

Figures for Simultaneous Clonal Subpopulation Inference from Single Cells and Bulk Sequencing Data

December 20, 2016

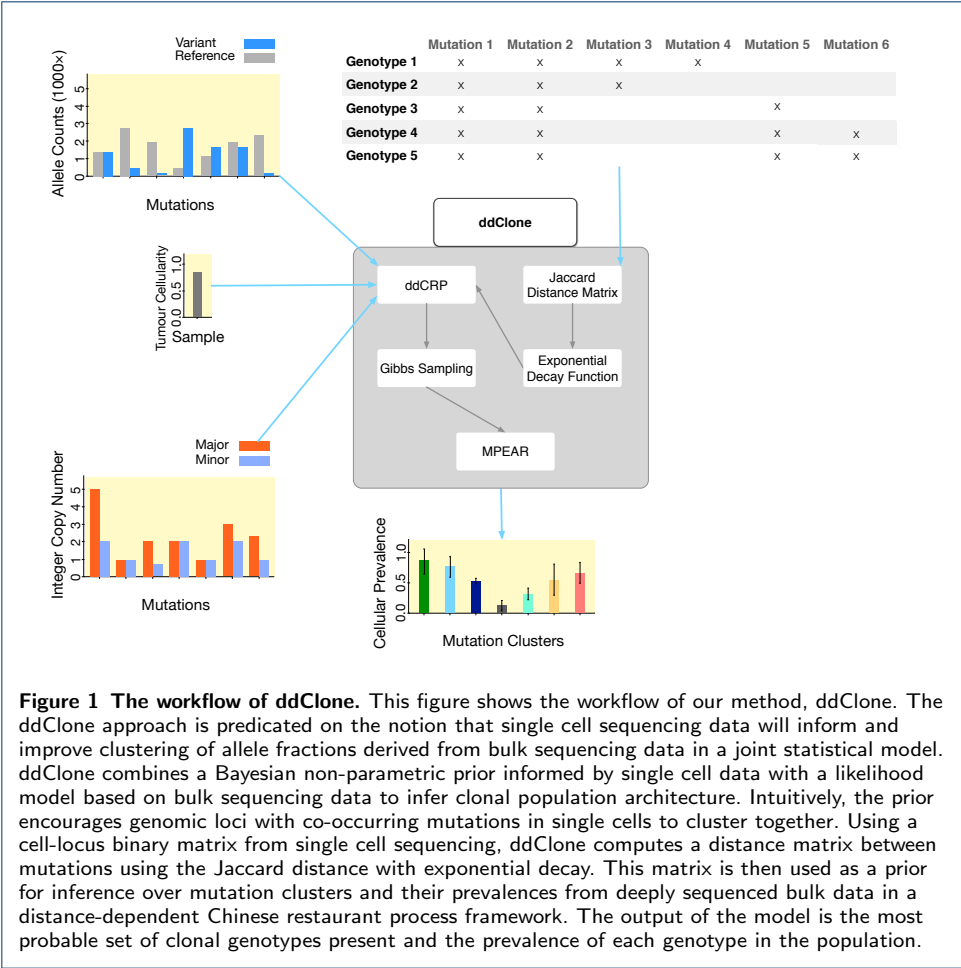


Figure 1 The workflow of ddClone. This figure shows the workflow of our method, ddClone. The ddClone approach is predicated on the notion that single cell sequencing data will inform and improve clustering of allele fractions derived from bulk sequencing data in a joint statistical model. ddClone combines a Bayesian non-parametric prior informed by single cell data with a likelihood model based on bulk sequencing data to infer clonal population architecture. Intuitively, the prior encourages genomic loci with co-occurring mutations in single cells to cluster together. Using a cell-locus binary matrix from single cell sequencing, ddClone computes a distance matrix between mutations using the Jaccard distance with exponential decay. This matrix is then used as a prior for inference over mutation clusters and their prevalences from deeply sequenced bulk data in a distance-dependent Chinese restaurant process framework. The output of the model is the most probable set of clonal genotypes present and the prevalence of each genotype in the population.

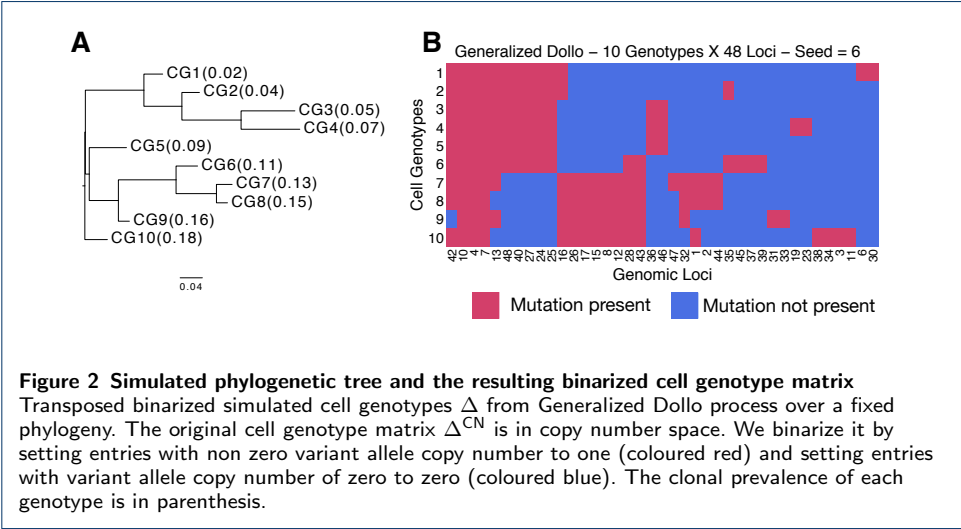
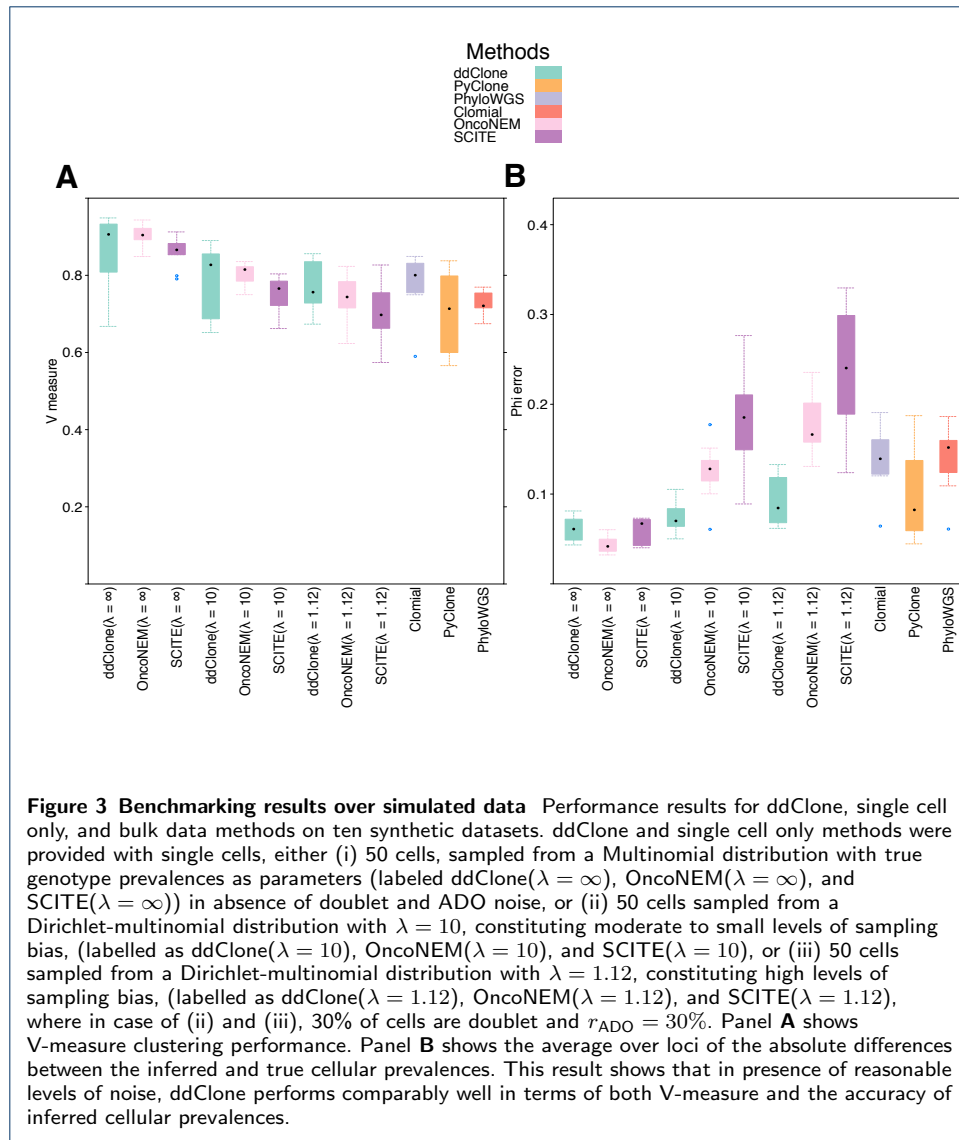
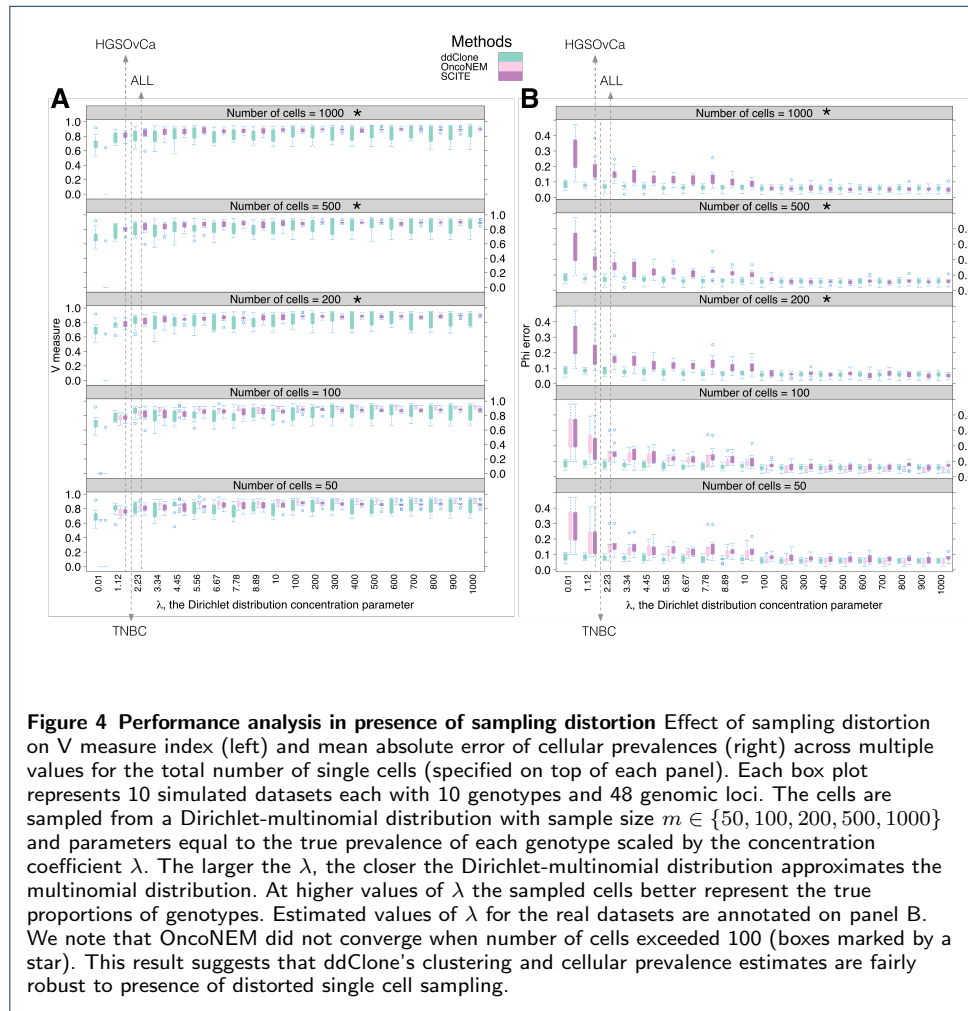


Figure 2 Simulated phylogenetic tree and the resulting binarized cell genotype matrix
Transposed binarized simulated cell genotypes Δ from Generalized Dollo process over a fixed phylogeny. The original cell genotype matrix Δ^{CN} is in copy number space. We binarize it by setting entries with non zero variant allele copy number to one (coloured red) and setting entries with variant allele copy number of zero to zero (coloured blue). The clonal prevalence of each genotype is in parenthesis.





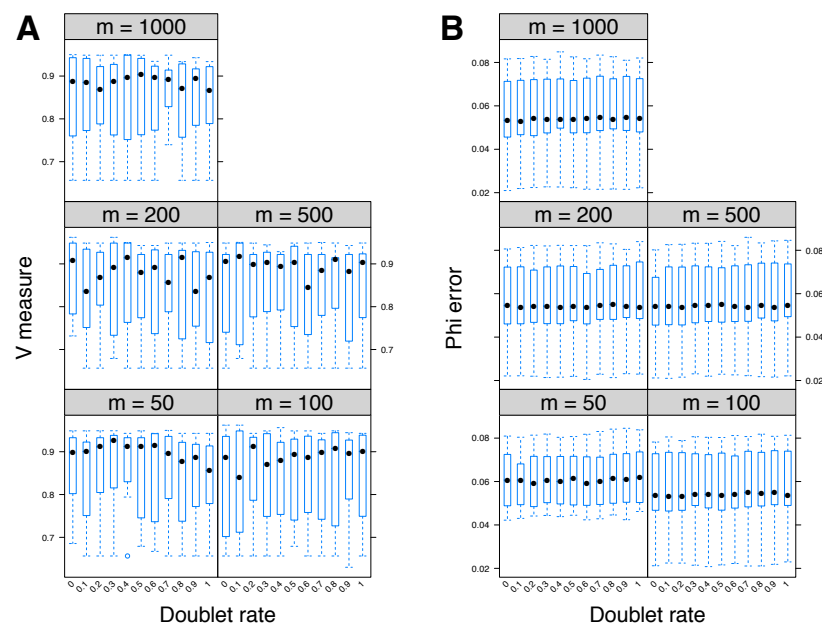


Figure 5 Performance analysis in presence of doublets Effect of presence of doublets on V measure index (left) and mean absolute error of cellular prevalences (right) across multiple values for the total number of single cells (specified as m on top of each panel). Each box plot represents 10 simulated datasets each with 10 genotypes and 48 genomic loci. The cells are sampled from a multinomial distribution with sample size of m and parameters equal to the true prevalence of each genotype. Progressively increasing percentage of doublet cells results in minor degrading performance in cellular prevalence estimate. Overall, this result suggests that ddClone's cellular prevalence estimates are robust to presence of uncorrected doublet noise.

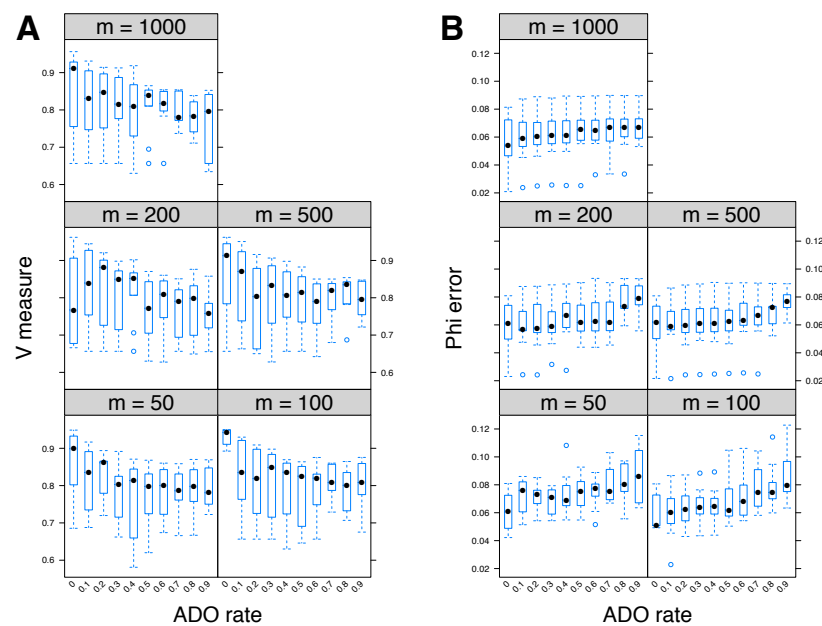
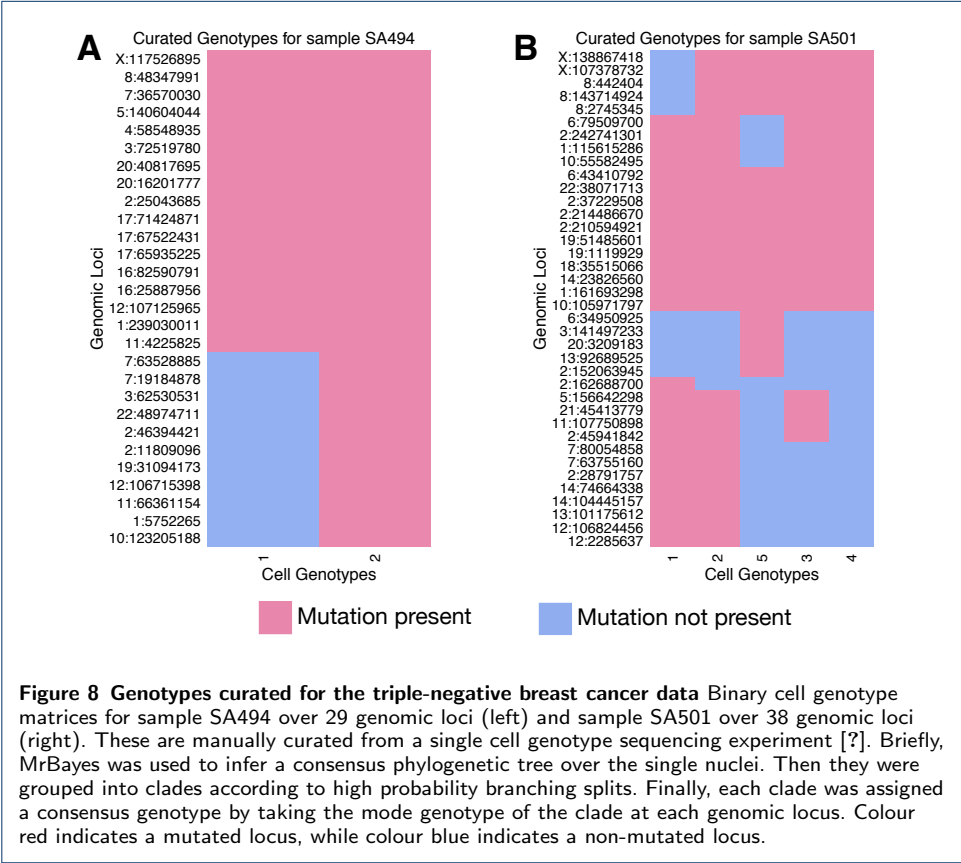
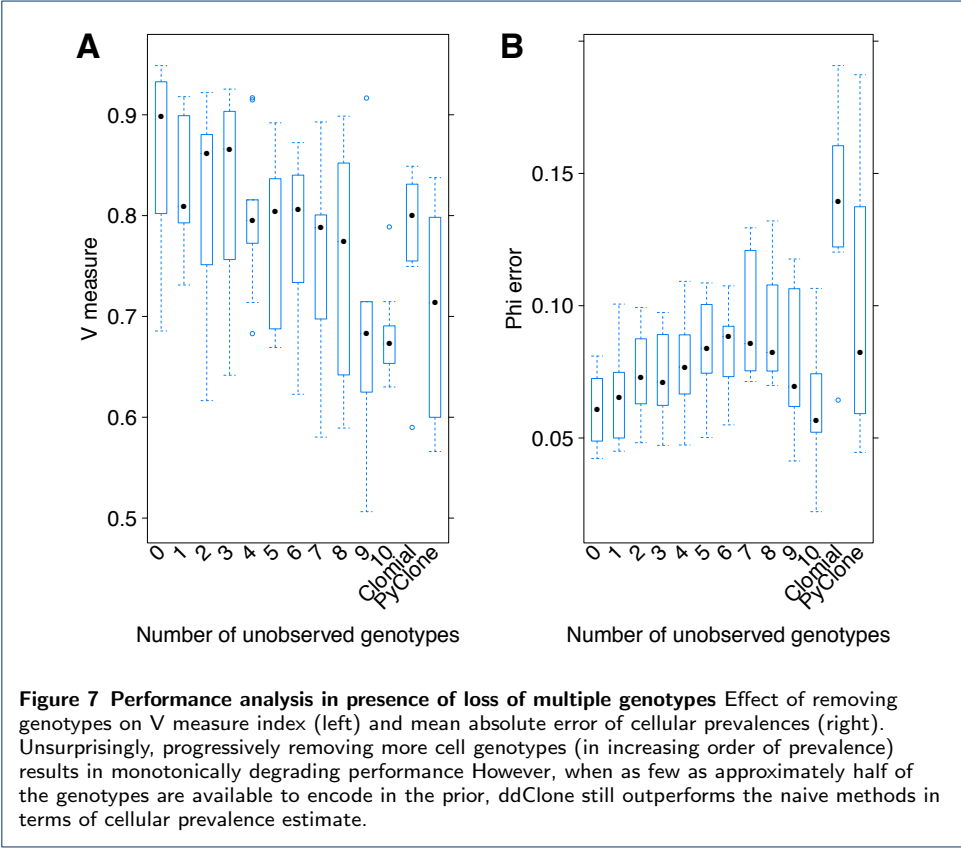


Figure 6 Performance analysis in presence of allele drop outs Effect of presence of allele drop outs (ADO) on V measure index (left) and mean absolute error of cellular prevalences (right) across multiple values for the total number of single cells (specified as m on top of each panel). Each box plot represents 10 simulated datasets each with 10 genotypes and 48 genomic loci. The cells are sampled from a multinomial distribution with sample size of m and parameters equal to the true prevalence of each genotype. As expected, progressively increasing the ADO rate results in degrading performance in both clustering and cellular prevalence estimates. The detrimental effect dampens as the number of sampled cells increases.



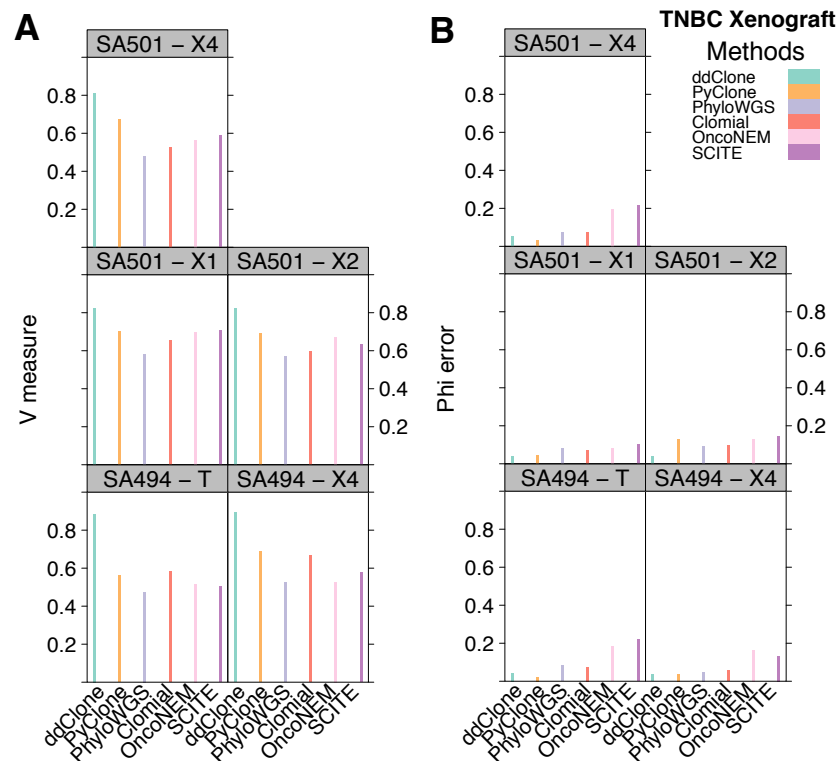
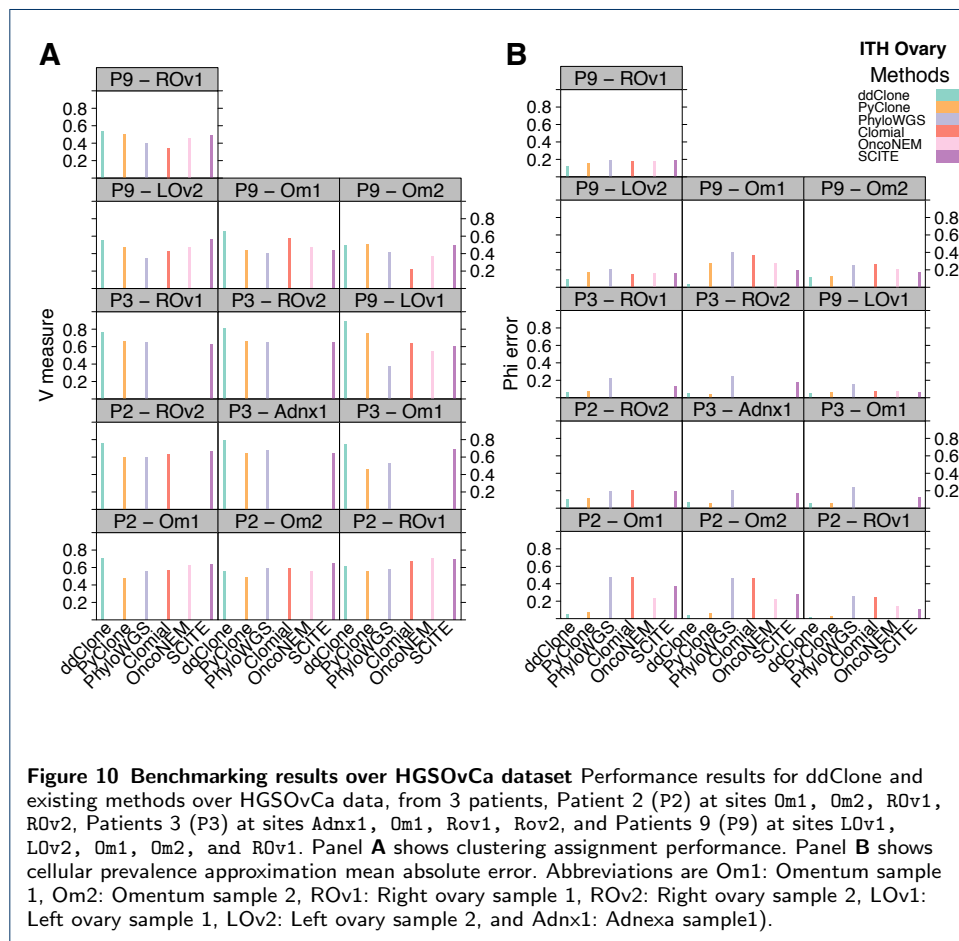


Figure 9 Benchmarking results over TNBC dataset Performance results for ddClone and existing methods over TNBC SA501 X1, X2, X4, and SA494 T, X4. Panel **A** shows clustering assignment performance. Panel **B** shows cellular prevalence approximation mean absolute error. Evaluated against multi-sample PyClone, ddClone outperforms the second best performing method (PyClone) in terms of V-measure (Wilcoxon rank sum test with p-value < 0.05) and performs as well (SA494, timepoint T) or better (all the other timepoints) than the second best performing method in terms of accuracy of inferred cellular prevalences.



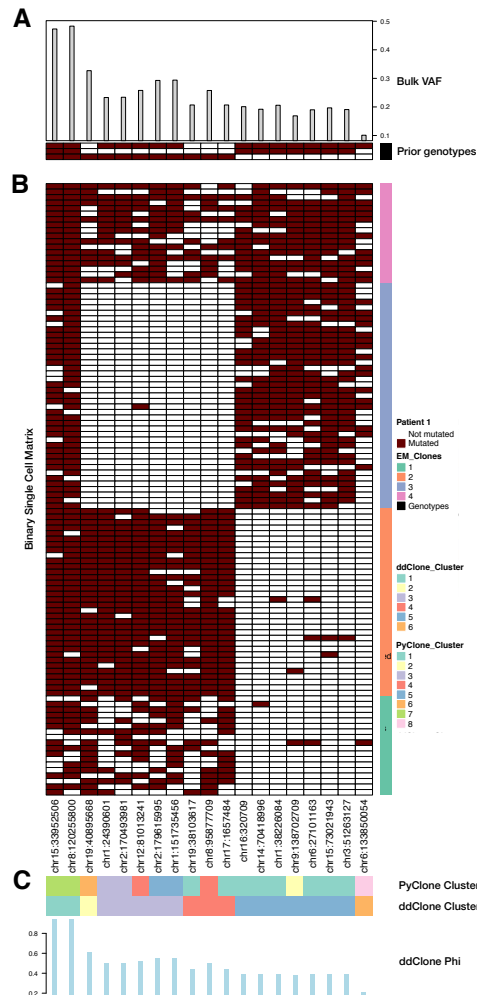


Figure 11 Analysis results of an acute lymphoblastic leukemia dataset [?] Analysis results of a patient with ALL (Patient 1) [?]. The variant allele frequencies from the bulk data (panel A, top) along with the consensus genotypes estimated from the binary cell matrix (panel A, bottom). These two constitute the input to the ddClone model. We note that the binary cell matrix (B) is displayed here for comparison and is not an input to ddClone. This binary cell matrix was used in [?] to cluster the cells into clones (vertical bar at the right side of the figure) and consensus genotypes (bottom part of panel A). ddClone clusters mutations into 6 groups (panel C, top) and estimates cellular prevalence (Φ) for each (panel C, bottom). ddClone's estimated Φ are highly correlated with the corrected bulk VAFs ($R^2 = 0.98$, also see Additional file 1) suggesting that it does not introduce unreasonable structure in the data. Furthermore, when there is evidence in the bulk, it can override its prior and splits clusters as necessary. For instance, even though locus chr19:40895668 has the same prior genotype as loci in cluster 4, its VAF in the bulk data is 1.5 times that of the mean of loci in cluster 4. This hints at a finer structure in cluster 4 and ddClone has automatically assigned chr19:40895668 to a separate cluster.

