

Supplementary Materials

| Method | | Clone1 | Clone2 | Clone3 | Clone4 | Clone5 | Clone6 | Clone7 | Clone8 | Clone9 | Clone10 | FP |
|--------|-----------|--------------|--------|--------|--------|--------|--------|--------|--------|--------|---------|-------|
| | | True Freq.,% | 50 | 25 | 12.5 | 6.25 | 3.125 | 1.56 | 0.78 | 0.39 | 0.19 | 0.097 |
| 2SNV | Match | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | × | 1 |
| | Freq., % | 51.8 | 23.7 | 12.5 | 6.4 | 2.3 | 1.2 | 0.7 | 0.3 | 0.1 | 0 | 1.9 |
| | 95 CI., % | | | | | | | | | | | |
| | Low | 50.0 | 22.5 | 11.2 | 6.2 | 2.2 | 1.1 | 0.6 | 0.3 | 0.09 | - | 1.7 |
| Upper | 52.2 | 23.8 | 12.4 | 6.6 | 2.4 | 1.3 | 0.8 | 0.4 | 0.14 | - | 4.0 | |
| PH | Match | ✓ | ✓ | ✓ | × | ✓ | × | ✓ | ✓ | × | × | 0 |
| | Freq.,% | 56.7 | 23.8 | 13.7 | 0 | 3.1 | 0 | 1.5 | 1.2 | 0 | 0 | 0 |
| | 95 CI.,% | | | | | | | | | | | |
| | Low | 49.2 | 22.3 | 12.6 | 0 | 3.5 | - | 1.7 | - | - | - | 0 |
| Upper | 57.0 | 23.0 | 13.3 | 8.0 | 4.3 | - | 2.9 | - | - | - | 7.4 | |

Table 1: Comparison of 2SNV and PredictHaplo on full data with bootstrapping. For all 33.5K reads, the sign “✓” (respectively, “×”) denotes fully matched (respectively, unmatched) true variant and the column FP reports the number of incorrectly predicted variants (false positives) and their total frequency. 95 CI showing results on 40 bootstraps, “-” means variant was never reconstructed.

| # of Reads | Method | Clones | Clone1 | Clone2 | Clone3 | Clone4 | Clone5 | Clone6 | Clone7 | Clone8 | Clone9 | Clone10 | FP |
|-------------|--------|--------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|-----|
| | | True Freq.,% | 50 | 25 | 12.5 | 6.25 | 3.125 | 1.56 | 0.78 | 0.39 | 0.19 | 0.097 | |
| 33.5K (all) | 2SNV | Match | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | × | 1 |
| | | Freq., % | 51.8 | 23.7 | 12.5 | 6.4 | 2.3 | 1.2 | 0.7 | 0.3 | 0.1 | 0 | 1.0 |
| | PH | Match | ✓ | ✓ | ✓ | × | ✓ | × | ✓ | ✓ | × | × | 0 |
| | | Freq., % | 56.7 | 23.8 | 13.7 | 0 | 3.1 | 0 | 1.5 | 1.2 | 0 | 0 | 0 |
| 16K | 2SNV | Match,% | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 0 | 0 | 0.2 |
| | | Freq., % | 52.4 | 23.7 | 12.5 | 6.4 | 2.3 | 1.1 | 0.7 | 0.3 | 0 | 0 | 0.6 |
| | PH | Match | 100 | 100 | 100 | 70 | 100 | 0 | 100 | 40 | 0 | 0 | 0.3 |
| | | Freq., % | 54.2 | 23.5 | 13.1 | 6.0 | 2.9 | 0 | 1.4 | 1.0 | 0 | 0 | 0.5 |
| 8K | 2SNV | Match,% | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 0 | 0 | 0 | 0 |
| | | Freq., % | 53.1 | 23.7 | 12.5 | 6.5 | 2.3 | 1.25 | 0.7 | 0 | 0 | 0 | 0 |
| | PH | Match,% | 100 | 100 | 100 | 0 | 100 | 0 | 100 | 20 | 0 | 0 | 0.2 |
| | | Freq., % | 58.1 | 24.0 | 12.7 | 0 | 3.1 | 0 | 1.6 | 1.3 | 0 | 0 | 0.5 |
| 4K | 2SNV | Match,% | 100 | 100 | 100 | 100 | 100 | 100 | 20 | 0 | 0 | 0 | 0 |
| | | Freq., % | 53.7 | 23.7 | 12.3 | 6.5 | 2.4 | 1.2 | 0.9 | 0 | 0 | 0 | 0 |
| | PH | Match,% | 100 | 100 | 100 | 0 | 70 | 0 | 10 | 0 | 0 | 0 | 0.3 |
| | | Freq., % | 60.1 | 23.9 | 12.8 | 0 | 3.5 | 0 | 2.5 | 0 | 0 | 0 | 0.5 |
| 2K | 2SNV | Match,% | 100 | 100 | 100 | 100 | 100 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | Freq., % | 55.2 | 23.4 | 12.5 | 6.9 | 2.0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | PH | Match,% | 100 | 100 | 90 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.3 |
| | | Freq., % | 60.4 | 24.3 | 15.6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.4 |
| 1K | 2SNV | Match,% | 100 | 100 | 100 | 100 | 10 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | Freq., % | 56.7 | 23.7 | 12.7 | 6.6 | 3.2 | 0 | 0 | 0 | 0 | 0 | 0 |
| | PH | Match,% | 90 | 90 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 |
| | | Freq., % | 72.8 | 26.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.4 |
| 0.5K | 2SNV | Match,% | 100 | 100 | 100 | 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | Freq., % | 62.0 | 23.7 | 12.8 | 7.6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | PH | Match,% | 50* | 50* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | Freq., % | 69.9 | 30.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Table 2: Comparison of 2SNV and PredictHaplo on full and sub-sampled data. For all 33.5K reads, the sign “✓” (respectively, “×”) denotes fully matched (respectively, unmatched) true variant and the column FP reports the number of incorrectly predicted variants (false positives) and their total frequency. For each sub-sample size (16K,...,0.5K), the table reports the percent of runs when a variant is completely matched and its average frequency. Similarly, the column FP reports the average number of false positive variants and their average total frequency.

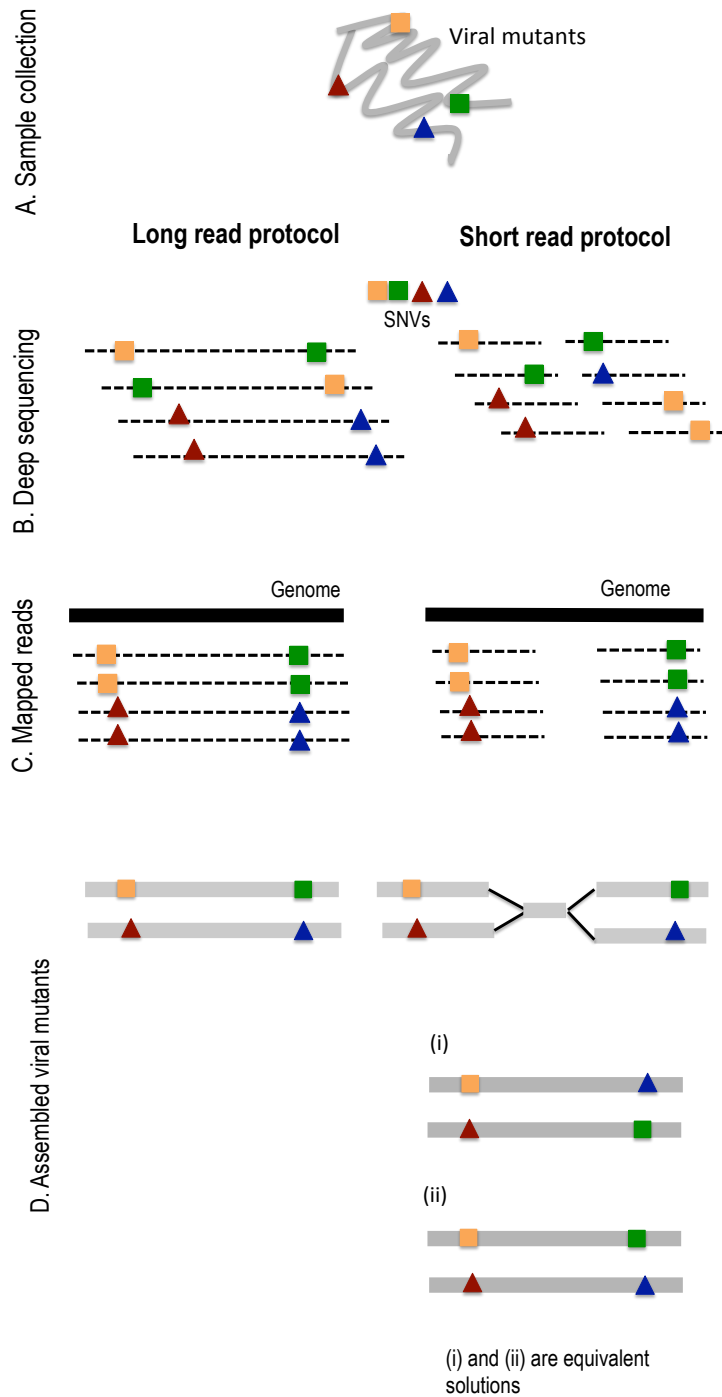


Figure 1: Overview of the long single-molecule sequencing protocol. (a) Extract the viral genomic DNA from the whole blood sample. (b) DNA material from the viral mutants is cleaved into sequence fragments using any suitable restriction enzyme. Amplified fragments are sequenced. (c) Long single-molecule reads are mapped to the reference genome. (d) SNVs are detected and assembled into the viral mutant variants. The short read protocol produces equivalent solutions.

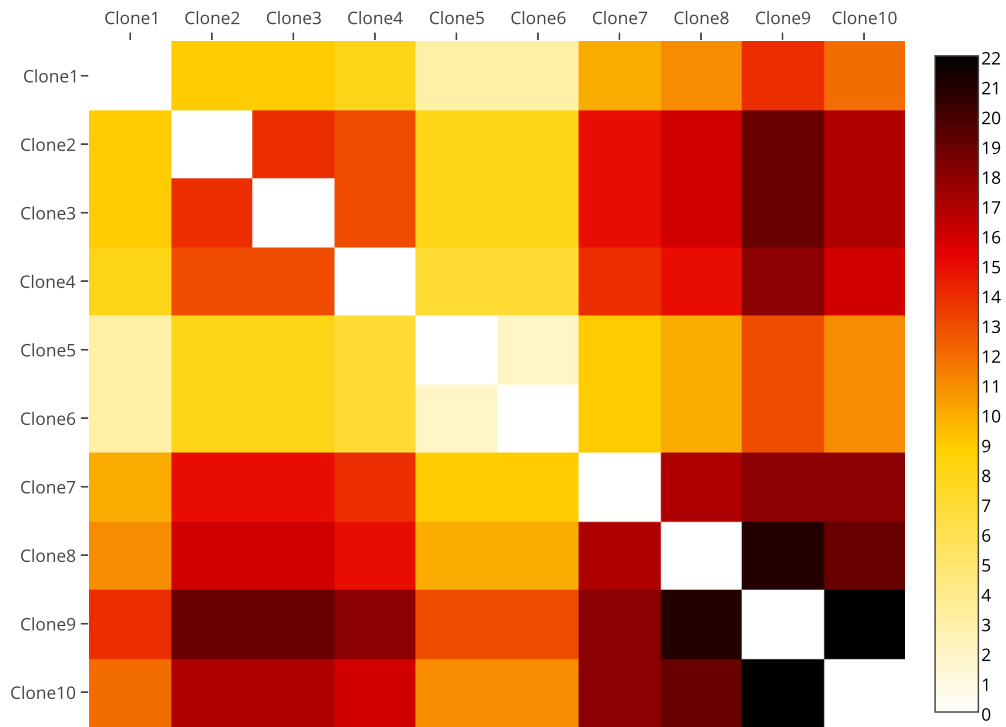


Figure 2: The heatmap representing pairwise edit distance between the 10 IAV clones.

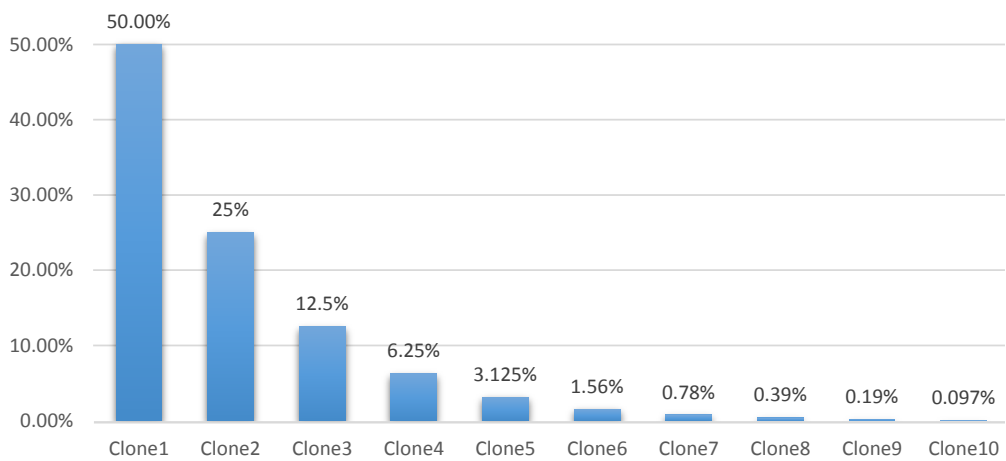


Figure 3: Frequency distribution of clones in the mixture.

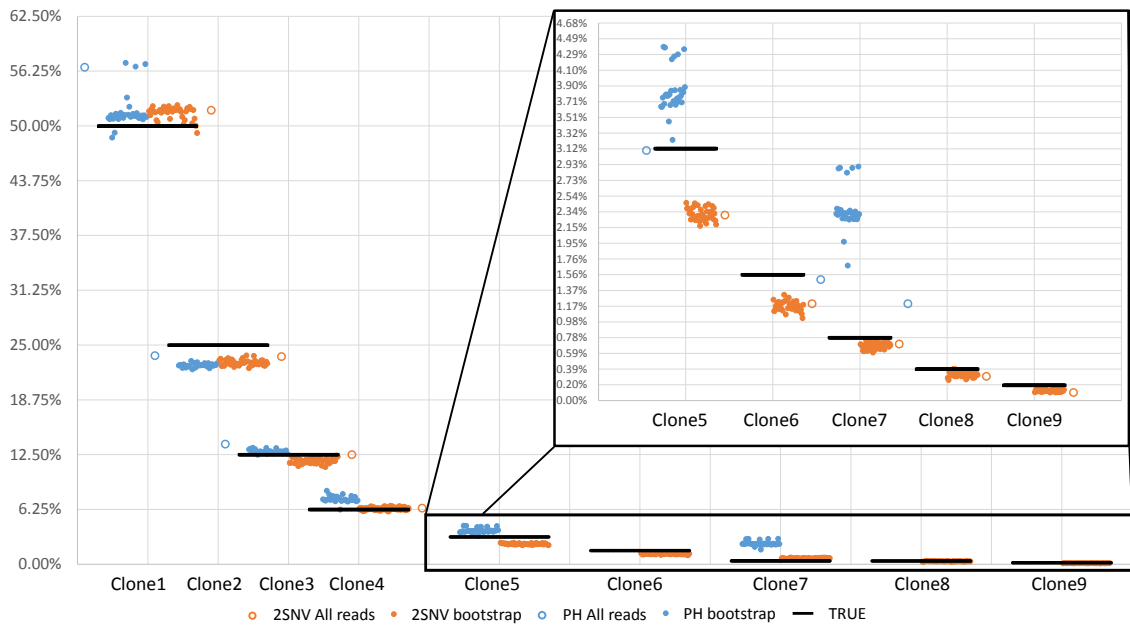


Figure 4: The results of running 2SNV and PredictHaplo (PH) on the original sample with all 33558 reads and on 40 bootstrapped samples (only 35 runs of PH were successful), y axis labels are clone frequencies and x axis labels are clone ids. Horizontal black bars are representing true clone frequency and colored dots are representing frequency reported by corresponding method in each of 40 runs. Clones 6 and 9 were never reconstructed by PH and clone 8 was reconstructed only on full data.

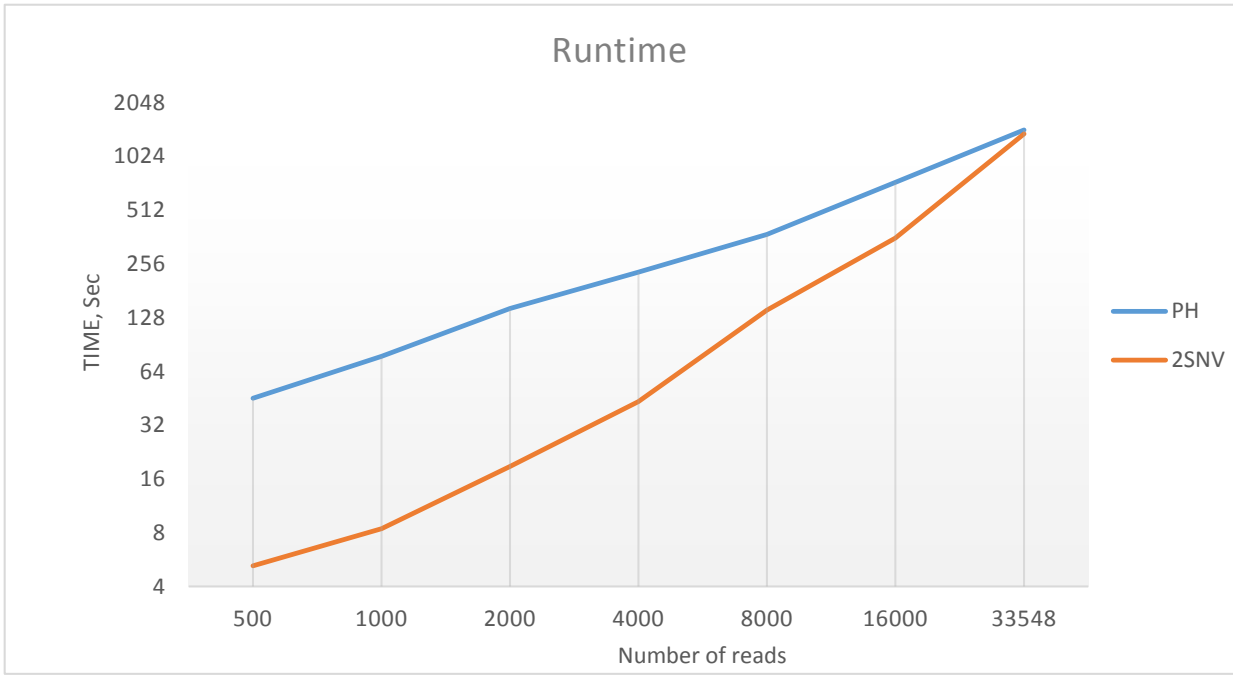


Figure 5: Runtime of PredictHaplo (PH) and 2SNV on datasets with different sizes. The runtime of 2SNV includes processing of alignment with b2w.

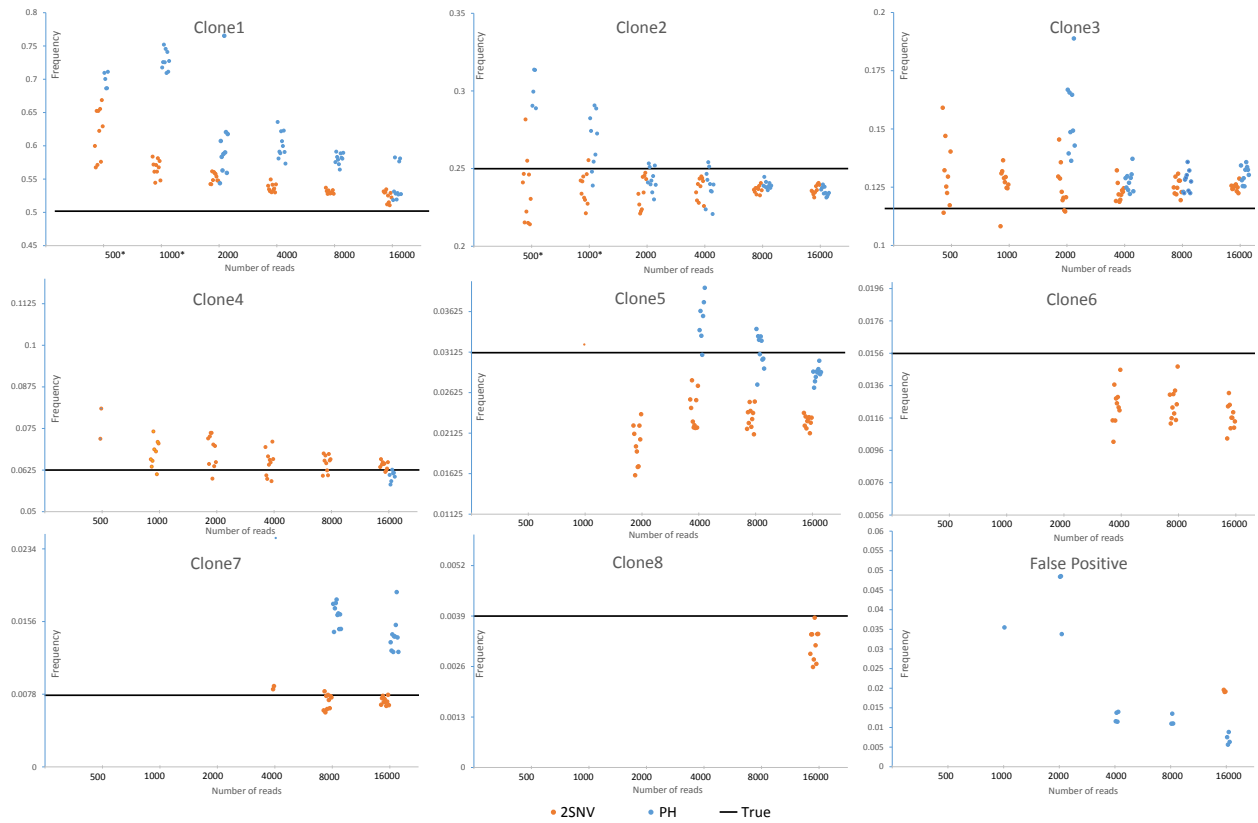


Figure 6: Dependency of accuracy on coverage represented by the number of reads. N reads ($N = 500, 1000, 2000, 4000, 8000, 16000$) were randomly selected 10 times from the original data and the both methods 2SNV and PredictHaplo were applied. For $N = 500^*$ and $N = 1000^*$ PredictHaplo gave results only in 5 and 9 runs, respectively. Each dot represents the reconstructed frequency of the clone in the respective runs.