

Legacy Data Confounds Modern Genomics Studies

Luke Anderson-Trocmé^{1,2}, Mathieu Bourgey^{1,2}, Fumihiko Matsuda³, Simon Gravel^{1,2}

*For correspondence:
simon.gravel@mcgill.ca (SG)

¹Department of Human Genetics, McGill University, Montreal, QC H3A 0G1, Canada;
²McGill University and Genome Quebec Innovation Centre, Montreal, QC H3A 0G1, Canada; ³Center for Genomic Medicine, Graduate School of Medicine, Kyoto University, Kyoto 606-8501, Japan

Abstract In investigating a population genetics signal, we noticed some discrepancies between the Japanese samples from the 1000 Genomes Project (1kGP) and a more recent high quality dataset. We found that the variants causing the signal were only present in individuals with low average sequencing quality. We identified 625 $p < 10^{-8}$ and 1048 $p < 10^{-6}$ single nucleotide polymorphisms (SNPs) that were associated to quality in the 1kGP dataset Japanese individuals. These variants are present in nearly all of the populations from the 1kGP. We then turned our attention to the rest of the 1000 Genomes populations and saw that there were similar batch effects in many of these populations. Some of these variants are being imputed onto genotype data and reach genome wide significance in recent publications. [*update values for all populations*](#)

Introduction

Reference Cohorts

The last 5 years have seen a drastic increase in the amount and quality of human genome sequence data. Hapmap, the 1000 Genomes Project, and the Simons Diversity project, for example, have made thousands of genomes publicly available for population and medical genetic analyses. Many more genomes are available indirectly through servers providing imputation services or variant frequency estimation.

Many of these large datasets have been assembled over many years. The first sequenced genomes in the 1000 Genomes project (1kGP), for example, were sequenced [*10*](#) years ago, at a time when sequencing technologies were still being . Because of the extraordinary value of freely available data, this early data from the 1kGP is still widely used as a reference panel for imputation, allele frequency estimations and a wide range of applications.

Yet this data is far from perfect [SG: cite papers that have pointed out issues](#). This raises the question of whether and how such legacy data should be included in contemporary analyses. Here we point out how large and previously unreported batch effects in the early phases of the 1000 Genomes Project still lead to incorrect genetic conclusions through population genetic analyses and indirect use through prominent imputation servers.

Motivations

Mutations can occur in the genome as a result of environmental inflictions as well as variation in genetic factors like DNA proofreading and repair. Many mutagens leave signatures that are specific to a genomic context. This means that the neighbouring nucleotides of a mutation can be used to gain insight into the cause of a mutation. Thymine dimers caused by UV radiation were detected

early on because they were highly specific and easy to study ? However, signatures caused by natural variations in the sequence of proteins involved in proofreading and repair are more difficult to study because they tend to occur less frequently and are more subtle . A large sample size can be used to increase statistical power to resolve some of these fine scale mutational signatures .

In 2015, Harris et al reported a **15%** overrepresentation of the TCC to TTC mutation in European populations. Harris and Pritchard published a follow up study comparing how various ethnic groups differ from one another with respect to their mutational spectrum. In this paper, they proposed that not only is there a distinct mutational spectrum difference in Europeans, but that many populations have their own mutational signatures. One result was of particular interest the heterogeneously distributed signal found in a subset of Japanese individuals. A mutational signature that is present unevenly across a population is unexpected because it suggested strong population structure; unlikely for a dense population living on an archipelago.

While trying to follow up on this observation, we were unable to reproduce this particular observation in a different dataset. Tracing back the source of the discrepancy between the two datasets, we identified a strong and previously unreported batch effect in the 1000 Genomes Project data, and found that this led to spurious results in a number of recently published studies.

Here we discuss the source of the batch effect and the ways in which it contaminated a number of recent studies. Fortunately, the main conclusions of most of these studies remain supported by other data. Yet this begs the question: When should we retire legacy data?

Many types of DNA sequencing machines have been available on the market over the past decade. Different types of sequencing technologies can have different technical limitations in quality control. This makes combining data an issue especially when multiple sequencing technologies are involved in the data production. Many of the legacy datasets produced using dated sequencing technologies have been known to contain a higher rate of false positives than their more recent counterparts. The errors in one dataset do not disappear when they are combined with another higher quality one. The technical biases caused by legacy sequencing technology is becoming increasingly relevant as newer technologies produce data with lower error rates, and more stringent quality controls.

Results

A peculiar mutational signature in Japan

If two groups of individuals are sampled from the same population, we expect there to be little deviation in allele frequencies between the two groups. Harris et al reported an excess of **AC→*CC* mutations in some Japanese individuals. While trying to follow up on this observations in a larger cohort of 886 Japanese individuals, **some info about cohort**, we did not find this particular signature. However, when comparing the allele frequencies between the Japanese individuals from the 1kGP and this larger dataset, we observed an unusually large number of private single nucleotide polymorphisms (SNPs), only found in one of the two groups. These mismatches were maintained after filtering for low-quality regions of the human genome and standard metrics such as Hardy-Weinberg equilibrium

Once mismatch sites were removed from the 1kGP data, the **AC→*CC* signal disappears [1](#), suggesting a technical artifact rather than a population structure effect. Regressions against different quality metrics provided by the 1kGP revealed that mean mapping quality per mapped base pair was an excellent correlate with prevalence of the **AC→*CC* mutational signature in 1kGP, with low-quality individuals showing elevated rates of the signature. ([SG: show supplementary Figure where you did the correlations](#))

Thus sequences with low mappability harbour mutations that reproduce poorly across studies and exhibit a particular mutational signature.

To identify SNPs that are associated with low quality, we performed a genome-wide association (GWA) study [1](#). GWA studies are commonly used to find genomic regions associated to biologically

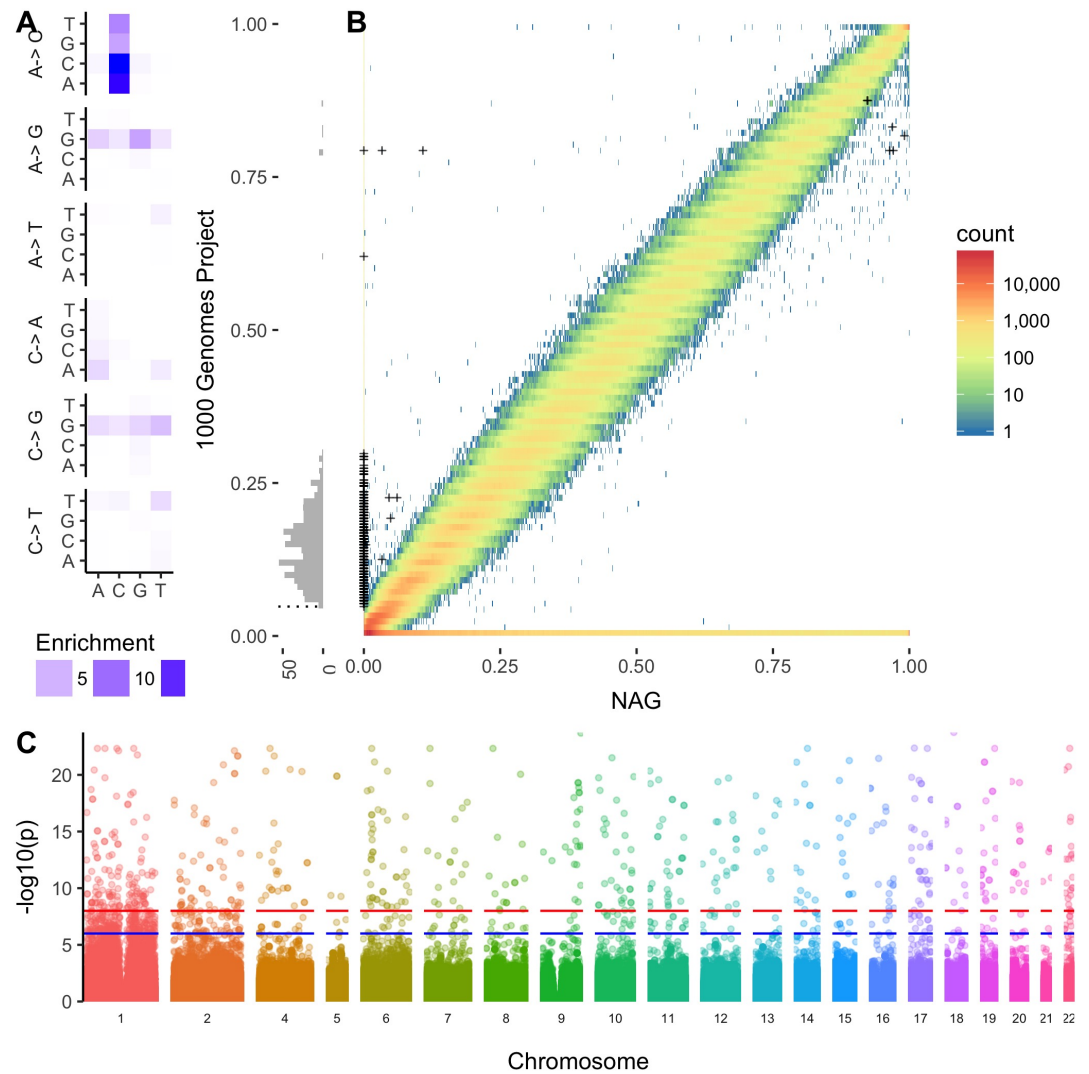


Figure 1. A Joint frequency spectrum plot of the Japanese from the 1000 Genomes Project and a more recent dataset. Crosses (+) are variants that reached genome wide significance in a GWA of sequencing quality. The histogram on the left of the plot is the distribution of significant variants. The lowest frequency variant able to reach genome wide significance is 5%. **B** Mutation spectrum of the 1048 variants that reached a genome wide significance with a p value less than in a GWA of sequencing quality. There is a **significant** enrichment in **AC→*CC* mutations. **C** Genome wide association of the average quality of mapped bases for the 104 Japanese individuals included in the 1000 Genomes Project. This GWA study identified 625 $p < 10^{-8}$ and 1048 $p < 10^{-6}$ SNPs that were associated to the average quality of SNPs mapped for an individual. The same analysis was performed independently for each of the populations in the 1000 Genomes Project.

91 relevant traits. In this case, we are using this analysis in a more unconventional way to identify
92 regions of the genome that are associated to biologically irrelevant traits like SNPs associated to
93 individuals with low quality.

94 Using the linear GWA study function offered by PLINK, we were able to identify 625 $p < 10^{-8}$ and
95 1048 $p < 10^{-6}$ SNPs that were associated to the average quality of SNPs mapped for an individual 1.
96 This GWA study included 104 individuals, *000* of them were sequenced in phase 1 while the rest
97 were sequenced in phase 3 of the 1000 Genomes Project. The variants that are associated to the
98 quality of mapped bases have an enrichment in *AC→*CC mutations.

99 **Alex: Having trouble following the first three sentences** Despite removing the variants identified
100 as being low quality, the signal persists. When we remove the variants significantly associated to
101 low quality, the enrichment in *AC→*CC mutations remains present, but *less* significant. When
102 we remove individuals with average quality of mapped SNPs below 30, the signal goes away. We
103 suspect that this has to do with the lack of power we have for more rare alleles. The lowest
104 frequency variants that are significantly associated to quality are at 5%. When we remove all the
105 variants above 1% in 1kGP but missing from the NAG data, the signal disappears. This suggests that
106 while the GWA approach can identify some of the low hanging bad apples, there are likely more of
107 these false positives nested inside these legacy cohorts.

108 One way to assess the validity of these unusual variants is to see if they are present in a higher
109 quality dataset. Upon comparing the Japanese 1000 Genomes cohort to a higher quality and
110 larger cohort, we identified *000* more variants that are beyond the expected frequency spectrum
111 deviation for individuals from the same population.

112 Sequencing quality over time

113 We turned our attention to the other populations in the 1kGP, we found that the sequencing done
114 in phase 1 was more variable and overall tended to include lower quality sequencing data 2. The
115 sequencing quality of individuals increases over time 2. By 2011 the sequencing quality seems to
116 level off, this also coincides with the phasing out of older sequencing technologies.

117 Overlap of significant SNPs

118 Comparing the results of each independent GWA study, we were able to identify over *0000* vari-
119 ants that were independently associated to low quality in multiple populations. This confirmation
120 using more than one GWAS is a strongly suggests that these variants might not be genuine 3.

121 Imputation

122 30% of the SNPs we identified as being associated with low quality were found to be imputed using
123 the Michigan Imputation Server. These should be removed from reference database.

124 Found to be included in other GWA studies

125 Once we identified SNPs that were clearly associated with low quality, we searched the literature for
126 any GWA studies that might have called these erroneous variants as being significantly correlated
127 with some biological trait. Using the NHGRI-EBI Catalog of published genome-wide association
128 studies we queried the rsIDs of the SNPs we identified as being low quality and found 6 recent
129 publications that had found at least one of the variants to have reached genome wide significance
130 in their study.

131 Five of these studies used the 1kGP as the reference database for imputation and one used the
132 1kGP cell cultures and sequence data. They used strict quality thresholds, including population
133 genetic statistical tests such as the Hardy-Weinberg equilibrium test, deviations in expected allele
134 frequency and sequencing data quality thresholds. They also removed rare alleles and alleles with
135 high degrees of missingness. Despite using the state of the art quality controls, these erroneous
136 variants managed not only to be imputed onto real genotype data, but they also reached genome
137 wide significance for biological traits.

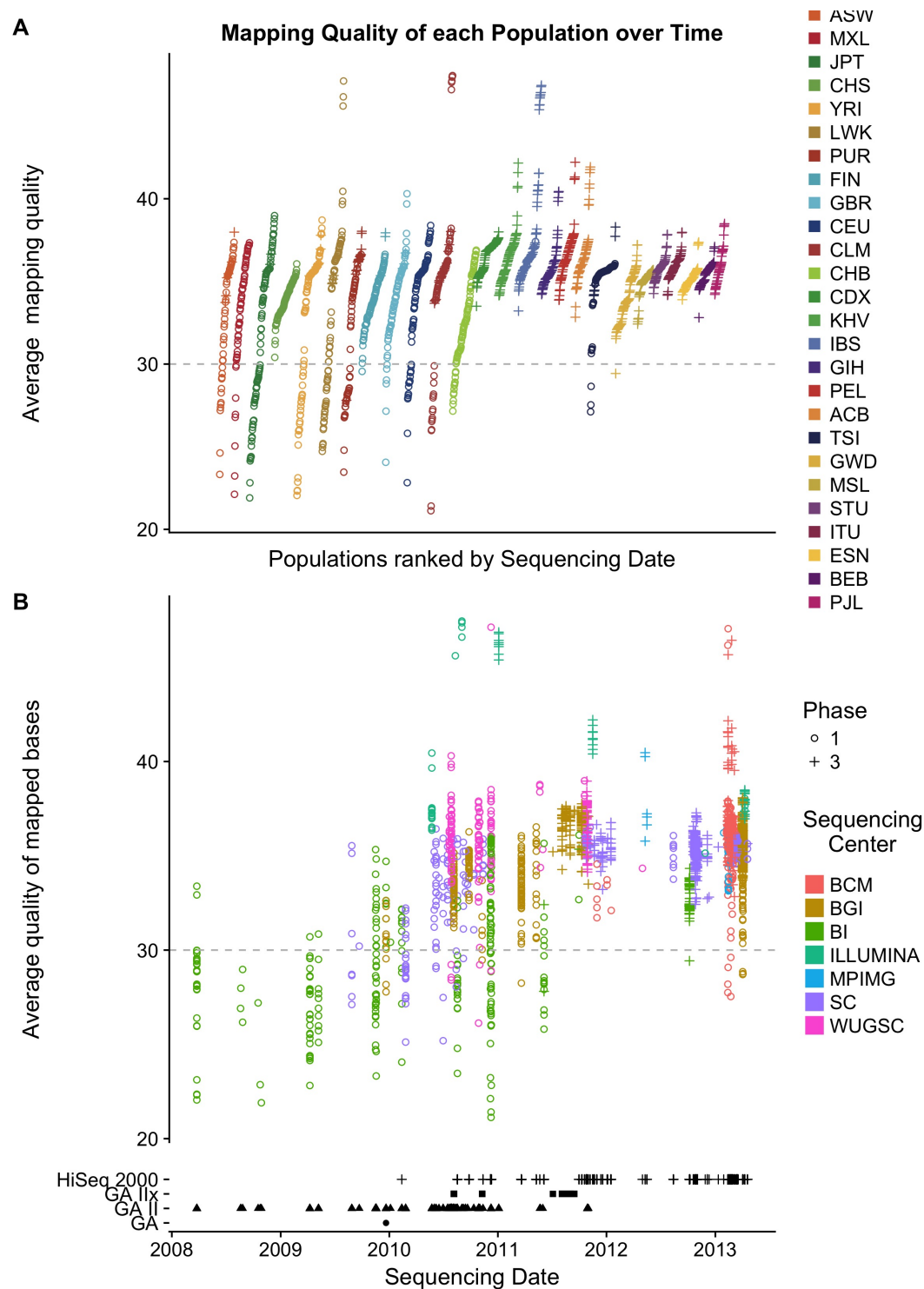


Figure 2. A The average mapping quality of each individual per population included in the 1000 Genomes sequencing project. The x-axis is ranked by populations with the least to the most variance, followed by average mapping quality per individual. **B** Same data as in **A** except the x-axis is sorted by sequencing date. The colors indicate the sequencing centers that produced the data for each individual.

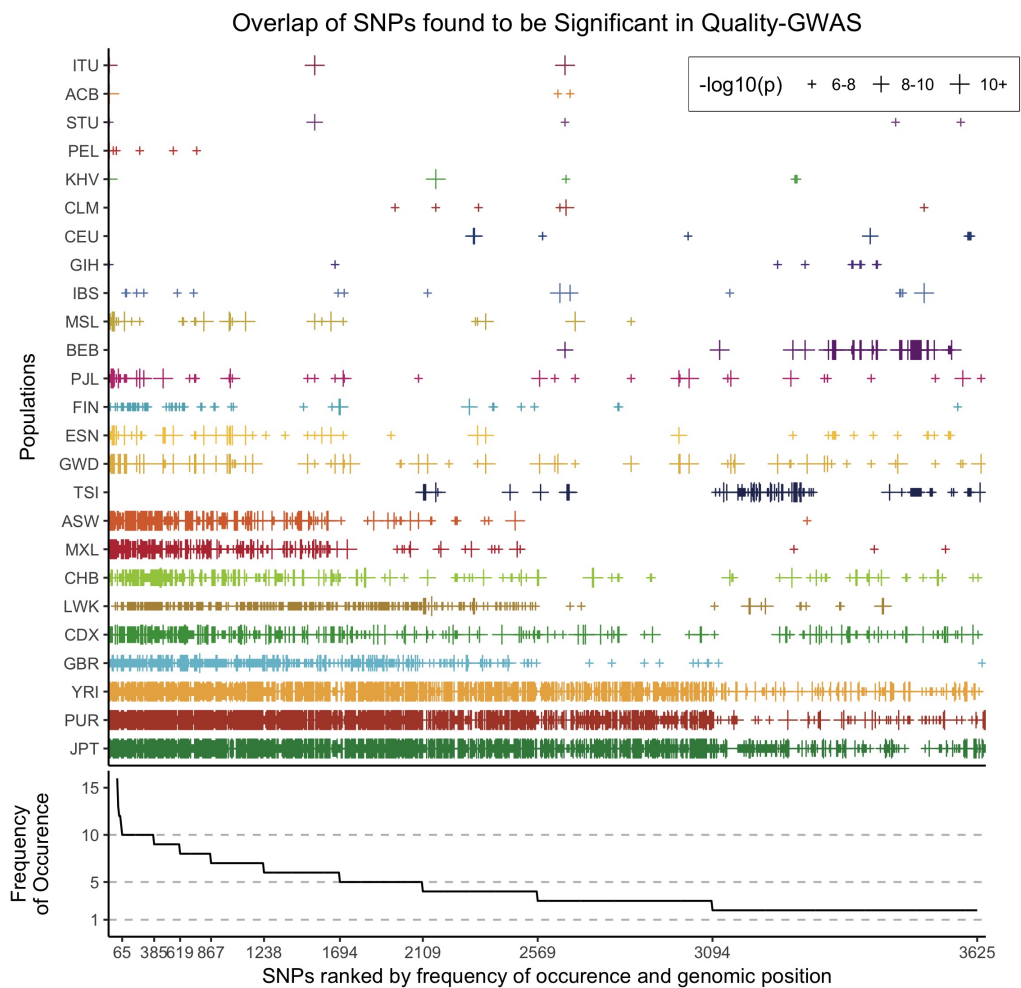


Figure 3. Overlap of SNPs identified independently to be associated with quality. The populations that have the most low-quality individuals also have the most low-quality variants. The same variants identified as being low quality independently in each population are found in other populations.

Discussion

Why do we care?

These SNPs matter because they reached genome wide significance for medically relevant traits. Including these variants in the GWAS catalogue could lead to false diagnosis at worst, or spurious correlations at least. Polygenic risk scores take into account all SNPs reaching significance, without much manual curation. It's likely that the variants we have identified are being included in these multi-locus risk scores.

Despite these variants reaching genome wide significance, the majority of researchers did not pursue these variants for further analysis. This is likely because these variants have no peak in LD that is characteristic of a biologically significant variant.

Recommendations

The most conservative approach would be to remove all individuals that don't meet the quality threshold as well as all the variants associated to low quality. In this case, we used a cut off of an average quality of mapped bases over 30. This threshold has been **previously used by studies**. It is the minimum requirements for **GATK variant calling** for them to have a minimum quality of 30.

Imputation

Imputation of genotype data is a probabilistic method that infers the bases of a given genome based on its similarity to a set of reference genomes. While on average, two humans differ in about 1/10,000 bases, this number is **lower** in closely related individuals, and **higher** in individuals from different continental origins. Modern chip sequencing will provide the genotype information for over 1 million bases of the genome. The unique combination of genotyped bases can be enough to identify haplotype blocks that are identical by descent in individuals from the reference database.

The accuracy of imputation depends on the size of the reference database; this varies significantly from one population to another. This means that individuals with ancestry that is less well represented in the reference database will have lower accuracy of imputation. To overcome this bias, reference databases are often combined to increase the sample size and in turn, the accuracy of imputation.

GWA studies from other papers

Since these variants are present in more than one population from the 1000 genomes project, they are more likely to be associated to biological traits as they would appear to be like any other variant that is shared among multiple populations. The only way to distinguish some of these more covert false positives is to use statistical tests associating the quality metrics of each position relative to each individual.

Luke: due to the temporal nature of the batch effects, and because entire populations were sequenced in one centre on one day, the false positives are more likely to cluster with population structure or case/control?

Conclusion

Our method identifies spurious mutations by correlating mutations with data quality metrics. We propose including our quality control methods to identify possible false positives in sequencing data. We have focused on the 1000 Genomes Project dataset as its quality metrics were freely available, however the issues of quality control are not limited to this one consortium. This study only used one dataset and one quality metric, but using this same approach can be used to identify more false positives in many more datasets.

As more large scale genotyping efforts are being imputed on the same legacy datasets, we must scrutinize the quality of the reference databases to avoid the propagation of false positives. These results bring forth many questions regarding the reliability of legacy datasets. Moreover, since there are so many broad applications of imputation, it frames the question for reference data turnover.

185 **Methods**

186 **Metadata**

187 The metadata used in this analysis was compiled from each of the index files from the 1000
188 Genomes file system. Average quality of mapped bases per sample was obtained from the BAS
189 files associated with each alignment file. Each BAS file has metadata regarding each sequencing
190 event for each sample. If a sample was sequenced more than once, we took the average of the
191 each quality score from each sequencing instance. The submission dates and sequencing centres
192 for each sample in the analysis was available in the sequence index files. This file also has multiple
193 entries per sample, however, we were unable to match the individual sequencing runs between the
194 bas files and the index file, which lead us to take the average of the quality scores and only kept
195 the earliest sequencing date per sample. The dates of the sequencing are only used to plot Figure.
196 **Average Quality of mapped bases: How was it calculated?**

197 **Data Availability**

198 Index of BAS files [available here](#).
199 Phase3 analysis sequence index file [available here](#)
200 **link to my compiled metadata file here**

201 **Quality Controls**

202 We reproduced the quality control pipelines used by Harris et. al as they applied the current
203 state of the art quality thresholds to remove questionable sequences especially for the high
204 standards for detecting population level differences. Several mask files were applied to remove
205 regions of the genome that might be lower quality, or might have very different mutation rates
206 or basepair complexity compared to the rest of the genome. The 1000 Genomes [strict mask](#)
207 was used to remove low quality regions of the genome , highly conserved regions were removed
208 using the [phastCons100way](#) mask file and highly repetitive regions were also removed using the
209 [NestedRepeats](#) mask file from RepeatMasker. Furthermore, only diallelic autosomal SNPs were
210 considered, with missingness below 0.01, MAF less than 0.1, and MAF greater than 0.9.

211 **Genome Wide Association**

212 Using PLINK v1.90b4.4 we ran a linear genome wide association study independently for each
213 population of the 1000 Genomes Project. We used the average quality of mapped bases per
214 individual as the phenotype for the analysis. We also controlled for population structure by
215 including the first 4 principle components of a PCA of each population using genotype data. **NOT
216 DONE YET**

217 **Mutation Spectrum**

218 We calculated the mutation spectrum for each list of significant SNPs for each population.
219 We also compared the mutation spectrum ratio between populations using a modified version
220 of the methods used in Harris et al. 2017.

221 **Imputation**

222 Using the Michigan Imputation Server, we imputed the genotype data from 1000 Genomes Project

223 **Code Availability**

224 **here is where I will put my git hub**

225 **Acknowledgments**

226 We would like to thank Kelly Harris for sharing her mutation spectrum pipelines.