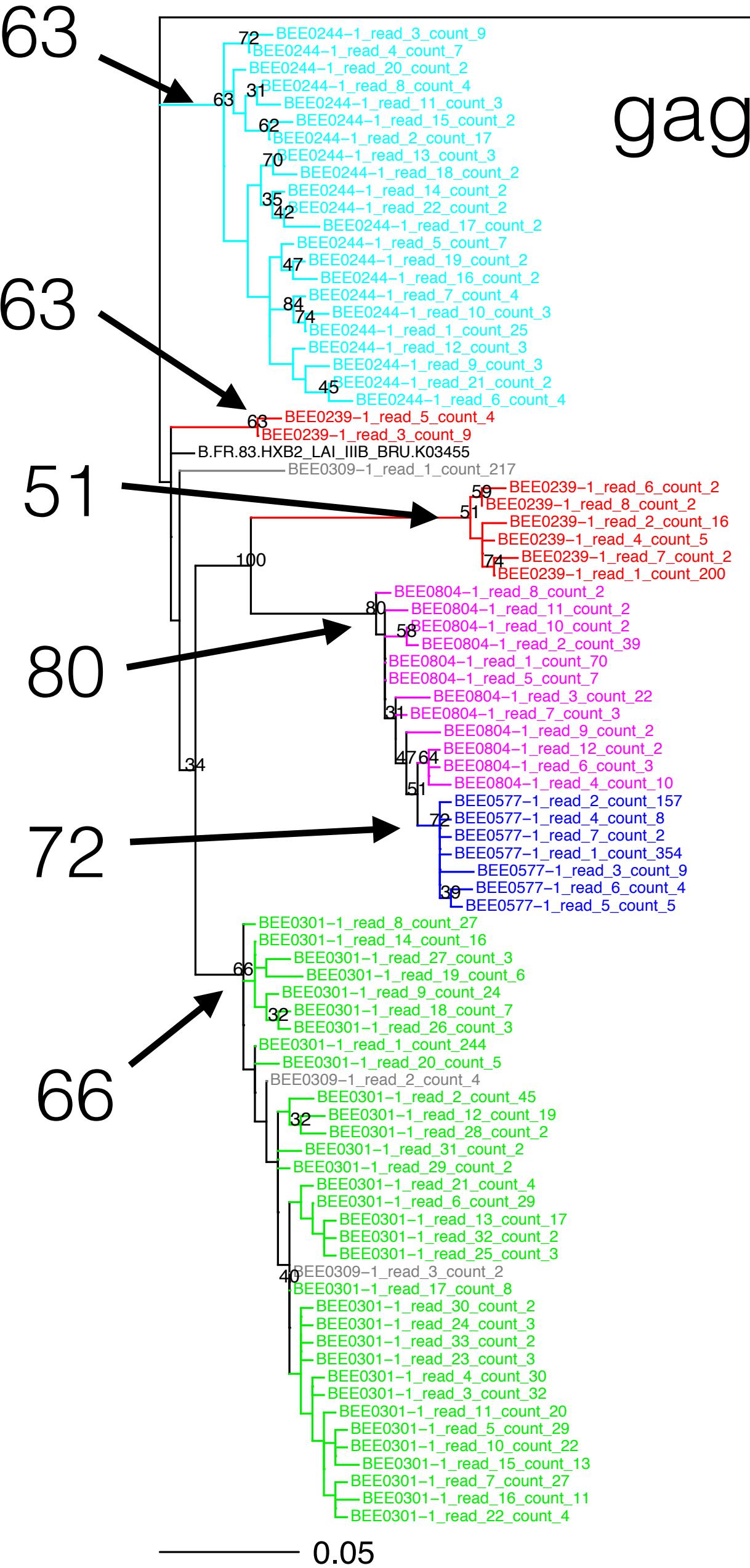
recipient?

donor?

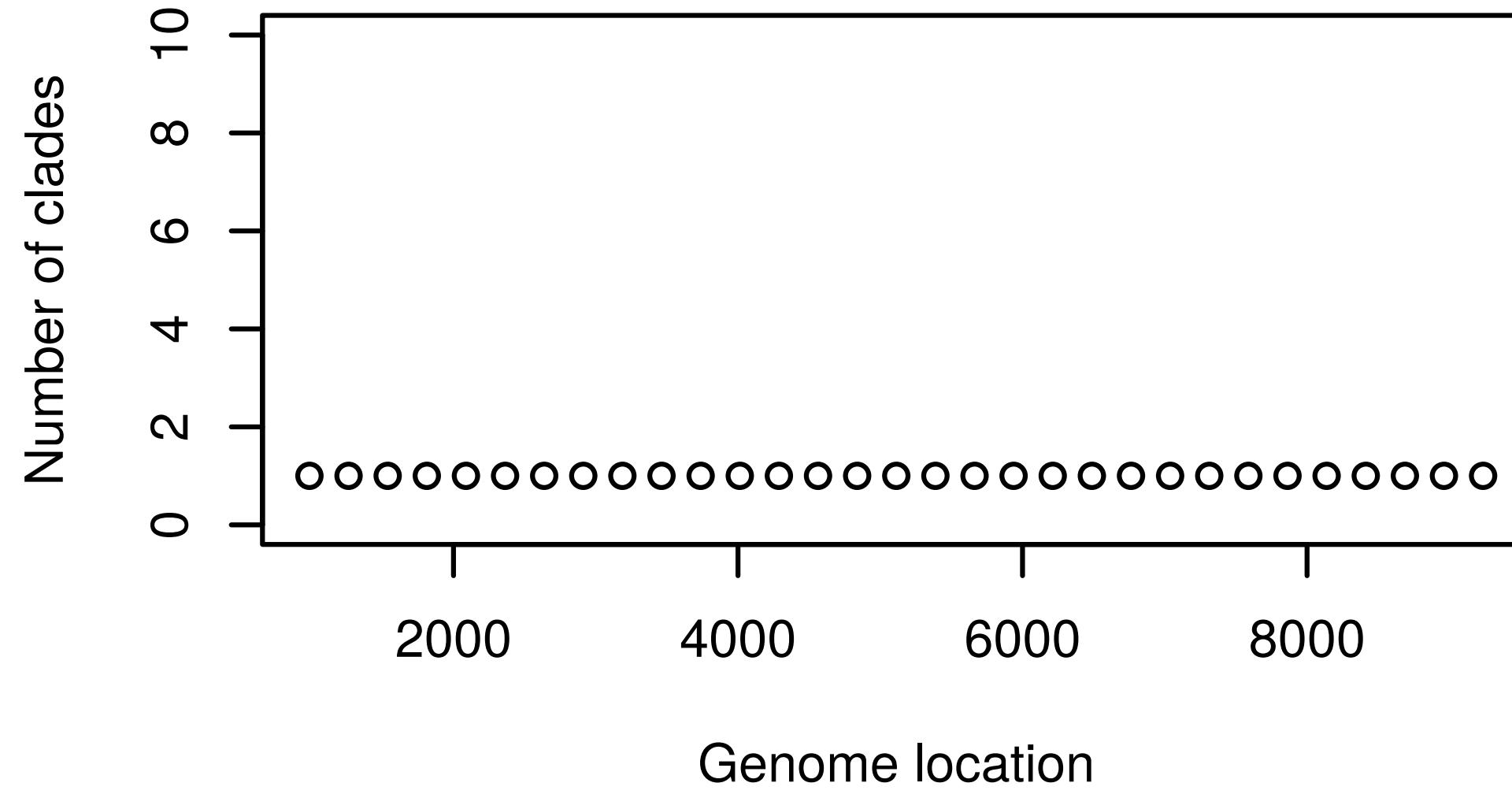
All trees excise the DRMs of
Johnson et al., Top. Anti. Med. 2011
and *Gatanaga et al., J. Biol. Chem. 2002* and the CTL-escape sites of
Carlson et al., J. Virology 2012

bootstraps (%)

values > 30 shown on nodes

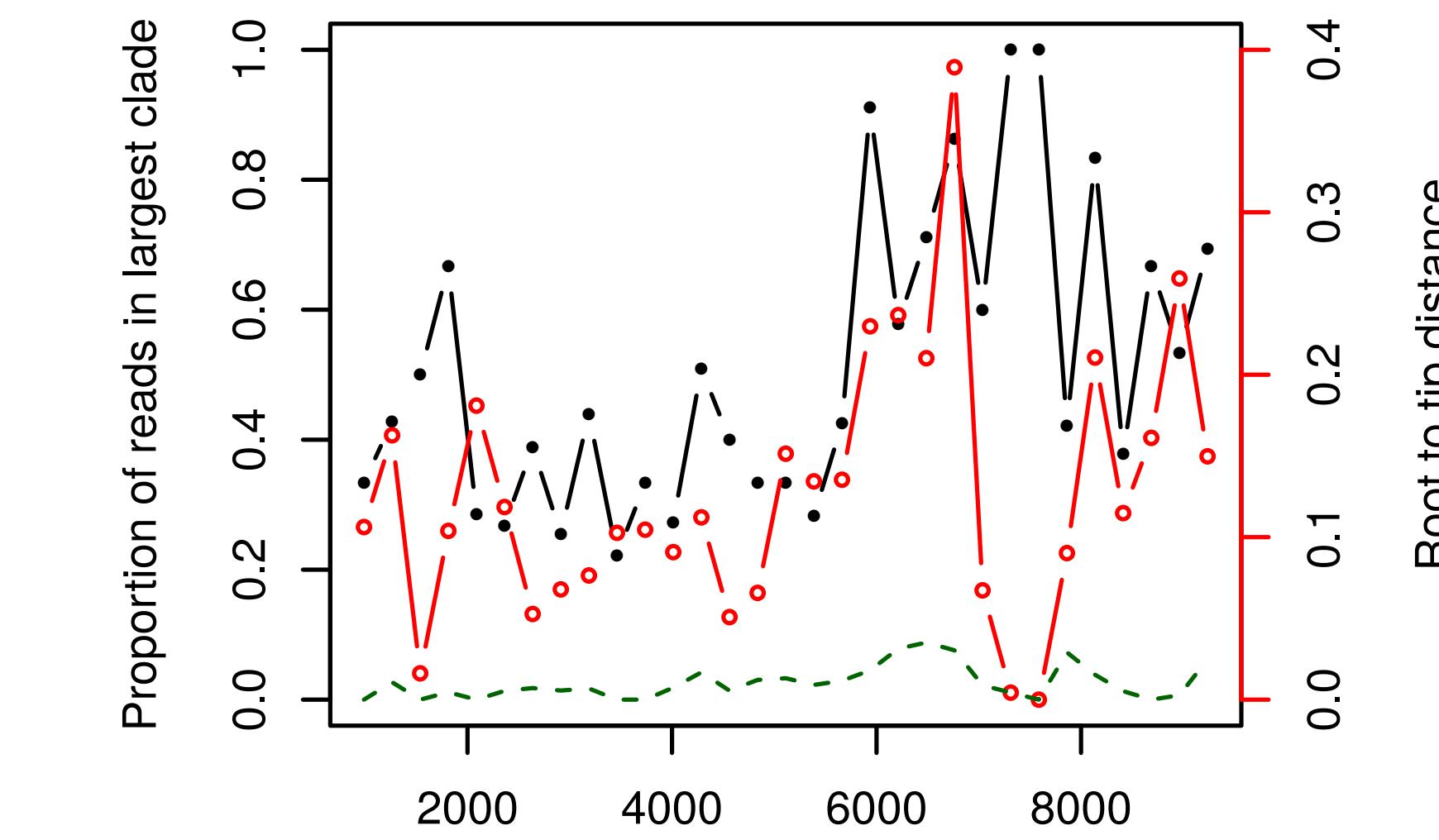
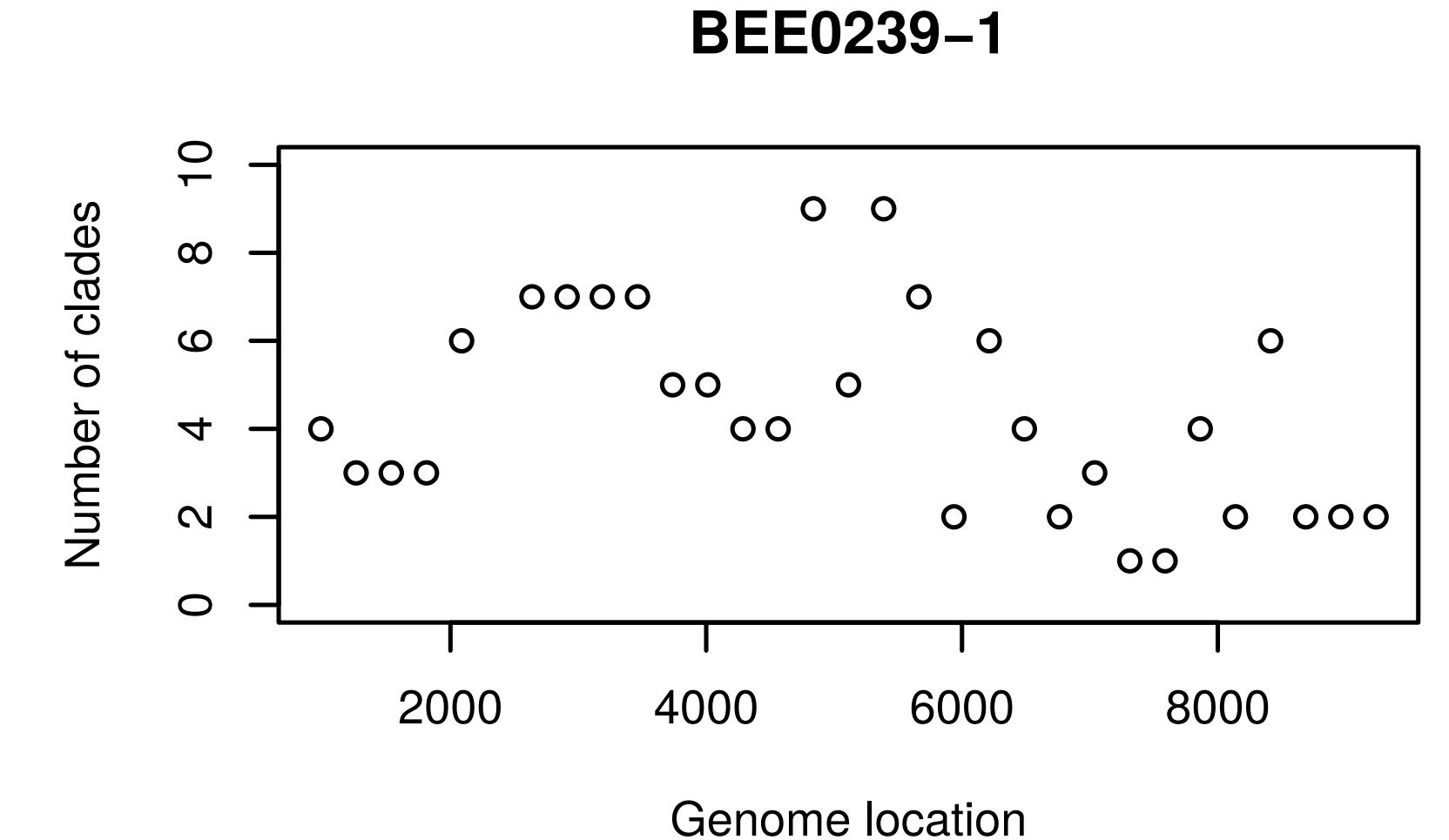


Per-patient summaries over the genome



single infection

clades = 1 throughout



dual infection

>1 clade, large root-tip, large proportion
of reads outside largest clade

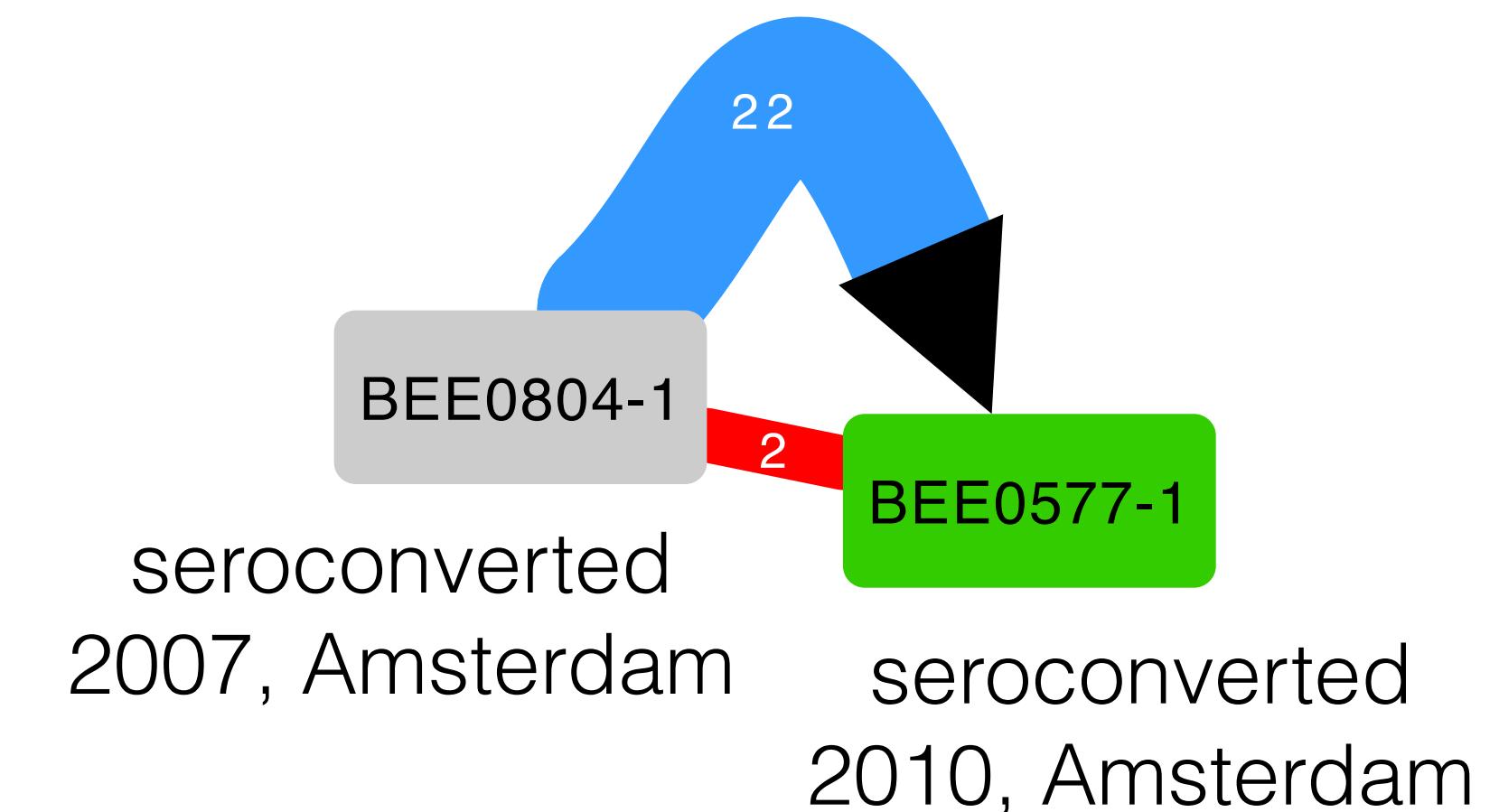
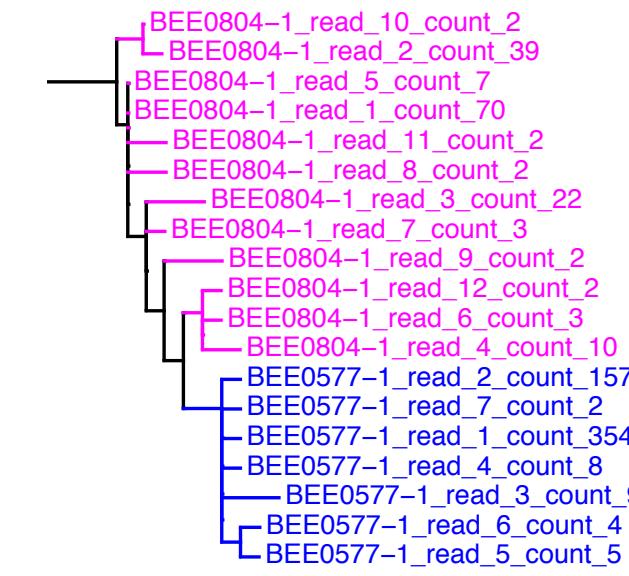
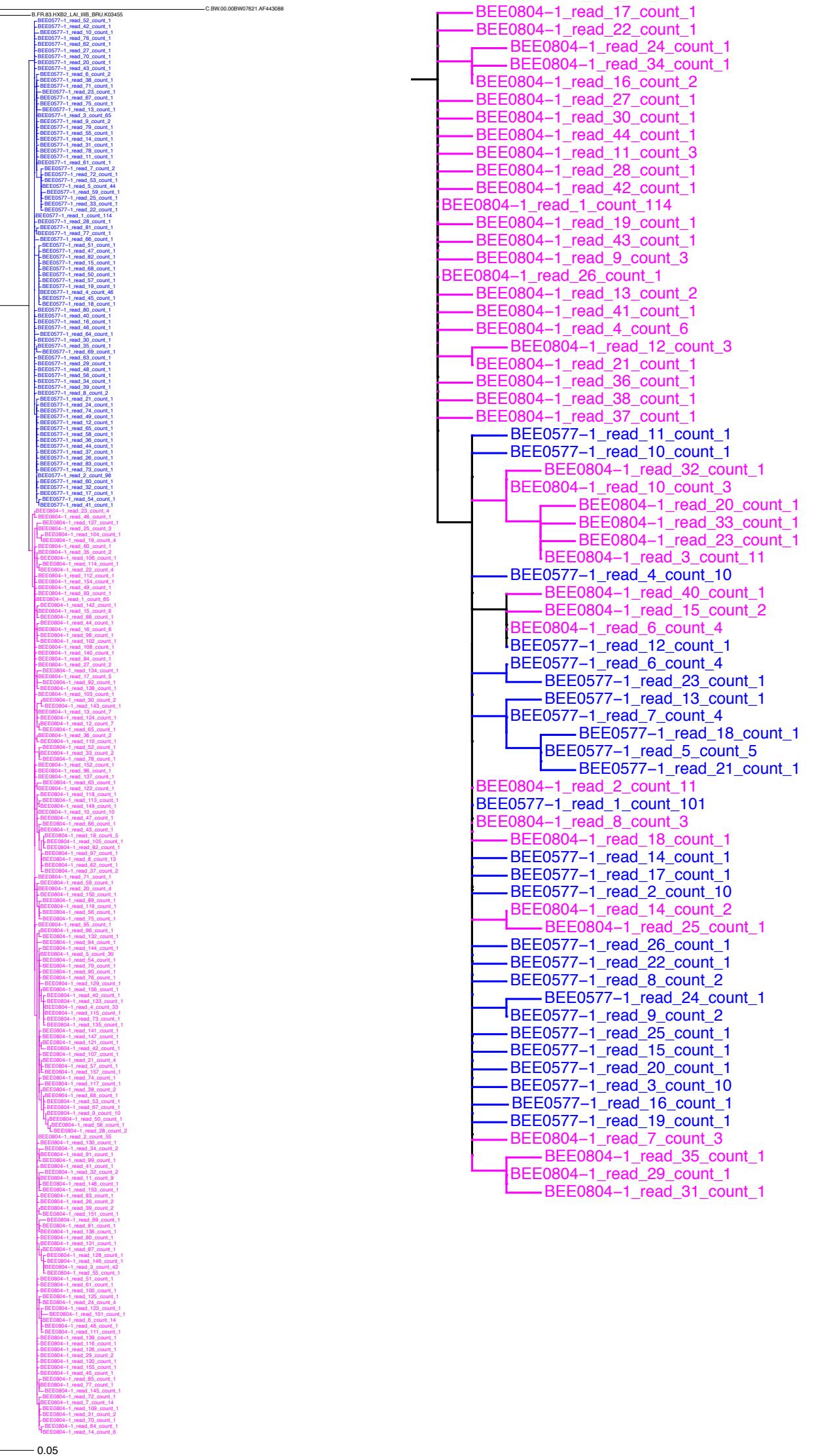
Two patients' reads can be...

separated, entangled, or clearly directed.

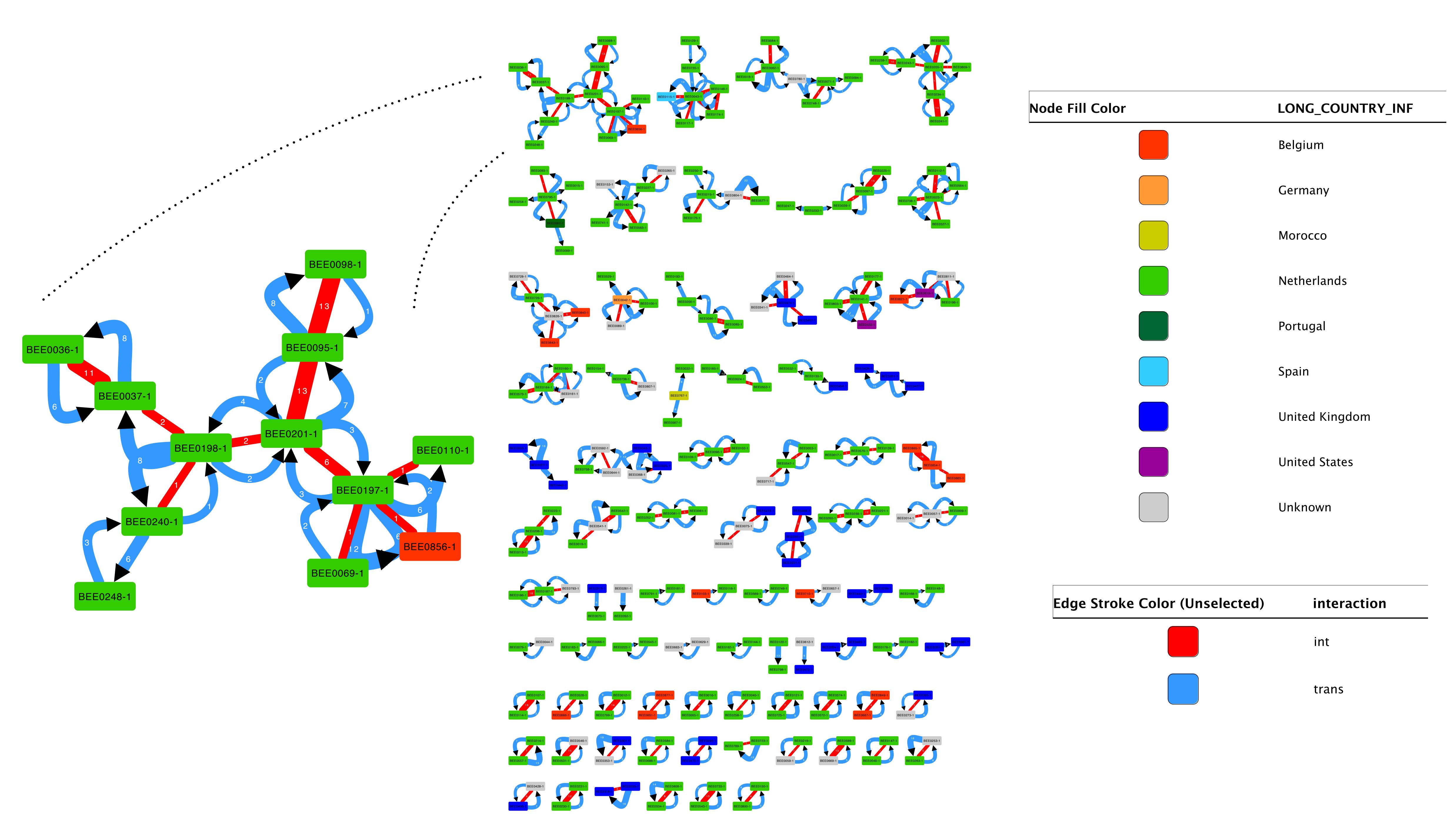
Count these over all windows (31 here):



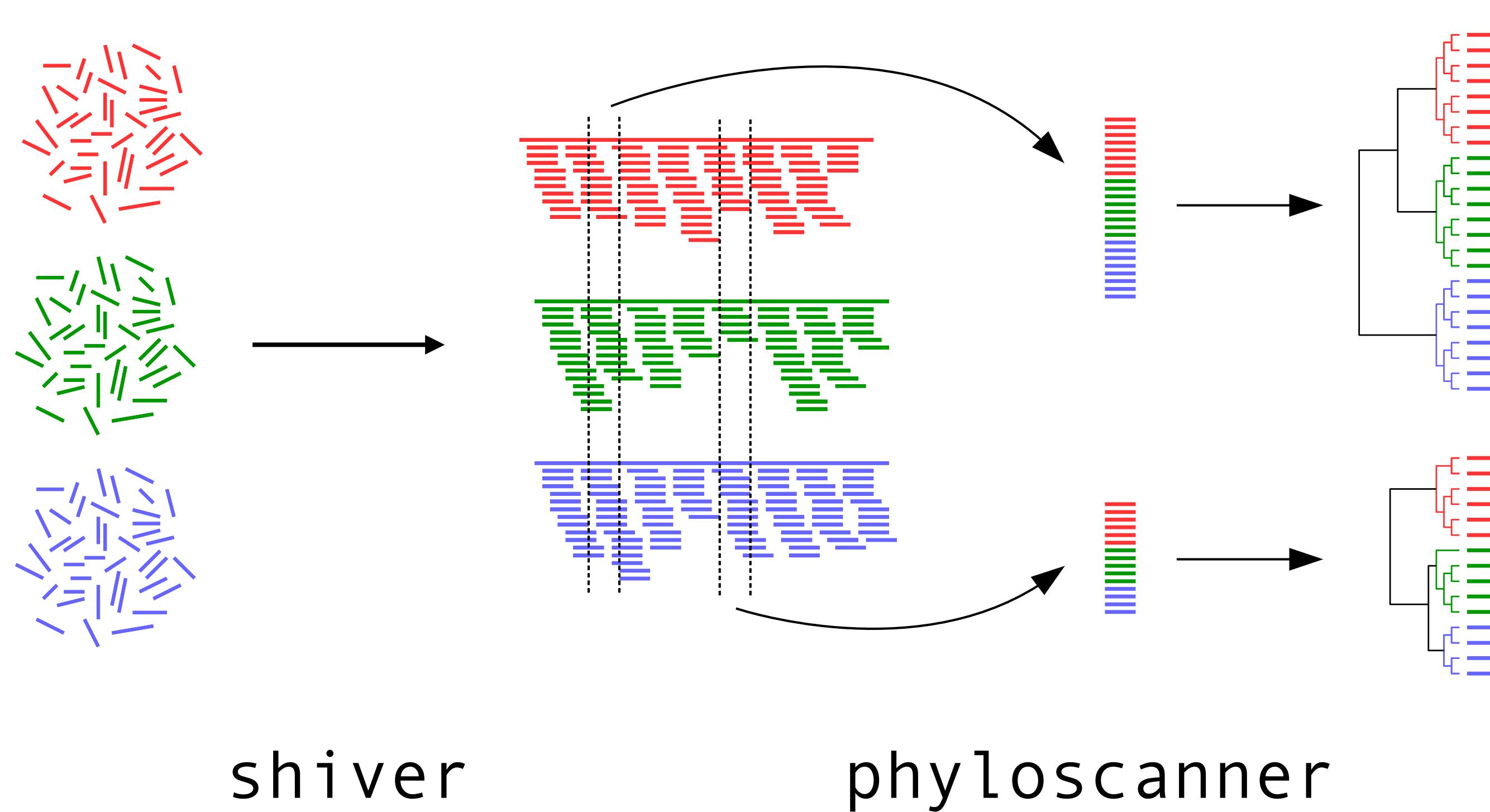
Matthew
Hall



See BEAST2 package *BEASTLIER* (github: *twoeventwo*)
c.f. Ypma et al., *Genetics* 2013;
Didelot et al., *Mol. Biol. Evol.* 2014;
Hall et al. *PLoS Comp. Biol.* 2015;
Kenah et al., arxiv 2016.
Romero-Severson et al. 2016



Summary



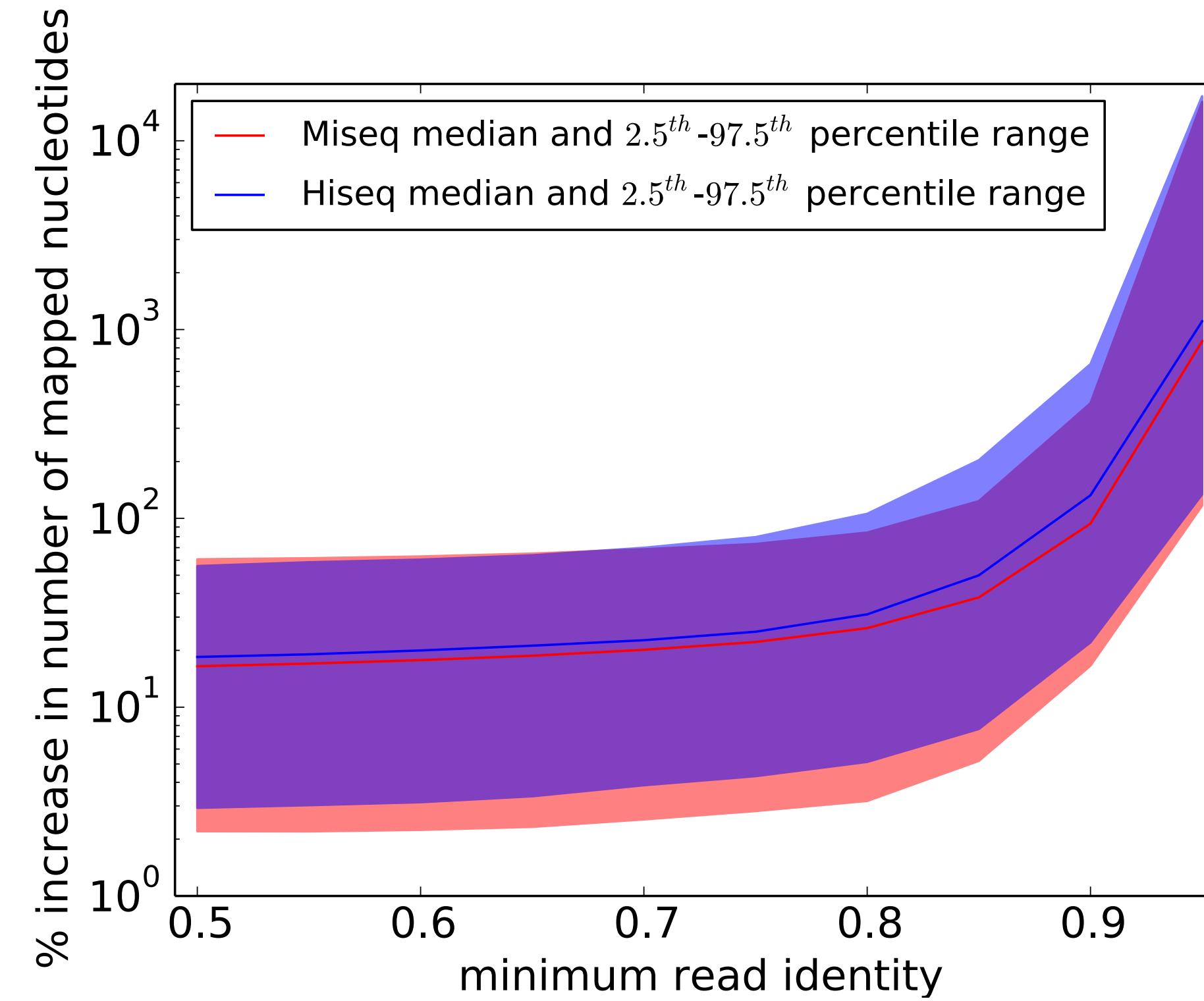
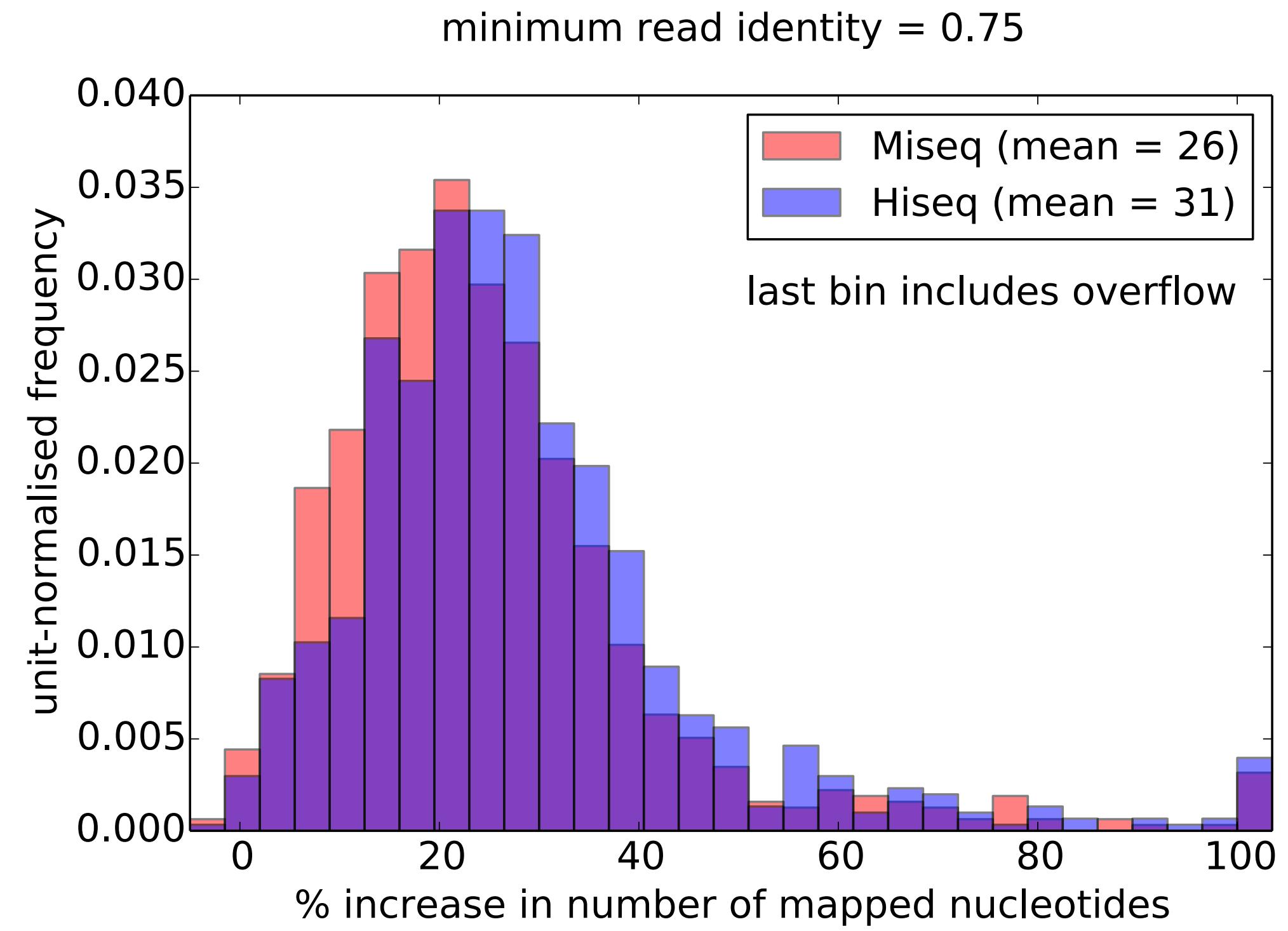
Analysis of phylogenetics between and within, along the genome, shows:

- **transmission**: one quasispecies as the descendant of another
- **dual infections**: two distinct quasispecies
- **recombination**: intermediates between two quasispecies
- **contamination**: repeated exact duplication of reads

Extra slides

Quantifying, not just arguing,
why you should use shiver

shiver maps more nucleotides

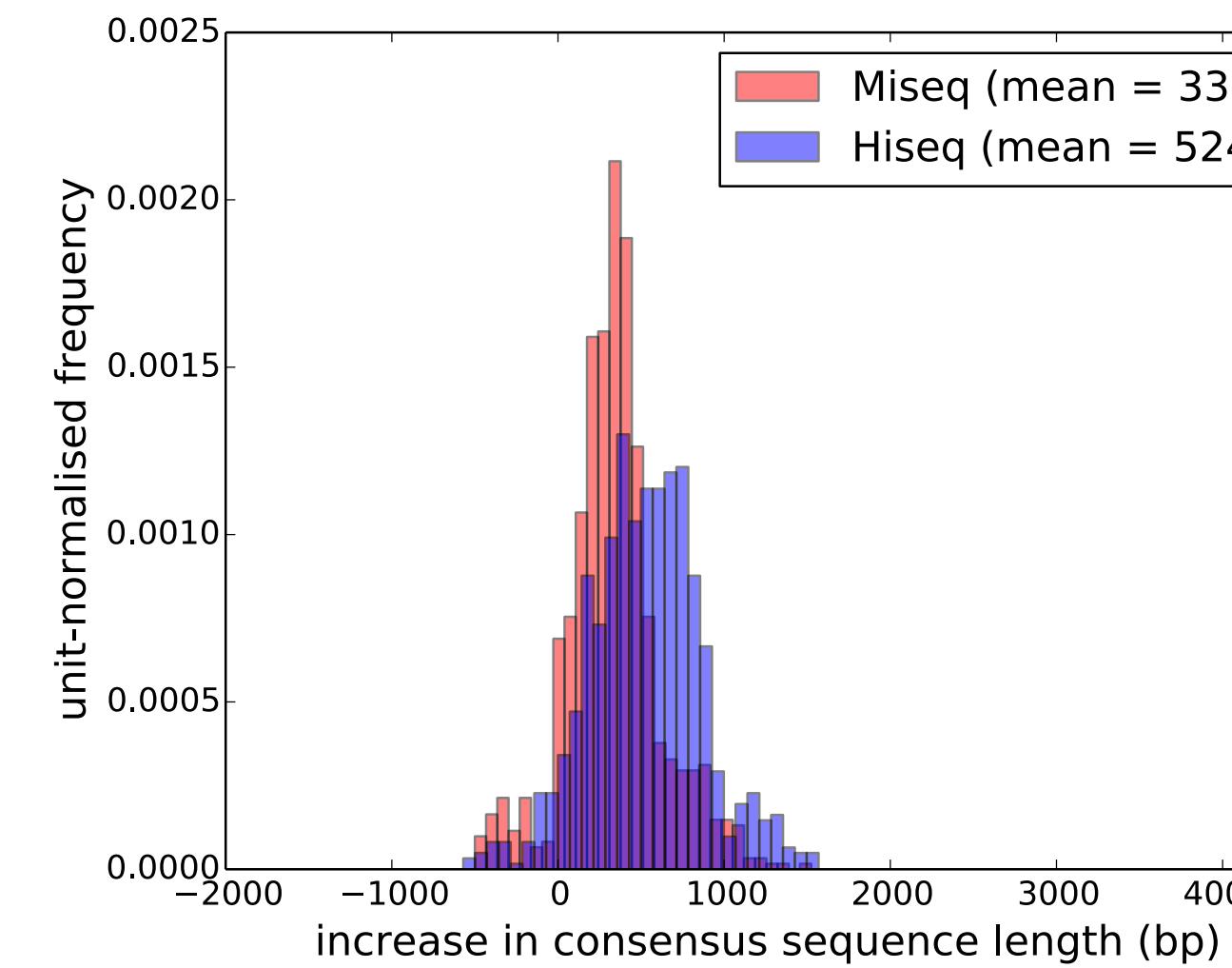


Map 904 Miseq and 864 Hiseq samples first to the closest of 160 existing references (chosen with Kraken, often HXB2) then to the custom reference constructed by shiver. Compare the number of nucleotides mapped. Repeat for different thresholds on read identity - the fraction of a read that agrees with the reference.

shiver reconstructs more of the genome

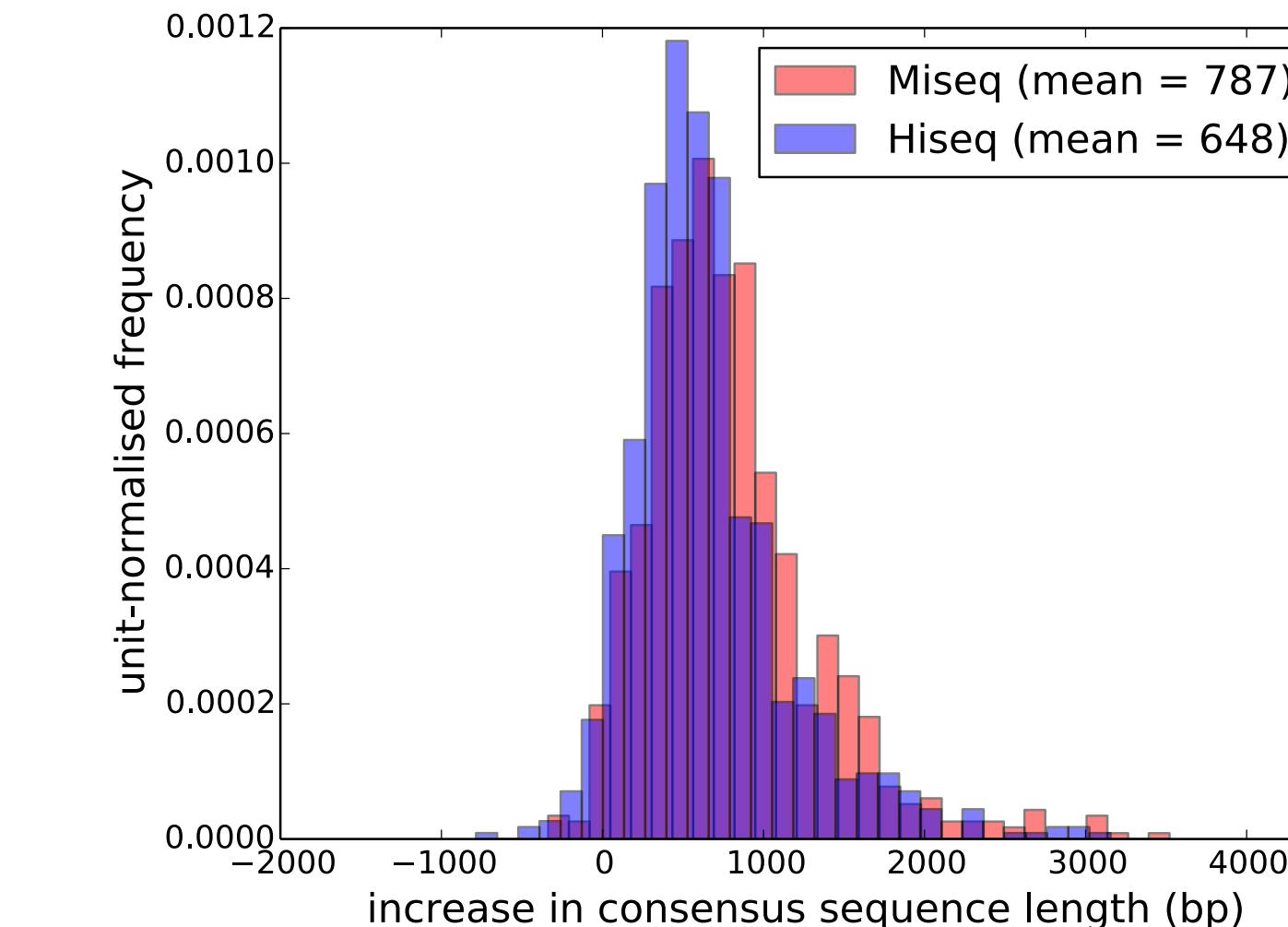
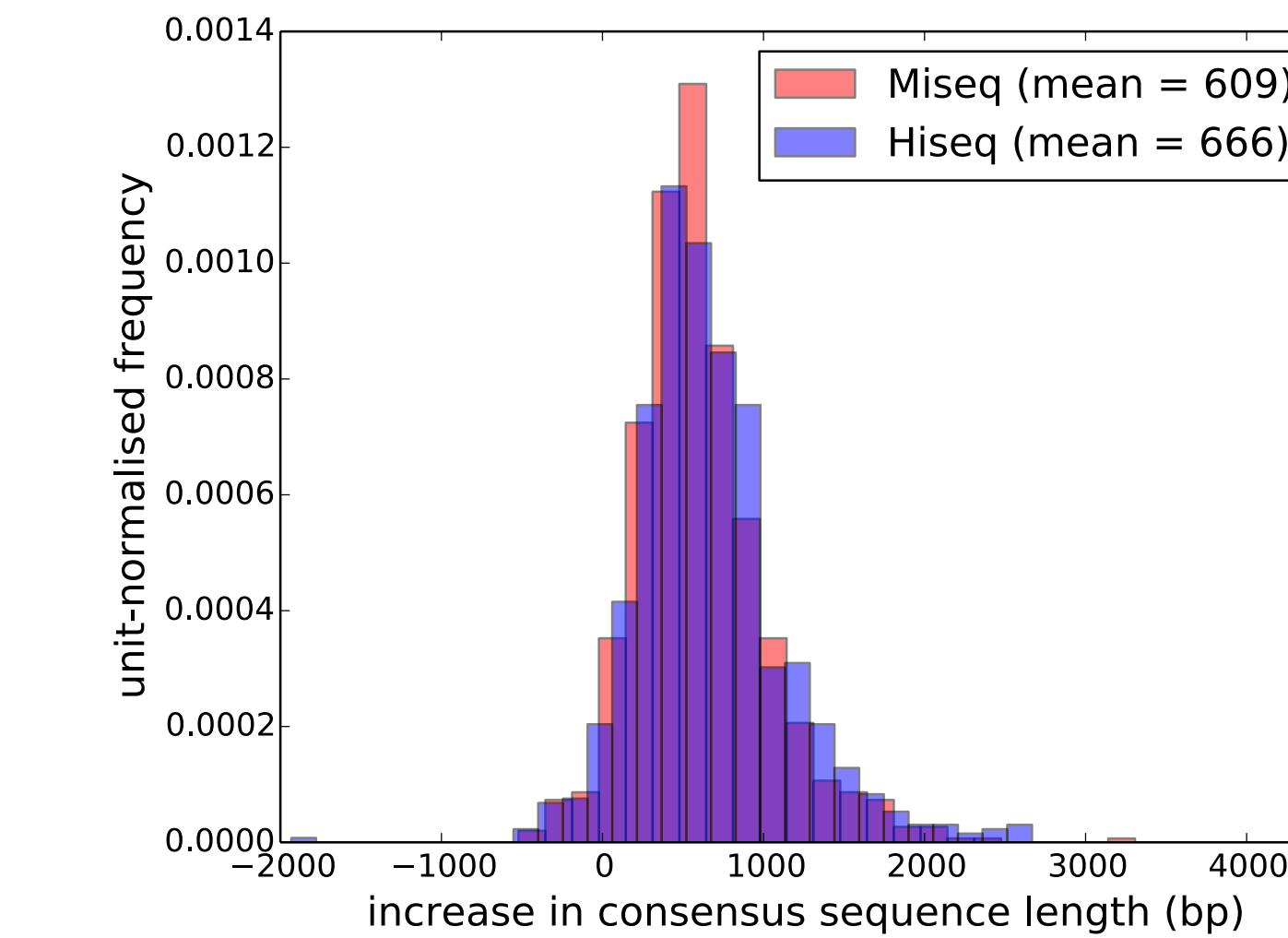
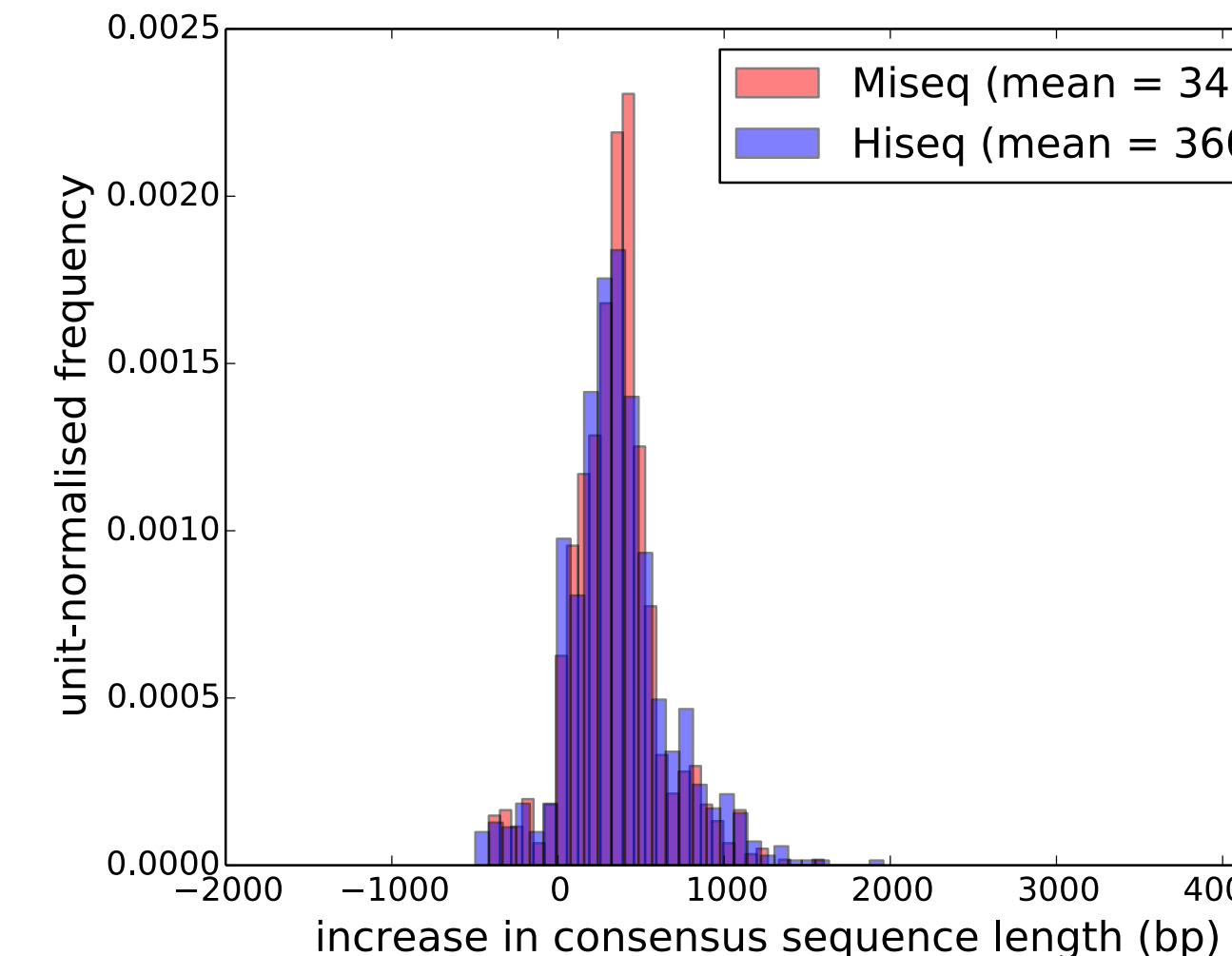
lenient on read identity (>50%)

lenient on
coverage
(>10)

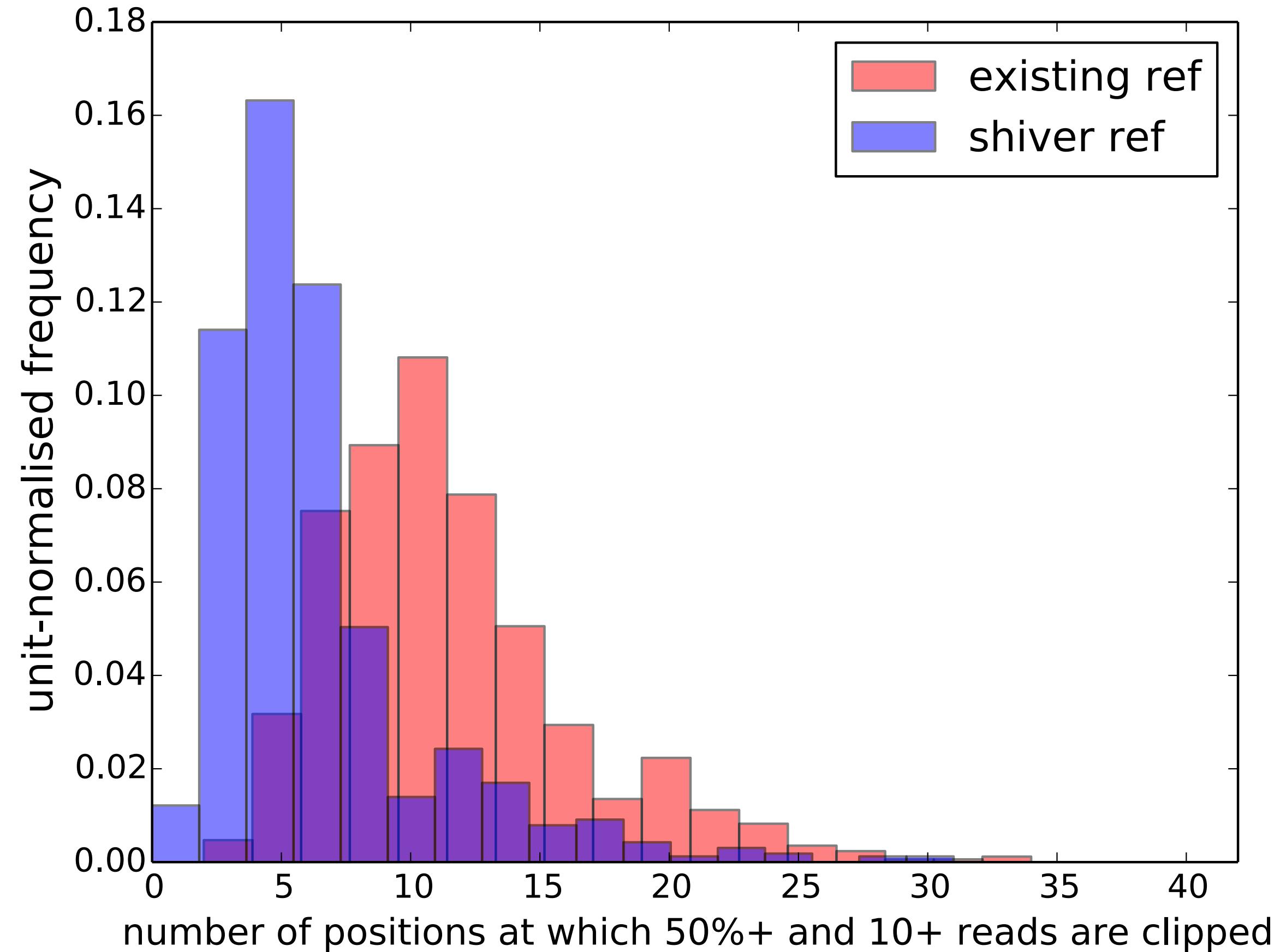


strict on read identity (>85%)

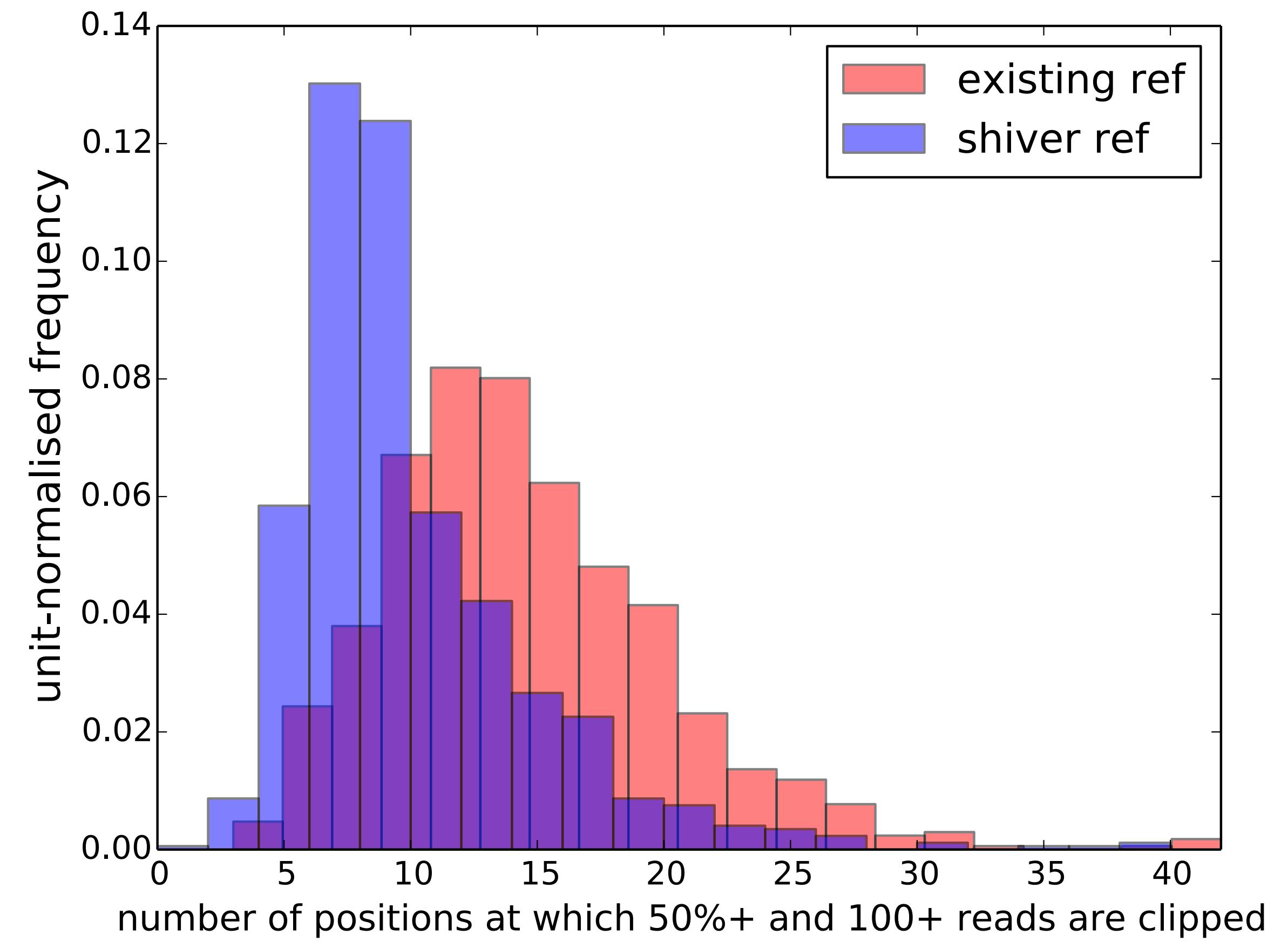
strict on
coverage



shiver decreases bias in the consensus



Miseq



Hiseq