

Rapid genotype imputation from sequence without reference panels

Robert W. Davies¹, Jonathan Flint², Simon Myers^{1, 3, 5}, Richard Mott^{1, 4, 5}

¹ Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, UK

² Center for Neurobehavioral Genetics, Semel Institute for Neuroscience and Human Behavior, University of California, Los Angeles, USA

³ Department of Statistics, University of Oxford, Oxford, UK

⁴ UCL Genetics Institute, University College London, London, UK

⁵ These authors contributed equally to this work

Abstract

Inexpensive genotyping methods are essential for genetic studies requiring large sample sizes. In human studies, array-based microarrays and high-density haplotype reference panels allow efficient genotype imputation for this purpose. However, these resources are typically unavailable in non-human settings. Here we describe a method (STITCH) for imputation based only on sequencing read data, without requiring additional reference panels or array data. We demonstrate its applicability even in settings of extremely low sequencing coverage, by accurately imputing 5.7 million SNPs at a mean r^2 of 0.98 in 2,073 outbred laboratory mice (0.15X sequencing coverage). In a sample of 11,670 Han Chinese (1.7X), we achieve accuracy similar to alternative approaches that require a reference panel, demonstrating that this approach can work for genetically diverse populations. Our method enables straightforward progression from low-coverage sequence to imputed genotypes, overcoming barriers that at present restrict the application of genome-wide association study technology outside humans.

Introduction

Over the last decade, genome-wide association studies (GWAS) have detected thousands of loci associated with complex traits in the human genome.¹ Generally, these involve genotyping 0.5-1M SNPs on DNA genotyping microarrays, and then employing externally generated haplotype reference panels such as those provided by the HapMap² and the 1000 Genomes Project³ to infer genotypes at tens of millions of additional sites, employing algorithms for phasing⁴ and imputation⁵⁻⁸.

In non-human species, large haplotype reference panels for fully genome-wide imputation are typically not available, and this fact has necessitated study designs incorporating high sample relatedness, and directed breeding, which can be leveraged to improve array-based genotype imputation⁹⁻¹². Moreover, inter-population differences within non-human species can further complicate genotyping array design and use. Arrays may work poorly when populations other than those used to design the chip are analyzed¹³, requiring the expensive development of either dense arrays, or many less dense, population specific arrays.

Given these issues, an attractive low-cost alternative to the use of arrays is to use low coverage next generation sequencing (LC-NGS) as a basis for imputation of complete genotypes¹⁴. Genotyping by LC-NGS could in principle be powerful:

even at modest sequencing depth LC-NGS reads sample the majority of segregating sites, including those specific to the population of interest.

However, to date, there exist no published genotype imputation methods specifically designed to use LC-NGS without additionally using genotyping microarrays or a haplotype reference panel. While methods such as Beagle (version 4)⁷ and findhap (version 4)¹² can be applied in this setting, they are tailored to work best with reference panels, and array data, which provide a framework of high-confidence genotypes.

Furthermore, the read-based nature of NGS provides useful phasing information on nearby variants from a single (paired) read, which is likely to be especially powerful in species with high SNP densities. While some approaches have started to use phase informative reads when phasing multiple heterozygous SNPs⁴, there are additional benefits to a fully read based imputation framework. First, in the absence of reference haplotypes, phasing with reads may help to initialize the phasing procedure. Second, at high SNP density, it is inaccurate to treat genotypes from the same read as being independent, which may help mitigate the influence of incorrectly mapped reads, *e.g.* near an indel, a cluster of false positive SNPs will contribute only once to the model, and not multiple times.

Here, we describe a genotype imputation algorithm STITCH (Sequencing To Imputation Through Constructing Haplotypes) suitable for population samples in any species sequenced at low coverage without requiring a haplotype

reference panel. The only requirement is a high-quality reference assembly for read-mapping. We demonstrate STITCH's utility on two datasets, from two species: first a set of 2,073 CrI:CFW(SW)-US_PO8 (CFW) mice sequenced to 0.15X, and second a set of 11,670 Han Chinese samples sequenced to 1.7X¹⁵.

Results

Overview of model

Our method, STITCH, models each chromosome in the population as a mosaic of K unknown founders or ancestral haplotypes. For fully outbred settings, these haplotypes can be thought of as informally capturing the set of distinct haplotypes within a region, so K may be large. We employ a hidden Markov model (HMM), whose parameters are sequentially updated using expectation maximization (EM), similar in spirit to the fastPHASE algorithm¹⁶. At each iteration of the EM, in the expectation step ancestral haplotype probabilities are generated for each sample, while in the maximization step ancestral haplotypes and other parameters are updated using sample haplotype membership. Both of these steps directly consider the underlying sequencing reads. An overview of the model used for imputation is presented in **Figure 1**.

Computationally, the algorithm has a per-iteration time complexity linear in the number of samples and SNPs, and when run in its standard “diploid” mode, has quadratic time complexity in the number K of ancestral haplotypes. Because the ability to model large K is essential for STITCH to handle human and other large

outbred populations, we developed an alternative mode, termed “pseudo-haploid”, with linear per-iteration time complexity in K . This is motivated by the observation that imputation in diploid individuals could be carried out with linear time complexity in K if the sequencing reads came with labels indicating their parental chromosomal origin (maternal, or paternal) – in other words, with phase information. In this setting, only reads mapping to the maternal chromosome would be required to impute mutations this chromosome carries, so the maternal and paternal chromosomes could be imputed separately. In the absence of such chromosome labels, in theory one could sample labels (e.g. by Gibbs sampling) for each read under our model. The relative contribution of a read to each chromosomes’ posterior likelihood would then depend on the probability it came from that chromosome. This sampling is in practice prohibitively computationally expensive. Therefore, in practice our pseudo-haploid makes several additional simplifying approximations (see the **Supplementary Note** for details, and discussion of other issues e.g. label-switching), to estimate, for every read, the probability it came from each chromosome. Given these probabilities, we can update the posterior ancestral haplotype probabilities for each chromosome separately, within the EM algorithm. This retains a common EM framework to both diploid and pseudo-haploid modes, thereby allowing the algorithm to switch between modes at any point. A full description of the diploid and pseudo-haploid modes is in the Methods and **Supplementary Note**, while guidance on parameter choice is in the **Supplementary Note**.

CFW outbred mice

We ran STITCH on low coverage sequence data (0.15X, paired end 100 base pair reads) from 2,073 outbred Crl:CFW(SW)-US_PO8 mice^{17,18}. These mice are thought to have descended from two outbred founders (i.e. K=4) about 100 generations ago. We imputed genotypes at 7.1 million single nucleotide polymorphic sites (SNPs) that either were polymorphic in the Mouse Genomes Project¹⁹ or passed VQSR²⁰ quality filtration¹⁷. Imputation accuracy was assessed in two ways – 44 mice genotyped on the Illumina MegaMUGA array (21,576 polymorphic SNPs) and four mice sequenced to 10X using an Illumina HiSeq. Correlations (r^2) between genotypes and imputed dosages were calculated either per-site for the array or aggregated across all SNPs in a given frequency range for the high coverage sequencing data.

We compared results from STITCH (K=4, diploid mode) to Beagle and findhap run without a reference panel^{7,12}. Genotypes across all frequencies from STITCH correlated highly with the Illumina MegaMUGA array (**Fig. 2a**, $r^2=0.972$) and 10X sequencing (**Fig. 2b**, $r^2=0.948$) (**Supplementary Table 1**). Filtering with an imputation info score > 0.4 (Methods) and Hardy-Weinberg Equilibrium (HWE) p-value > 10^{-6} improved accuracy further, to r^2 of 0.981 and 0.974, with 5.72M SNPs (81%) retained. In general, imputation performance was good across all allele frequencies, except for a slight decrease at low frequency (<5%) SNPs (**Fig. 2**) that are expected to be challenging for low-coverage sequencing. Beagle under default conditions achieved r^2 's of 0.080 and 0.219 compared to 10X

sequencing and array without QC filtering, respectively, while findhap achieved 0.58 and 0.55 (**Fig 2, Supplementary Table 1**).

We performed additional analyses to explore parameter choices. We found that optimal results were achieved with K=4 for STITCH (**Supplementary Table 2**), as expected from the population's ancestry. Results for Beagle did not differ appreciably when the number of iterations, window size, and model scale factor were changed, while the results reported above for findhap were the best observed when varying parameters of the method over a range of values suggested by the findhap documentation (**Supplementary Table 2**). Of the 3 methods, findhap was approximately 12 times faster than Beagle and 38 times faster than STITCH, although if parallelized by chromosome, imputation for all samples for any of the methods could be performed in less than 48 hours on a modest computational server. Ignoring phase information from reads in applying STITCH (*i.e.* treating each variant in a read as independent) reduced accuracy considerably, from $r^2=0.97$ to 0.87 with the Illumina MegaMUGA array (**Supplementary Table 2**).

CONVERGE study

To explore performance in human data, we ran STITCH on low coverage sequence data (1.7X, paired end 83 base pair reads) from 11,670 Han Chinese women¹⁵. Details of read mapping and variant calling are as detailed previously¹⁵. We used the first 10 Mbp of chromosome 20 to test the imputation algorithms and compared our predictions with genotypes from 72 individuals

genotyped on the Illumina HumanOmniZhongHua-8 array and 9 individuals sequenced at 10X coverage¹⁵.

Following preliminary testing (**Supplementary Table 3**), we applied STITCH with $K=40$ “founder” haplotypes, with 40 rounds of updating to estimate parameters and perform imputation. The first 38 rounds were in the faster “pseudo-haploid” mode and the final 2 in the slower but more accurate “diploid” mode. STITCH achieved close correspondence to Illumina array results (**Fig. 3A**, $r^2=0.920$, **Supplementary Table 4**) and 10X sequencing (**Fig. 3B**, $r^2=0.949$), results that improved when SNPs were filtered ($\text{info} > 0.4$, HWE $p\text{-value} > 10^{-6}$) (array $r^2=0.939$, 10X $r^2=0.960$). Accuracy declined for $K < 40$, and was marginally improved for $K > 40$. Running additional slower “diploid” mode iterations also improved accuracy only marginally, and fully diploid imputation became computationally prohibitive beyond $K=30$. Results were essentially unchanged when STITCH was run ignoring read information, reflecting the low SNP density in humans. Beagle without a reference panel achieved reduced r^2 values of 0.886 and 0.930 for sequencing and array without QC filtering, respectively, while the best parameter settings we identified for findhap achieved 0.414 and 0.550 (**Fig 3, Supplementary Table 4**).

We next compared STITCH to applying Beagle with additional reference panel information. In this setting, Beagle is modestly more accurate than STITCH, at the cost of run time (**Fig. 3C, Fig. 3D, Supplementary Table 5, Supplementary Table 6**). For example, Beagle achieved an r^2 of 0.943 versus 0.922 compared to the array for STITCH before SNP QC, although it took 7.3X as long.

We then repeated the imputation strategy for Beagle used in the original CONVERGE study¹⁵ of first imputing all sites without a reference panel, then imputing the subset of variants with a reference panel, and replacing SNPs in the former with the latter when they existed. We compared these results to those from STITCH, run without a reference panel on the entire set of SNPs. Results between these two strategies were essentially the same between STITCH and Beagle, (**Supplementary Table 7**) with STITCH achieving an r^2 of 0.972 (array) and Beagle 0.968 under the most stringent QC scenario, which retained 75% of common sites (>5% minor allele frequency). Results for STITCH were generated 5.3X faster than Beagle under this strategy.

Effect of sample size and coverage on imputation

We next examined the consequences of altering sample size and sequence coverage (**Fig. 4**). For the CFW mice, for the full 0.15X coverage using STITCH, sample size above 500 has little impact on performance, while at down-sampled lower coverage, increasing sample size to 2,073 leads to substantially increased performance. Surprisingly, even at 0.06X for the full sample of 2,073 animals, results are only marginally poorer than using 0.15X. For the CONVERGE samples using STITCH, sample size has less of an influence across the range of sequencing coverages considered, although results did consistently improve with increasing sequencing depth. STITCH outperforms Beagle without a reference panel over the range of low coverages considered here (0.3-1.7X).

Effect of variant filtration on imputation performance

Methods of genotyping from next generation sequencing typically employ an initial step of variant filtration, to reject any newly discovered sites whose quality control metrics differ from those at known variant sites. One such method is the GATK Variant Quality Score Recalibrator²⁰. Since we developed STITCH to be applicable to populations in which a catalogue of variant sites was unavailable, we investigated whether prior variant filtration was necessary, or whether STITCH itself could be used to filter SNPs directly. We compared imputation at filtered variant sites in the CFW population, as defined using VQSR and known variable sites, to a two-step strategy where no prior variant catalogue is used. As a first step, all discovered sites in the sample are imputed without filtration. In the second step, only those variants that pass quality control (QC) filters from the first step are re-imputed. Results indicate this to be a viable strategy (**Supplementary Table 8**). For the one step strategy with variant filtration (the original study design), 152K SNPs on chromosome 19 were imputed, with 122K SNPs passing QC at an r^2 of 0.968. For the two-stage approach, 355K variants were imputed in the first round of imputation, with 136K passing QC. In the second round of imputation, these 136K were re-imputed, with 128K of them passing QC at an r^2 of 0.952. Overlap between the two approaches was 116K, with marginally better r^2 in the overlap from the two-step approach, with results specific to either set having lower r^2 . These results indicate that prior knowledge of variable sites is not needed to impute accurately using STITCH.

Discussion

Inexpensive genotyping microarrays and imputation with large reference panels have made genome-wide association studies tractable in humans, but these resources are unavailable in many species, and are not ideal for human populations in parts of the world where appropriate reference populations have not yet been deeply sequenced. Our method alleviates this bottleneck, by imputing high quality genotypes directly from low coverage sequencing data. The method delivered highly accurate imputation at a depth of only 0.15X in the CFW mouse population. In a higher coverage situation of 1.7X in humans, STITCH performed similarly to a method using a reference panel, without requiring such a panel. This simplifies the imputation pipeline, and allows application in populations where no reference is available. We also introduce an approximation that achieves linear as opposed to quadratic time scaling with the number of founder haplotypes with very little loss of accuracy, making the method suitable for the analysis of very large and ancestrally complex populations.

Importantly, imputation results were better when using direct phase information, especially in the CFW mice, while the two-stage CFW imputation procedure showed that careful filtering of candidate SNPs based on prior variation is not essential. This can simplify the analysis, by alleviating the need for running separate SNP-filtering procedures, *e.g.* the VQSR²⁰.

The differences in imputation performance we observe in the mouse and human samples reflect their different genetic histories. Our method involves two alternating processes – reconstructing founder haplotypes, and determining in an individual sample which pair of founder haplotypes it is most similar to at each locus. Because the CFW mouse population was founded about 100 generations ago from just two progenitors, physical distances between haplotype switches are large, making it relatively easy to identify which of the small number of founder haplotypes an individual carries, even at low coverage. By contrast, in the CONVERGE Han Chinese population sample, haplotypic diversity is far greater, and consequently haplotypic switches occur much more frequently. This explains why imputation in humans is less accurate for a given level of sequence coverage and why increasing K – the modeled number of founder haplotypes – had little influence on performance in mice, but increased accuracy in humans. Because these datasets represent relatively extreme scenarios in terms of haplotypic diversity, we expect that STITCH will work well in intermediate settings, without haplotype reference panels.

Our method delivers the greatest accuracy improvements for populations with recent strong bottlenecks, such as those studied in agricultural or plant genomics^{21,22,23}. While poorer quality reference assemblies than those available for mice and humans will impact the performance of STITCH in other species, in future the decreasing cost of constructing high quality reference assemblies using single molecule long read technologies and optimal mapping techniques may mitigate this issue²⁴.

Although STITCH out-performed findhap for imputation using low-coverage sequencing data in the scenarios evaluated here, in cases where additional genotyping array data is also available, findhap may perform well. Specifically, if additional microarray data is available for a set of samples drawn from the same population as those sequenced at low coverage, findhap obtained comparable accuracy to our STITCH runs that used the read unaware option¹², and offers a speed advantage.

For human samples, at 1-1.5X sample coverage, STITCH accurately imputes all common variation, making it suitable for any population that lacks a reference panel, or one with an incomplete variant catalogue. We imagine that our method might be particularly appropriate for ethnic groups so far not subject to GWAS, population isolates, and for ancient humans, where low coverage sequencing is common.

Acknowledgements

R. W. Davies is supported by a grant from the Wellcome Trust 097308/Z/11/Z. S. R. M is supported by Investigator Award 098387/Z/12/Z. This work was funded by the Wellcome Trust (WT090532/Z/09/Z, WT083573/Z/07/Z, WT089269/Z/09/Z, WT098387/Z/12/Z)

Author Contributions

R.W.D, S.M, and R.M. developed the method. R.W.D. wrote the algorithm and performed analyses. All authors contributed to study design, and reviewed and contributed to the final manuscript.

Correspondence should be addressed to R.W.D. (rwdavies@well.ox.ac.uk)

References for main text

1. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42**, D1001–D1006 (2014).
2. Frazer, K. A. *et al.* A second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**, 851–861 (2007).
3. The 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
4. Delaneau, O., Zagury, J.-F. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat. Methods* **10**, 5–6 (2013).
5. Howie, B. N., Donnelly, P. & Marchini, J. A Flexible and Accurate Genotype Imputation Method for the Next Generation of Genome-Wide Association Studies. *PLoS Genet* **5**, e1000529 (2009).
6. Li, Y., Willer, C. J., Ding, J., Scheet, P. & Abecasis, G. R. MaCH: Using Sequence and Genotype Data to Estimate Haplotypes and Unobserved Genotypes. *Genet. Epidemiol.* **34**, 816–834 (2010).
7. Browning, S. R. & Browning, B. L. Rapid and Accurate Haplotype Phasing and Missing-Data Inference for Whole-Genome Association Studies By Use of Localized Haplotype Clustering. *Am. J. Hum. Genet.* **81**, 1084–1097 (2007).

8. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G. R. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat. Genet.* **44**, 955–959 (2012).
9. Swarts, K. *et al.* Novel Methods to Optimize Genotypic Imputation for Low-Coverage, Next-Generation Sequence Data in Crop Plants. *Plant Genome* **7**, 0 (2014).
10. Huang, B. E. & George, A. W. R/mpMap: A computational platform for the genetic analysis of multi-parent recombinant inbred lines. *Bioinformatics* **27**, 727-729 (2011).
11. Sargolzaei, M., Chesnais, J. P. & Schenkel, F. S. A new approach for efficient genotype imputation using information from relatives. *BMC Genomics* **15**, 478 (2014).
12. VanRaden, P. M., Sun, C. & O'Connell, J. R. Fast imputation using medium or low-coverage sequence data. *BMC Genet.* **16**, 82 (2015).
13. Didion, J. P. *et al.* Discovery of novel variants in genotyping arrays improves genotype retention and reduces ascertainment bias. *BMC Genomics* **13**, 34 (2012).
14. Pasaniuc, B. *et al.* Extremely low-coverage sequencing and imputation increases power for genome-wide association studies. *Nat. Genet.* **44**, 631–635 (2012).
15. CONVERGE Consortium. Sparse whole-genome sequencing identifies two loci for major depressive disorder. *Nature*, **523**, 588–591 (2015).
16. Scheet, P. & Stephens, M. A Fast and Flexible Statistical Model for Large-Scale Population Genotype Data: Applications to Inferring Missing Genotypes and Haplotypic Phase. *Am. J. Hum. Genet.* **78**, 629–644 (2006).

17. Nicod, J. *et al.* Genome-wide association of multiple complex traits in outbred mice by ultra low-coverage sequencing. *Nat. Genet.* In press
18. Yalcin, B. *et al.* Commercially Available Outbred Mice for Genome-Wide Association Studies. *PLoS Genet* **6**, e1001085 (2010).
19. Keane, T. M. *et al.* Mouse genomic variation and its effect on phenotypes and gene regulation. *Nature* **477**, 289–294 (2011).
20. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
21. Freedman, A. H. *et al.* Genome Sequencing Highlights the Dynamic Early History of Dogs. *PLoS Genet* **10**, e1004016 (2014).
22. The Bovine HapMap Consortium. Genome-Wide Survey of SNP Variation Uncovers the Genetic Structure of Cattle Breeds. *Science* **324**, 528–532 (2009).
23. Daetwyler, H. D. *et al.* Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle. *Nat. Genet.* **46**, 858–865 (2014).
24. VanBuren, R. *et al.* Single-molecule sequencing of the desiccation-tolerant grass *Oropetium thomaeum*. *Nature* **527**, 508–511 (2015).

Figure Legends

Figure 1. Overview of STITCH

After initializing various parameters (left), represented here by the ancestral haplotypes, 40 EM iterations are performed (middle). Each iteration involves i) determining hidden haplotype states (going down, left side) using current parameters and sample reads, and ii) parameter updates (going up, right side) using sample reads and haplotype probabilities (hidden states). Once the expectation-maximization iterations are completed, imputed genotypes are generated using the haplotype probabilities and ancestral haplotypes from the final iteration (right). This example uses real data from the CFW mice with K=4 founder haplotypes for approximately 3,000 base pairs on chromosome 19 containing 20 imputed SNPs. Each of the SNPs in the 4 reconstructed haplotypes are shown as a vertical bar split proportionally to the probability of emitting the reference (black) or alternate (grey). Sample reads are similarly coloured.

Figure 2. Performance of STITCH on CFW mice compared to external validation

Validation dataset is the Illumina MegaMUGA array (a) and 10X Illumina sequencing (b). Results are shown for STITCH (K=4, diploid mode), Beagle (default) and findhap (maxlen=10000, minlen=100, steps=3, iters=4) genome-wide for n=2,073 mice featuring 7.07M SNPs before QC and 5.72M after QC. STITCH is run using K=4, diploid method, 40 iterations. Post-QC is SNPs with info>0.4 and HWE p-value > 1×10^{-6} .

Figure 3. Performance of STITCH on CONVERGE humans compared to external validation

Validation dataset is the Illumina HumanOmniZhongHua-8 array (a, c) and 10X sequencing (b, d). Results are shown for STITCH (K=40, 38 pseudo-haploid iterations, 2 diploid iterations), Beagle (default (a,b), 3 iterations with reference panel (c,d)), and findhap (maxlen=50000, minlen=500, steps=3, iters=4) for the first 10 Mbp of chromosome 20 for n=11,670 Han Chinese samples, either for all SNPs (a,b), or for SNPs also present in the 1000 Genomes ASN reference panel (c,d). Post-QC is SNPs with info>0.4 and HWE p-value > 1×10^{-6} .

Figure 4. Effects of reduced sequence coverage

Results are shown for CFW mice (a) and CONVERGE humans using STITCH (b) and Beagle run without a reference panel (c). Validation is using array data, with each value representing the average for common SNPs (allele frequency 5–95%), without correction for post-imputation QC. Downsampling of samples and reads, as shown in the legends, was performed at random, except that samples necessary for accuracy assessment were always retained. STITCH settings are the same as for the full CFW, CONVERGE datasets. Colours representing downsampling sequence depth are the same for STITCH and Beagle.

Online methods

Overview and simulation under the model

Here we outline the model by describing how one would simulate (read) data from it, given knowledge of the underlying parameters. In the following section, we then more rigorously lay out how we infer parameters and perform inference of genotypes. We describe technical details of the EM procedure and parameter updating in the **Supplementary Note**.

We consider a population of individuals that can be approximated as having been founded G generations ago from K unknown ancestral founding haplotypes. Consider a haplotype from a single chromosomal region with T SNPs from a present day individual drawn from our model. A starting state (haplotype) k is chosen with probability π_k . Let d_t and p_t be the physical distance and average recombination rate between SNPs t and $t+1$, respectively. Therefore $\sigma_t = d_t p_t$ is the recombination distance between SNPs t and $t+1$ in one generation, so in the G generations since founding the probability of recombination between these SNPs is $1 - \exp\{-G\sigma_t\}$.

Conditional on recombination locations, we sample an ancestral haplotype for each non-recombining interval. We allow genetic drift to play a sizeable role in the proportions of the ancestral haplotypes in each short genomic interval. As such, we model the probability of choosing ancestral haplotype k to the right

(from SNP $t+1$), given a recombination between SNP t and SNP $t+1$ as $\alpha_{t,k}$. This choice is made independently of the state at SNP t .

Finally, the reads are sampled conditional on the local haplotype background. Gene conversion, de novo mutation, read-mapping errors and other issues mean that not all chromosomes and reads descended from ancestral haplotype k will be an exact match to the ancestral sequence. We therefore model that for each read, for SNP t and ancestral haplotype k , that the alternate base will be drawn with probability $\theta_{t,k}$, and the reference base with probability $1-\theta_{t,k}$. Inherent in this is the assumption that different reads are emitted from different samplings of $\theta_{t,k}$; in reality there would be a simple sampling of $\theta_{t,k}$ for each haplotype, and reads sampled conditional on these real underlying bases. This assumption is necessary for computational reasons, and has reduced impact for low coverage sequencing data.

Consider sampling the r^{th} read, R_r . To do this, first choose read boundaries and determine $u_{r,j}$ the set of indices of the SNPs in the read for $j=1,\dots,J_r$, where J_r is the number of SNPs in the r^{th} read R_r . Paired end reads can be easily accommodated in this way by allowing discontinuous $u_{r,j}$ within read r . We make the assumption that recombinations are infrequent enough that reads have constant haplotype state over their length; as such, each read has a central SNP, call it c_r , and state membership over the read is drawn from the central SNP. Therefore, the underlying “real” bases for the sequencing read are sampled according to $\theta_{t,k}$ for t in $u_{r,j}$ where $k=k_{t'}$ for $t'=c_r$. To sample “observed” bases, we sample $b_{r,j}$, the set of

base qualities of the SNPs in the read – in practice these qualities are externally provided - and then sample observed bases $s_{r,j}$ according to the real bases and the base qualities.

Expectation and hidden state determination

In the HMM, for the haploid model, let q_t be the hidden state at SNP t , *i.e.* $q_t \in \{1, \dots, K\}$. For the diploid model, let $q_t = (k_{t,1}, k_{t,2})$ be the hidden states at SNP t . Let $\lambda = \{\pi, \sigma, \alpha, \theta\}$ be the parameters of the model. The pseud-haploid model is described in the **Supplementary Note**.

Initial haploid state probabilities for the $k=1, \dots, K$ different states are defined as π_k . Diploid initial state probabilities are taken by multiplying together haploid state probabilities.

For state transitions, with probability $\exp\{-G\sigma_t\}$, no recombination occurs between SNPs t and $t+1$, while with probability $1-\exp\{-G\sigma_t\}$, a recombination occurs and a new state q_{t+1} is chosen at SNP t according to $\alpha_{t,k'}$ for $k'=k_{t+1}$. This gives the haploid transition matrix

$$P(q_{t+1} = k_{t+1} | q_t = k_t, \lambda) = \begin{cases} e^{-G\sigma_t} + (1 - e^{-G\sigma_t})\alpha_{t,k_{t+1}} & \text{if } k_{t+1} = k_t \\ (1 - e^{-G\sigma_t})\alpha_{t,k_{t+1}} & \text{if } k_{t+1} \neq k_t \end{cases}$$

Assuming independence between the two chromosomes then the diploid transition probability from state $q_t=(k_{t,1}, k_{t,2})$ at SNP t to $q_{t+1}=(k_{t+1,1}, k_{t+1,2})$ at SNP $t+1$ is:

$$P\left(q_{t+1} = (k_{t+1,1}, k_{t+1,2}) \mid q_t = (k_{t,1}, k_{t,2})\right) = P(q_{t+1} = k_{t+1,1} \mid q_t = k_{t,1}) \times P(q_{t+1} = k_{t+1,2} \mid q_t = k_{t,2})$$

For the emission of reads, for read R_r , let c_r be the index of the most central SNP in that read, choosing at random when a read intersects exactly two SNPs. Reads that don't intersect any SNPs are removed as they are uninformative. Consider the probability of an observation of a set of reads whose central SNP is t , or in other words $O_t=\{R_r \mid c_r=t\}$. For SNP j in read R_r , $s_{r,j}$ is the observed sequencing read (0 = reference, 1 = alternate), and $b_{r,j}$ is the Phred scaled base quality, *i.e.* the log probability that the base is called erroneously, so let $\epsilon_{r,j}=10^{(-b_{r,j}/10)}$. Then, given the underlying (unobserved) genotype of this read at this SNP is g

$$P(s_{r,j} \mid g) = \begin{cases} 1 - \epsilon_{r,j} & \text{if } s_{r,j} = g \\ \frac{1}{3} \epsilon_{r,j} & \text{if } s_{r,j} \neq g \end{cases}$$

For convenience, set $\phi^i_{r,j}=P(s_{r,j} \mid g=i)$. We disregard sequenced bases which are not the reference or alternate base. Paired end reads are handled as the indices of the SNPs in the read $u_{r,j}$, are allowed to be discontinuous. Given there are J_r SNPs in read R_r , the probability of drawing read R_r from haplotype k is the product of the contribution of each SNP $j=1, \dots, J_r$ in that read. For the j^{th} SNP, this probability is the probability the read contained the alternate base $\phi^1_{r,j}$ times the probability $\theta_{t,k}$ for $t=u_{r,j}$ that ancestral haplotype k emitted the alternate base, added to the equivalent probability for the reference base. Taken together, this yields

$$P(R_r|q_t = k, \lambda) = \prod_{j=1}^{J_r} (\theta_{u_r,j,k} \phi_{r,j}^1 + (1 - \theta_{u_r,j,k}) \phi_{r,j}^0).$$

Let O_t be the set of reads with central SNP t . In the haploid model, the probability of the observations at locus t is

$$P(O_t|q_t = k_t, \lambda) = \prod_{R_r \in O_t} P(R_r|q_t = k_t, \lambda)$$

In the diploid model, each read is equally likely to come from either the maternal or paternal chromosome, giving

$$P(R_r|q_t = (k_{t,1}, k_{t,2}), \lambda) = \frac{1}{2}P(R_r|q = k_{t,1}, \lambda) + \frac{1}{2}P(R_r|q = k_{t,2}, \lambda)$$

For every SNP t , the probability of the observations at that locus is

$$P(O_t|q_t = (k_{t,1}, k_{t,2}), \lambda) = \prod_{R_r \in O_t} P(R_r|q_t = (k_{t,1}, k_{t,2}), \lambda)$$

Finally, note that for SNPs which are not covered by reads, we set $P(O_t|q_t=k_t,\lambda)=1$ for all k_t .

CFW mouse sequencing

Full details on the Crl:CFW(SW)-US_PO8 (CFW) mice, including sample acquisition, age, sex and sequencing are provided elsewhere¹⁷. CFW mice are from a commercial outbred colony¹⁸. Sample pre-processing was done in accordance with best practice recommendations²⁰. Sequencing reads from low coverage samples were mapped to mm10 using bwa²⁵, remapped using

Stampy²⁶, PCR duplicates were marked using Picard, files were merged using Picard, indel realignment was performed using the GATK²⁷, and base quality score recalibration was performed using the GATK. Variant calling was done using the GATK UnifiedGenotyper and filtered by the GATK VQSR, using as training data a set of variants from the Mouse Genomes Project¹⁹ and a sensitivity threshold of 80%. Sites in the training set which failed VQSR were nonetheless retained. In total, 7.07M SNPs were called on the autosomes and chromosome X. 4 mice were additionally sequenced at 10X. Genotypes for these mice were generated using the GATK UnifiedGenotyper using the genotype given alleles option. For comparisons with low coverage imputation, individual genotypes from the high coverage samples were set to missing if the read depth was less than 5 or more than 25, or if the genotype quality was less than 10.

CFW MegaMUGA Array Genotyping

48 of the 2,073 mice were sent to Neogen and genotyped using the Mega Mouse Universal Genotyping Array (MegaMUGA), an array built upon the Illumina Infinium platform with 77,808 SNPs (Neogen, Lincoln, Nebraska, USA). Genotype calling was performed by Neogen using GenCall.

After genotyping, recorded genders were compared to X and Y chromosome marker information, revealing no gender mismatches on the arrays. Samples were further compared to imputation and array QC metrics (call rate and 10% GC score); this revealed 4 of 48 samples had poorly performing array metrics. These 4 samples were subsequently removed from further analysis.

For the 77,808 SNPs for which we had genotypes, 144 were not mappable from mm9 to mm10 using liftOver, and out of the remaining sites, 29,694 intersected between imputation and MegaMUGA. Out of those, we removed sequentially for the following reasons: 17 for allele disagreements between sequencing and the array; 3,819 monomorphic array sites; 3,160 SNPs with an imputed SNP within 25 bp of the array target SNP (as off-target variation can affect microarray genotyping)¹³ ; 56 sites with an array Hardy-Weinberg Equilibrium p-value of less than 1×10^{-10} . Subsequent comparisons between CFW imputed dosages and array genotypes were made for the remaining 21,576 sites.

CONVERGE

Full details for the processing of the CONVERGE data, including low coverage, high coverage, Illumina array data, ethics committees and informed consent have been published elsewhere¹⁵. In brief, 11,670 low coverage (1.7X) Han Chinese samples were called using the GATK to yield a set of 20.5M variants. 9 samples were sequenced to 10X and used independently to call variants (5.9M). For our analysis, high coverage sample genotypes with a read depth of lower than 5, read depth greater than 25, or genotype quality of less than 10 were masked out. 72 samples were genotyped using the Illumina HumanOmniZhongHua-8 (v1.0B) BeadChip. Of the 21057 sites present on chromosome 20 on the array and used for imputation, we removed 292 sites with >5% missingness, 7642 sites with probes within 25bp of another site in the imputed dataset, and 0 sites with an

array Hardy-Weinberg Equilibrium p-value of less than 1×10^{-10} . This left 13,123 sites used for assessing accuracy.

Beagle

For both the CFW and CONVERGE samples, Beagle version 4.0 (beagle.r1399.jar) was run using default parameters, unless otherwise noted⁷. Genotype likelihoods from VCF files were used as inputs. For the CONVERGE study, the 1000 Genomes Asian reference panel was used.

findhap

For both the CFW and CONVERGE samples, findhap version 4 was run with default parameters, unless otherwise noted¹². For both CFW and CONVERGE, allele depth information from VCF files was used to construct input files for findhap. Since pedigree information was not available for the CFW or CONVERGE study, input pedigree files were made with missing values for maternal and paternal inheritance. Results for CFW used maxhap=10000 and CONVERGE used maxhap=25000. Both methods used default options of overlap=10, lowdense=0.07, and errrate=0.01.

Correlation between dosages and validation

For the arrays, correlations (r^2) are generated per-SNP between array genotypes and imputed dosages. When reported by frequency, values are averaged over all SNPs in that frequency bin; otherwise averaging is over all SNPs. For sequencing, correlations are between genotypes and imputed dosages across all samples. When reported by frequency, values are generated using only SNPs in that frequency bin, otherwise averaging is over all SNPs. When aggregating genotypes across sites, to remove an upward bias in correlations due to genotype encoding, dosages were recoded from the default of 0 as the reference allele and 1 the alternate allele to 0 for homozygous major allele and 1 homozygous for the minor allele.

Software

STITCH v1.0.0 was used for all analyses in this paper. STITCH is available from <http://www.stats.ox.ac.uk/~myers/>

Methods-only references

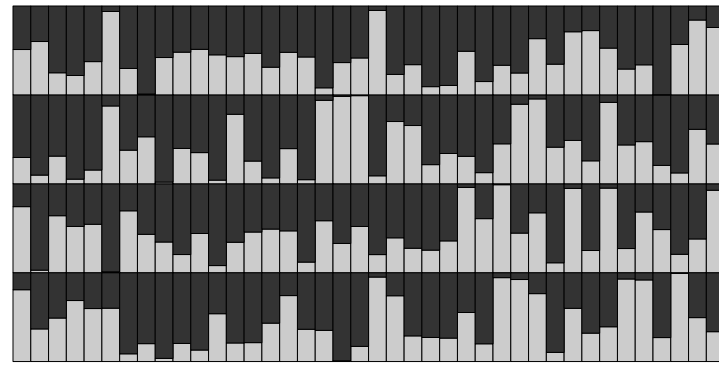
25. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
26. Lunter, G. & Goodson, M. Stampy: A statistical algorithm for sensitive and fast mapping of Illumina sequence reads. *Genome Res.* **21**, 936–939 (2011).

27. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).

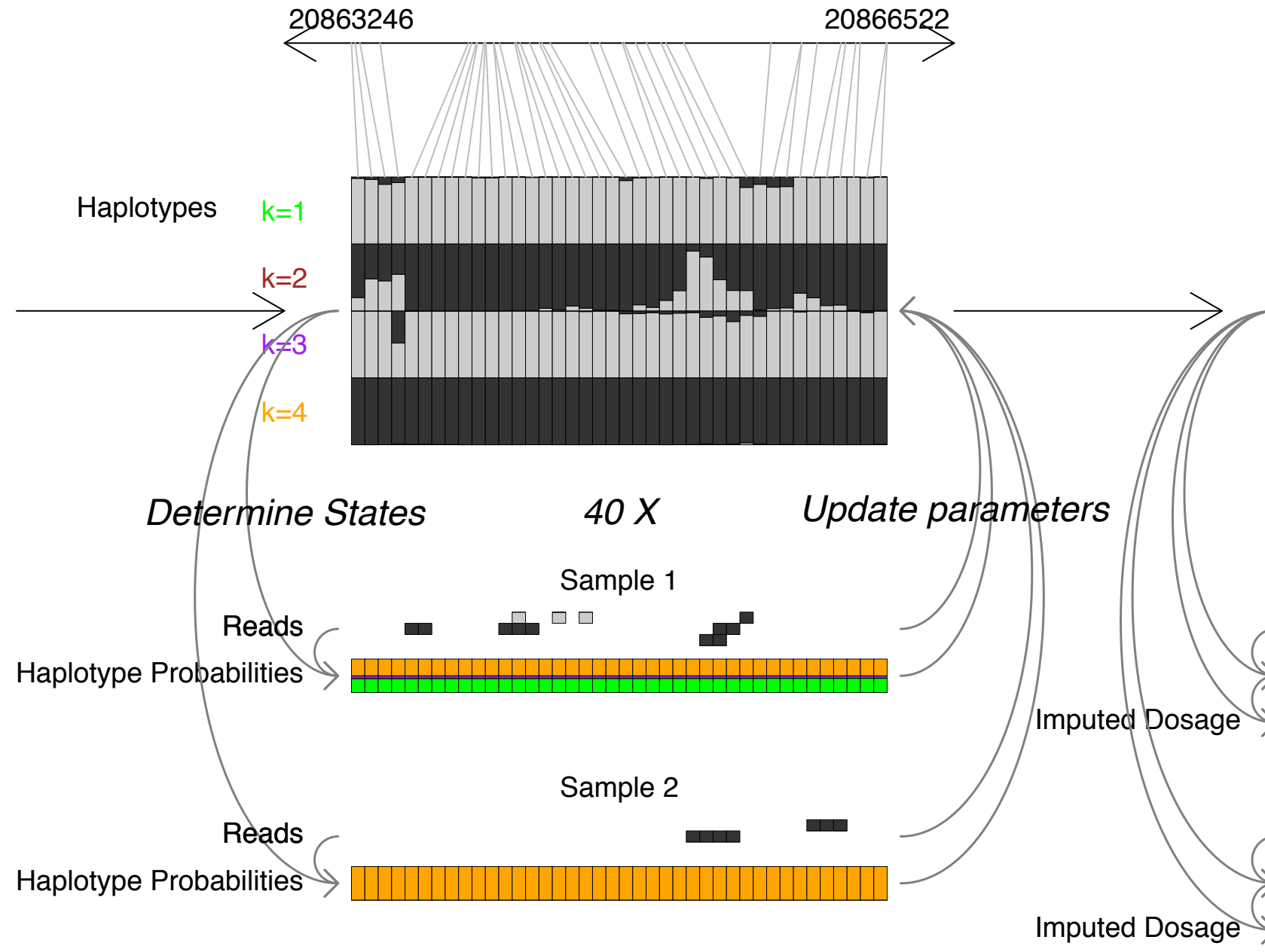
Competing financial interests

The authors declare no competing financial interests.

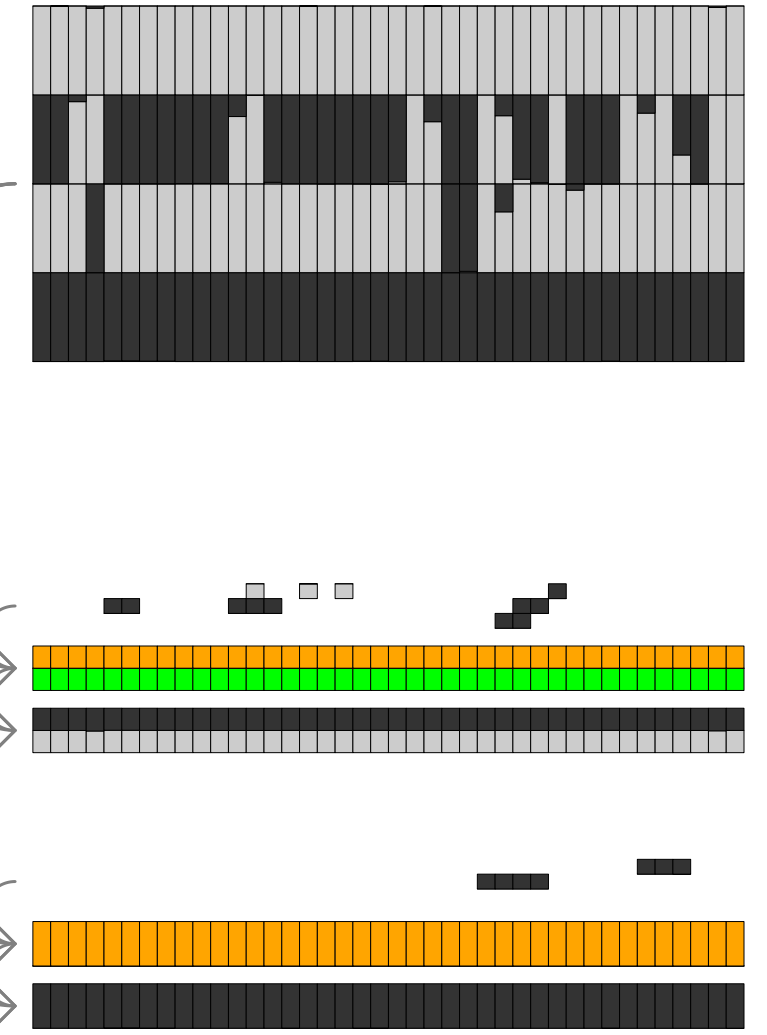
Initialize

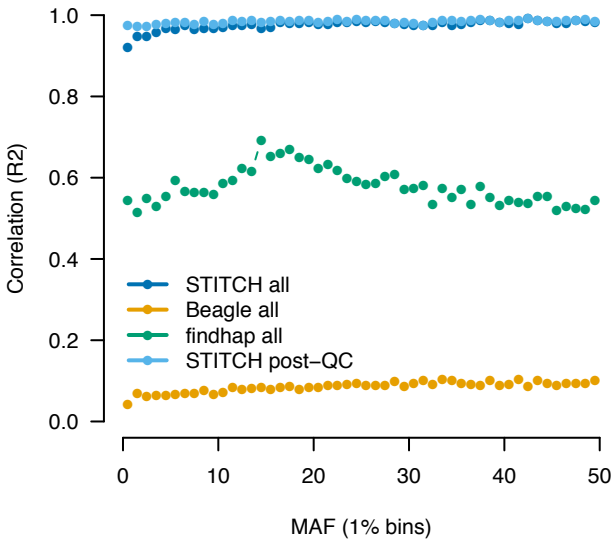
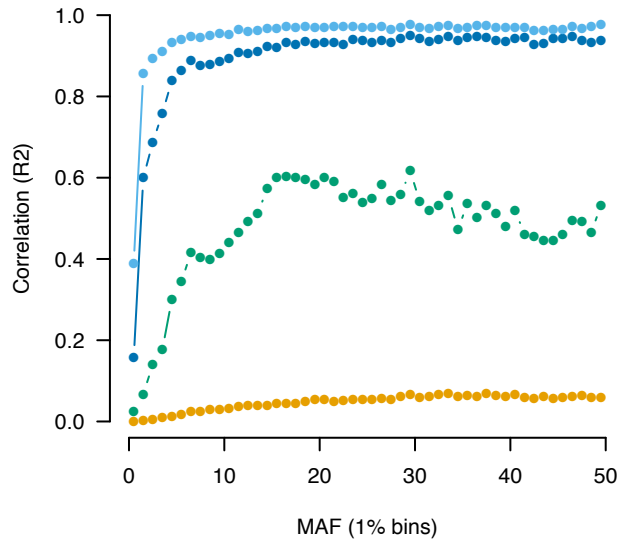


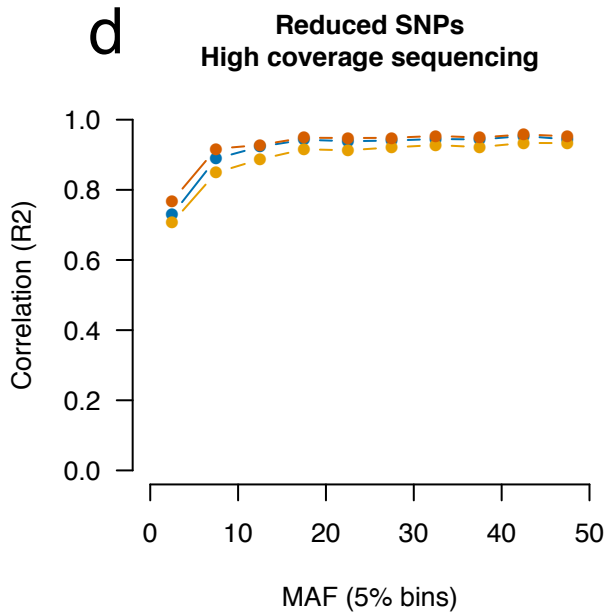
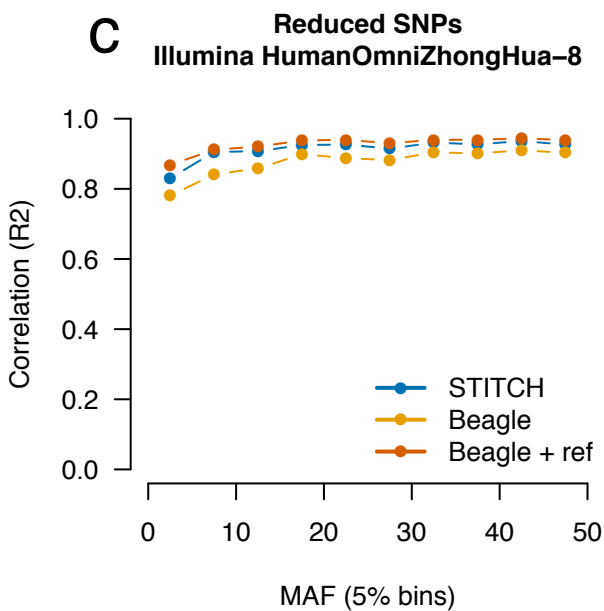
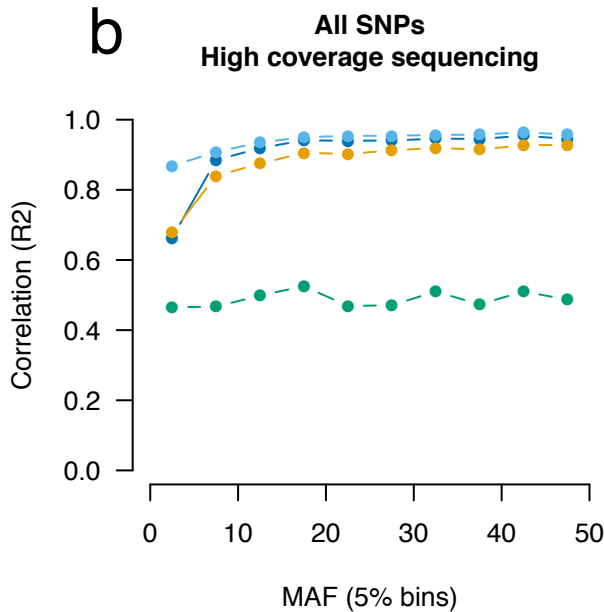
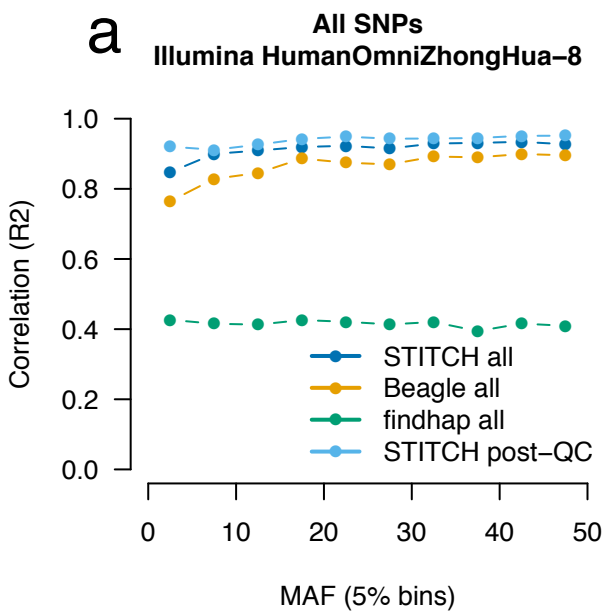
Iterate

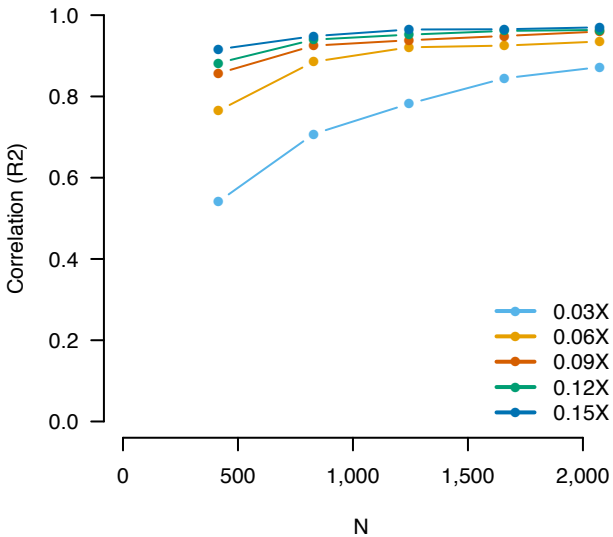
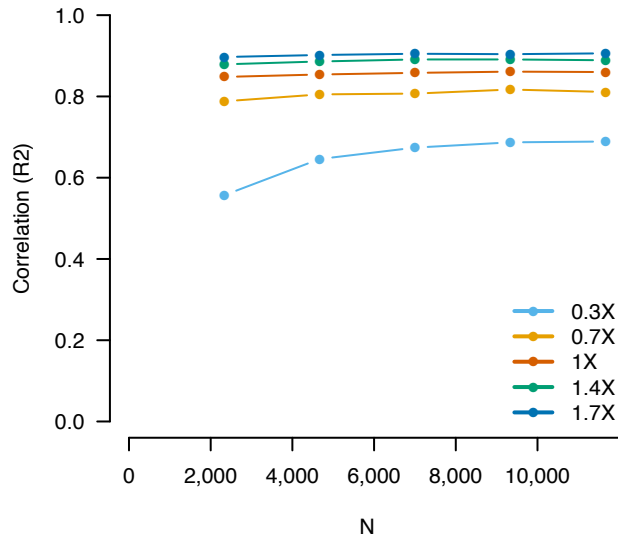
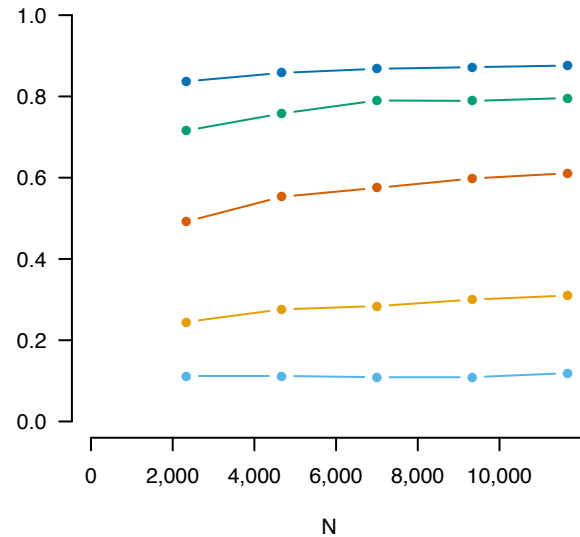


Impute



a**Illumina MegaMUGA****b****High coverage sequencing**



a**CFW – Outbred mice – STITCH****b****CONVERGE – Human – STITCH****c****CONVERGE – Human – Beagle**

Supplementary Information for “Rapid genotype imputation from sequence without reference panels”

Davies *et al.*

May 11, 2016

Contents

| | | |
|----------|---|-----------|
| 1 | Supplementary Note | 2 |
| 1.1 | Pseudo-haploid model | 3 |
| 1.2 | Maximization and parameter updating | 6 |
| 1.2.1 | Useful Variables | 7 |
| 1.2.2 | Haploid model | 8 |
| 1.2.3 | Pseudo-haploid model | 10 |
| 1.2.4 | Diploid model | 11 |
| 1.3 | Efficient calculation of forward backward variables | 14 |
| 1.4 | Initialization | 16 |
| 1.5 | Parameter bounding | 16 |
| 1.6 | Heuristics | 16 |
| 1.7 | Guidance behind parameter options | 17 |
| 2 | Supplementary Tables | 19 |

1 Supplementary Note

Definitions for commonly used variables

| Symbol | Definition |
|----------------|---|
| K | Number of ancestral or founder haplotypes |
| T | Number of SNPs in region |
| N | Number of sample individuals |
| G | Number of generations since population founding |
| R_r | Read with index r which spans J_r SNPs, with SNP indices u_r , sequenced bases s_r and base qualities b_r , or $R_r = \{u_r, s_r, b_r\}$ |
| J_r | Number of SNPs spanned by read R_r |
| c_r | Central SNP for read R_r |
| O_t | Set of reads with central SNP t , $O_t = \{R_r c_r = t\}$ |
| O | Set of observations for each SNP t on the chromosome, $O = \{O_t t = 1, \dots, T\}$ |
| $u_{r,j}$ | For SNP j in read R_r , its index with respect to the chromosomal listing of SNPs (<i>e.g.</i> If the physical position of SNP t in the region is L_t for $t = 1, \dots, T$, then SNP j in read R_r has physical position $L_{u_{r,j}}$) |
| $s_{r,j}$ | Sequencing base for SNP j in read R_r , with $s_{r,j} = 1$ for the alternate base and 0 for the reference base |
| $b_{r,j}$ | Base quality for SNP j in read R_r |
| $R_{r,j}$ | Subset of read R_r for SNP j , or $R_{r,j} = \{u_{r,j}, s_{r,j}, b_{r,j}\}$ |
| $\phi_{r,j}^i$ | Probability of SNP j from read R_r coming from an underlying genotype i , or $P(s_{r,j} g = i)$ |
| I_t | Variable counting the number of recombinations that take place between SNPs t and $t + 1$ |
| H_r^j | Variable that takes value 1 if SNP j from read R_r is the alternate base and value 0 if it is the reference base |
| H_r | Variable that takes value 1 if read R_r comes from the maternal haplotype and 2 if it comes from the paternal haplotype |
| π_k | Probability of starting in state k at the first SNP |
| σ_t | Recombination distance between SNPs t and $t + 1$ |
| $\alpha_{t,k}$ | Probability of switching into state k between SNPs t and $t + 1$ |
| $\theta_{t,k}$ | Probability that haplotype k emits the alternate base at SNP t |
| λ | Parameters of the model $\lambda = \{\pi, \sigma, \alpha, \theta\}$ |

In the main text, we described the model by showing one would simulate it, and more formally laid out the details necessary to generate probabilities under the model. Here, we further specify the model by describing how the expectation of the complete data likelihood can be used in an EM framework to provide updated parameters λ^{i+1} which guarantee no decrease in the likelihood of the observed data. Doing this requires state space augmentation and calculating expectations over hidden states in the Markov model. Here we show how these

expectations are calculated for the haploid, diploid and pseudo-haploid cases. Later, initialization, bounding, and heuristics of the model are given as well.

First, we give a brief review of notation (and see list above). We consider a genomic region with T SNPs, and sequencing reads from N individuals. For each individual, we index their reads with r so we speak of read R_r . We define a central SNP c_r for read R_r , so that for each SNP in the region, we observe a set of reads, $O_t = \{R_r | c_r = t\}$. Read R_r consists of a triplet of vectors: u_r , the indices; s_r the reference (0) or alternate (1) bases; and the base qualities b_r . From this we use $\phi_{r,j}^i$, the probability SNP j in read R_r has underlying genotype i .

We model our population as having been founded G generations ago with K ancestral haplotypes. Sampling a set of observations for an individual can be thought of as 1) choosing an initial haplotype k according to π_k , the prior probability of starting in state k ; 2) choosing where to switch states according to σ_t , the genetic distance between SNPs t and $t+1$; 3) choosing which haplotype k to sample at recombination breakpoints according to $\alpha_{t,k}$, the local probability of switching into haplotype k at SNP $t+1$; and 4) sampling reads by i) choosing read breakpoints and determining $u_{r,j}$, the indices of the SNPs in the read; ii) obtaining $b_{r,j}$, the base qualities of the SNPs in the read; iii) choosing the real bases of the SNPs in the read according to $\theta_{u_{r,j},k}$, the probability that haplotype k emits the alternate base at SNP $u_{r,j}$; iv) observing sequenced bases $s_{r,j}$ according to $b_{r,j}$ and the real bases.

In the unaugmented hidden state space, the haploid model corresponds to a set of $k_t \in 1, \dots, K \quad \forall t = 1, \dots, T$. For the diploid model, this consists of a set of pairs of states $(k_{t,1}, k_{t,2})$, while for the pseudo-haploid mode, it is two hidden states $k_{t,1}$ and $k_{t,2}$. In the augmented hidden state space, we further consider knowledge of: how many recombinations occur between SNPs t and $t+1$, defined by variable I_t ; whether base j of read R_r is a reference or alternate base, defined by variable H_r^j ; and whether read R_r comes from the maternal or paternal haplotype, defined by variable H_r . Utilization of the augmented hidden state space is necessary for updating parameters, as explained below.

1.1 Pseudo-haploid model

The diploid model presented here and used in fastPHASE and other similar algorithms suffers from a quadratic computational complexity due to the need to sum over K^2 possible diploid states at each site. With sequencing reads, the observed data fundamentally comes from either the first (*e.g.* maternal) or second (*e.g.* paternal) haplotype. If we had labels for each read as to whether they came from the maternal or paternal haplotype, we would have separable likelihoods, and could use the maternal reads to infer the maternal states, and likewise for the paternal reads and paternal states, which would have computational cost proportional to $2 \times K$ as opposed to K^2 .

In the diploid EM algorithm, we use the current set of parameters to generate the posterior probability of the pair of hidden states given the observations, and use these to generate a new set of parameters that maximize the

likelihood. An alternative approach is to average over sampled hidden states realized through a hypothetical Gibbs sampler that i) samples labels conditional on states, observations, and parameters, and ii) samples states conditional on labels, observations and parameters. Implementing such a Gibbs sampler in reality would be computationally unwise, as it would likely take at least as long as the original diploid EM. However, with certain assumptions about the posterior distribution of the labels, we can approximate the posterior distribution of the hidden states quickly.

Let q_1 be the full hidden state for haplotype 1, the maternal haplotype. Let H_r be the label for read r with $H_r = 1$ corresponding to the maternal haplotype and $H_r = 2$ corresponding to the paternal haplotype. Let $O = \{R_r\}$ be the set of all reads, with $|O|$ reads in total, and let H correspond to an assignment of labels $H \in \mathcal{H} = \{1, 2\}^{|O|}$. Let $\mathcal{R}_h = \{R_r | H_r = h\}$ be the set of reads with label h . Then we have

$$P(q_1|O, \lambda) = \sum_{H \in \mathcal{H}} P(q_1, H|O, \lambda) \quad (1)$$

$$= \sum_{H \in \mathcal{H}} P(q_1|H, O, \lambda)P(H|O, \lambda) \quad (2)$$

$$= \sum_{H \in \mathcal{H}} P(q_1|H, O, \lambda) \prod_{r=1}^{|O|} P(H_r|O, \lambda) \quad (3)$$

where the last equality requires the approximation that the probability of the labels are independent of each other. Now, the probability of a state given labels and reads can be further written as

$$P(q_1|H, O, \lambda) = \frac{P(O|H, q_1, \lambda)P(q_1|H, \lambda)}{P(O|H, \lambda)} \quad (4)$$

$$= \frac{\left(\prod_{r: H_r=1} P(R_r|q_1, \lambda)\right)P(\mathcal{R}_2|\text{hap2}, \lambda)P(q_1|\lambda)}{P(\mathcal{R}_1|\text{hap1}, \lambda)P(\mathcal{R}_2|\text{hap2}, \lambda)} \quad (5)$$

where we use $P(q_1|H, \lambda) = P(q_1|\lambda)$, since labels don't affect state probabilities without observations, and where $(\mathcal{R}_1|\text{hap1}, \lambda)$ is the probability of observing the set of reads labeled as coming from haplotype 1, conditional on their having come from haplotype 1. If we further approximate $P(\mathcal{R}_1|\text{hap1}, \lambda) = \prod_{r: H_r=1} P(R_r|\text{hap1}, \lambda)$, and approximate $P(R_r|\text{hap1}, \lambda) = P(R_r|\lambda)$, we get that

$$P(q_1|H, O, \lambda) = P(q_1|\lambda) \prod_{r: H_r=1} (P(R_r|q_1|\lambda))/(P(R_r|\lambda)) \quad (6)$$

This gives us that

$$P(q_1|O, \lambda) = \left(\sum_{H \in \mathcal{H}} P(q_1|\lambda) \left(\prod_{r: H_r=1} \frac{P(R_r|q_1, \lambda)}{P(R_r|\lambda)} \right) \right) \left(\prod_{r=1}^{|O|} P(H_r|O, \lambda) \right) \quad (7)$$

$$= P(q_1|\lambda) \sum_{H \in \mathcal{H}} \prod_{r=1}^{|O|} \left(P(H_r|O, \lambda) \left(\mathcal{I}\{H_r = 1\} \frac{P(R_r|q_1, \lambda)}{P(R_r|\lambda)} + \mathcal{I}\{H_r = 2\} 1 \right) \right) \quad (8)$$

$$= P(q_1|\lambda) \prod_{r=1}^{|O|} \left(P(H_r = 1|O, \lambda) \frac{P(R_r|q_1, \lambda)}{P(R_r|\lambda)} + P(H_r = 2|O, \lambda) \right) \quad (9)$$

Therefore, we get that read r contributes $P(H_r = 1|O, \lambda)P(R_r|q_1, \lambda) + P(H_r = 2|O, \lambda)P(R_r|\lambda)$ to the likelihood, after multiplying by the constant $P(R_r|\lambda)$, as opposed to $P(R_r|q_1, \lambda)$ as it would under a fully seperable model. When testing on real data, we found that we achieved marginally but consistently better performance using $P(H_r = 1|O, \lambda)P(R_r|q_1, \lambda) + P(H_r = 2|O, \lambda)P(R_r|\text{hap}2, \lambda)$ instead, so this equation was used when calculating the state probabilities.

To use this, we need an estimate of the probability of a label given the data. To do this, consider a read R_r , with lead SNP c_r , and label H_r . Then we can calculate the following

$$P(H_r = 1|O, \lambda) = \sum_{q_1, q_2} P(H_r|q_1, q_2, O, \lambda)P(q_1, q_2|O, \lambda) \quad (10)$$

$$= \sum_{q_1, q_2} P(H_r|q_1, q_2, R_r, \lambda)P(q_1, q_2|O, \lambda) \quad (11)$$

$$= \sum_{q_1, q_2} \frac{P(R_r|q_1, \lambda)}{P(R_r|q_1, \lambda) + P(R_r|q_2, \lambda)} P(q_1, q_2|O, \lambda) \quad (12)$$

$$= \mathbb{E}_{q_1, q_2} \left[\frac{P(R_r|q_1, \lambda)}{P(R_r|q_1, \lambda) + P(R_r|q_2, \lambda)} | O, \lambda \right] \quad (13)$$

$$\approx \frac{\mathbb{E}_{q_1} [P(R_r|q_1, \lambda) | O, \lambda]}{\sum_{h=1}^2 \mathbb{E}_{q_h} [P(R_r|q_h, \lambda) | O, \lambda]} \quad (14)$$

This uses a prior probability on labels of $P(H_r = 1) = P(H_r = 2) = \frac{1}{2}$. We also use the approximation that the expectation of ratios is equivalent to the ratio of expectations, to avoid a calculation with computational complexity of order K^2 . To perform this calculation we use

$$P(R_r|\text{hap}h, \lambda) = \mathbb{E}_{q_h} [P(R_r|q_h, \lambda) | O, \lambda] \approx \sum_{k=1}^K P(R_r|q_h = k, \lambda') P(q_k | O, \lambda') \quad (15)$$

where λ' are the parameters from the previous iteration.

Therefore, in calculating the complete data probability for the pseudo-haploid model for haplotype $H = 1$, we use the probability of the observation at SNP t

given state $q_t = k_t$ and parameters λ as

$$\begin{aligned}
P_{H=h}(O_t|q_t = k_t, \lambda) &= \prod_{j=1}^{J_r} P_{H=h}(R_r|q_t = k_t, \lambda) \\
&= \prod_{j=1}^{J_r} P(H_r = h|O, \lambda)P(R_r|q_1, \lambda) + P(H_r \neq h|O, \lambda)P(R_r|\text{hap}h, \lambda)
\end{aligned} \tag{16}$$

where $P(H_r = h|O, \lambda)$ is from Equation 14, $P(R_r|q_1, \lambda)$ is as defined in the main text, $P(H_r \neq h|O, \lambda) = 1 - P(H_r = h|O, \lambda)$, and $P(R_r|\text{hap}h, \lambda)$ is from Equation 15.

1.2 Maximization and parameter updating

In the EM algorithm, one defines a ‘‘complete dataset’’ D including the observed data (O , the reads), as well as the hidden parameters (Q , the hidden states). Given a set of parameters λ , one defines the log-likelihood of the complete data as $L(\lambda) = \log(l(\lambda|D)) = \log(l(\lambda|D = (O, Q)))$. Given a current set of parameters λ^i , we generate a new set of parameters λ^{i+1} to maximize the expectation of $l(\lambda^{i+1})$ with respect to the distribution of hidden parameters obtained by λ^i

$$\begin{aligned}
U(\lambda^{i+1}, \lambda^i) &= \mathbb{E}[l(\lambda^{i+1})|O, \lambda^i] \\
&= \sum_Q P(Q|O, \lambda^i) \log(P(O, Q|\lambda^{i+1}))
\end{aligned} \tag{17}$$

Standard theory implies that by choosing λ^{i+1} to maximize $U(\lambda^{i+1}, \lambda^i)$, we also increase the likelihood of the observed data, $l(\lambda^{i+1}|O) > l(\lambda^i|O)$ [1].

In applying the EM algorithm, we first initialize with a set of parameters λ^0 . For each subsequent iteration $i = 1, 2, \dots$, we then iteratively alternate between the ‘‘Expectation’’ phase, where we calculate $U(\lambda^{i+1}, \lambda^i)$, and the ‘‘Maximization’’ phase, where we calculate λ^{i+1} to maximize $U(\lambda^{i+1}, \lambda^i)$. In the Expectation phase, the crucial component is calculating the state probabilities $P(Q|O, \lambda^i)$ - these are calculated using the forward and backward algorithms. To calculate the updates in the Maximization stage, we must further augment the latent space to model how many recombinations occur between SNPs, whether emissions were due to occurrences of an alternate base or a reference base, and whether observed reads were from the maternal or paternal haplotype. To calculate the updates in the Maximization stage, we must further augment the latent space to model whether transitions occur due to recombinations or not, and whether emissions were due to occurrences of an alternate base or a reference base. In this new augmented latent space, for some fixed set of hidden parameters for the N samples, consider some sums that can be calculated. Let n_k^1 be the number of sample haplotypes in state k at the first SNP, n_{stay}^t be the number of sample haplotypes which do not recombine between SNPs t and $t + 1$, $n_{\text{switch},k}^t$ be the number of sample haplotypes which switch into ancestry k

between SNPs t and $t+1$, and $n_{k,s}^t$ be the number of reads that have a reference $s = 0$ or alternate $s = 1$ base for SNP t that are in state k for their central SNP. Then the complete data log likelihood is

$$\begin{aligned}
l(\lambda) &= \log(P(O, Q|\lambda)) \\
&= \sum_{k=1}^K n_k^1 \log(\pi_k) \\
&+ \sum_{t=1}^{T-1} n_{\text{stay}}^t \log(e^{-G\sigma t}) + \sum_{t=1}^{T-1} \sum_{k=1}^K n_{\text{switch},k}^t \log((1 - e^{-G\sigma t})\alpha_{t,k}) \\
&+ \sum_{t=1}^T \sum_{k=1}^K n_{k,1}^t \log(\theta_{t,k}) + \sum_{t=1}^T \sum_{k=1}^K n_{k,0}^t \log((1 - \theta_{t,k})) \tag{18}
\end{aligned}$$

Calculating updates for a parameter is done by taking the derivative of $U(\lambda^{i+1}, \lambda^i)$ with respect to that parameter, setting it equal to 0 and solving. Employing the notation $\mathbb{E}[x|O, \lambda] = \mathbb{E}_\lambda[x]$, it is easy to calculate the following updates for $\lambda^{i+1} = (\pi^{i+1}, \theta^{i+1}, \alpha^{i+1}, \sigma^{i+1})$

$$\pi_k^{i+1} = \frac{\mathbb{E}_{\lambda^i}[n_k^1]}{\sum_{j=1}^K \mathbb{E}_{\lambda^i}[n_j^1]} \tag{19}$$

$$\theta_{t,k}^{i+1} = \frac{\mathbb{E}_{\lambda^i}[n_{k,1}^t]}{\mathbb{E}_{\lambda^i}[n_{k,0}^t] + \mathbb{E}_{\lambda^i}[n_{k,1}^t]} \tag{20}$$

$$\alpha_{t,k}^{i+1} = \frac{\mathbb{E}_{\lambda^i}[n_{\text{switch},k}^t]}{\sum_{j=1}^K \mathbb{E}_{\lambda^i}[n_{\text{switch},j}^t]} \tag{21}$$

$$\sigma_t^{i+1} = \frac{1}{-G} \log \left(\frac{\sum_{k=1}^K \mathbb{E}_{\lambda^i}[n_{\text{switch},k}^t]}{\sum_{k=1}^K \mathbb{E}_{\lambda^i}[n_{\text{switch},k}^t] + \mathbb{E}_{\lambda^i}[n_{\text{stay}}^t]} \right) \tag{22}$$

1.2.1 Useful Variables

We use a standard forward backward HMM implementation with a set of parameters λ . Recall that q_t is the hidden state at SNP t . We use the following notations for states k_t at SNP t and k_{t+1} at SNP $t+1$

$$\begin{aligned}
\alpha_t(k_t) &= P(O_1 O_2 \dots O_t, q_t = k_t | \lambda) \\
\beta_t(k_t) &= P(O_{t+1} O_{t+2} \dots O_T | q_t = k_t, \lambda) \\
\gamma_t(k_t) &= P(q_t = k_t | O, \lambda) = \frac{\alpha_t(k_t) \beta_t(k_t)}{P(O | \lambda)} \\
\xi_t(k_t, k_{t+1}) &= P(q_t = k_t, q_{t+1} = k_{t+1} | O, \lambda) \\
&= \frac{\alpha_t(k_t) P(q_{t+1} = k_{t+1} | q_t = k_t, \lambda) \beta_{t+1}(k_{t+1}) P(O_{t+1} | q_{t+1} = k_{t+1}, \lambda)}{P(O | \lambda)}
\end{aligned}$$

The diploid version of these equations, where we go from state $(k_{t,1}, k_{t,2})$ at SNP t to state $(k_{t+1,1}, k_{t+1,2})$ at SNP $t + 1$ is

$$\begin{aligned}
\alpha_t(k_{t,1}, k_{t,2}) &= P(O_1 O_2 \dots O_t, q_t = (k_{t,1}, k_{t,2}) | \lambda) \\
\beta_t(k_{t,1}, k_{t,2}) &= P(O_{t+1} O_{t+2} \dots O_T | q_t = (k_{t,1}, k_{t,2}), \lambda) \\
\gamma_t(k_{t,1}, k_{t,2}) &= P(q_t = (k_{t,1}, k_{t,2}) | O, \lambda) = \frac{\alpha_t(k_{t,1}, k_{t,2}) \beta_t(k_{t,1}, k_{t,2})}{P(O | \lambda)} \\
\xi_t\left((k_{t,1}, k_{t,2}), (k_{t+1,1}, k_{t+1,2})\right) &= P(q_t = (k_{t,1}, k_{t,2}), q_{t+1} = (k_{t+1,1}, k_{t+1,2}) | O, \lambda) \\
&= \frac{1}{P(O | \lambda)} \alpha_t(k_{t,1}, k_{t,2}) P(q_{t+1} = (k_{t+1,1}, k_{t+1,2}) | q_t = (k_{t,1}, k_{t,2}), \lambda) \times \\
&\quad \beta_{t+1}(k_{t+1,1}, k_{t+1,2}) P(O_{t+1} | q_t = (k_{t,1}, k_{t,2}), \lambda)
\end{aligned}$$

1.2.2 Haploid model

Initial probabilities

To update the prior parameters, we need the expectation of n_k^1 , which we define as the number of sample haplotypes in state k at the first SNP. Denote the probability that the sample is in the first state at SNP t by $\gamma_t(k)$. Let $\gamma_{n,t}(k)$ be $\gamma_t(k)$ for sample n . We can therefore calculate the required expectation from the main text as

$$\mathbb{E}_\lambda[n_k^1] = \sum_{n=1}^N \gamma_{n,1}(k) \tag{23}$$

Transition matrix probabilities

To update the transition parameters, we use an augmented state space where we have knowledge of how many recombinations occurred between two SNPs. Define a variable I_t as the count of the number of recombinations between SNPs t and $t + 1$; in the haploid model, this takes value 0 or 1. This will allow us to calculate the expectation of n_{stay}^t , the number of sample haplotypes that do not recombine between SNPs t and $t + 1$, and $n_{\text{switch},k}^t$, the number that switch into state k between SNPs t and $t + 1$.

We extend our transition probability to include I_t as follows

$$P(q_{t+1} = k_{t+1}, I_t | q_t = k_t, \lambda) = \begin{cases} e^{-G\sigma_t} & \text{if } k_t = k_{t+1} \text{ and } I_t = 0 \\ 0 & \text{if } k_t \neq k_{t+1} \text{ and } I_t = 0 \\ (1 - e^{-G\sigma_t}) \alpha_{t,k_{t+1}} & \text{if } I_t = 1 \end{cases}$$

Recall that $\xi_t(k_t, k_{t+1})$ is

$$\xi_t(k_t, k_{t+1}) = \frac{\alpha_t(k_t) P(q_{t+1} = k_{t+1} | q_t = k_t, \lambda) \beta_{t+1}(k_{t+1}) P(O_{t+1} | q_{t+1} = k_{t+1}, \lambda)}{P(O | \lambda)} \tag{24}$$

Denote the probability given the observed data O that across SNP t , the sample has states k_t, k_{t+1} and indicator I_t by $\xi_t(k_t, k_{t+1}, I_t)$. Then

$$\xi_t(k_t, k_{t+1}, I_t) = \frac{\alpha_t(k_t)P(q_{t+1} = k_{t+1}, I_t | q_t = k_t, \lambda)\beta_{t+1}(k_{t+1})P(O_{t+1} | q_{t+1} = k_{t+1}, \lambda)}{P(O | \lambda)} \quad (25)$$

Let $\xi_t(k_t, k_{t+1}, I_t)$ be $\xi_{n,t}(k_t, k_{t+1}, I_t)$ for sample n . We can therefore calculate expectations as

$$\mathbb{E}_\lambda[n_{\text{stay}}^t] = \sum_{n=1}^N \sum_{k=1}^K \xi_{n,t}(k, k, I_t = 0) \quad (26)$$

$$\mathbb{E}_\lambda[n_{\text{switch},k}^t] = \sum_{n=1}^N \sum_{i=1}^K \xi_{n,t}(i, k, I_t = 1) \quad (27)$$

and since

$$\mathbb{E}_\lambda[n_{\text{stay}}^t] = 1 - \sum_{n=1}^N \sum_{i=1}^K \sum_{k=1}^K \xi_{n,t}(i, k, I_t = 1) = 1 - \sum_{k=1}^K \mathbb{E}_\lambda[n_{\text{switch},k}^t] \quad (28)$$

it is therefore sufficient to calculate $\mathbb{E}_\lambda[n_{\text{switch},k}^t]$ to perform the EM updating from the main text.

Emission matrix probabilities

To update the emission parameters, we use an augmented state space where we have knowledge of whether emissions were due to the alternate or reference base. Recall that $\phi_{r,j}^i$ is the probability SNP j in read R_r came from a read with underlying genotype i . Denote by H_r^j a variable which takes value 1 if the underlying base is the alternate base and 0 if it is the reference base. We will use this to calculate the expectation of $n_{k,s}^t$, the number of reads with a base at SNP t that contain the alternate ($s = 1$) or reference ($s = 0$) base where the sample was in state k at the central SNP of the read.

Recall that the original definition of the probability of read R_r given hidden state k at SNP t and parameters λ is

$$P(R_r | q_t = k, \lambda) = \prod_{j=1}^{J_r} P(R_{r,j} | q_t = k, \lambda) = \prod_{j=1}^{J_r} (\phi_{r,j}^1 \theta_{u_{r,j},k} + \phi_{r,j}^0 (1 - \theta_{u_{r,j},k})) \quad (29)$$

We extend our emission probability to include H_r^j as follows

$$P(R_r, H_r^j | q_t = k_t, \lambda) = \begin{cases} \left[\prod_{i \neq j} P(R_{r,i} | q_t = k_t, \lambda) \right] \phi_{r,j}^1 \theta_{u_{r,j},k} & \text{if } H_r^j = 1 \\ \left[\prod_{i \neq j} P(R_{r,i} | q_t = k_t, \lambda) \right] \phi_{r,j}^0 (1 - \theta_{u_{r,j},k}) & \text{if } H_r^j = 0 \end{cases} \quad (30)$$

For read R_r with central SNP c_r , the probability of the observation (set of reads)

at SNP $t = c_r$ and H_r^j becomes

$$P(O_t, H_r^j | q_t = k_t, \lambda) = \begin{cases} P(O_t | q_t = k_t, \lambda) \frac{\phi_{r,j}^1 \theta_{u_{r,j},k}}{\phi_{r,j}^1 \theta_{u_{r,j},k} + \phi_{r,j}^0 (1 - \theta_{u_{r,j},k})} & \text{if } H_r^j = 1 \\ P(O_t | q_t = k_t, \lambda) \frac{\phi_{r,j}^0 (1 - \theta_{u_{r,j},k})}{\phi_{r,j}^1 \theta_{u_{r,j},k} + \phi_{r,j}^0 (1 - \theta_{u_{r,j},k})} & \text{if } H_r^j = 0 \end{cases} \quad (31)$$

We expand $\gamma_t(k_t)$ as

$$\begin{aligned} \gamma_t(k_t) &= \frac{\alpha_t(k_t) \beta_t(k_t)}{P(O|\lambda)} \\ &= \frac{[\sum_{l=1}^K \alpha_{t-1}(l) P(q_t = k_t | q_{t-1} = l, \lambda)] P(O_t | q_t = k, \lambda) \beta_t(k_t)}{P(O|\lambda)} \end{aligned} \quad (32)$$

where we note that for $t = 1$, we substitute π_k for $[\sum_{l=1}^K \alpha_{t-1}(l) P(q_t = k_t | q_{t-1} = l, \lambda)]$. Denote the probability that for SNP j in read R_r with central SNP $t = c_r$, the sample has a hidden state k_t and has indicator H_r^j given observed data O and parameters λ by $\gamma_t(k_t, H_r^j)$. Then

$$\begin{aligned} \gamma_t(k_t, H_r^j) &= \frac{[\sum_{l=1}^K \alpha_{t-1}(l) P(q_t = k_t | q_{t-1} = l, \lambda)] P(O_t, H_r^j | q_t = k, \lambda)}{P(O|\lambda)} \\ &= \begin{cases} \gamma_t(k_t) \frac{\phi_{r,j}^1 \theta_{u_{r,j},k}}{\phi_{r,j}^1 \theta_{u_{r,j},k} + \phi_{r,j}^0 (1 - \theta_{u_{r,j},k})} & \text{if } H_r^j = 1 \\ \gamma_t(k_t) \frac{\phi_{r,j}^0 (1 - \theta_{u_{r,j},k})}{\phi_{r,j}^1 \theta_{u_{r,j},k} + \phi_{r,j}^0 (1 - \theta_{u_{r,j},k})} & \text{if } H_r^j = 0 \end{cases} \end{aligned} \quad (33)$$

Let $\gamma_{n,t}(k_t, H_r^j)$ be $\gamma_t(k_t, H_r^j)$ for sample n , and let A_n be the complete set of SNPs j from reads R_r for sample n such that $u_{r,j} = t$. We can therefore calculate the required expectations from the main text as

$$\mathbb{E}_\lambda[n_{k,1}^t] = \sum_{n=1}^N \sum_{(r,j) \in A_n} \gamma_{n,c_r}(k, H_r^j = 1) \quad (34)$$

$$\mathbb{E}_\lambda[n_{k,0}^t] = \sum_{n=1}^N \sum_{(r,j) \in A_n} \gamma_{n,c_r}(k, H_r^j = 0) \quad (35)$$

1.2.3 Pseudo-haploid model

In the pseudo-haploid model, the only changes to the likelihood occur through the emissions; as such, we need to re-calculate Equations 34 and 35. To update the emission parameters for the pseudo-haploid model, we use an augmented state space where we have knowledge of whether emissions were due to the alternate or reference base, and further have knowledge of whether the read came from the maternal or paternal haplotype. Recall that $\phi_{r,j}^i$ is the probability that observed base j in read R_r came from a read with underlying genotype i .

Recall that H_r^j is an indicator variable which takes value 1 if the underlying base is the alternate base and 0 if it is the reference base. Let H_r take value 1 if the read came from the maternal haplotype and 2 if it came from the paternal haplotype. We will use these to calculate the expectation of $n_{k,s}^t$, the number of reads that emit the alternate base ($s = 1$) or reference base ($s = 0$) given they are in state k at the central SNP of the read.

Recall that for each individual, we make two forward backward passes of the algorithm, once for the maternal haplotype ($h = 1$), and a second time for the paternal haplotype ($h = 2$). We also attempt to probabilistically infer for each read which haplotype it came from. Let H refer to the haplotype we are currently modelling (maternal or paternal).

First, recall that the original definition of the probability while modelling haplotype h of read R_r given hidden state k at SNP t and parameters λ is

$$P_{H=h}(R_r|q_t = k_t, \lambda) = P(R_r|q_t = k_t, \lambda)P(H_r = h|O, \lambda) + P(R_r|H_r \neq h, \lambda)P(H_r \neq h|O, \lambda) \quad (36)$$

For notational convenience, set $F_{r,j,h} = P(H_r = h|O, \lambda) \left[\prod_{i \neq j} P(R_{r,i}|q_t = k, \lambda) \right]$. We therefore expand the emission probability to include H_r^j and H_r as follows

$$P_{H=h}(R_r, H_r^j, H_r|q_t = k, \lambda) = \begin{cases} F_{r,j,h} \theta_{u_{r,j},k} \phi_{r,j}^1 & \text{if } H_r^j = 1, H_r = h \\ F_{r,j,h} (1 - \theta_{u_{r,j},k}) \phi_{r,j}^0 & \text{if } H_r^j = 0, H_r = h \\ P(H_r \neq h|O, \lambda) P(R_r|H_r \neq h, \lambda) & \text{if } H_r \neq h \end{cases}$$

Denote the probability that haplotype h of the sample is in state k_t at SNP t with H_r^j and H_r given observed data O and parameters λ by $\gamma_{t,h}(k_t, H_r^j, H_r)$. Then, we get that

$$\gamma_{t,h}(k_t, H_r^j, H_r) = \begin{cases} \gamma_{t,h}(k_t) \frac{F_{r,j,h} \theta_{u_{r,j},k} \phi_{r,j}^1}{P_{H=h}(R_r|q_t=k_t, \lambda)} & \text{if } H_r^j = 1, H_r = h \\ \gamma_{t,h}(k_t) \frac{F_{r,j,h} (1 - \theta_{u_{r,j},k}) \phi_{r,j}^0}{P_{H=h}(R_r|q_t=k_t, \lambda)} & \text{if } H_r^j = 0, H_r = h \end{cases}$$

Let $\gamma_{n,t,h}(k_t, H_r^j, H_r)$ be $\gamma_{t,h}(k_t, H_r^j, H_r)$ for sample n , and let A_n be the complete set of SNPs j from reads R_r for sample n such that $u_{r,j} = t$. We can therefore calculate the required expectations from the main text as

$$\mathbb{E}_\lambda[n_{k,1}^t] = \sum_{n=1}^N \sum_{(r,j) \in A_n} \sum_{h=1}^2 \gamma_{n,c_r,h}(k, H_r^j = 1, H_r = h) \quad (37)$$

$$\mathbb{E}_\lambda[n_{k,0}^t] = \sum_{n=1}^N \sum_{(r,j) \in A_n} \sum_{h=1}^2 \gamma_{n,c_r,h}(k, H_r^j = 0, H_r = h) \quad (38)$$

1.2.4 Diploid model

Initial probabilities

To update the prior parameters, we need the expectation of n_k^1 , which we define as the number of sample haplotypes in state k at the first SNP. Denote the probability that sample n is in pairs of states $(k_{t,1}, k_{t,2})$ at SNP t given observed data O by $\gamma_{n,t}(k_{t,1}, k_{t,2})$. We can therefore calculate the required expectation from the main text as

$$\mathbb{E}_\lambda[n_k^1] = \sum_{n=1}^N \sum_{j=1}^K (\gamma_{n,1}(k, j) + \gamma_{n,1}(j, k)) \quad (39)$$

Transition probabilities

To update the transition parameters for the diploid model, we use an augmented state space where we have knowledge of how many recombinations occurred between two SNPs. Here we define a variable I_t which counts the number of recombinations that occur between SNPs t and $t+1$ for the two haplotypes of the diploid sample, and takes values 0, 1 or 2. This will allow us to calculate the expectation of n_{stay}^t , the number of sample haplotypes that do not recombine between SNPs t and $t+1$, and $n_{\text{switch},k}^t$, the number of sample haplotypes that switch into state k between SNPs t and $t+1$.

We can therefore extend the diploid transition probability to include I_t by multiplying the haploid transition probabilities as follows

$$P(q_{t+1} = (k_{t+1,1}, k_{t+1,2}), I_t | q_t = (k_{t,1}, k_{t,2}), \lambda) = \begin{cases} e^{-2G\sigma_t} & \text{if } I_t = 0 \text{ and } k_{t+1,1} = k_{t,1} \text{ and } k_{t+1,2} = k_{t,2} \\ e^{-G\sigma_t}(1 - e^{-G\sigma_t})\alpha_{t,k_{t+1,1}} & \text{if } I_t = 1 \text{ and } k_{t+1,1} \neq k_{t,1} \text{ and } k_{t+1,2} = k_{t,2} \\ e^{-G\sigma_t}(1 - e^{-G\sigma_t})\alpha_{t,k_{t+1,2}} & \text{if } I_t = 1 \text{ and } k_{t+1,1} = k_{t,1} \text{ and } k_{t+1,2} \neq k_{t,2} \\ e^{-G\sigma_t}(1 - e^{-G\sigma_t})(\alpha_{t,k_{t+1,1}} + \alpha_{t,k_{t+1,2}}) & \text{if } I_t = 1 \text{ and } k_{t+1,1} = k_{t,1} \text{ and } k_{t+1,2} = k_{t,2} \\ (1 - e^{-G\sigma_t})^2 \alpha_{t,k_{t+1,1}} \alpha_{t,k_{t+1,2}} & \text{if } I_t = 2 \\ 0 & \text{otherwise} \end{cases} \quad (40)$$

Denote the probability under the diploid model that the sample is in states $(k_{t,1}, k_{t,2})$ at SNP t and states $(k_{t+1,1}, k_{t+1,2})$ at SNP $t+1$ and has indicator variable I_t given observed data O and parameters λ by $\xi_t((k_{t,1}, k_{t,2}), (k_{t+1,1}, k_{t+1,2}), I_t)$. Then

$$\xi_t((k_{t,1}, k_{t,2}), (k_{t+1,1}, k_{t+1,2}), I_t) = \frac{1}{P(O|\lambda)} \alpha_t(k_{t,1}, k_{t,2}) \beta_{t+1}(k_{t+1,1}, k_{t+1,2}) P(O_{t+1} | q_t = (k_{t,1}, k_{t,2}), \lambda) \times P(q_{t+1} = (k_{t+1,1}, k_{t+1,2}), I_t | q_t = (k_{t,1}, k_{t,2}), \lambda) \quad (41)$$

Let $m_{\text{switch},k}^t$ be the number of haplotypes of the sample that switch into state k between SNPs t and $t+1$. We can calculate $\mathbb{E}_\lambda[n_{\text{switch},k}^t]$, and from this $\mathbb{E}_\lambda[n_{\text{stay}}^t]$, by summing across $\mathbb{E}_\lambda[m_{\text{switch},k}^t]$ for all N samples, and so we can calculate the required expectations from the main text by performing the calculations below. Note that we simplify the summation to give a formulation that enables

quadratic versus linear computational complexity in K . A similar approach is done for the haploid model to achieve linear versus quadratic computational complexity (not shown).

$$\begin{aligned}
\mathbb{E}_\lambda[m_{\text{switch},k}^t] &= \sum_{k_1=1}^K \sum_{k_2=1}^K \sum_{k_3=1}^K \sum_{j=0}^2 j \times \left(\xi_t\left((k_1, k_2), (k, k_3), I_t = j\right) + \xi_t\left((k_1, k_2), (k_3, k), I_t = j\right) \right) \\
&= \sum_{k_1=1}^K \sum_{k_2=1}^K 1 \times \left(\xi_t\left((k_1, k_2), (k, k_2), I_t = 1\right) + \xi_{n,t}\left((k_1, k_2), (k_1, k), I_t = 1\right) \right) \\
&+ \sum_{k_1=1}^K \sum_{k_2=1}^K \sum_{k_3=1}^K 2 \times \left(\frac{1}{2} \xi_t\left((k_1, k_2), (k, k_3), I_t = 2\right) + \frac{1}{2} \xi_t\left((k_1, k_2), (k_3, k), I_t = 2\right) \right) \\
&\tag{42}
\end{aligned}$$

$$\begin{aligned}
&= \sum_{k_1=1}^K \sum_{k_3=1}^K 2 \times \xi_t\left((k_1, k_3), (k, k_3), I_t = 1\right) \\
&+ \sum_{k_1=1}^K \sum_{k_2=1}^K \sum_{k_3=1}^K 2 \times \xi_t\left((k_1, k_2), (k, k_3), I_t = 2\right) \\
&\tag{43}
\end{aligned}$$

$$\begin{aligned}
&= 2 \sum_{k_1=1}^K \sum_{k_3=1}^K \frac{\alpha_t(k_1, k_3) \beta_{t+1}(k, k_3) P(O_{t+1}|q_t = (k, k_3), \lambda) \alpha_{t,k} (1 - e^{-G\sigma_t}) e^{-G\sigma_t}}{P(O|\lambda)} \\
&+ 2 \sum_{k_1=1}^K \sum_{k_2=1}^K \sum_{k_3=1}^K \frac{\alpha_t(k_1, k_2) \beta_{t+1}(k, k_3) P(O_{t+1}|q_t = (k, k_3), \lambda) \alpha_{t,k} \alpha_{t,k_3} (1 - e^{-T\sigma_t})^2}{P(O|\lambda)} \\
&= \frac{2\alpha_{t,k}}{P(O|\lambda)} \sum_{k_3=1}^K \left((1 - e^{-G\sigma_t}) e^{-G\sigma_t} \left[\sum_{k_1=1}^K \alpha_t(k_1, k_3) \right] \right. \\
&\left. + \alpha_{k_3}^t (1 - e^{-G\sigma_t})^2 \left[\sum_{k_1=1}^K \sum_{k_2=1}^K \alpha_t(k_1, k_2) \right] \right) \beta_{t+1}(k, k_3) P(O_{t+1}|q_t = (k, k_3), \lambda) \\
&\tag{44}
\end{aligned}$$

Emission probabilities

To update the emission parameters for the diploid model, we use an augmented state space as in for the pseudo-haploid model where we have knowledge of whether emissions were due to the alternate or reference base, and further have knowledge of whether the read came from the maternal or paternal haplotype. Recall that: $\phi_{r,j}^i$ is the probability that observed base j in read R_r came from a read with underlying genotype i ; H_r^j is a variable which takes value 1 if the underlying base is the alternate base and 0 if it is the reference base; and H_r is a variable that takes value 1 if the read came from the maternal haplotype and 2 if from the paternal haplotype. We will use these to calculate the expectation of $n_{k,s}^t$, the number of reads that emit the alternate base ($s = 1$) or reference base ($s = 0$) given they are in state k at their central SNP.

Recall from the main text that the probability of an observation (set of reads) at SNP t in the diploid model is

$$P(O_t|q_t = (k_{t,1}, k_{t,2}), \lambda) = \frac{1}{2}P(R_r|q_t = k_{t,1}, \lambda) + \frac{1}{2}P(R_r|q_t = k_{t,2}, \lambda) \quad (45)$$

For notational convenience set

$$F_{r,j,H_r} = \frac{\frac{1}{2}P(R_r|q_t = k_{H_r}, \lambda)}{\frac{1}{2}P(R_r|q_t = k_{t,1}, \lambda) + \frac{1}{2}P(R_r|q_t = k_{t,2}, \lambda)} \left(\frac{1}{\theta_{t,k_{t,H_r}}\phi_{r,j}^1 + (1 - \theta_{t,k_{t,H_r}})\phi_{r,j}^0} \right) \quad (46)$$

We can therefore calculate the probability that SNP j in read R_r with central SNP c_r has indicator variable H_r^j and H_r and observation for SNP $t = c_r$ of O_t given the pair of hidden states $(k_{t,1}, k_{t,2})$ and parameters λ as

$$P(O_t, H_r^j, H_r|q_t = (k_{t,1}, k_{t,2}), \lambda) = \begin{cases} P(O_t, |q_t = (k_{t,1}, k_{t,2}), \lambda)F_{r,j,H_r}\theta_{t,k_{t,H_r}}\phi_{r,j}^1 & \text{if } H_r^j = 1 \\ P(O_t, |q_t = (k_{t,1}, k_{t,2}), \lambda)F_{r,j,H_r}(1 - \theta_{t,k_{t,H_r}})\phi_{r,j}^0 & \text{if } H_r^j = 0 \end{cases}$$

Denote the probability for SNP j in read R_r that at the central SNP of the read $t = c_r$ is in the pair of states $(k_{t,1}, k_{t,2})$ given the observed data O and parameters λ by $\gamma_t(k_{t,1}, k_{t,2}, H_r^j, H_r)$. Then

$$\gamma_t(k_{t,1}, k_{t,2}, H_r^j, H_r) = \begin{cases} \gamma_t(k_{t,1}, k_{t,2})F_{r,j,H_r}\theta_{t,k_{t,H_r}}\phi_{r,j}^1 & \text{if } H_r^j = 1 \\ \gamma_t(k_{t,1}, k_{t,2})F_{r,j,H_r}(1 - \theta_{t,k_{t,H_r}})\phi_{r,j}^1 & \text{if } H_r^j = 0 \end{cases}$$

Let $\gamma_{n,t}(k_{t,1}, k_{t,2}, H_r^j, H_r)$ be $\gamma_t(k_{t,1}, k_{t,2}, H_r^j, H_r)$ for sample n , and let A_n be the complete set of SNPs j and reads R_r for sample n such that $u_{r,j} = t$. We can calculate the required expectations from the main text as

$$\mathbb{E}_\lambda[n_{k,s}^t] = \sum_{n=1}^N \sum_{(r,j) \in A_n} \sum_{i=1}^K \left(\gamma_{n,c_r}(k, i, H_r^j = s, H_r = 1) + \gamma_{n,c_r}(i, k, H_r^j = s, H_r = 2) \right) \quad (47)$$

1.3 Efficient calculation of forward backward variables

We take the time here to write out the forward backwards calculations that we used for the diploid case, as symmetries in the transition matrix allow us to make the calculation in quadratic, rather than quartic time with respect to K . Similar calculations (not shown) are used for the haploid model to ensure linear versus quadratic computational complexity in K . We note that these calculations are not original and are given in very similar form in the original fastPHASE paper [2], but we reproduce them here as they represent important simplifications for computational reasons

$$\begin{aligned}
\alpha_{t+1}(k_3, k_4) &= \left[\sum_{k_1=1}^K \sum_{k_2=1}^K \alpha_t(k_1, k_2) P(q_{t+1} = (k_3, k_4) | q_t = (k_1, k_2), \lambda) \right] P(O_{t+1} | q_{t+1} = (k_3, k_4), \lambda) \\
&= \left[\alpha_t(k_3, k_4) (e^{-G\sigma_t})^2 + \sum_{k=1}^K e^{-G\sigma_t} (1 - e^{-G\sigma_t}) \alpha_{t,k_3} \alpha_t(k, k_4) + \right. \\
&\quad \left. \sum_{k=1}^K e^{-G\sigma_t} (1 - e^{-G\sigma_t}) \alpha_{t,k_4} \alpha_t(k_3, k) + \right. \\
&\quad \left. \sum_{k_1=1}^K \sum_{k_2=1}^K (1 - e^{-G\sigma_t})^2 \alpha_{t,k_3} \alpha_{t,k_4} \alpha_t(k_1, k_2) \right] P(O_{t+1} | q_{t+1} = (k_3, k_4)) \\
&= \left[\alpha_t(k_3, k_4) (e^{-G\sigma_t})^2 + \alpha_{t,k_3} A_{t,1}(k_4) + \alpha_{t,k_4} A_{t,2}(k_3) + \alpha_{t,k_3} \alpha_{t,k_4} B_t \right] \times \\
&\quad P(O_{t+1} | q_{t+1} = (k_3, k_4), \lambda)
\end{aligned}$$

where

$$A_{t,1}(k_4) = e^{-G\sigma_t} (1 - e^{-G\sigma_t}) \sum_{k=1}^K \alpha_t(k, k_4) \quad (48)$$

$$A_{t,2}(k_3) = e^{-G\sigma_t} (1 - e^{-G\sigma_t}) \sum_{k=1}^K \alpha_t(k_3, k) \quad (49)$$

$$B_t = (1 - e^{-G\sigma_t})^2 \sum_{k_1=1}^K \sum_{k_2=1}^K \alpha_t(k_1, k_2) \quad (50)$$

As such, the forward calculation can be done in quadratic time with respect to the number of ancestral haplotypes K .

Similarly, for the backward calculation we get that

$$\begin{aligned}
\beta_t(k_1, k_2) &= \sum_{k_3=1}^K \sum_{k_4=1}^K P(q_{t+1} = (k_3, k_4) | q_t = (k_1, k_2), \lambda) P(O_{t+1} | q_{t+1} = (k_3, k_4), \lambda) \beta_{t+1}(k_3, k_4) \\
&= (e^{-G\sigma_t})^2 P(O_{t+1} | q_{t+1} = (k_1, k_2), \lambda) \beta_{t+1}(k_1, k_2) + \\
&\quad (e^{-G\sigma_t}) (1 - e^{-G\sigma_t}) \left(\sum_{k=1}^K \alpha_{t,k} P(O_{t+1} | q_{t+1} = (k, k_2), \lambda) \beta_{t+1}(k, k_2) + \right. \\
&\quad \left. \sum_{k=1}^K \alpha_{t,k} P(O_{t+1} | q_{t+1} = (k_1, k), \lambda) \beta_{t+1}(k_1, k) \right) + \\
&\quad (1 - e^{-G\sigma_t})^2 \sum_{k_3=1}^K \sum_{k_4=1}^K \alpha_{t,k_3} \alpha_{t,k_4} P(O_{t+1} | q_{t+1} = (k_3, k_4), \lambda) \beta_{t+1}(k_3, k_4) \\
&= (e^{-G\sigma_t})^2 P(O_{t+1} | q_{t+1} = (k_1, k_2), \lambda) \beta_{t+1}(k_1, k_2) + E_{t,1}(k_2) + E_{t,2}(k_1) + F_t
\end{aligned}$$

where

$$E_{t,1}(k_2) = (e^{-G\sigma_t})(1 - e^{-G\sigma_t}) \sum_{k=1}^K \alpha_{t,k} P(O_{t+1}|q_{t+1} = (k, k_2), \lambda) \beta_{t+1}(k, k_2) \quad (51)$$

$$E_{t,2}(k_1) = (e^{-G\sigma_t})(1 - e^{-G\sigma_t}) \sum_{k=1}^K \alpha_{t,k} P(O_{t+1}|q_{t+1} = (k_1, k), \lambda) \beta_{t+1}(k_1, k) \quad (52)$$

$$F_t = (1 - e^{-G\sigma_t})^2 \sum_{k_3=1}^K \sum_{k_4=1}^K \alpha_{t,k_3} \alpha_{t,k_4} P(O_{t+1}|q_{t+1} = (k_3, k_4), \lambda) \beta_{t+1}(k_3, k_4) \quad (53)$$

1.4 Initialization

Haploid probabilities π_k are initialized with equal weights $\pi_k = \frac{1}{K}$, as are diploid priors $\pi_{k_1, k_2} = \frac{1}{K \times K}$. The state probabilities $\alpha_{t,k}$ are also initialized with equal weights $\alpha_{t,k} = \frac{1}{K}$. The recombination distance is initialized assuming a constant recombination rate multiplied by the physical distance between SNPs, for example assuming $\sigma_t = d_t \times 0.5\text{cM/Mb}$ where d_t is the physical distance between SNPs t and $t + 1$. Finally, given a lower bound δ on emission probabilities, for example $\delta = 0.0001$, $\theta_{t,k}$ are sampled from a uniform distribution with minimum value δ and maximum value $1 - \delta$. Note that G is left as a user set parameter, which can be approximated for outbred populations using external estimates of N_e with $G = \frac{4N_e}{K}$.

1.5 Parameter bounding

After parameter updating, newly calculated parameters are bounded with default but user tunable parameters. Prior probabilities π_k , new state parameters $\alpha_{t,k}$, and emission probabilities $\theta_{t,k}$ (and $1 - \theta_{t,k}$) whose values are less than a threshold are set equal to that threshold, and then probabilities re-normalized as appropriate to have sum 1. Under default conditions this bound is 1×10^{-4} . For the recombination distance, values of σ_t that exceed implied upper (default 100 cM/Mb) and lower (default 0.1 cM/Mb) bounds are reset to the bound value.

1.6 Heuristics

Since emission probabilities θ are initialized at random, STITCH can get stuck in local minima, for which two heuristics are employed at various (default) iterations. First, to help overcome unnecessary switches between ancestral haplotype backgrounds, at iterations 4, 8, 12 and 16, pairs of haplotype states are calculated for each sample between pairs of nearby SNPs (starting at SNP 51, then every further 100th SNP) by multiplying their marginal ancestral haplotype

probabilities. If, across all samples, for each pair of nearby SNPs, there exists a re-ordering of ancestral haplotype states that minimizes switching, then that switch, or switches, is performed, and local SNPs (plus or minus 20 from the break) are reset with θ from a $U(0, 1)$ distribution. Second, to help fill unused ancestral haplotypes, and to overcome superimposed ancestral haplotypes, at iterations 6, 10, 14 and 18, ancestral haplotype usage in the most recent iteration is discretized by averaging over 100 SNP intervals, and every continuous interval of infrequently used ancestral haplotype ($< 0.5\%$) is identified. Values of θ over each interval are then refilled for that ancestral haplotype by copying from another sampled ancestral haplotype chosen with sampling probability proportional to ancestral haplotype usage over that interval. θ is then reset using 80% of these filled values and 20% noise from a $U(0, 1)$ distribution.

1.7 Guidance behind parameter options

STITCH contains many parameter options that can be modified by the user, for example upper and lower bounds on recombination rate. However, most of these are reasonable for the majority of anticipated applications of STITCH. For the analyses presented here for the CFW and CONVERGE populations, we varied: K (option `K`), the number of ancestral haplotypes; whether the diploid or pseudo-haploid method was used (option `method`); the number of pseudo-haploid iterations (option `switchModelIteration`); the number of generations when the population was founded (or can be so approximated) G (option `nGen`) (which we set as 100 for the CFW analyses and $\frac{4 \times 20000}{K}$ for the CONVERGE studies). We also, for model evaluation purposes only, invoked a flag on whether reads were split into new reads containing one SNP each (option `readAware`), the number of computer cores available to the process (option `nCores`), and whether the process is running in a server or cluster environment (option `environment`).

We anticipate that in using STITCH, the majority of users will achieve desired results, both in terms of accuracy and computational speed, through varying K , G , the method (diploid or pseudo-haploid), and the number of pseudo-haploid iterations.

In terms of selecting K , the diploid or pseudo-haploid method, and the number of pseudo-haploid iterations, we recommend imputing a small region of the genome, such as a chromosome, using the diploid mode with a range of K , and then evaluate performance. We recommend that to evaluate imputation performance, users obtain validation data, using either genotyping microarrays or higher coverage sequencing (like 10X). In the absence of external validation data, we recommend the info score distribution or its average. If, for the diploid method and a choice of K , results start to deteriorate, then choose the diploid mode and K that gave optimal performance. If results do not deteriorate but become computationally impractical, we recommend applying the pseudo-haploid method for a range of pseudo-haploid and diploid iterations (as was done here for CONVERGE), and choosing the combination that gives optimal results under the given computational constraints.

For G (or `nGen`), we recommend setting this to a reasonable *a priori* esti-

mate, like was available for the CFW mice, or to use $\frac{4 \times N_e}{K}$, when the population is wild or has not been through a strong bottleneck. We note that STITCH should be fairly robust to this parameter choice. Users may also increase the minimum and maximum allowed recombination rates if they are less certain about this parameter.

Finally, while we do not give specific guidance on study design strategies and sequencing depths, we note that in designing low coverage sequencing only studies, users should try to ensure adequate population sequencing coverage to ensure the ancestral haplotypes are well reconstructed, particularly in the case when the founding structure is well known. For example, if a population was founded with $K = 8$ haplotypes, then to achieve a given level of per-ancestral haplotype coverage (e.g. 30X), while sequencing each sample at a given level (e.g. 0.2X), one should consider sequencing in excess of $\frac{30 \times K}{0.2} = 1200$ samples. Drift in the population (*i.e.* non-equal ancestral haplotype usage in the population) would require additional samples or depth for reconstruction of rare haplotypes in the population.

2 Supplementary Tables

Supplementary Table 1A: Genotype concordance for CFW using STITCH (K=4, diploid) at all SNPs Results give genotype concordance stratified by genotype class and allele frequency. Discrete genotype calls are generated for imputation as the the genotype with the maximum genotype posterior probability. Results are given genome-wide (autosome and chromosome X). Allele freqs = allele frequencies are the frequency of the minor allele. Type is either High Cov = high coverage (10X) sequencing (4 samples) or Array = MegaMuga (44 samples). Columns contain either Num = Number of non-missing genotypes considered (samples times SNPs for sequencing or array), or Per = Percent of imputed best guess genotypes that match sequencing or array genotypes. Hom major = homozygous for the major allele, Het = heterozygous, Hom Minor = homozygous for the minor allele. Note that truth (sequencing or array) genotypes contain some missing data

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|-----------|---------|---------------|---------------|
| [0,0.01) | High Cov | 1,139,724 | 99.98 | 3,958 | 17.08 | 14 | 0 |
| [0.01,0.02) | High Cov | 1,016,641 | 99.84 | 20,516 | 64.72 | 124 | 4.84 |
| [0.02,0.05) | High Cov | 3,756,581 | 99.71 | 213,186 | 81.62 | 3,407 | 52.51 |
| [0.05,0.1) | High Cov | 4,072,071 | 99.57 | 552,747 | 89.85 | 22,958 | 75.8 |
| [0.1,0.2) | High Cov | 3,781,946 | 99.23 | 1,164,122 | 92.84 | 117,126 | 90.83 |
| [0.2,0.3) | High Cov | 1,973,117 | 98.69 | 1,274,072 | 94.86 | 204,474 | 93.69 |
| [0.3,0.4) | High Cov | 1,201,299 | 98.08 | 1,296,792 | 95.83 | 328,621 | 95.31 |
| [0.4,0.5] | High Cov | 730,043 | 97.29 | 1,417,056 | 96.59 | 452,854 | 96.13 |
| [0,0.01) | Array | 3,101 | 99.97 | 106 | 56.6 | 3 | 33.33 |
| [0.01,0.02) | Array | 19,788 | 99.93 | 803 | 84.43 | 20 | 50 |
| [0.02,0.05) | Array | 135,504 | 99.9 | 9,386 | 93.51 | 312 | 73.08 |
| [0.05,0.1) | Array | 161,620 | 99.86 | 24,850 | 95.99 | 1,238 | 81.18 |
| [0.1,0.2) | Array | 163,416 | 99.76 | 55,438 | 97.59 | 5,965 | 94.25 |
| [0.2,0.3) | Array | 82,595 | 99.57 | 54,880 | 98.27 | 9,586 | 97.83 |
| [0.3,0.4) | Array | 45,709 | 99.33 | 49,416 | 98.36 | 14,094 | 98.24 |
| [0.4,0.5] | Array | 33,823 | 99.21 | 53,605 | 98.79 | 22,887 | 98.93 |

Supplementary Table 1B: Genotype concordance for CFW using Beagle (default) at all SNPs

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|-----------|---------|---------------|---------------|
| [0,0.01) | High Cov | 1,139,728 | 95 | 3,958 | 19.45 | 10 | 0 |
| [0.01,0.02) | High Cov | 1,016,641 | 91.05 | 20,516 | 14.9 | 124 | 0 |
| [0.02,0.05) | High Cov | 3,756,615 | 90.2 | 213,186 | 16.04 | 3,373 | 0.18 |
| [0.05,0.1) | High Cov | 4,072,282 | 84.46 | 552,747 | 23.3 | 22,747 | 0.32 |
| [0.1,0.2) | High Cov | 3,782,741 | 65.12 | 1,164,122 | 45.7 | 116,331 | 0.87 |
| [0.2,0.3) | High Cov | 1,973,436 | 23.14 | 1,274,072 | 84.16 | 204,155 | 2.16 |
| [0.3,0.4) | High Cov | 1,196,914 | 14.13 | 1,296,792 | 90.73 | 333,006 | 3.85 |
| [0.4,0.5] | High Cov | 718,518 | 13.02 | 1,417,056 | 90.02 | 464,379 | 6.1 |
| [0,0.01) | Array | 3,101 | 91.36 | 106 | 17.92 | 3 | 0 |
| [0.01,0.02) | Array | 19,788 | 93.68 | 803 | 15.44 | 20 | 0 |
| [0.02,0.05) | Array | 135,504 | 90.06 | 9,386 | 17.45 | 312 | 0 |
| [0.05,0.1) | Array | 161,581 | 84.58 | 24,850 | 23.57 | 1,277 | 0.23 |
| [0.1,0.2) | Array | 163,393 | 63.19 | 55,438 | 48.4 | 5,988 | 0.78 |
| [0.2,0.3) | Array | 82,553 | 21.1 | 54,880 | 88.32 | 9,628 | 1.65 |
| [0.3,0.4) | Array | 45,562 | 15.71 | 49,416 | 92.11 | 14,241 | 1.94 |
| [0.4,0.5] | Array | 33,002 | 15.13 | 53,605 | 91.53 | 23,708 | 3.28 |

Supplementary Table 1C: Genotype concordance for CFW using findhap (maxlen=10000, minlen=100, steps=3, iters=4) at all SNPs

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|-----------|---------|---------------|---------------|
| [0,0.01) | High Cov | 1,138,593 | 96.5 | 3,958 | 30.19 | 1,145 | 12.4 |
| [0.01,0.02) | High Cov | 1,012,728 | 92.28 | 20,516 | 52.67 | 4,037 | 9.14 |
| [0.02,0.05) | High Cov | 3,739,212 | 89.72 | 213,186 | 81.45 | 20,776 | 13.59 |
| [0.05,0.1) | High Cov | 4,040,317 | 89.29 | 552,747 | 86.98 | 54,712 | 15.61 |
| [0.1,0.2) | High Cov | 3,730,302 | 91.25 | 1,164,122 | 79.28 | 168,770 | 34.79 |
| [0.2,0.3) | High Cov | 1,966,686 | 93.17 | 1,274,072 | 71.31 | 210,905 | 51.8 |
| [0.3,0.4) | High Cov | 1,198,042 | 90.45 | 1,296,792 | 62.92 | 331,878 | 64.81 |
| [0.4,0.5] | High Cov | 722,465 | 87.28 | 1,417,056 | 57.79 | 460,432 | 74.34 |
| [0,0.01) | Array | 3,101 | 97.84 | 106 | 37.74 | 3 | 0 |
| [0.01,0.02) | Array | 19,745 | 93.78 | 803 | 69.12 | 63 | 7.94 |
| [0.02,0.05) | Array | 135,174 | 90.06 | 9,386 | 86.74 | 642 | 19.47 |
| [0.05,0.1) | Array | 160,288 | 87.06 | 24,850 | 89.34 | 2,570 | 18.6 |
| [0.1,0.2) | Array | 161,093 | 89.2 | 55,438 | 82.78 | 8,288 | 33.69 |
| [0.2,0.3) | Array | 82,595 | 91.56 | 54,880 | 73.4 | 9,586 | 52.22 |
| [0.3,0.4) | Array | 45,673 | 87.74 | 49,416 | 68 | 14,130 | 61.78 |
| [0.4,0.5] | Array | 33,145 | 83.72 | 53,605 | 63.99 | 23,565 | 69.89 |

Supplementary Table 1D: Genotype concordance for CFW using STITCH (K=4, diploid) at all SNPs (post QC)

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|-----------|---------|---------------|---------------|
| [0,0.01) | High Cov | 1,139,724 | 99.98 | 3,958 | 17.08 | 14 | 0 |
| [0.01,0.02) | High Cov | 1,016,641 | 99.84 | 20,516 | 64.72 | 124 | 4.84 |
| [0.02,0.05) | High Cov | 3,756,581 | 99.71 | 213,186 | 81.62 | 3,407 | 52.51 |
| [0.05,0.1) | High Cov | 4,072,071 | 99.57 | 552,747 | 89.85 | 22,958 | 75.8 |
| [0.1,0.2) | High Cov | 3,781,946 | 99.23 | 1,164,122 | 92.84 | 117,126 | 90.83 |
| [0.2,0.3) | High Cov | 1,973,117 | 98.69 | 1,274,072 | 94.86 | 204,474 | 93.69 |
| [0.3,0.4) | High Cov | 1,201,299 | 98.08 | 1,296,792 | 95.83 | 328,621 | 95.31 |
| [0.4,0.5] | High Cov | 730,043 | 97.29 | 1,417,056 | 96.59 | 452,854 | 96.13 |
| [0,0.01) | Array | 3,101 | 99.97 | 106 | 56.6 | 3 | 33.33 |
| [0.01,0.02) | Array | 19,788 | 99.93 | 803 | 84.43 | 20 | 50 |
| [0.02,0.05) | Array | 135,504 | 99.9 | 9,386 | 93.51 | 312 | 73.08 |
| [0.05,0.1) | Array | 161,620 | 99.86 | 24,850 | 95.99 | 1,238 | 81.18 |
| [0.1,0.2) | Array | 163,416 | 99.76 | 55,438 | 97.59 | 5,965 | 94.25 |
| [0.2,0.3) | Array | 82,595 | 99.57 | 54,880 | 98.27 | 9,586 | 97.83 |
| [0.3,0.4) | Array | 45,709 | 99.33 | 49,416 | 98.36 | 14,094 | 98.24 |
| [0.4,0.5] | Array | 33,823 | 99.21 | 53,605 | 98.79 | 22,887 | 98.93 |

Supplementary Table 2: Performance of CFW study under different programs and options Results are given for chromosomes 18 and 19. All STITCH results are for the diploid model with 40 iterations. Program options are as follows. For STITCH, RU refers to read unaware (i.e. split each read spanning multiple SNPs into sub-reads spanning one read each). For Beagle, shown are the number of iterations (i.e. burnin-its, phase-its, and impute-its to this value), window is the window size, and msf is the (singlescale) model scale factor. For findhap, options correspond directly to parameter options. Note that times for STITCH do not include the generation of input data from BAMs, which took about 1-1.5 hours per chromosome for chromosomes 18 and 19, irrespective of other program options. Similarly, times for findhap do not include conversion time from VCF to the findhap input format. Av r² is the average r^2 for SNPs on the Illumina MegaMUGA array, with no filtration for QC for any method. Time is the average time in hours for chromosomes 18 and 19, where all programs were run on 1 core on 2.60 GHz Intel E5-2650 chips

| Program | Options | Time | Av r ² |
|---------|---|------|-------------------|
| STITCH | K=2 | 7.2 | 0.622 |
| STITCH | K=3 | 11.3 | 0.957 |
| STITCH | K=4 | 18.6 | 0.972 |
| STITCH | K=5 | 25.3 | 0.97 |
| STITCH | K=6 | 37.5 | 0.966 |
| STITCH | K=7 | 49 | 0.964 |
| STITCH | K=8 | 59 | 0.967 |
| STITCH | K=4, RU | 18.8 | 0.873 |
| Beagle | its=5, window=50000, msf=1 | 6.1 | 0.074 |
| Beagle | its=5, window=1000000, msf=1 | 4.7 | 0.073 |
| Beagle | its=10, window=50000, msf=1 | 17.2 | 0.085 |
| Beagle | its=20, window=50000, msf=1 | 34.1 | 0.109 |
| Beagle | its=5, window=50000, msf=0.4 | 72.4 | 0.088 |
| Beagle | its=5, window=50000, msf=0.6 | 7.7 | 0.079 |
| Beagle | its=5, window=50000, msf=0.8 | 6.6 | 0.073 |
| Beagle | its=5, window=50000, msf=1.0 | 5.3 | 0.072 |
| Beagle | its=5, window=50000, msf=1.2 | 4.9 | 0.071 |
| Beagle | its=5, window=50000, msf=1.4 | 5.7 | 0.071 |
| Beagle | its=5, window=50000, msf=1.6 | 5.2 | 0.071 |
| Beagle | its=5, window=50000, msf=1.8 | 5.2 | 0.071 |
| Beagle | its=5, window=50000, msf=2.0 | 5.1 | 0.071 |
| findhap | maxlen=100000, minlen=1000, steps=3, iters=4 | 0.6 | 0.225 |
| findhap | maxlen=100000, minlen=1000, steps=2, iters=6 | 0.7 | 0.226 |
| findhap | maxlen=100000, minlen=1000, steps=5, iters=10 | 2.2 | 0.15 |
| findhap | maxlen=10000, minlen=100, steps=3, iters=4 | 0.5 | 0.523 |
| findhap | maxlen=50000, minlen=500, steps=3, iters=4 | 0.6 | 0.281 |
| findhap | maxlen=200000, minlen=2000, steps=3, iters=4 | 0.5 | 0.169 |

Supplementary Table 3: Performance of CONVERGE study under different programs and options with no reference panel Results are given for the first 10 Mbp of chromosome 20, run in 0.5 Mbp regions with 0.1 Mbp buffers. Program options are as follows. For STITCH, all options were run using 40 EM iterations, split into either diploid (D) or pseudo-haploid (PH) iterations, while RU refers to read unaware (i.e. split each read spanning multiple SNPs into sub-reads spanning one read each). For Beagle, shown are the number of iterations (i.e. burnin-its, phase-its, and impute-its to this value). For findhap, options correspond directly to parameter options. Note that times for STITCH do not include the generation of input data from BAMs, which took about 30 minutes per region, irrespective of other program options. Similarly, times for findhap do not include conversion time from VCF to the findhap input format. Av r2 is the average r^2 for SNPs on the Illumina HumanOmniZhongHua-8 array for common (MAF 5% to 95%) variants, with no filtration for QC for any method. Time is the average in hours for each 0.5Mbp region, where all programs were run on 4 cores on 2.60 GHz Intel E5-2650 chips.

| Program | Options | Time | Av r2 |
|---------|---|------|-------|
| STITCH | K=20, its=40D | 24.5 | 0.922 |
| STITCH | K=20, its=40PH | 8.0 | 0.875 |
| STITCH | K=20, its=34PH;6D | 10.6 | 0.920 |
| STITCH | K=20, its=35PH;5D | 9.9 | 0.919 |
| STITCH | K=20, its=36PH;4D | 9.6 | 0.918 |
| STITCH | K=20, its=37PH;3D | 9.3 | 0.917 |
| STITCH | K=20, its=38PH;2D | 8.8 | 0.911 |
| STITCH | K=20, its=39PH;1D | 8.4 | 0.898 |
| STITCH | K=20, its=38PH;2D, RU | 9.4 | 0.910 |
| STITCH | K=30, its=40D | 52.2 | 0.927 |
| STITCH | K=30, its=38PH;2D | 12.4 | 0.917 |
| STITCH | K=40, its=38PH;2D | 16.5 | 0.920 |
| STITCH | K=60, its=38PH;2D | 27.7 | 0.923 |
| STITCH | K=80, its=38PH;2D | 42.2 | 0.925 |
| STITCH | K=100, its=38PH;2D | 61.1 | 0.927 |
| Beagle | its=5 | 12.5 | 0.874 |
| findhap | maxlen=100000, minlen=1000, steps=3, iters=4 | 0.4 | 0.437 |
| findhap | maxlen=100000, minlen=1000, steps=2, iters=6 | 0.4 | 0.437 |
| findhap | maxlen=100000, minlen=1000, steps=5, iters=10 | 1.4 | 0.426 |
| findhap | maxlen=10000, minlen=100, steps=3, iters=4 | 0.3 | 0.434 |
| findhap | maxlen=50000, minlen=500, steps=3, iters=4 | 0.4 | 0.448 |
| findhap | maxlen=200000, minlen=2000, steps=3, iters=4 | 0.5 | 0.414 |

Supplementary Table 4A: Genotype concordance for CONVERGE using STITCH (K=40, 38 PH iterations, 2 D iterations) (without a reference panel) at all SNPs Results give genotype concordance stratified by genotype class and allele frequency. Discrete genotype calls are generated for imputation as the the genotype with the maximum genotype posterior probability. Results are given for the first 10 Mbp region of chromosome 20, run in 20 0.5 Mbp regions with 0.1 Mbp buffers. Allele freqs = allele frequencies are the frequency of the minor allele. Type is either High Cov = high coverage (10X) sequencing (9 samples) or Array = HumanOmniZhongHua-8 (72 samples). Columns contain either Num = Number of non-missing genotypes considered (samples times SNPs for sequencing or array), or Per = Percent of imputed best guess genotypes that match sequencing or array genotypes. Hom major = homozygous for the major allele, Het = heterozygous, Hom Minor = homozygous for the minor allele. Note that truth (sequencing or array) genotypes contain some missing data

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 16,234 | 99.98 | 879 | 27.19 | 14 | 0 |
| [0.01,0.02) | High Cov | 3,973 | 99.72 | 483 | 62.11 | 6 | 16.67 |
| [0.02,0.05) | High Cov | 11,325 | 99.59 | 1,725 | 83.65 | 31 | 19.35 |
| [0.05,0.1) | High Cov | 14,634 | 99.33 | 2,931 | 91.23 | 116 | 65.52 |
| [0.1,0.2) | High Cov | 28,004 | 99.13 | 10,059 | 96.01 | 1,183 | 84.53 |
| [0.2,0.3) | High Cov | 18,880 | 98.59 | 12,285 | 96.81 | 2,134 | 90.63 |
| [0.3,0.4) | High Cov | 15,750 | 97.69 | 15,356 | 97.49 | 4,300 | 94.12 |
| [0.4,0.5] | High Cov | 10,469 | 96.89 | 16,280 | 97.61 | 7,445 | 96.15 |
| [0,0.01) | Array | 9,691 | 99.95 | 100 | 70 | 0 | NA |
| [0.01,0.02) | Array | 4,519 | 99.91 | 156 | 73.08 | 0 | NA |
| [0.02,0.05) | Array | 11,883 | 99.82 | 834 | 88.13 | 12 | 83.33 |
| [0.05,0.1) | Array | 18,317 | 99.42 | 3,003 | 91.44 | 114 | 73.68 |
| [0.1,0.2) | Array | 29,193 | 98.92 | 10,165 | 93.38 | 866 | 85.68 |
| [0.2,0.3) | Array | 20,139 | 97.89 | 12,835 | 94.39 | 2,201 | 89.41 |
| [0.3,0.4) | Array | 14,833 | 97.2 | 15,879 | 95.44 | 4,171 | 93.02 |
| [0.4,0.5] | Array | 10,172 | 96.29 | 16,290 | 95.67 | 6,766 | 94.8 |

Supplementary Table 4B: Genotype concordance for CONVERGE using Beagle (default) (without a reference panel) at all SNPs

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 16,234 | 100 | 879 | 56.09 | 14 | 0 |
| [0.01,0.02) | High Cov | 3,973 | 100 | 483 | 64.6 | 6 | 0 |
| [0.02,0.05) | High Cov | 11,325 | 99.97 | 1,725 | 73.74 | 31 | 6.45 |
| [0.05,0.1) | High Cov | 14,634 | 99.64 | 2,931 | 84.61 | 116 | 50.86 |
| [0.1,0.2) | High Cov | 28,004 | 99.49 | 10,059 | 90.55 | 1,183 | 82.25 |
| [0.2,0.3) | High Cov | 18,880 | 98.98 | 12,285 | 93.09 | 2,134 | 88.71 |
| [0.3,0.4) | High Cov | 15,750 | 97.96 | 15,356 | 94.27 | 4,300 | 93.37 |
| [0.4,0.5] | High Cov | 10,469 | 97.05 | 16,280 | 95.12 | 7,445 | 96.55 |
| [0,0.01) | Array | 9,691 | 100 | 100 | 54 | 0 | NA |
| [0.01,0.02) | Array | 4,519 | 100 | 156 | 57.69 | 0 | NA |
| [0.02,0.05) | Array | 11,883 | 99.91 | 834 | 76.5 | 12 | 75 |
| [0.05,0.1) | Array | 18,317 | 99.77 | 3,003 | 81.22 | 114 | 69.3 |
| [0.1,0.2) | Array | 29,193 | 99.47 | 10,165 | 86.54 | 866 | 83.03 |
| [0.2,0.3) | Array | 20,139 | 98.46 | 12,835 | 89.44 | 2,201 | 86.37 |
| [0.3,0.4) | Array | 14,833 | 97.5 | 15,879 | 91.91 | 4,171 | 91.63 |
| [0.4,0.5] | Array | 10,172 | 96.23 | 16,290 | 92.98 | 6,766 | 94.66 |

Supplementary Table 4C: Genotype concordance for CONVERGE using findhap (maxlen=50000, minlen=500, steps=3, iters=4) (without a reference panel) at all SNPs

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 13,485 | 99.69 | 562 | 48.22 | 10 | 0 |
| [0.01,0.02) | High Cov | 3,701 | 99.08 | 453 | 60.71 | 5 | 20 |
| [0.02,0.05) | High Cov | 10,771 | 97.96 | 1,623 | 61.06 | 27 | 3.7 |
| [0.05,0.1) | High Cov | 14,328 | 94.45 | 2,875 | 68.49 | 115 | 18.26 |
| [0.1,0.2) | High Cov | 26,632 | 93.43 | 9,443 | 69 | 1,083 | 42.84 |
| [0.2,0.3) | High Cov | 17,884 | 87.82 | 11,566 | 66.72 | 1,964 | 50.61 |
| [0.3,0.4) | High Cov | 15,229 | 84.5 | 14,665 | 67.23 | 4,212 | 60.73 |
| [0.4,0.5] | High Cov | 9,766 | 77.63 | 15,546 | 67.34 | 7,050 | 67.48 |
| [0,0.01) | Array | 8,690 | 99.57 | 94 | 46.81 | 0 | NA |
| [0.01,0.02) | Array | 3,894 | 98.02 | 133 | 46.62 | 0 | NA |
| [0.02,0.05) | Array | 11,302 | 97.34 | 773 | 56.4 | 8 | 37.5 |
| [0.05,0.1) | Array | 17,884 | 93.63 | 2,934 | 63.53 | 113 | 20.35 |
| [0.1,0.2) | Array | 27,294 | 91.11 | 9,529 | 62.09 | 813 | 31.12 |
| [0.2,0.3) | Array | 18,845 | 85.16 | 12,052 | 63.57 | 2,046 | 40.27 |
| [0.3,0.4) | Array | 14,260 | 78.72 | 15,331 | 65.34 | 4,068 | 49.68 |
| [0.4,0.5] | Array | 9,903 | 72.2 | 15,787 | 66.08 | 6,603 | 58.47 |

Supplementary Table 4D: Genotype concordance for CONVERGE using STITCH (K=40, 38 PH iterations, 2 D iterations) (without a reference panel) at all SNPs (that pass QC)

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 6,266 | 99.97 | 241 | 83.82 | 1 | 0 |
| [0.01,0.02) | High Cov | 2,725 | 99.67 | 323 | 85.45 | 2 | 50 |
| [0.02,0.05) | High Cov | 10,001 | 99.59 | 1,541 | 90.53 | 29 | 17.24 |
| [0.05,0.1) | High Cov | 13,916 | 99.35 | 2,760 | 94.28 | 109 | 68.81 |
| [0.1,0.2) | High Cov | 27,327 | 99.19 | 9,773 | 97.24 | 1,155 | 85.97 |
| [0.2,0.3) | High Cov | 18,291 | 98.9 | 11,938 | 97.76 | 2,064 | 92.34 |
| [0.3,0.4) | High Cov | 15,294 | 98.12 | 14,869 | 98.14 | 4,170 | 95.4 |
| [0.4,0.5] | High Cov | 10,174 | 97.67 | 15,802 | 98.15 | 7,260 | 96.85 |
| [0,0.01) | Array | 6,763 | 99.94 | 76 | 90.79 | 0 | NA |
| [0.01,0.02) | Array | 3,898 | 99.9 | 132 | 84.09 | 0 | NA |
| [0.02,0.05) | Array | 11,068 | 99.83 | 787 | 92.12 | 11 | 90.91 |
| [0.05,0.1) | Array | 17,968 | 99.49 | 2,925 | 92.89 | 110 | 76.36 |
| [0.1,0.2) | Array | 27,902 | 99.06 | 9,643 | 95.49 | 809 | 90.36 |
| [0.2,0.3) | Array | 18,709 | 98.36 | 11,832 | 96.49 | 2,047 | 92.82 |
| [0.3,0.4) | Array | 14,263 | 97.69 | 15,252 | 96.48 | 4,000 | 95.15 |
| [0.4,0.5] | Array | 9,573 | 97.49 | 15,336 | 96.77 | 6,376 | 96.86 |

Supplementary Table 5: Performance of CONVERGE study under different programs and options with a reference panel Results are given for the first 10 Mbp of chromosome 20, run in 0.5 Mbp regions with 0.1 Mbp buffers. Program options are as follows. For STITCH, all options were run using 40 EM iterations, split into either diploid (D) or pseudo-haploid (PH) iterations. For Beagle, shown are the number of iterations (i.e. burnin-its, phase-its, and impute-its to this value). Note that times for STITCH do not include the generation of input data from BAMs, which took about 30 minutes per region, irrespective of other program options. Av r2 is the average r^2 for SNPs on the Illumina HumanOmniZhongHua-8 array for common (MAF 5% to 95%) variants, with no filtration for QC for any method. Time is the average in hours for each 0.5Mbp region, where all programs were run on 4 cores on 2.60 GHz Intel E5-2650 chips.

| Program | Options | Time | Av r2 |
|---------|---------------------|-------|-------|
| STITCH | K=20, its=38PH;2D | 5.4 | 0.911 |
| STITCH | K=40, its=38PH;2D | 10.2 | 0.922 |
| STITCH | K=60, its=38PH;2D | 16.6 | 0.925 |
| Beagle | its=5, no ref panel | 7.8 | 0.886 |
| Beagle | its=4 | 114.4 | 0.946 |
| Beagle | its=3 | 74.5 | 0.943 |
| Beagle | its=2 | 39.7 | 0.939 |
| Beagle | its=1 | 12.0 | 0.930 |

Supplementary Table 6A: Genotype concordance for CONVERGE using STITCH (K=40, 38 PH iterations, 2 D iterations) (without a reference panel) at reference panel SNPs (1000G ASN) Results give genotype concordance stratified by genotype class and allele frequency. Discrete genotype calls are generated for imputation as the the genotype with the maximum genotype posterior probability. Results are given for the first 10 Mbp region of chromosome 20, run in 20 0.5 Mbp regions with 0.1 Mbp buffers. Allele freqs = allele frequencies are the frequency of the minor allele. Type is either High Cov = high coverage (10X) sequencing (9 samples) or Array = HumanOmniZhongHua-8 (72 samples). Columns contain either Num = Number of non-missing genotypes considered (samples times SNPs for sequencing or array), or Per = Percent of imputed best guess genotypes that match sequencing or array genotypes. Hom major = homozygous for the major allele, Het = heterozygous, Hom Minor = homozygous for the minor allele. Note that truth (sequencing or array) genotypes contain some missing data

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 13,968 | 99.96 | 596 | 35.91 | 10 | 0 |
| [0.01,0.02) | High Cov | 3,798 | 99.74 | 460 | 63.04 | 6 | 0 |
| [0.02,0.05) | High Cov | 11,259 | 99.72 | 1,714 | 85.3 | 31 | 16.13 |
| [0.05,0.1) | High Cov | 14,634 | 99.33 | 2,931 | 91.88 | 116 | 63.79 |
| [0.1,0.2) | High Cov | 27,980 | 99.1 | 10,054 | 95.91 | 1,183 | 85.88 |
| [0.2,0.3) | High Cov | 18,875 | 98.56 | 12,282 | 96.88 | 2,133 | 90.53 |
| [0.3,0.4) | High Cov | 15,737 | 97.59 | 15,315 | 97.43 | 4,298 | 94.04 |
| [0.4,0.5] | High Cov | 10,463 | 96.82 | 16,261 | 97.72 | 7,434 | 95.96 |
| [0,0.01) | Array | 8,975 | 99.96 | 97 | 64.95 | 0 | NA |
| [0.01,0.02) | Array | 4,519 | 99.82 | 156 | 80.13 | 0 | NA |
| [0.02,0.05) | Array | 11,740 | 99.74 | 833 | 88.36 | 12 | 83.33 |
| [0.05,0.1) | Array | 18,317 | 99.45 | 3,003 | 91.81 | 114 | 76.32 |
| [0.1,0.2) | Array | 29,144 | 98.93 | 10,142 | 93.7 | 866 | 86.95 |
| [0.2,0.3) | Array | 20,139 | 97.85 | 12,835 | 94.66 | 2,201 | 90.37 |
| [0.3,0.4) | Array | 14,833 | 97.01 | 15,879 | 95.54 | 4,171 | 92.78 |
| [0.4,0.5] | Array | 10,172 | 96.06 | 16,290 | 95.75 | 6,766 | 94.77 |

Supplementary Table 6B: Genotype concordance for CONVERGE using Beagle (default) (without a reference panel) at reference panel SNPs (1000G ASN)

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 13,968 | 100 | 596 | 55.54 | 10 | 0 |
| [0.01,0.02) | High Cov | 3,798 | 100 | 460 | 66.09 | 6 | 0 |
| [0.02,0.05) | High Cov | 11,259 | 99.95 | 1,714 | 75.61 | 31 | 12.9 |
| [0.05,0.1) | High Cov | 14,634 | 99.64 | 2,931 | 85.77 | 116 | 53.45 |
| [0.1,0.2) | High Cov | 27,980 | 99.41 | 10,054 | 91.87 | 1,183 | 83.94 |
| [0.2,0.3) | High Cov | 18,875 | 98.95 | 12,282 | 93.93 | 2,133 | 89.45 |
| [0.3,0.4) | High Cov | 15,737 | 97.99 | 15,315 | 95.09 | 4,298 | 93.9 |
| [0.4,0.5] | High Cov | 10,463 | 97.24 | 16,261 | 95.79 | 7,434 | 96.7 |
| [0,0.01) | Array | 8,975 | 100 | 97 | 56.7 | 0 | NA |
| [0.01,0.02) | Array | 4,519 | 100 | 156 | 58.33 | 0 | NA |
| [0.02,0.05) | Array | 11,740 | 99.88 | 833 | 78.63 | 12 | 75 |
| [0.05,0.1) | Array | 18,317 | 99.72 | 3,003 | 83.25 | 114 | 71.05 |
| [0.1,0.2) | Array | 29,144 | 99.32 | 10,142 | 88.55 | 866 | 84.06 |
| [0.2,0.3) | Array | 20,139 | 98.38 | 12,835 | 90.7 | 2,201 | 87.6 |
| [0.3,0.4) | Array | 14,833 | 97.57 | 15,879 | 92.88 | 4,171 | 92.14 |
| [0.4,0.5] | Array | 10,172 | 96.26 | 16,290 | 93.49 | 6,766 | 95.26 |

Supplementary Table 6C: Genotype concordance for CONVERGE using Beagle (its=3) (with a reference panel) at reference panel SNPs (1000G ASN)

| Allele freqs | Type | Num Hom Major | Per Hom Major | Num Het | Per Het | Num Hom Minor | Per Hom Minor |
|--------------|----------|---------------|---------------|---------|---------|---------------|---------------|
| [0,0.01) | High Cov | 13,968 | 99.96 | 596 | 66.95 | 10 | 0 |
| [0.01,0.02) | High Cov | 3,798 | 99.63 | 460 | 81.09 | 6 | 33.33 |
| [0.02,0.05) | High Cov | 11,259 | 99.86 | 1,714 | 91.54 | 31 | 19.35 |
| [0.05,0.1) | High Cov | 14,634 | 99.55 | 2,931 | 95.19 | 116 | 67.24 |
| [0.1,0.2) | High Cov | 27,980 | 99.32 | 10,054 | 97.13 | 1,183 | 88.33 |
| [0.2,0.3) | High Cov | 18,875 | 98.95 | 12,282 | 97.44 | 2,133 | 92.45 |
| [0.3,0.4) | High Cov | 15,737 | 98.2 | 15,315 | 97.54 | 4,298 | 95.23 |
| [0.4,0.5] | High Cov | 10,463 | 97.69 | 16,261 | 97.82 | 7,434 | 97.19 |
| [0,0.01) | Array | 8,975 | 99.98 | 97 | 81.44 | 0 | NA |
| [0.01,0.02) | Array | 4,519 | 99.96 | 156 | 83.33 | 0 | NA |
| [0.02,0.05) | Array | 11,740 | 99.8 | 833 | 93.28 | 12 | 83.33 |
| [0.05,0.1) | Array | 18,317 | 99.6 | 3,003 | 94.21 | 114 | 84.21 |
| [0.1,0.2) | Array | 29,144 | 99.22 | 10,142 | 95.5 | 866 | 91.11 |
| [0.2,0.3) | Array | 20,139 | 98.72 | 12,835 | 95.96 | 2,201 | 94 |
| [0.3,0.4) | Array | 14,833 | 97.9 | 15,879 | 96.52 | 4,171 | 95.35 |
| [0.4,0.5] | Array | 10,172 | 97.53 | 16,290 | 96.62 | 6,766 | 97.04 |

Supplementary Table 7: Performance of STITCH on CONVERGE study original imputation Results are over the first 10 Mbp of chromosome 20. Beagle methodology was the same as done in the original CONVERGE paper and as explained in the text. STITCH results are for $K = 40$, 38 pseudo-haploid iterations, 2 diploid iterations. All sites with removal of SNPs failing QC also removed SNPs with Hardy-Weinberg p-value less than 10^{-6} . Av r2 is the average r^2 for SNPs on the Illumina HumanOmniZhongHua-8 array for high frequency (MAF 5% to 95%) variants.

| Method | SNP set | % SNPs | Av r2 |
|--------|----------|--------|-------|
| Beagle | All | 100 | 0.933 |
| STITCH | All | 100 | 0.92 |
| Beagle | info>0.4 | 90 | 0.939 |
| STITCH | info>0.4 | 90 | 0.939 |
| Beagle | info>0.9 | 78 | 0.968 |
| STITCH | info>0.9 | 75 | 0.972 |

Supplementary Table 8: Effect of filtering on imputation performance Results are given for chromosome 19. QC is defined per-run and reflects $\text{info} > 0.4$ and HWE p-value $> 1 \times 10^{-6}$. r^2 values are against the 4 10X mice.

| Set | Description | SNPs | Number of SNPs | Ti/Tv | VQSR r2 | No VQSR r2 |
|-----|-----------------------------|---------|----------------|-------|---------|------------|
| 1 | VQSR | All | 152,486 | 2.07 | 0.937 | |
| 2 | VQSR | Post-QC | 122,878 | 2.21 | 0.968 | |
| 3 | No VQSR, Round 1 | All | 355,123 | 1.48 | | 0.745 |
| 4 | No VQSR, Round 1 | Post-QC | 136,164 | 2.08 | | 0.945 |
| 5 | No VQSR, Round 2 | All | 136,164 | 2.08 | | 0.938 |
| 6 | No VQSR, Round 2 | Post-QC | 128,054 | 2.14 | | 0.952 |
| 7 | Intersect Set 2 and Set 6 | | 115,567 | 2.22 | 0.967 | 0.969 |
| 8 | Present Set 2, absent Set 6 | | 7,311 | 2.13 | 0.915 | |
| 9 | Present Set 6, absent Set 2 | | 12,487 | 1.55 | | 0.930 |

References

- [1] Dempster AP, Laird NM, Rubin DB. Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society Series B (Methodological)*. 1977;39(1):1–38. Available from: <http://www.jstor.org/stable/2984875>.
- [2] Scheet P, Stephens M. A Fast and Flexible Statistical Model for Large-Scale Population Genotype Data: Applications to Inferring Missing Genotypes and Haplotypic Phase. *The American Journal of Human Genetics*. 2006 Apr;78(4):629–644. Available from: <http://www.sciencedirect.com/science/article/pii/S000292970763701X>.