



UK Health
Security
Agency

Consensus & Variant Calling

Dr. David Bibby,
Genomics and Clinical Virology,
21st March 2024

Overview

Consensus & Variant – what are they?

- How to build a consensus and define variants

What can be done with variant analysis - examples

- Features: Drug resistance, epitopes, species/strain identification
- Phyletics: Linkage, Dual infections, Transmission, Quasispecies reconstruction

Technical pitfalls - examples

- Virus – Laboratory – Bioinformatics

Validation, validation, validation

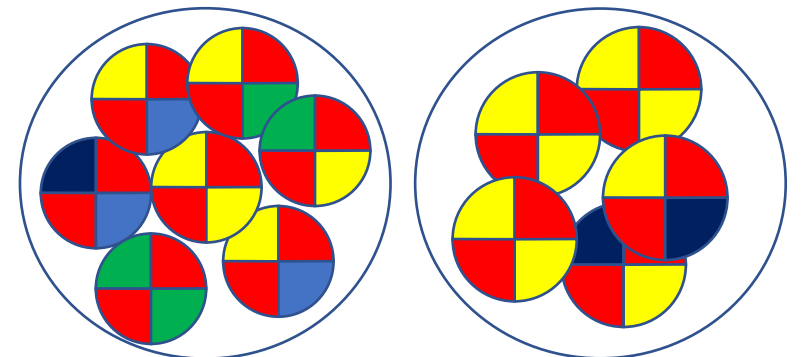
- Reproducibility
- Standardised materials
- EQA schemes
- Clinical validation?

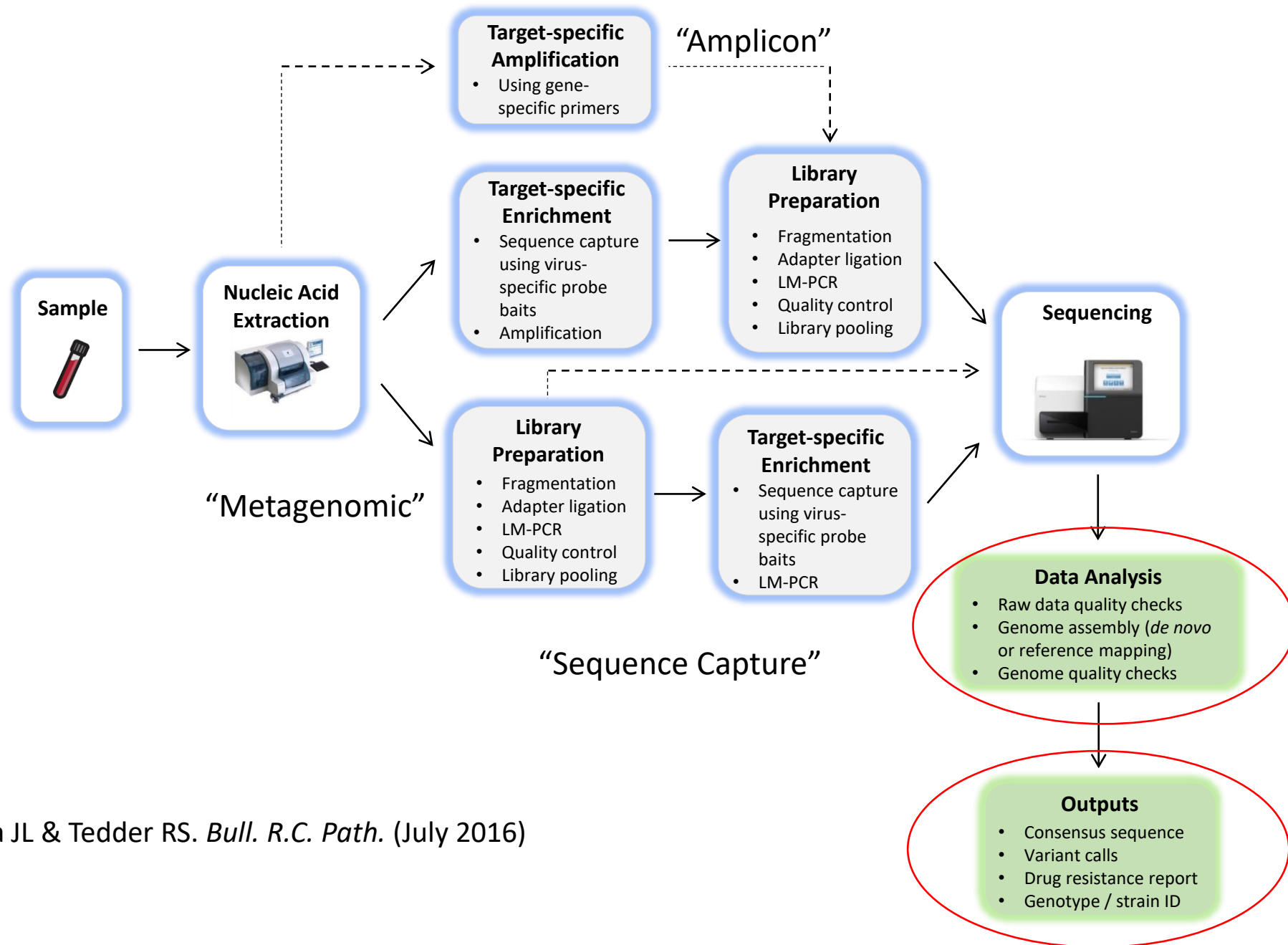
Consensus and Variant

Consensus: "The sequence of the most frequent nucleotides at each position"

Variants: "Differences between a test sequence and a reference"

- Viruses often exist in populations of related sequences, i.e. 'quasi-species'
- A consensus of a viral sequence may often contain mixed bases, incorporating the variants above a set frequency
 - e.g. 15-20% to mimic Sanger detection





Mbisa JL & Tedder RS. *Bull. R.C. Path.* (July 2016)

How to build a consensus

Sequencer output:

- Giant file containing all sequences from all samples (and controls)
- Each read has an adapter sequence added during the sample library prep
- These enable the reads to be 'binned' according to sample ID

The bins are FASTQ files

- Paired end – Forward and Reverse (often R1 & R2 files)
- Adapters usually trimmed before further analysis

Reads are e.g. Reference Mapped → SAM file

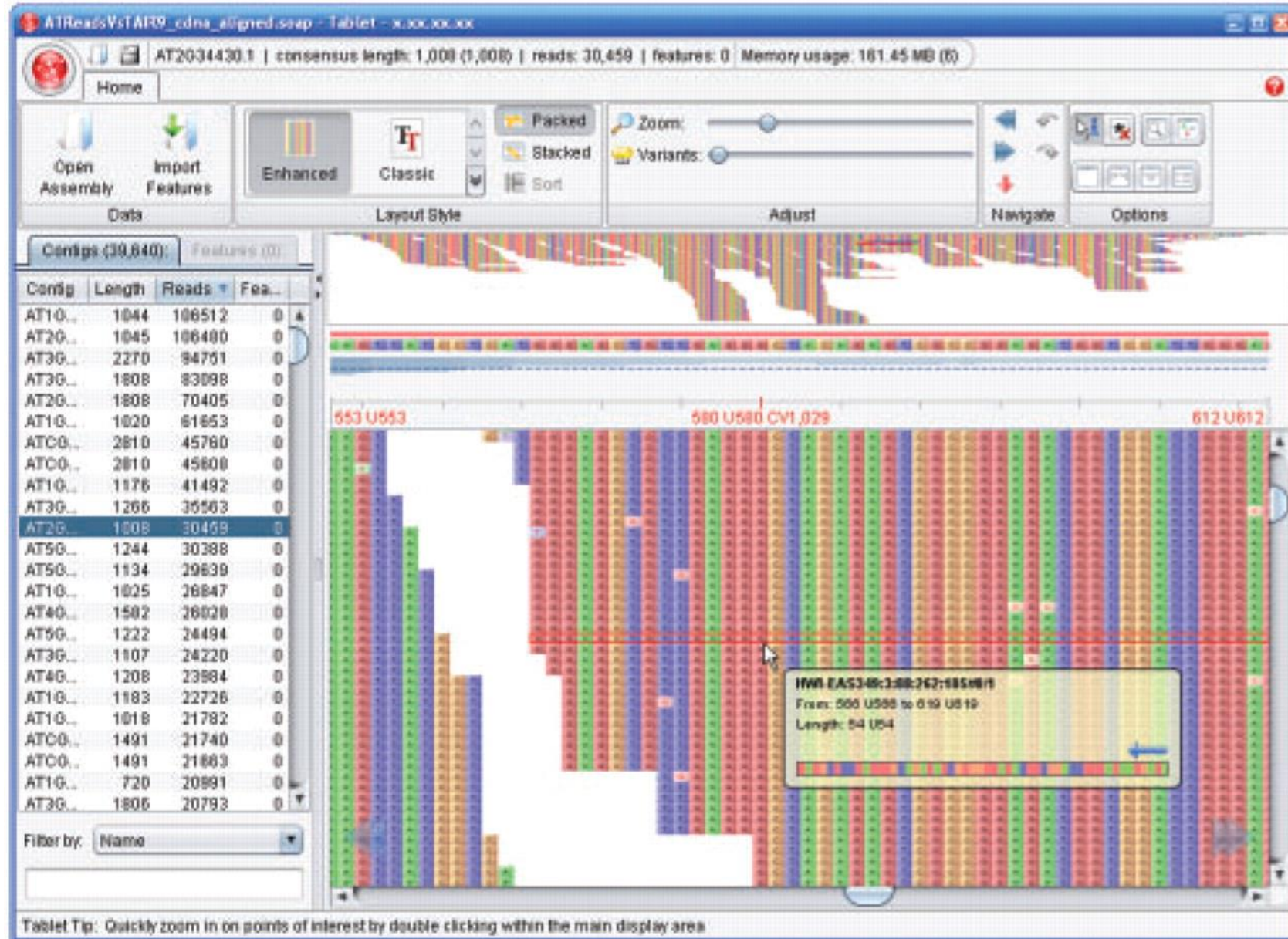
```
bwa mem my_virus_ref.fasta sample1_R1.fastq sample1_R2.fastq > sample1.sam
```

- SAM files are converted to BAM files

```
samtools view -Sbhu sample1.sam | samtools sort > sample1.bam
```

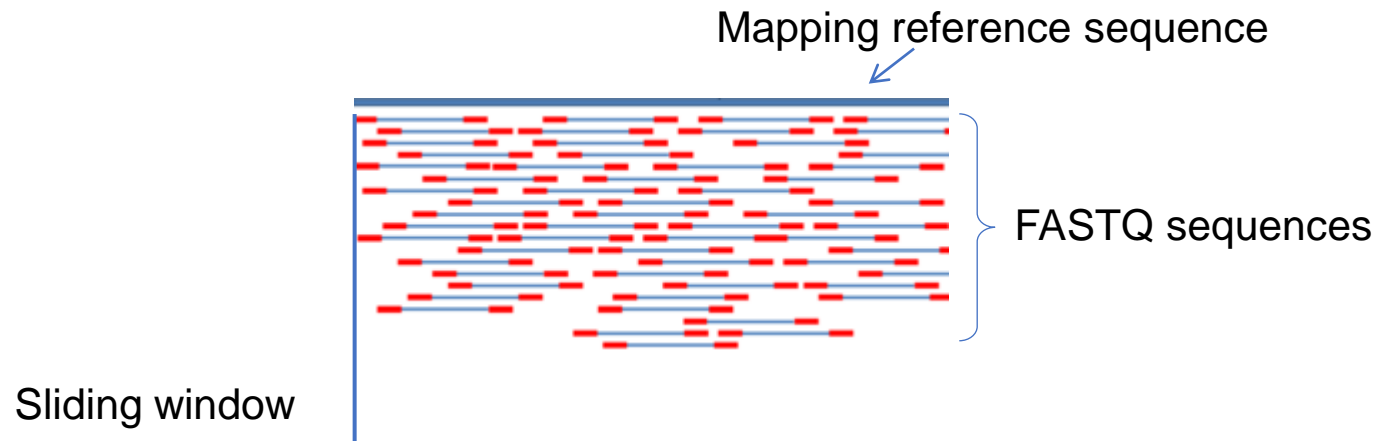
```
samtools index sample1.bam
```

- These can be viewed in Tablet / Genious / IGV etc.



How to build a consensus

- Several tools are available to derive consensus sequences from a SAM/BAM file:
 mpileup
 V-Phaser
 QuasiBAM
- Slide along the sequence, interrogating all reads covering each position



How to build a consensus

Considerations

- Quality of bases within a read
 - Phred score exclusion thresholds (usually 30, sometimes 20)
 - Quality of the read mapping
 - Map Quality exclusion thresholds
 - i.e. where the degree of homology to the reference sequence is low
 - Are these contaminants or rare sequence motif(s)?
 - Handling of insertions / deletions / variants
- Can be very dependent upon choice of mapping software
Its parameters, and/or reference sequence**

Variant calling – mpileup

mpileup (samtools) <http://www.htslib.org/doc/samtools.html>

1. Iterates through each position in a reference (i.e. one row per position)
2. Iterates through each read covering that position and adds a match type symbol...

| | | |
|--------|---------------------|---|
| . | , | Match to reference (forward & reverse respectively) |
| ^ | \$ | Start and finish of a read respectively |
| ACGTN, | acgtn | Mismatch to reference (fwd & rev respectively) |
| [+-] | [0-9]+[ACGTNacgtn]+ | Insertions / deletions |

3. ...and a Quality symbol (Phred Score)

Variant calling – mpileup

```
1 215906528 G 21 ,,,,,,,,,,,,,,,,,,,,,, ;=?./:??>>;=7?>>@A?==:
1 215906529 A 18 ,,,,,,,,,,,,,,,,,,,,,, D>AA:@A>9>?;;?>>@=
[...]
```

| | | | | | |
|---|-----------|---|----|--------------------------|--------------------|
| 1 | 215906547 | C | 15 | gGg\$,GggGG,,.... | <;80;><9=86=C>= |
| 1 | 215906548 | G | 19 | c\$,ccC.,,,,,,,,,,,,,,^. | ;58610=7=>75=7<46; |

```
[...]
```

| | | | | | |
|---|-----------|---|----|---------------------|-----------------|
| 1 | 215906555 | G | 16 | .\$aaaaaA.AAAaAAA^A | 2@>?8?;<:335?:> |
|---|-----------|---|----|---------------------|-----------------|

apprize.info

Variant calling – VCF

VCF “Variant Call Format”

<http://vcftools.sourceforge.net>

- Developed for human genome annotations by 1,000 Genomes project
- Useful for sparse variation in long, multi-chromosome genomes
- Lists variations from a reference in a tabular format
 - One row per variant
 - (At least) 8 columns:

| CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO |
|-------|-----|----|-----|-----|------|--------|------|
|-------|-----|----|-----|-----|------|--------|------|

CHROM = Chromosome

POS = Position

REF = Reference

ALT = Alternative (variant)

Variant calling – VCF

```
##fileformat=VCFv4.3
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:.,.
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ 0|0:49:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTC G,GTCT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP 0/1:35:4 0/2:17:2 1/1:40:3
```

samtools.github.io

Variant calling – V-phaser

V-Phaser & V-Phaser 2

- Developed by the Broad Institute
- Considers read position, strand bias, quality scores, dinucleotide frequency, forward & reverse read, and phasing (linkage)
- Reports variant frequency and absolute read numbers by forward & reverse read
- Similar to VCF, but for viral populations
 - One row per variant
 - Seven columns:

| Ref_Pos | Var | Cons | Strd_bias_pval | Type | Var_perc | SNP_or_LP_Profile |
|---------|-----|------|----------------|------|----------|-------------------|
|---------|-----|------|----------------|------|----------|-------------------|

Macalalad AR *et al.* PLoS Computational Biology 2012 8(3):e1002417

Yang X, *et al.* BMC Genomics 2013 14:674

Variant calling – V-phaser

| # | Ref_Pos | Var | Cons | Strd_bias | Type | Var_perc | SNP_or_LP_Profile | | |
|----------|---------|-----|------|-----------|------|----------|-------------------|---------|---------|
| # | ----- | | | | | | | | |
| | 1448 | G | C | 0.2919 | snp | 7.477 | C:53:46 | G:2:6 | |
| | 1462 | T | A | 1 | snp | 6.604 | A:49:47 | G:1:2 | T:4:3 |
| | 1476 | C | T | 1 | snp | 7.273 | A:0:1 | C:4:4 | T:50:51 |
| | 1480 | T | C | 0.6589 | snp | 11.21 | A:0:1 | C:45:47 | G:1:1 |
| | 1481 | A | G | 1 | snp | 7.273 | A:4:4 | C:1:1 | G:49:51 |
| | 1488 | C | T | 0.8233 | snp | 9.91 | C:5:6 | G:3:0 | T:46:51 |
| | 1568 | T | C | 1 | snp | 7.865 | C:37:45 | T:4:3 | |
| | 1872 | A | G | 1 | snp | 8.14 | A:3:4 | G:37:42 | |
| | 3473 | A | G | 1 | snp | 2.857 | A:2:1 | G:56:46 | |
| | 3481 | T | C | 1 | snp | 2.913 | C:56:44 | T:2:1 | |
| | 3511 | C | T | 1 | snp | 2.885 | C:2:1 | T:52:49 | |
| | 3514 | T | C | 1 | snp | 2.83 | A:0:1 | C:52:50 | T:2:1 |
| | 3527 | A | T | 1 | snp | 3.061 | A:2:1 | T:49:46 | |
| | 3530 | G | A | 1 | snp | 3.125 | A:46:47 | G:2:1 | |
| | 3532 | C | T | 1 | snp | 3.125 | C:2:1 | T:46:47 | |
| | 3559 | G | C | 1 | snp | 3.75 | C:34:43 | G:2:1 | |
| | 3570 | C | A | 1.127 | snp | 4.878 | A:35:43 | C:2:2 | |
| | 3574 | T | C | 1.127 | snp | 4.762 | A:0:1 | C:36:43 | T:2:2 |
| | 3577 | C | T | 1.127 | snp | 4.878 | C:2:2 | T:37:41 | |
| | 3592 | T | C | 1.127 | snp | 4.819 | C:40:39 | T:2:2 | |
| | 3601 | A | G | 1 | snp | 3.614 | A:1:2 | G:39:41 | |
| | 3605 | C | A | 1 | snp | 3.704 | A:37:41 | C:1:2 | |
| | 3616 | G | A | 0.9257 | snp | 14.29 | A:31:35 | G:6:5 | |
| | 6583 | T | C | 0.2925 | snp | 8.654 | C:51:44 | T:7:2 | |
| | 6882 | T | A | 0.8081 | snp | 21.21 | A:32:46 | T:8:13 | |
| | 6895 | G | A | 0.7799 | snp | 39.81 | A:24:37 | G:17:24 | T:0:1 |
| | 7150 | T | C | 1.004 | snp | 9.639 | C:31:42 | G:2:0 | T:5:3 |
| | 7387 | G | A | 0.5027 | snp | 10.71 | A:35:40 | G:6:3 | |
| | 9176 | A | G | 0.7282 | snp | 24.29 | A:5:12 | G:18:35 | |
| # | ----- | | | | | | | | |
| # Summar | LPV: 0 | | | | | | | | |

Variant calling – QuasiBAM

- **QuasiBAM** (UKHSA)
- Produces a table of nucleotide & codon frequencies for an entire reference
- One row per nucleotide position, i.e. more like mpileup
- 14 Columns:
 - 1-3 **Position, Reference nucleotide, Depth**
 - 4-9 **A / C / G / T / Gap / Insertion frequencies**
 - 10 **Insertion sequences and their frequencies**
 - 11-12 **Reference Amino Acid, Depth**
 - 13-14 **Codon / Amino Acid frequencies**
- Can be parameterized
 - Strandedness
 - Gap-masking
 - Primer-mediated error filtering

Variant calling – QuasiBAM

| Pos | Ref_N | Depth | A | C | G | T | Gap | Ins | I_Desc | Ref_AA | AA_depth | Cod | AA |
|------|-------|-------|--------|--------|--------|--------|-------|-------|-------------|--------|----------|--|--|
| 4492 | C | 23097 | 0 | 99.753 | 0 | 0 | 0 | | | P | 22244 | CCC:21813:98.062 CCT:295:1.326 | P:22120:99.443 |
| 4493 | C | 23048 | 0 | 99.683 | 0 | 0 | 0 | | | P | 22465 | CCT:22064:98.215 CTT:302:1.344 | P:22081:98.291 L:302:1.344 |
| 4494 | C | 23623 | 0 | 98.650 | 0 | 1.300 | 0 | | | L | 22991 | CTG:22574:98.186 TTG:306:1.331 | L:22963:99.878 |
| 4495 | T | 23607 | 0 | 0 | 0 | 99.924 | 0 | | | C | 22904 | TGC:22744:99.301 | C:22796:99.528 |
| 4496 | G | 23547 | 0 | 0 | 99.643 | 0 | 0 | | | A | 22876 | GCT:22701:99.235 | A:22729:99.357 |
| 4497 | C | 23323 | 0 | 99.734 | 0 | 0 | 0 | | | L | 22800 | CTT:22587:99.066 | L:22708:99.596 |
| 4498 | T | 23511 | 0 | 0 | 0 | 99.860 | 0 | | | L | 22928 | TTA:22635:98.722 | L:22671:98.879 |
| 4499 | T | 23515 | 0 | 0 | 0 | 99.468 | 0 | | | * | 22688 | TAA:19794:87.244 TAG:2636:11.618 | *:22442:98.916 |
| 4500 | A | 23389 | 99.376 | 0 | 0 | 0 | 0 | | | K | 22584 | AAG:19401:85.906 AGG:2630:11.645 AA:-278:1.231 | K:19526:86.459 R:2637:11.676 X:395:1.749 |
| 4501 | A | 23320 | 87.414 | 0 | 12.543 | 0 | 0 | | | R | 22621 | AGG:19329:85.447 GGG:2820:12.466 A-G:269:1.189 | R:19369:85.624 G:2824:12.484 X:275:1.216 |
| 4502 | G | 24246 | 0 | 0 | 98.189 | 0 | 1.192 | | | G | 23482 | GGG:22955:97.756 -GG:276:1.175 | G:22998:97.939 X:281:1.197 |
| 4503 | G | 24131 | 0 | 0 | 99.731 | 0 | 0 | | | G | 23371 | GGG:23209:99.307 | G:23264:99.542 |
| 4504 | G | 24352 | 0 | 0 | 99.782 | 0 | 0 | | | G | 23454 | GGG:23275:99.237 | G:23339:99.51 |
| 4505 | G | 24122 | 0 | 0 | 99.718 | 0 | 0 | | | G | 23214 | GGG:22957:98.893 | G:23083:99.436 |
| 4506 | G | 24106 | 0 | 0 | 99.722 | 0 | 0 | | | G | 22025 | GGA:21713:98.583 | G:21853:99.219 |
| 4507 | G | 23894 | 0 | 0 | 99.456 | 0 | 0 | | | E | 21459 | GAA:21207:98.826 | E:21208:98.83 |
| 4508 | A | 22601 | 99.345 | 0 | 0 | 0 | 0 | | | K | 21422 | AAG:21237:99.136 | K:21261:99.248 |
| 4509 | A | 22591 | 99.708 | 0 | 0 | 0 | 0 | | | R | 21844 | AGG:17994:82.375 AGA:3758:17.204 | R:21762:99.625 |
| 4510 | G | 22841 | 0 | 0 | 99.764 | 0 | 0 | | | G | 21640 | GGC:17818:82.338 GAC:3733:17.25 | G:17845:82.463 D:3736:17.264 |
| 4511 | G | 22795 | 17.043 | 0 | 82.843 | 0 | 0 | | | A | 21368 | GCA:17613:82.427 ACA:3698:17.306 | A:17618:82.45 T:3698:17.306 |
| 4512 | C | 22402 | 0 | 99.853 | 0 | 0 | 0 | | | H | 21443 | CAC:21350:99.566 | H:21399:99.795 |
| 4513 | A | 22367 | 99.978 | 0 | 0 | 0 | 0 | | | T | 21502 | ACC:21381:99.437 | T:21443:99.726 |
| 4514 | C | 22774 | 0 | 99.750 | 0 | 0 | 0 | | | P | 21968 | CCT:21806:99.263 | P:21851:99.467 |
| 4515 | C | 22582 | 0 | 99.703 | 0 | 0 | 0 | | | L | 21981 | CTC:21824:99.286 | L:21871:99.5 |
| 4516 | T | 22841 | 0 | 0 | 0 | 99.781 | 0 | | | S | 22318 | TCA:22131:99.162 | S:22220:99.561 |
| 4517 | C | 22753 | 0 | 99.780 | 0 | 0 | 0 | | | H | 21775 | CAT:21427:98.402 | H:21430:98.416 |
| 4518 | A | 22928 | 99.603 | 0 | 0 | 0 | 0 | 1.396 | T:320:1.396 | I | 21820 | ATT:21364:97.91 | I:21521:98.63 |
| 4519 | T | 22553 | 0 | 0 | 0 | 99.056 | 0 | | | F | 22155 | TTT:21779:98.303 | F:21786:98.334 |
| 4520 | T | 22721 | 0 | 0 | 0 | 99.080 | 0 | | | F | 22039 | TTT:21816:98.988 | F:21829:99.047 |
| 4521 | T | 22804 | 0 | 0 | 0 | 99.961 | 0 | | | F | 22041 | TTT:21893:99.329 | F:22011:99.864 |
| 4522 | T | 22430 | 0 | 0 | 0 | 99.911 | 0 | | | F | 21689 | TTT:21545:99.336 | F:21547:99.345 |
| 4523 | T | 22551 | 0 | 0 | 0 | 99.463 | 0 | | | L | 21655 | TTG:21466:99.127 | L:21636:99.912 |
| 4524 | T | 22159 | 0 | 0 | 0 | 99.973 | 0 | | | C | 21555 | TGC:21400:99.281 | C:21480:99.652 |

Uses of variant analysis

Features

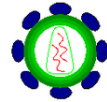
- Typing
- Resistance
- Epitopes

Quasispecies reconstruction

- Linkage
- Dual infections
- Transmission

Uses of variant analysis - Features

Here, the consensus can be submitted to 'conventional' tools for interpretation



geno2pheno[ngs-freq]

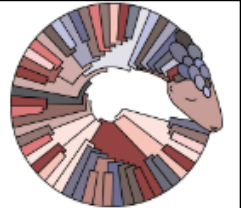


DENGUE, ZIKA & CHIKUNGUNYA
VIRUSES TYPING TOOL

Pangolin COVID-19

Lineage Assigner

Phylogenetic Assignment of Named
Global Outbreak LINEages

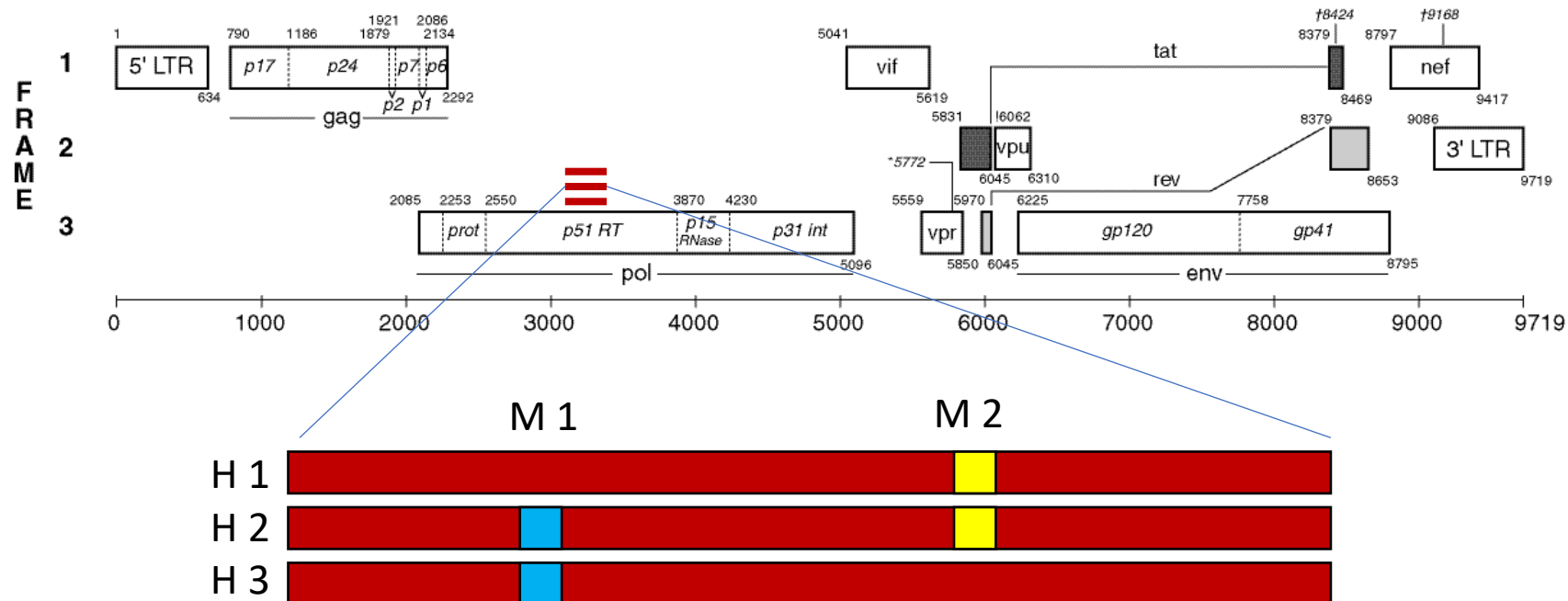


- Use the variant calling tool(s) to produce consensus at different mixed-base thresholds to interrogate minority variants.

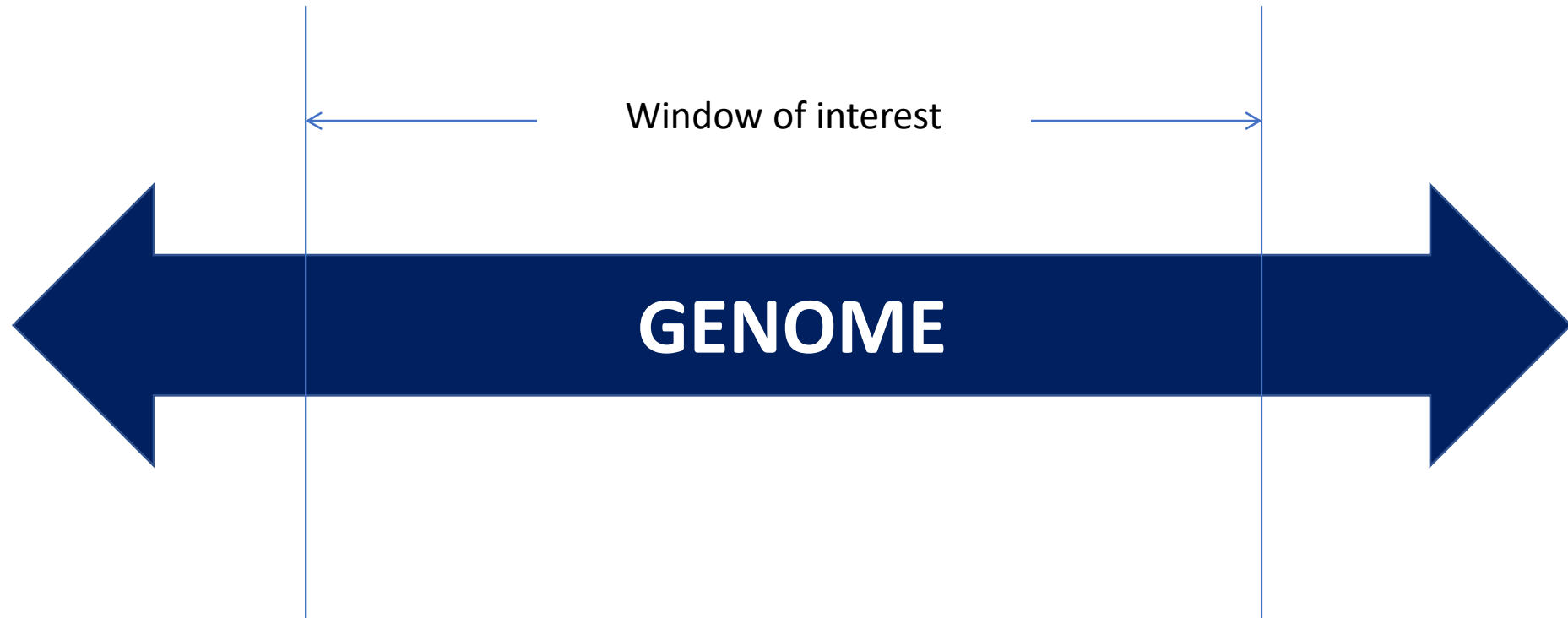
Validation, validation, validation

Uses of variant analysis - Quasispecies

- Each read derives from an individual virus genome molecule
- Linkage of variants on reads enables binning of haplotypes
- Examine all reads that map across a short, specified region of the genome:



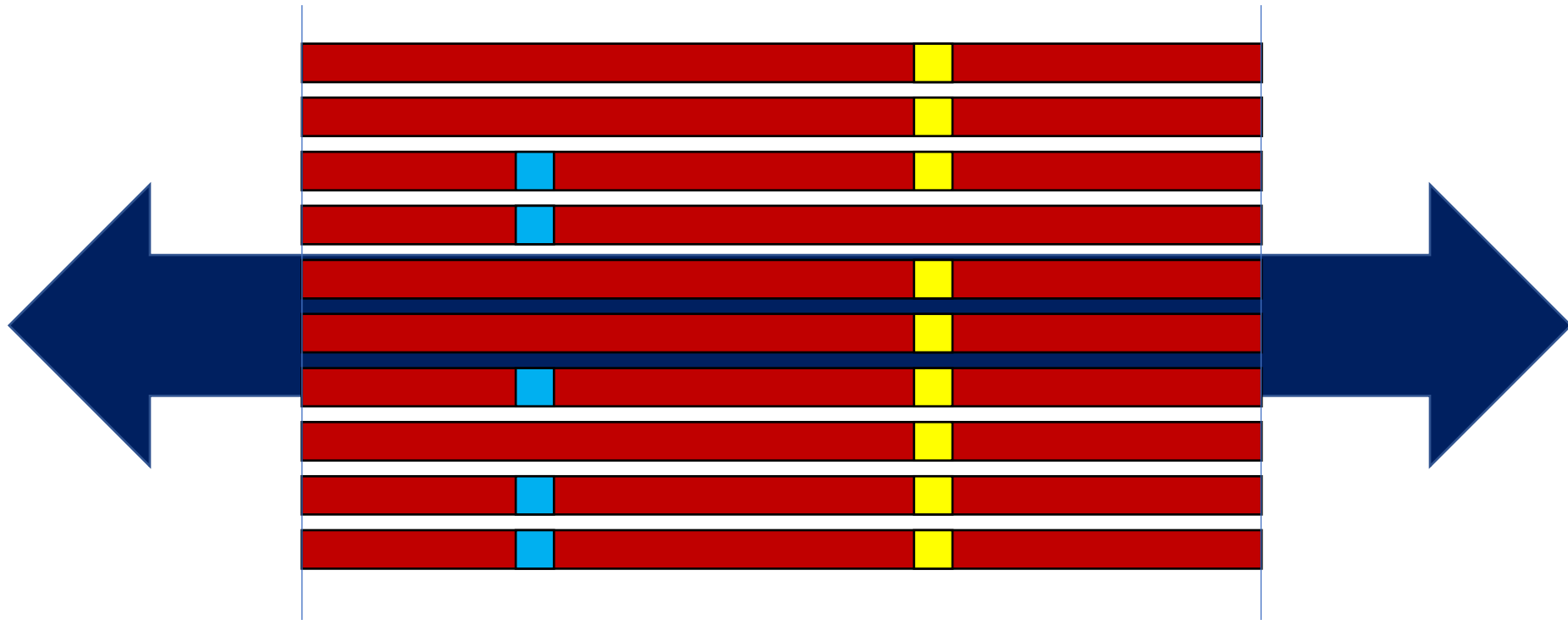
Uses of variant analysis - Quasispecies



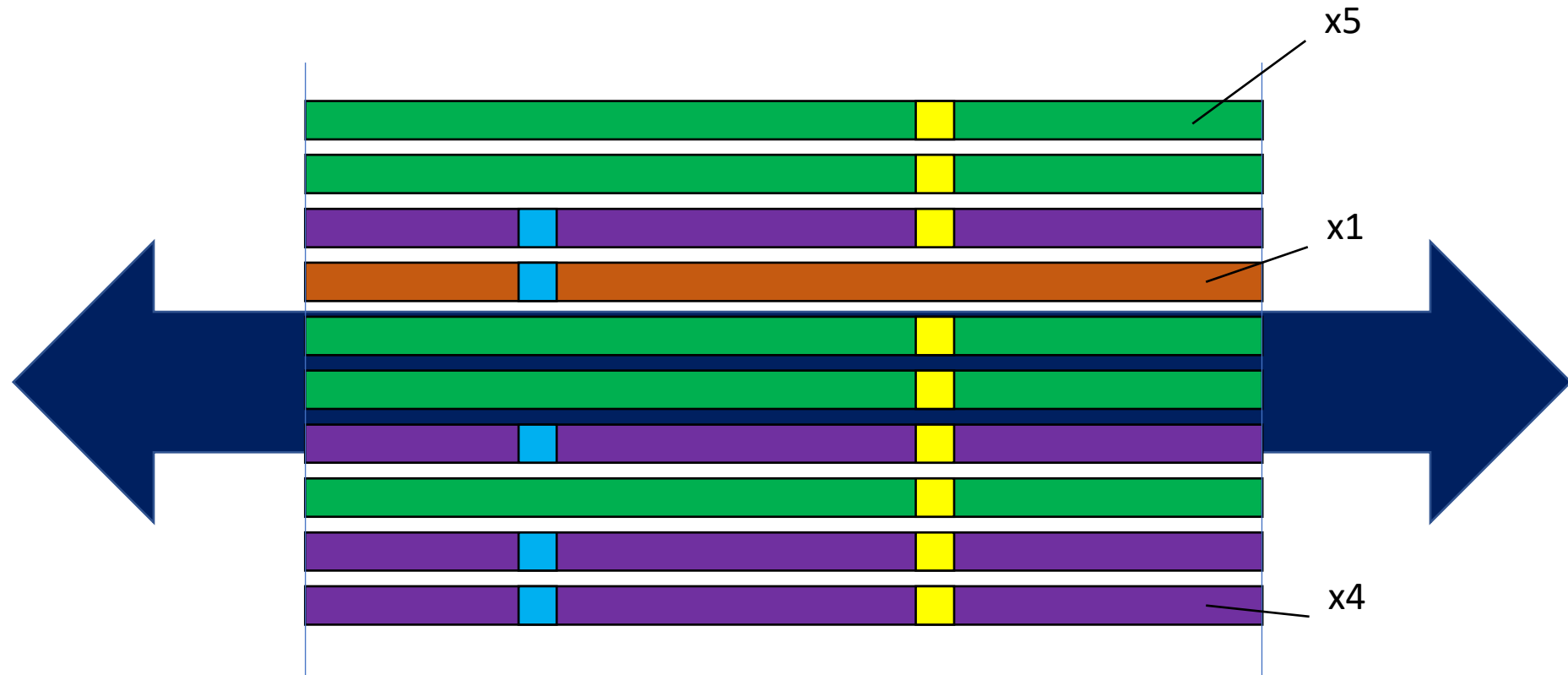
Uses of variant analysis - Quasispecies



Uses of variant analysis - Quasispecies



Uses of variant analysis - Quasispecies



Uses of variant analysis - Quasispecies

- Local data can be expanded to generate longer haplotypes

Haploclique <https://github.com/cbg-ethz/haploclique>

QuasiRecomb <https://github.com/cbg-ethz/QuasiRecomb>

QuRe <https://sourceforge.net/projects/quire>

PredictHaplo <http://bmda.cs.unibas.ch/software.html>

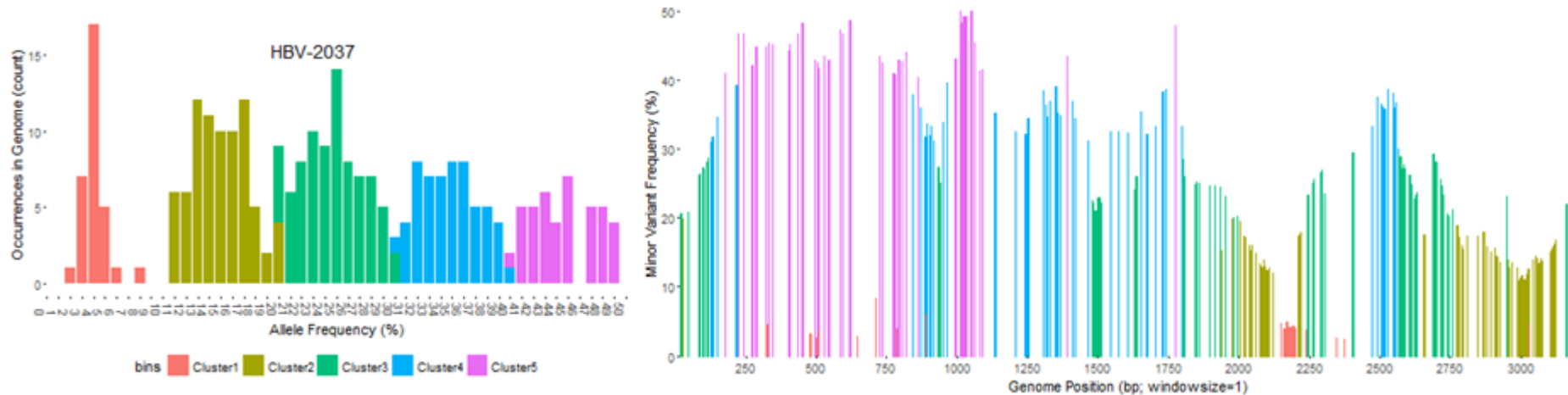
Efficiency of reconstruction is “varied”!

Beerenwinkel N *et al.* Front Microbiol. 2012 3:329

Prosperi MCF *et al.* Sci Rep. 2013 3:2837

Uses of variant analysis - Quasispecies

- Correlate mutation frequencies across the genome



Mathew Beale

Uses of variant analysis - Quasispecies

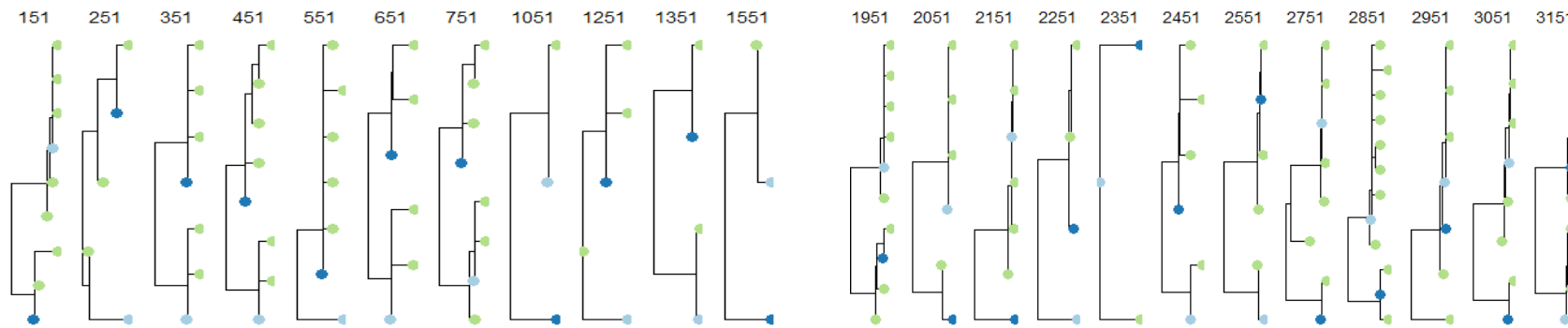
- Correlate mutation frequencies across the genome



Mathew Beale

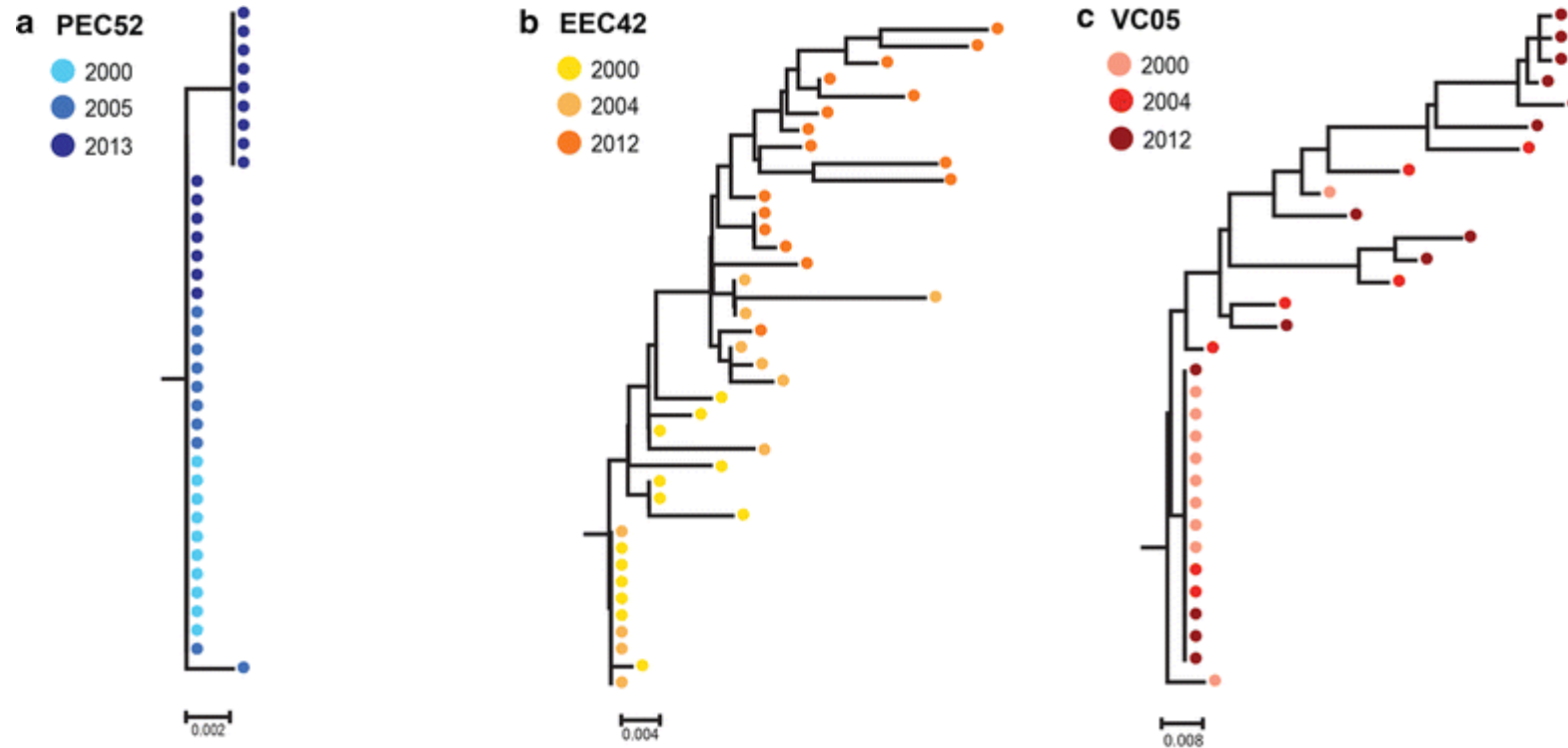
Uses of variant analysis - Quasispecies

- Correlate mutation frequencies across the genome



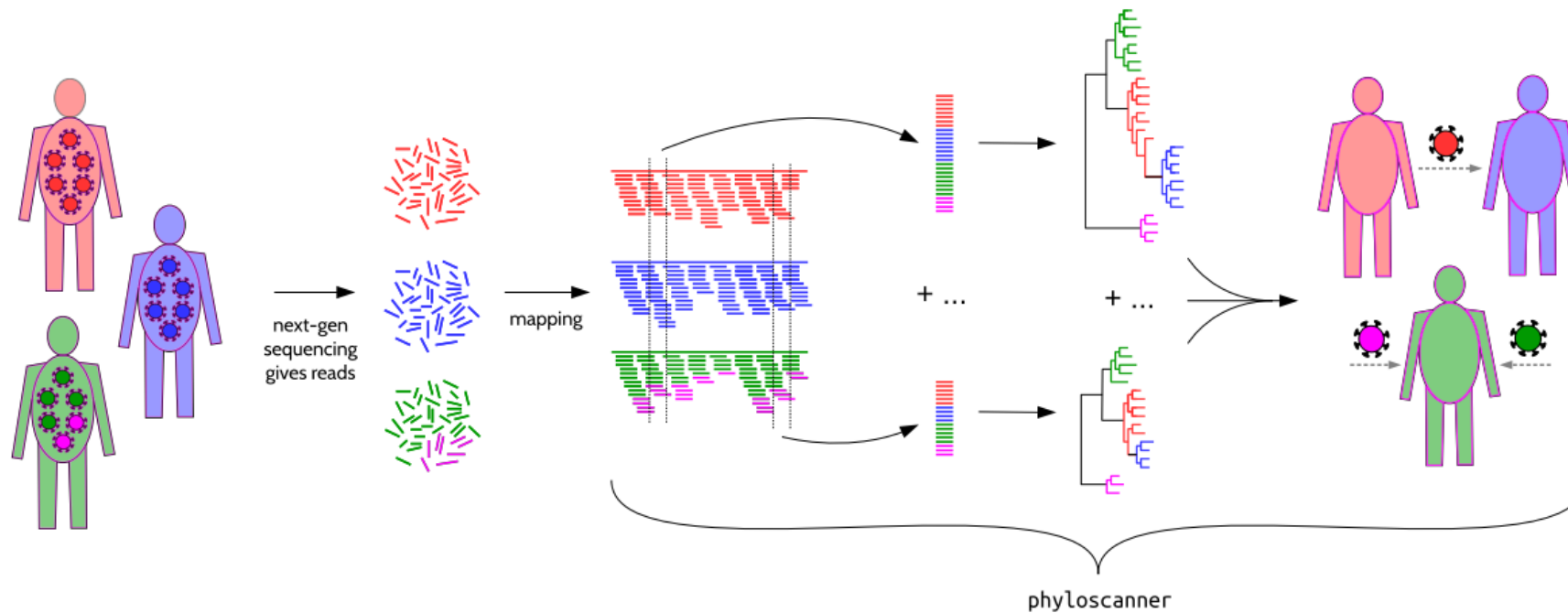
Mathew Beale

Uses of variant analysis - Quasispecies



de Azevedo SSD *et al.* Retrovirology 2017 14:29

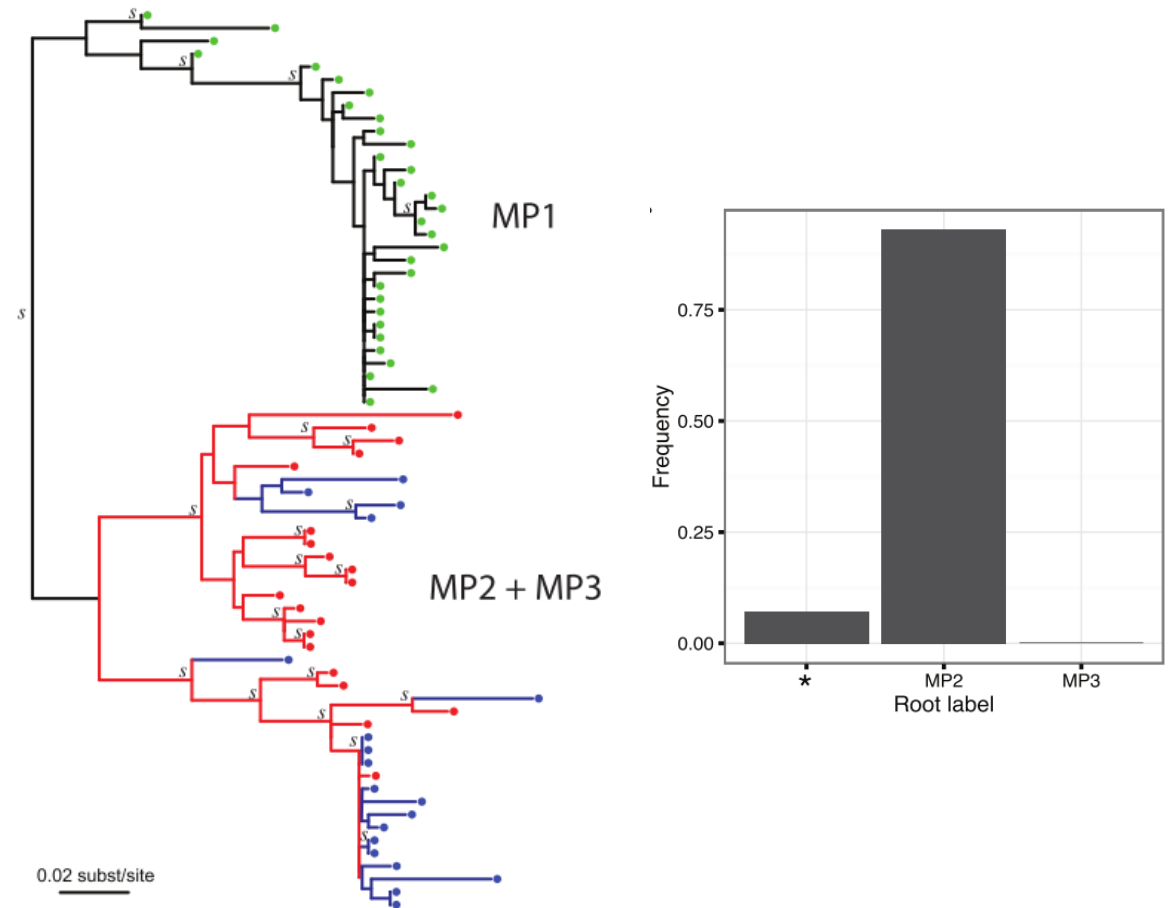
Uses of variant analysis - Quasispecies



Wymant C *et al.* Mol Biol Evol 2017

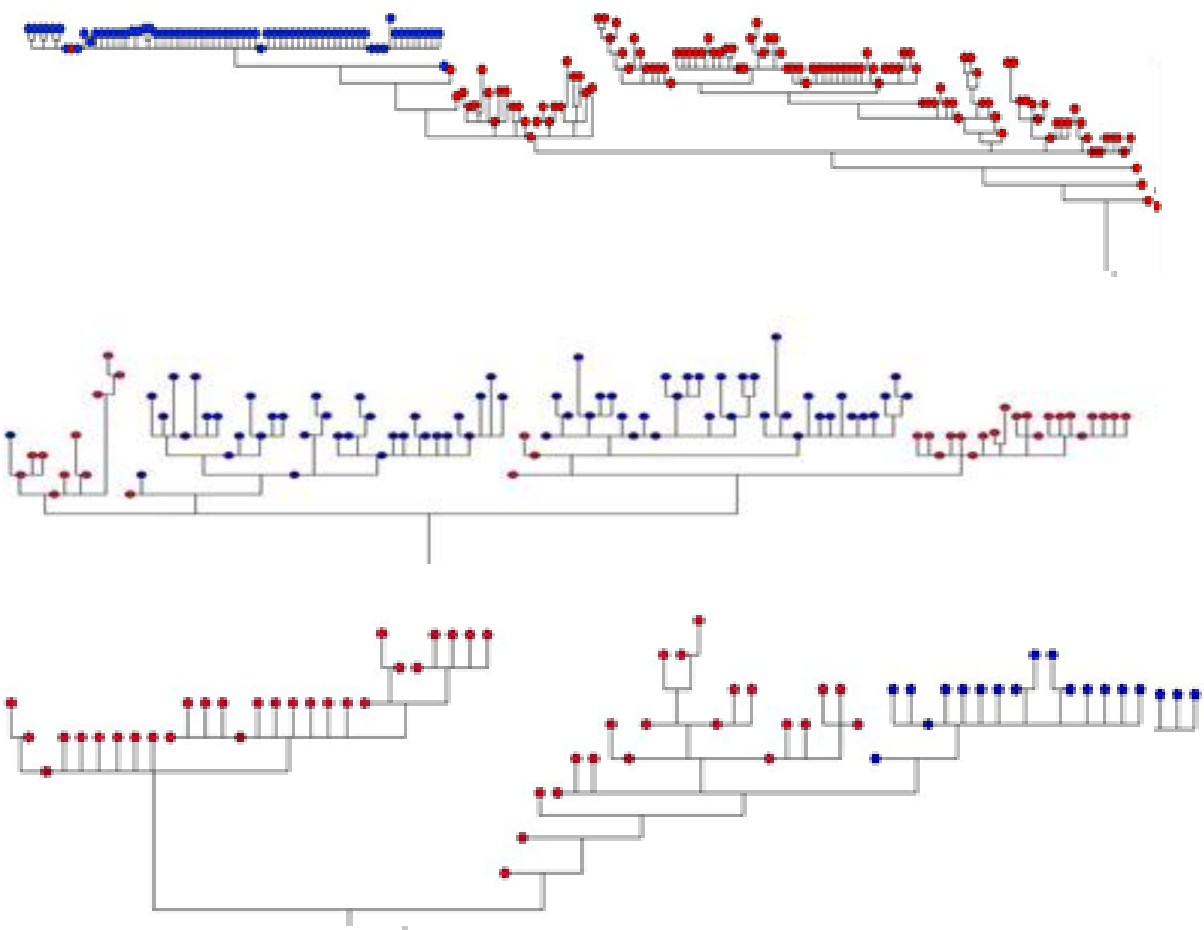
Uses of variant analysis - Quasispecies

- Three patients – MP1, 2 & 3
- **Who infected whom?**
- MP1 is independent from the cluster
- MP2 (red) → MP3 (blue)?
- Or *vice versa*?



Romero-Severson EO *et al.* Genetics 2017 207(3):1089

Uses of variant analysis - Quasispecies



| Transmission | |
|--------------|----|
| MSM | x3 |
| HET | x5 |
| MtCT | x2 |

| Subtype | |
|-----------|---------|
| B | x5 |
| C | x2 |
| G, 01, 02 | x1 each |

• Three patterns of sampling (🧴) & transmission window

(1) (2) (3)

TIME →

| Pair | A → B | A ↔ B | B → A | A ? B | Transmission & sampling pattern | Expected result |
|------|-------|-------|-------|-------|---------------------------------|-----------------|
| | | | | | | |
| 1 | | 2 | 8 | | MSM 1 | |
| 2 | | 3 | 6 | 1 | HET 2 | |
| 3 | 10 | | | | HET 1 | |
| 4 | 1 6 | 1 | 1 | 1 | HET 3 | ? |
| 5 | | | 5 5 | | HET 3 | ? |
| 6 | 3 | 5 | 2 | | MSM 2 | |
| 7 | 3 | 6 | 1 | | MSM 2 | |
| 8 | 4 | | | 6 | HET 3 | ? |
| 9 | | 10 | | | MtCT 1 | |
| 10 | 6 3 | 1 | | | MtCT 1 | |

Table 1. Summary of tree topologies from 10 most populous tiles for the ten linked pairs.

Colours describe the relationship with known / unknown transmission histories:

Consistent

Inconsistent

Suggestive

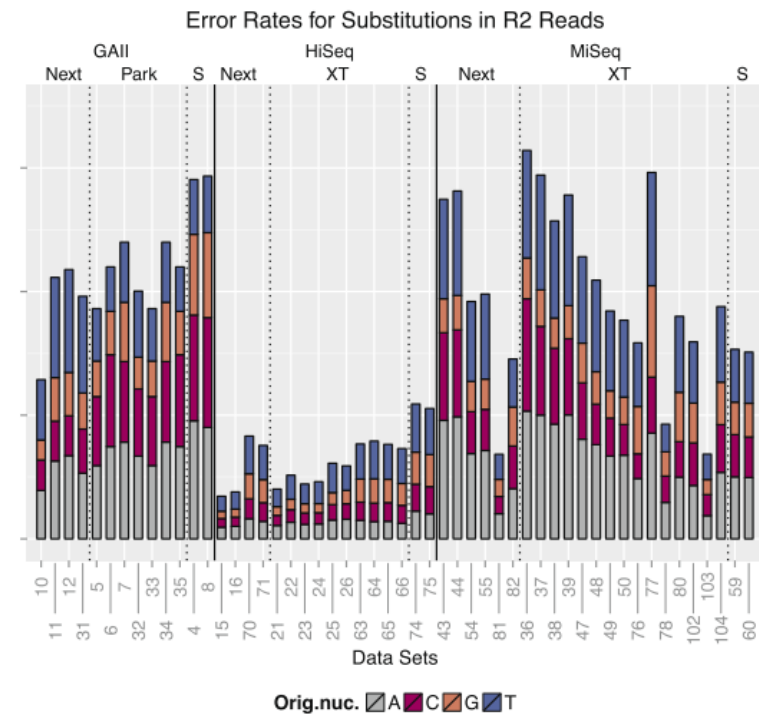
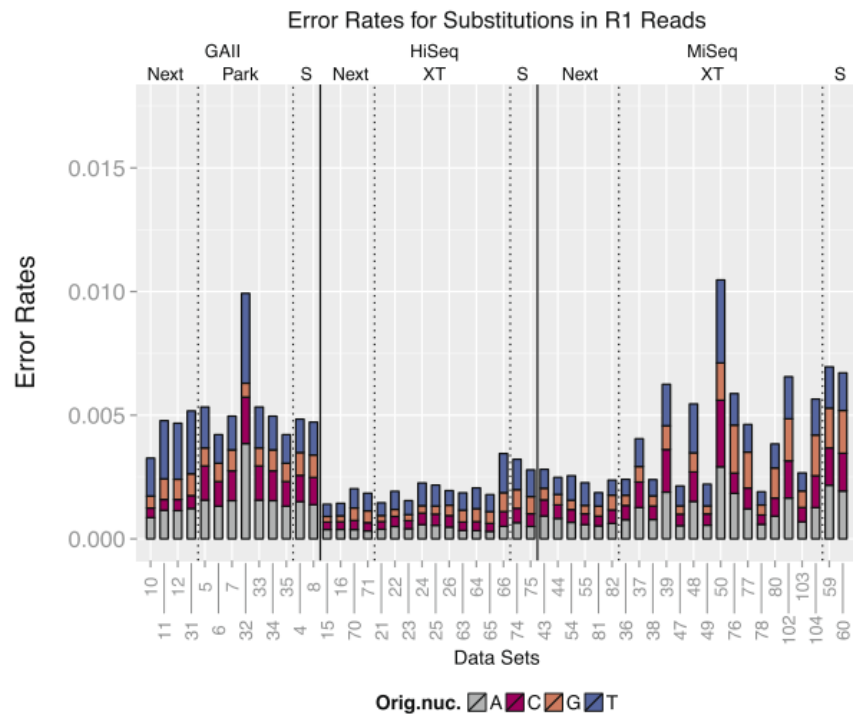
Bibby DF *et al.* HIV Dynamics & Evolution 2017

Technical pitfalls

Frequency of minor variant detection limited by experimental noise
Many sources of error:

1. Sequencing
2. Amplicon-based sequencing
3. Nucleotide content
4. Hexamer priming
5. Product degradation
6. Contamination
7. Bioinformatics

Technical pitfalls – Sequencing



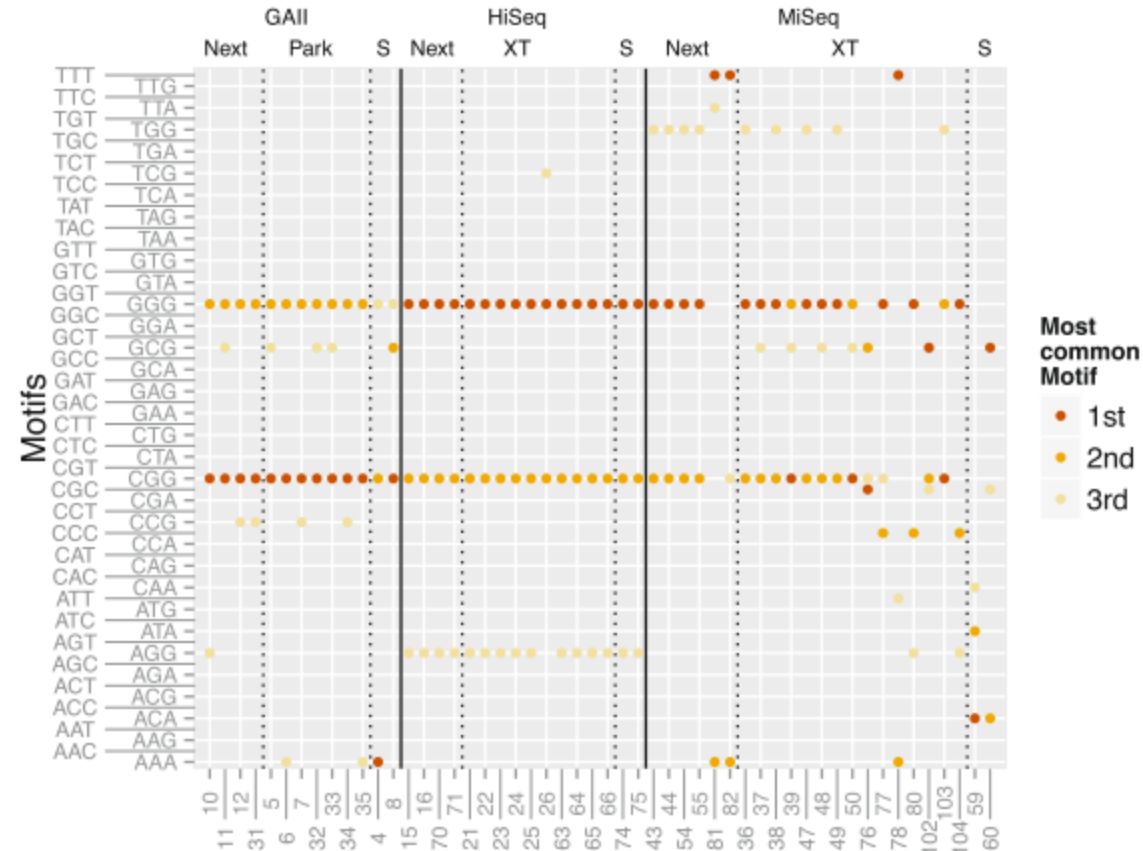
Schirmer M *et al.* BMC Bioinformatics 17(1):125

Technical pitfalls – Sequencing

Substitutions

xGG motif

GGG = CGG > AGG > TGG



Schirmer M *et al.* BMC Bioinformatics 17(1):125

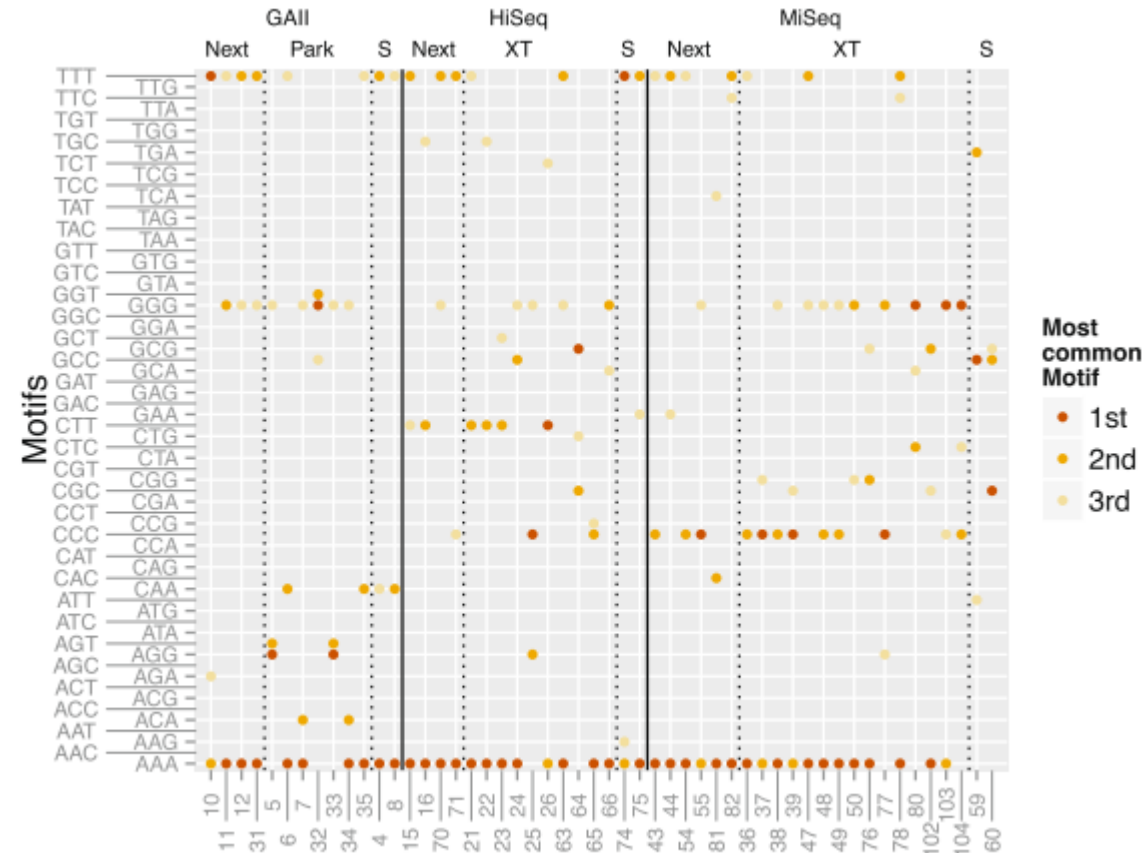
Technical pitfalls – Sequencing

Indels

Homopolymeric tracts

AAA > CCC = GGG = TTT

High Quality Scores



Schirmer M *et al.* BMC Bioinformatics 17(1):125

Technical pitfalls – Amplicons

Table 1. Error rate of *Taq* DNA polymerase.

| Amplicon | Substitution rate | Deletion rate | Insertion rate | Total error rate | Total bases |
|--|------------------------------|-----------------------------|-----------------------------|----------------------|-------------|
| <i>Sanger (dideoxy)</i> | | | | | |
| LacZ-1 | 1.2×10^{-4} (98.8%) | 1.6×10^{-6} (1.2%) | - (0.0%) | 1.3×10^{-4} | 323,802 |
| <i>Pacific Biosciences RSII</i> | | | | | |
| LacZ-1 | 1.7×10^{-4} (97.3%) | 4.7×10^{-6} (2.6%) | 1.8×10^{-7} (0.1%) | 1.8×10^{-4} | 35,879,784 |
| LacZ-2 | 1.7×10^{-4} (96.1%) | 5.1×10^{-6} (2.9%) | 1.8×10^{-6} (1.0%) | 1.8×10^{-4} | 15,857,446 |
| DNA-1 | 1.4×10^{-4} (97.2%) | 3.9×10^{-6} (2.8%) | 1.2×10^{-7} (0.1%) | 1.4×10^{-4} | 18,680,811 |
| DNA-2 | 1.4×10^{-4} (97.5%) | 3.4×10^{-6} (2.4%) | 1.5×10^{-7} (0.1%) | 1.4×10^{-4} | 27,978,748 |

Reported error rates are per base per doubling as detailed in Materials and Methods. Numbers in parentheses are percentages of the total error rate.

Table 6. PCR-mediated recombination rate by *Taq* DNA polymerase.

| Template pair | N_{re}^a | N_{total}^b | Recombination rate c | Strands with at least 1 recombination event |
|---------------|------------|---------------|-------------------------|---|
| DNA-1:DNA-1x | 19,943 | 77,725,936 | 9.6×10^{-5} | 23% |
| DNA-2:DNA-2x | 14,687 | 44,271,304 | 1.3×10^{-4} | 28% |

^a Number of recombination events.

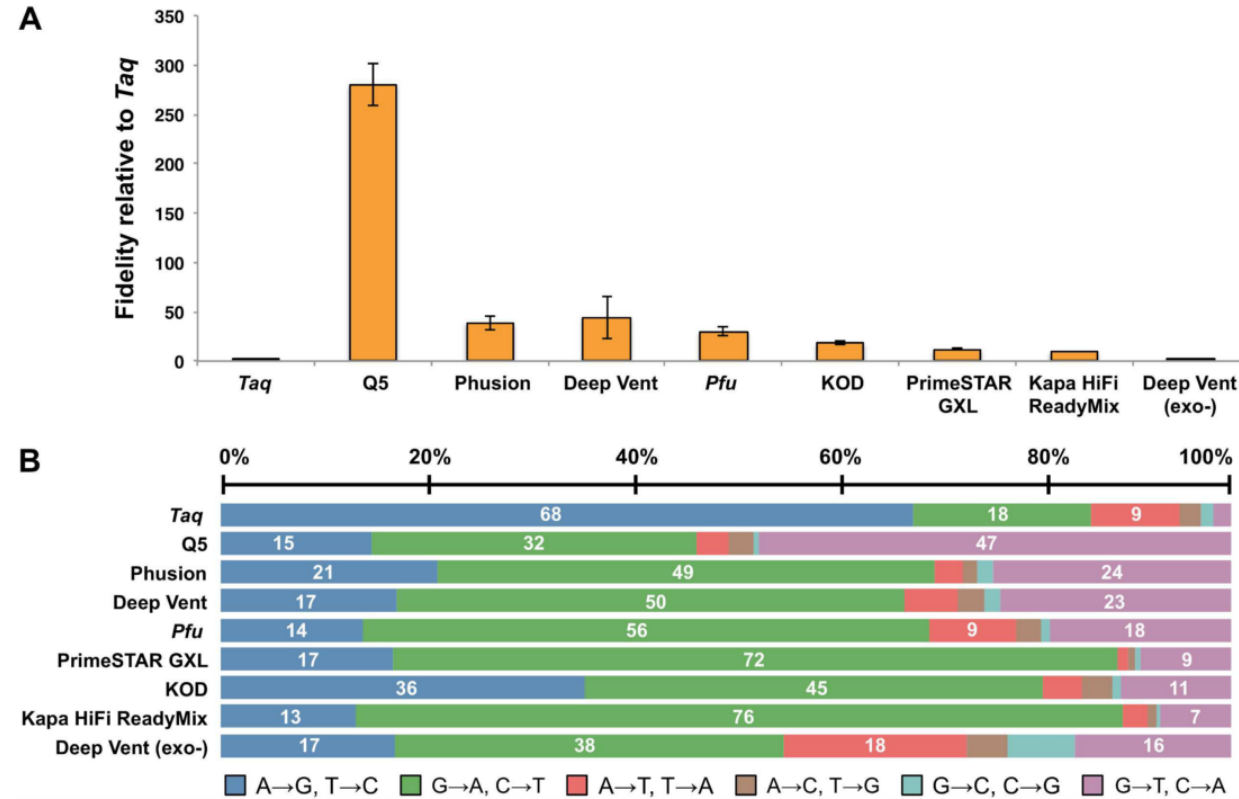
^b Total number of analyzed sequenced bases.

^c Recombination rate is per base per doubling. Recombination rate is doubled to account for “cryptic” recombination events.

1kb, 16x cycles

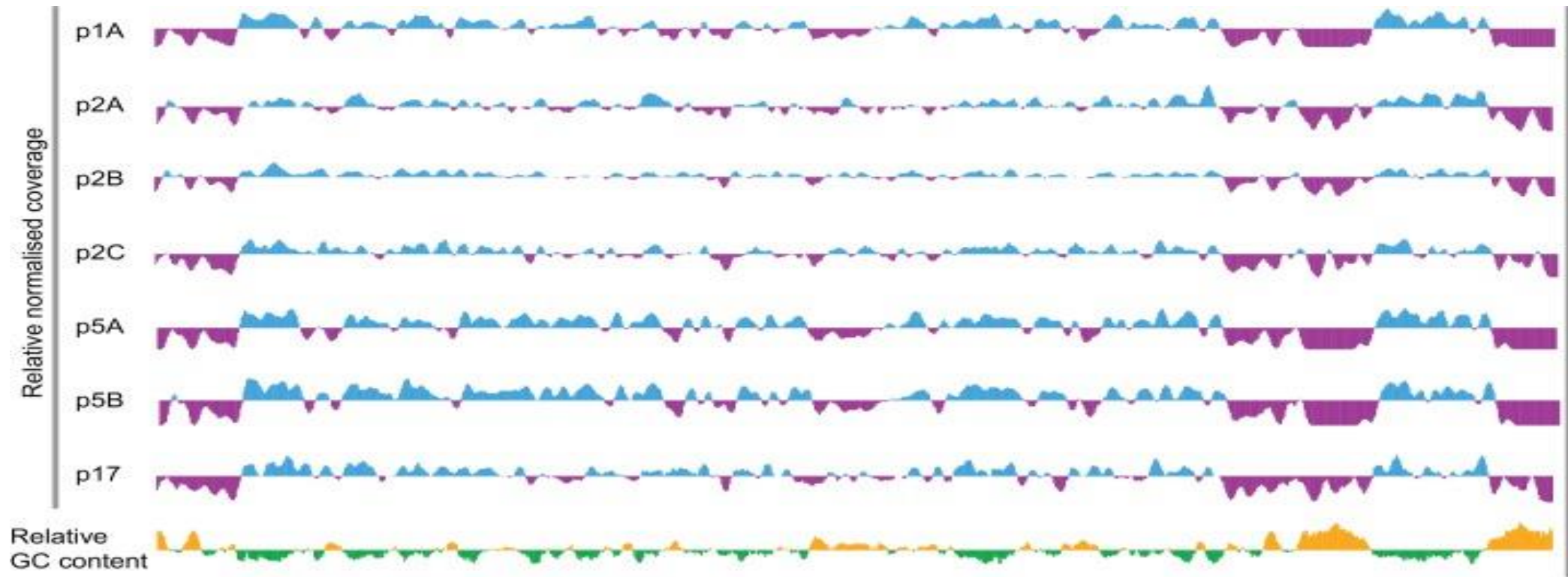
Potapov V *et al.* PLoS One 2017 12(1): e0169774

Technical pitfalls – Amplicons



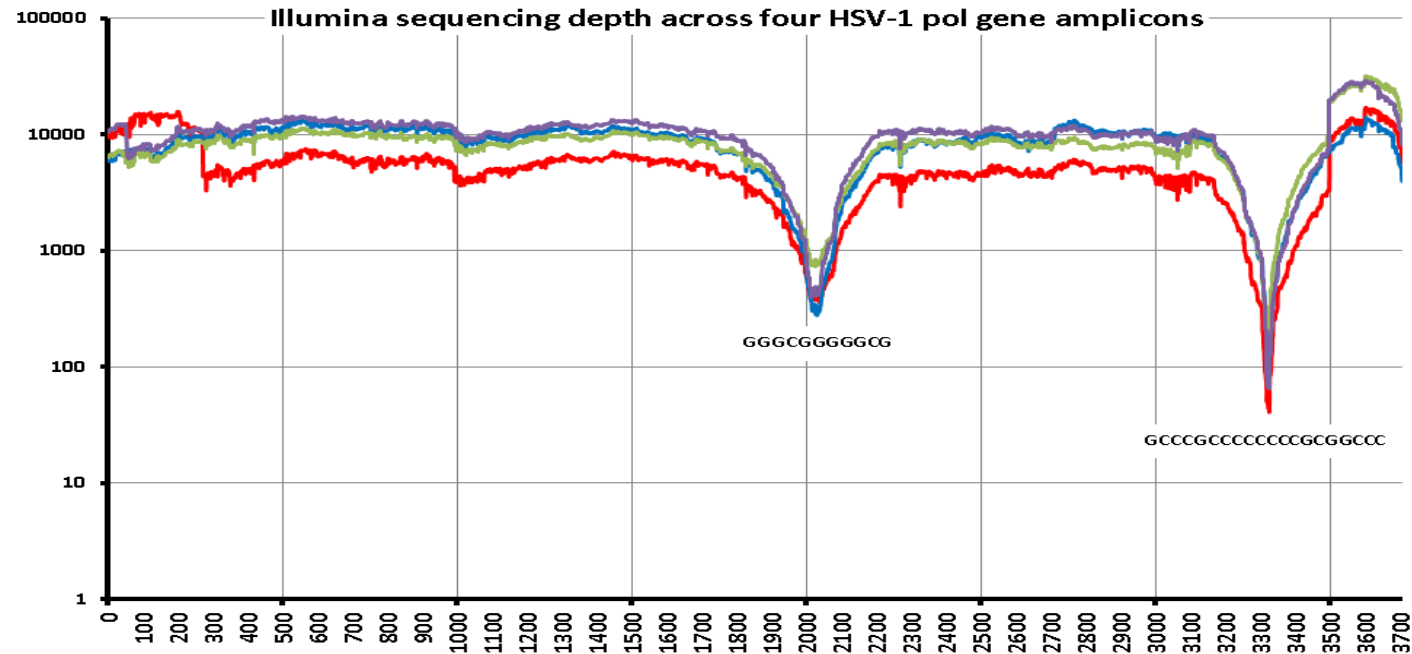
Potapov V *et al.* PLoS One 2017 12(1): e0169774

Technical pitfalls – Nucleotide content



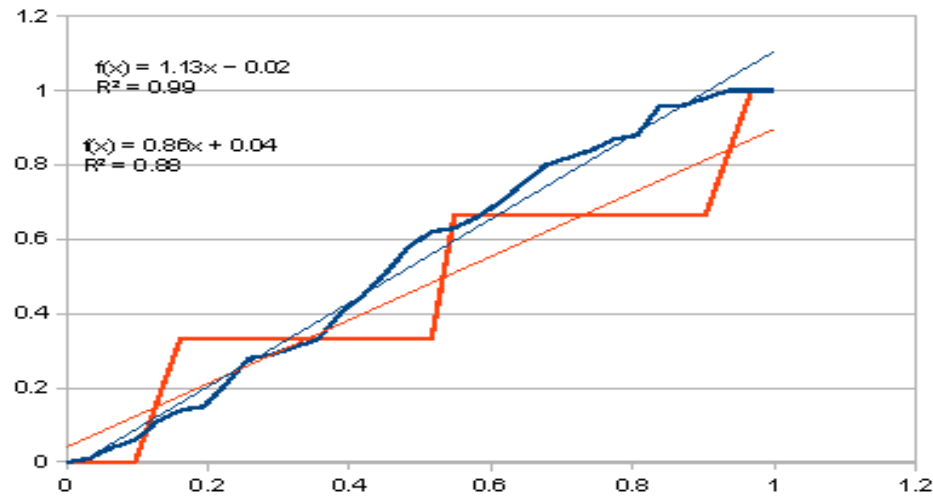
Karamitros T *et al.* PLoS One 2016 11(6):e0157600

Technical pitfalls – Nucleotide content

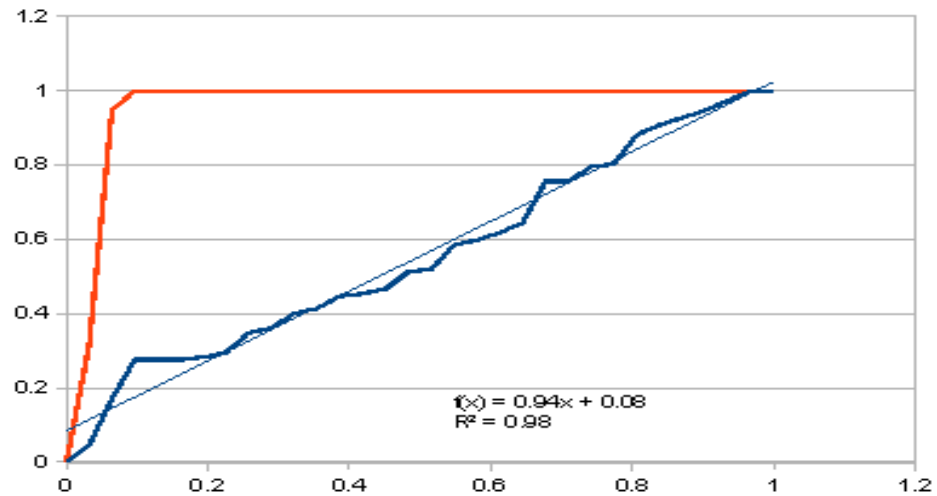


UKHSA (unpublished)

Technical pitfalls – Hexamer priming



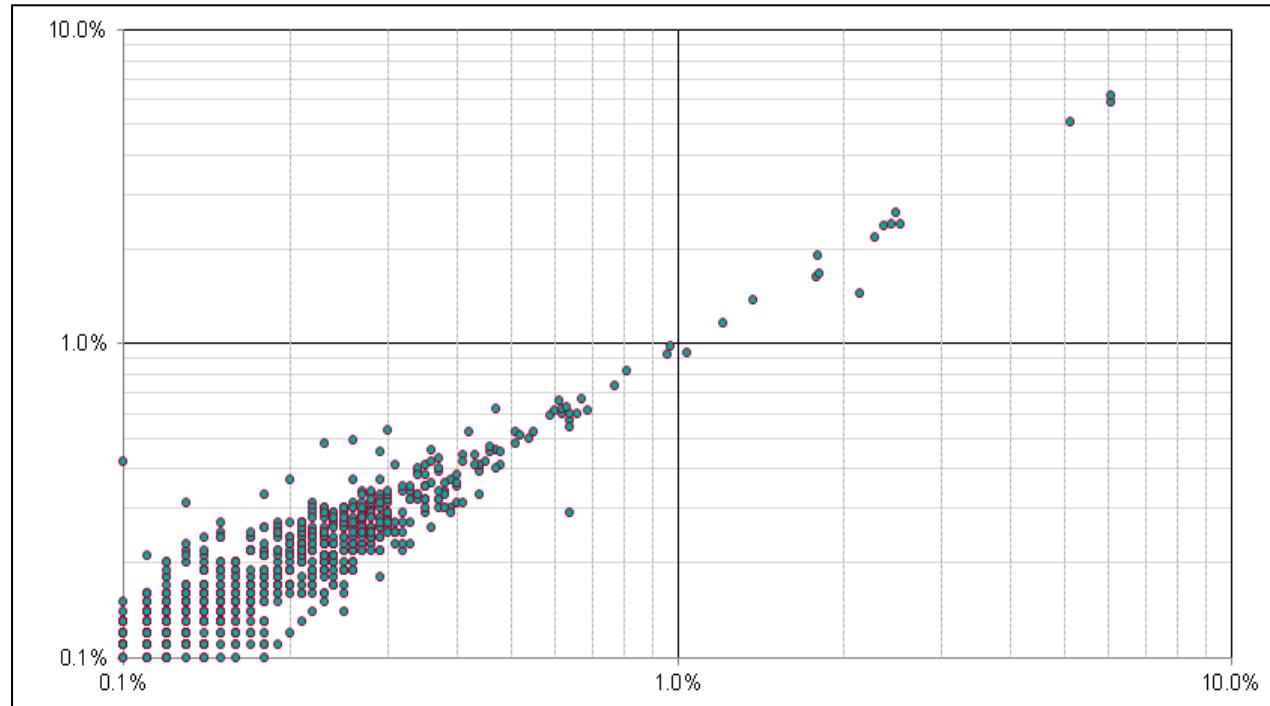
The position of the minority variant nucleotide (red line) is evenly distributed along the read lengths (as is the majority variant in blue)



100% of the minority variants are within 7% (10nt) of a read terminus – artefact from insert-priming

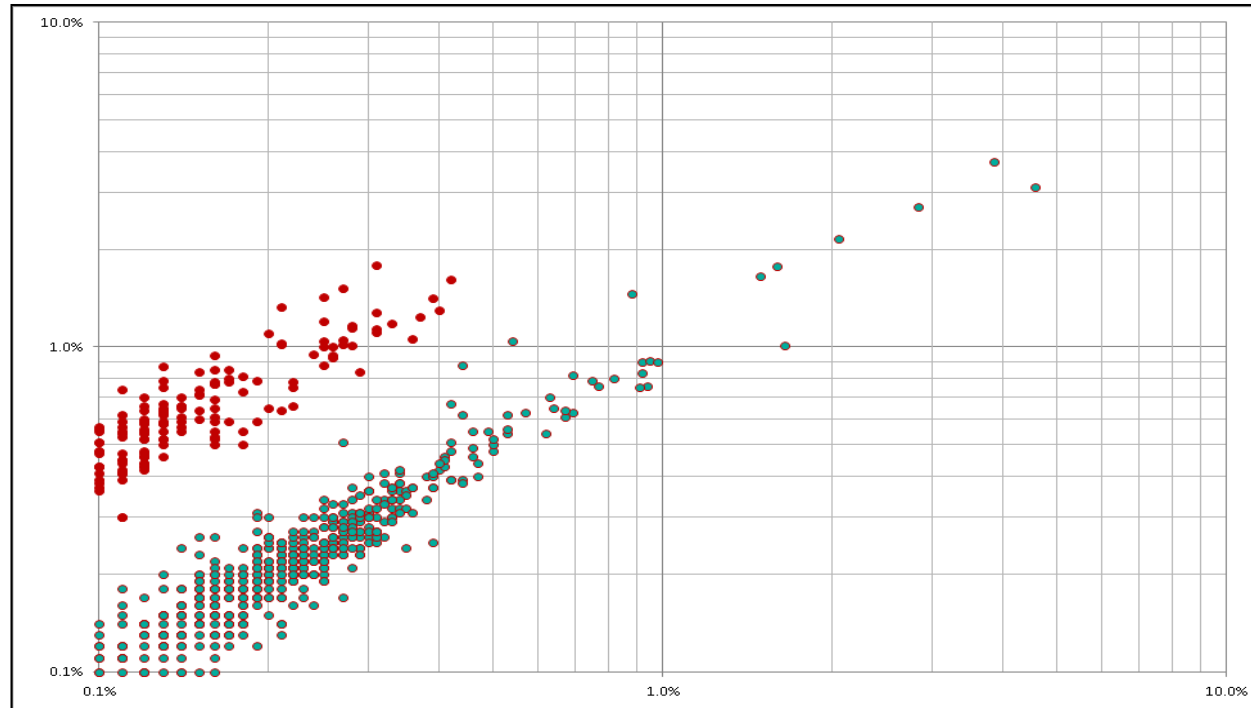
UKHSA (unpublished)

Technical pitfalls – Product degradation



UKHSA (unpublished)

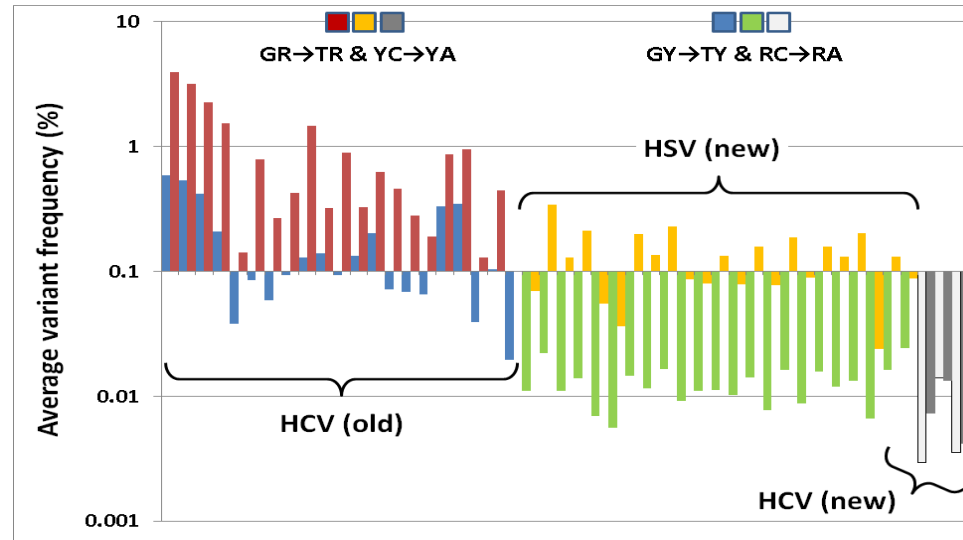
Technical pitfalls – Product degradation



UKHSA (unpublished)

Technical pitfalls – Product degradation

- Much investigation revealed context-specific conversion of dinucleotides
 $YC \rightarrow YA$
 $GR \rightarrow TR$ } previously only seen in sonicated fragments – Costello *et al.* NAR 2013
- The frequency of converted bases is proportional to the time spent at 4°C



UKHSA (unpublished)

Technical pitfalls – Contamination

“Sequences not belonging to that sample present in the FASTQ set”

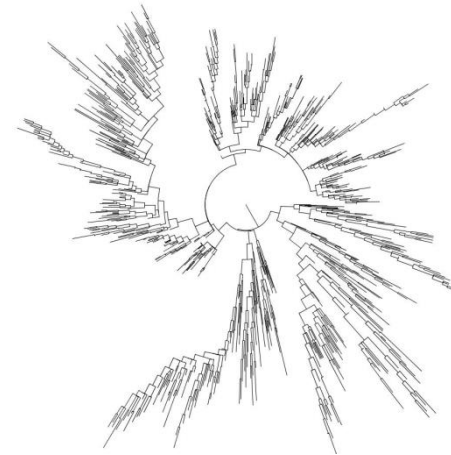
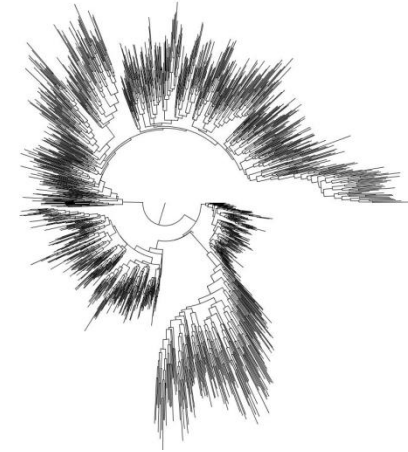
- Index-switching
 - Occurs during library prep / sequencing
- Laboratory contamination
 - Similar to PCR
 - Spatial and temporal separation of work areas
 - Rotation of adapters
 - Robotics
 - Rotation of control positions
 - Alternation of template types

Technical pitfalls – Bioinformatics

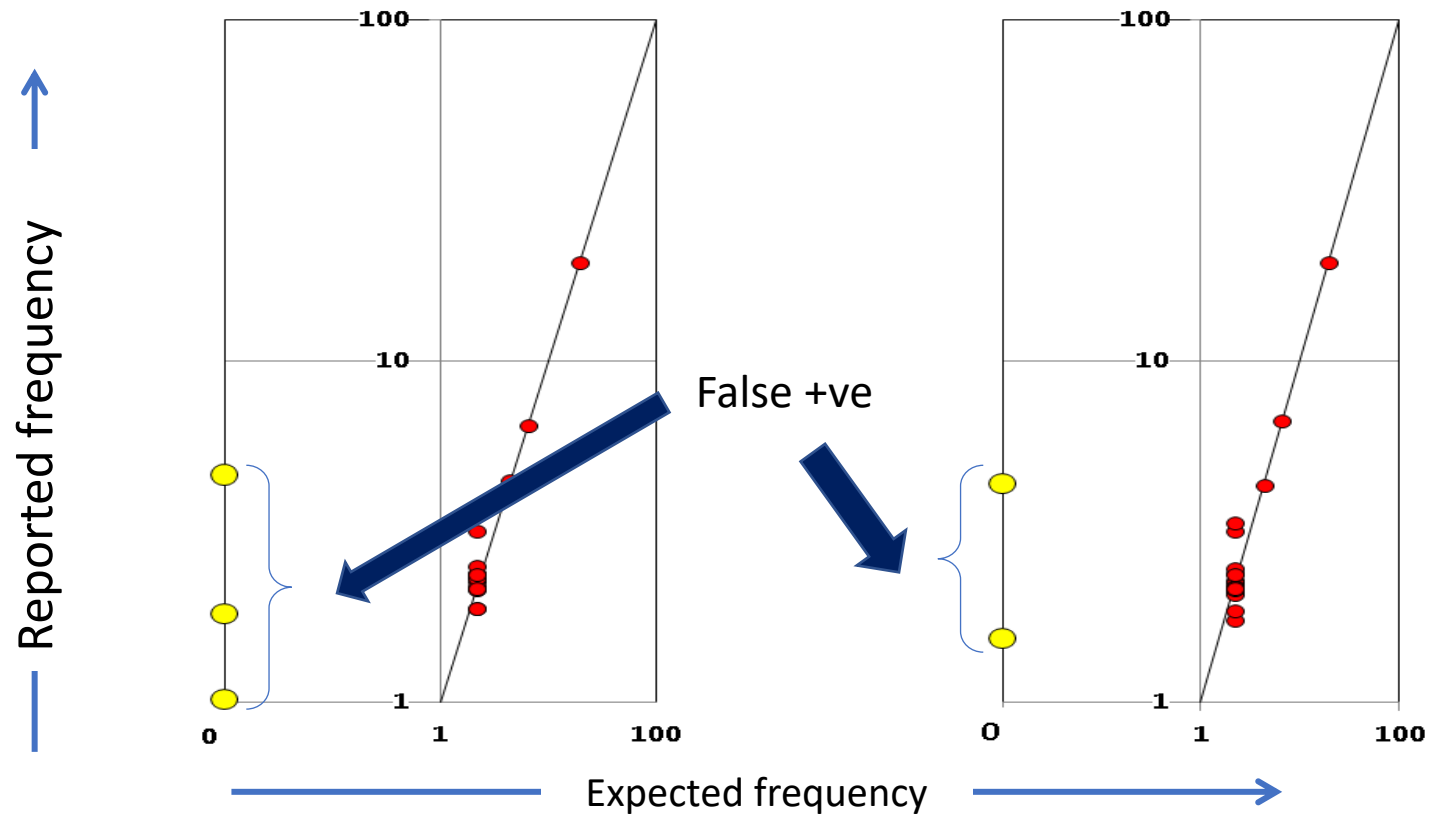
- Synthetic FASTQ datasets
 - HIV quasi-species generation

Adapted from Pandit A & de Boer R, Retrovirology 2014 11(1):56
 - FASTQ generation using empirical error profiles / quality scores

ART - Huang W, *et al.* Bioinformatics 2012 28(4):593-4
- Two pre-production pipelines tested

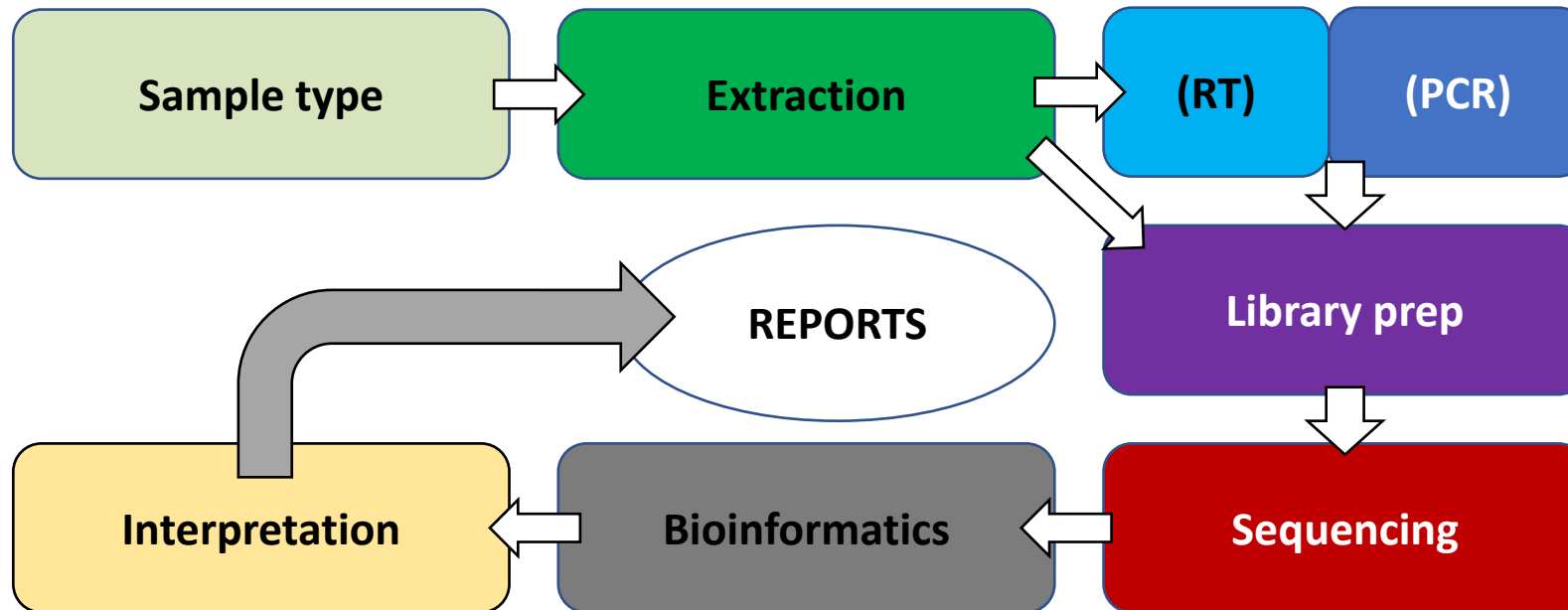


Technical pitfalls – Bioinformatics



UKHSA (unpublished)

Validation



Validation – How low can you go?

- For HIV & HCV, Sanger at 15-20% == resistant
 - What does 10% mean?
 - 5%?
- Reproducibility, repeatability, accuracy and precision is critical
 - Validating against Sanger is relatively straightforward, both clinically and technically
 - Validating lower frequencies is quite the opposite
- Clinical utility of lower frequency variants unproven








Cut-offs vary considerably between assays

Validation – How low can you go?

J Antimicrob Chemother 2023; **78**: 656–664
<https://doi.org/10.1093/jac/dkac430> Advance Access publication 4 February 2023

**Journal of
Antimicrobial
Chemotherapy**

Frequency matters: comparison of drug resistance mutation detection by Sanger and next-generation sequencing in HIV-1

Suraj Balakrishna^{1,2*}, Tom Loosli^{1,2}, Maryam Zaheri^{2,3}, Paul Frischknecht¹, Michael Huber^{2,3}, Katharina Kusejko ^{1,2}, Sabine Yerly⁴, Karoline Leuzinger⁵, Matthieu Perreau⁶, Alban Ramette ⁷, Chris Wymant ⁸, Christophe Fraser^{8,9}, Paul Kellam¹⁰, Astrid Gall¹¹, Hans H. Hirsch ¹², Marcel Stoeckle¹², Andri Rauch¹³, Matthias Cavassini ¹⁴, Enos Bernasconi¹⁵, Julia Notter¹⁶, Alexandra Calmy¹⁷, Huldrych F. Günthard ^{1,2}, Karin J. Metzner^{1,2} and Roger D. Kouyos ^{1,2}

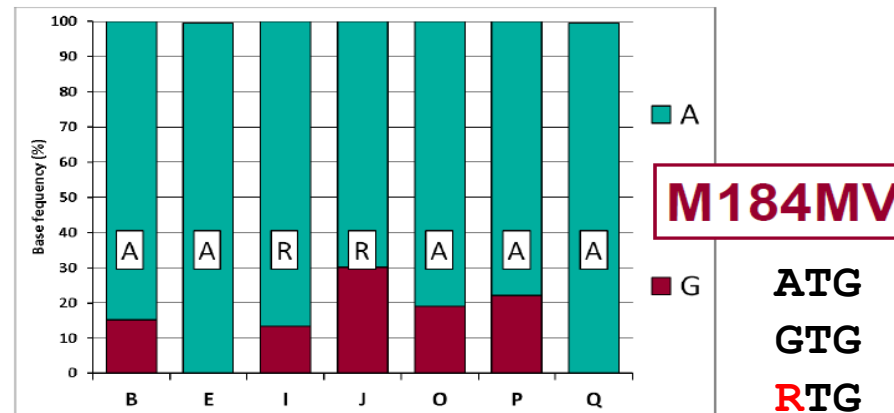
Conclusions: We found high concordance between SS and NGS but also a substantial number of low-abundance HIV-DRMs detected only by NGS at lower variant-calling thresholds. Our findings suggest that a substantial fraction of the low-abundance HIV-DRMs detected at thresholds <3% may represent sequencing errors and hence should not be overinterpreted in clinical practice.

Validation – Copy number & variant frequency

How reliable is a variant frequency call?

- When the depth of coverage (i.e. reads covering that position) is high/low?
- What levels constitute 'high' and 'low' depths for PCR, sequence capture and metagenomic approaches?

Sample: A Domain: PR/RT Position: 847 Consensus: R

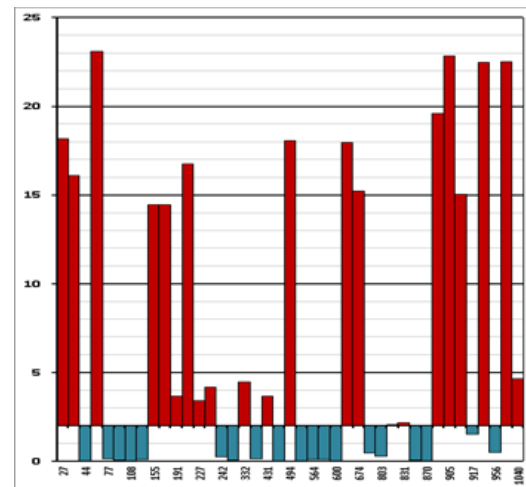


UKHSA (unpublished)

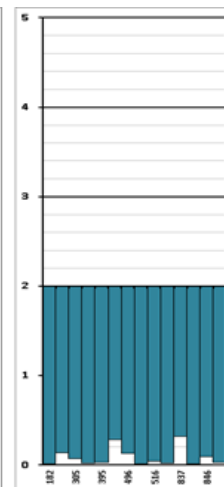
Validation – Copy number & variant frequency

How reliable is a variant frequency call?

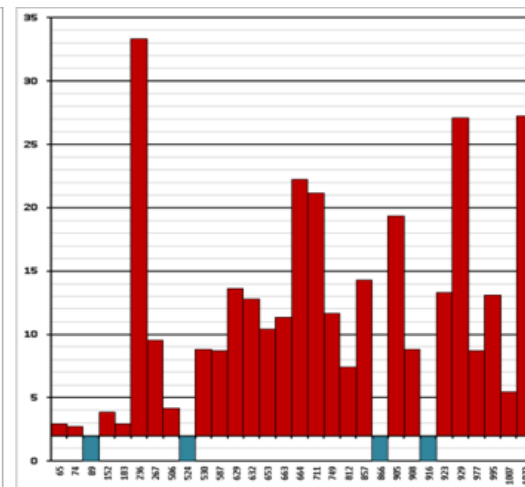
- When the depth of coverage (i.e. reads covering that position) is high/low?
- What levels constitute 'high' and 'low' depths for PCR, sequence capture and metagenomic approaches?



Sample 1



Sample 2



Sample 3

UKHSA (unpublished)

Validation – Copy number & variant frequency

How reliable is a variant frequency call?

- PCR produces large amounts of material
- How representative of the starting population is the amplicon mix?
 - Depth is not correlated with reliability!
- If the starting virus copy number is 100,000 copies per library, 10% = 10,000
 - But at 1,000 copies, how reliable is 5% (50 viruses)?
 - Reverse transcription (to generate cDNA) is notoriously inefficient and error-prone – how many viruses are represented?
 - There are multiple PCR cycles in the library prep too...
- Often, the amount of starting material / viral load is unknown

Validation – Copy number & variant frequency



Perspective

Fact and Fiction about 1%: Next Generation Sequencing and the Detection of Minor Drug Resistant Variants in HIV-1 Populations with and without Unique Molecular Identifiers

Shuntai Zhou ^{1,*} and Ronald Swanstrom ^{1,2}

[Viruses \(2020\) 12\(8\): 850](#)

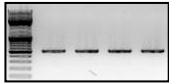
Validation – Standardised materials



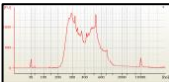
Sample – haplotypes mixed at precise frequencies (e.g. 1, 2, 5, 10 & 20%)



DNA / RNA extract



PCR product



Library – storage issues



FASTQ datasets - artificial quasispecies and synthetic FASTQs

5' ATGACGTGGGA3'
3' TACTGCACCCCT5'

Consensus sequences – to test interpretation mechanisms

Summary

Consensus & variant calling

- Many tools for mapping (BWA, Bowtie, smalt, Tanoti)
- Several tools for variant calling (QuasiBAM, V-Phaser)

Choice of reference & user-defined parameters is CRITICAL

Using the consensus

- Submit to usual tools
- Different mixed-base thresholds to incorporate minority variants

BEWARE – Interpretations may not be validated on NGS-derived data...

- Phyletic analysis is coming

Summary

Technical pitfalls

- Experimental approach influences the result in unpredictable ways
 - Low-frequency variation can arise through diverse processes
- Reproducibility experiments are essential
 - Across a range of conditions and samples

Validation

ESSENTIAL – especially when detecting ‘new’ data, e.g. <20% minority

- All components of the assay need independent investigation
- Look at all the data
- ESPECIALLY THE BIOINFORMATICS!