

# A reference haplotype panel for genome-wide imputation of short tandem repeats

Shubham Saini<sup>1,2</sup>, Ileena Mitra<sup>2,3</sup>, Nima Mousavi<sup>1,2,4</sup>, Stephanie Feupe Fotsing<sup>1,2,5</sup>, Melissa Gymrek<sup>1,2,\*</sup>

<sup>1</sup> Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA USA

<sup>2</sup> Department of Medicine, University of California San Diego, La Jolla, CA USA

<sup>3</sup> Bioinformatics and Systems Biology Program, University of California, San Diego, La Jolla, CA, USA.

<sup>4</sup> Department of Electrical and Computer Engineering, University of California San Diego, La Jolla, CA USA

<sup>5</sup> Department of Biomedical Informatics, University of California, San Diego, La Jolla, CA, USA.

\* Correspondence should be addressed to [mgymrek@ucsd.edu](mailto:mgymrek@ucsd.edu)

## Abstract

Short tandem repeats (STRs) are involved in dozens of Mendelian disorders and have been implicated in a variety of complex traits. However, existing technologies focusing on single nucleotide polymorphisms (SNPs) have not allowed for systematic STR association studies. Here, we leverage next-generation sequencing data from 479 families to create a SNP+STR reference haplotype panel for genome-wide imputation of STRs into SNP data. Imputation achieved an average of 97% concordance between genotyped and imputed STR genotypes in an external dataset compared to 63% expected under a random model. Performance varied widely across STRs, with near perfect concordance at bi-allelic STRs vs. 70% at highly polymorphic forensics markers. We demonstrate that imputation increases power over individual SNPs to detect STR associations using simulated phenotypes and gene expression data. This resource will enable the first large-scale STR association studies using existing SNP datasets, and will likely yield new insights into complex traits.

## Introduction

Genome-wide association studies (GWAS) have become increasingly successful at identifying genetic loci significantly associated with complex traits in humans, largely due to the enormous growth in available sample sizes<sup>1-3</sup>. Hundreds of thousands of individuals have been genotyped using commodity genotyping arrays. These arrays take advantage of the correlation structure between nearby variants induced by linkage disequilibrium (LD), which allows genome-wide imputation based on genotypes of only a small subset of loci<sup>4</sup>. However, GWAS based on single nucleotide polymorphism (SNP) associations face important limitations. Even with sample sizes of up to 100,000 individuals, common SNPs still fail to explain the majority of heritability for many complex traits<sup>1,5</sup>.

One compelling hypothesis explaining the “missing heritability” dilemma is that complex variants, such as multi-allelic repeats not in strong LD with common SNPs are important drivers of complex traits but are largely invisible to current analyses. Indeed, dissection of the strongest schizophrenia association, located in the major histocompatibility complex, revealed a poorly tagged polymorphic copy number variant (CNV) to be the causal variant<sup>6</sup>. The signal could not be localized to a single SNP and could only be explained after deep characterization of the underlying CNV. This and subsequent discoveries<sup>7,8</sup> highlight the importance of considering alternative variant classes.

Short tandem repeats (STRs), consisting of repeated motifs of 1-6bp in tandem, comprise more than 3% of the human genome<sup>9</sup>. Multiple lines of evidence support a role of STRs in complex traits<sup>10-12</sup>, particularly in neurological and psychiatric phenotypes. Due to their rapid mutation rates<sup>13</sup>, STRs exhibit high rates of heterozygosity<sup>14</sup> and likely contribute more *de novo* mutations per generation than all other known sources of genetic variation. Furthermore, STRs have been shown to play a significant role in regulating gene expression<sup>15,16</sup>, splicing<sup>17-19</sup>, and DNA methylation<sup>16</sup>. Intriguingly, more than 30 Mendelian disorders are caused by STR expansions via a range of mechanisms, including polyglutamine aggregation (Huntington’s Disease, ataxias<sup>20</sup>), hypermethylation (Fragile X Syndrome<sup>21</sup>), and RNA toxicity (ALS/FTD<sup>22</sup>). Furthermore, causal STRs driving existing GWAS signals have already been identified<sup>23</sup>.

Existing technologies have not allowed for systematic STR association studies. Next-generation sequencing (NGS) can be used to directly genotype short STRs, but NGS is still too expensive

to perform on sufficiently large cohorts for GWAS of most complex traits. An alternative approach is to impute STRs into existing SNP array datasets. Previous studies have demonstrated that STRs are often in significant LD with nearby SNPs<sup>24–26</sup> and found that STRs and SNPs provide complementary information about the evolutionary history of a genomic region. Despite widespread SNP-STR LD, statistical phasing of STRs and SNPs is challenging for several reasons: SNP-STR LD is notably weaker than SNP-SNP LD<sup>24</sup> due to the rapid mutation rates<sup>27,28</sup> and high prevalence of recurrent mutations in STRs. As a result, the relationship between STR repeat number and SNP haplotype can be complicated and nonlinear, with the same STR allele present on multiple SNP haplotypes and vice versa. Furthermore, LD patterns at STRs vary widely as a function of properties of the repeat, such as the repeat unit length, mutation rate, and mutation step size<sup>24</sup>. Finally, STRs are prone to genotyping errors induced during PCR amplification<sup>29,30</sup>, further ambiguating phase information.

Sequencing related samples allows haplotype resolution by directly tracing inheritance patterns. The recent generation of deep NGS using PCR-free protocols for hundreds of nuclear families in combination with accurate tools for genotyping STRs from NGS<sup>31</sup> now enables applying this technique genome-wide. Here, we profiled STRs in 479 families and used pedigree information to phase STR genotypes onto SNP haplotypes to create a genome-wide reference for imputation. We used this panel to impute STRs into an external dataset of similar ethnic background with average 97% concordance with observed STR genotypes. Imputation accuracy varied across STRs, ranging from nearly perfect concordance at bi-allelic STRs to around 70% for highly polymorphic forensics markers. We show that STR imputation achieves greater power than individual SNPs to detect underlying STR associations and demonstrate the utility of our panel by detecting novel STRs associated with gene expression. Finally, we imputed genotypes at STRs previously implicated in human disorders and found that we could accurately identify specific SNP haplotypes associated with long normal alleles most at risk for expansion.

To facilitate use by the community, we have released a phased SNP+STR haplotype panel for samples genotyped as part of the 1000 Genomes Project (see **Data availability**). This resource will enable the first large-scale studies of STR associations in hundreds of thousands of available SNP datasets, and will likely yield significant new insights into complex traits.

## Results

## ***A catalog of STR variation in 479 families***

We first generated the deepest catalog of STR variation to date in a large cohort of families included in the Simons Simplex Collection (SSC) (see **URLs**). We focused on 1,916 individuals from 479 family quads (parents and two children) that were sequenced to an average depth of 30x using Illumina's PCR-free protocol. Based on comparison to 1000 Genomes Project samples, we estimated the cohort to consist primarily of Europeans (83%), with 2.0%, 9.0%, and 3.6% of East Asian, South Asian, and African ancestry respectively (**Supplementary Figure 1**). We used HipSTR<sup>31</sup> to profile autosomal STRs in each sample. HipSTR takes aligned reads and a reference set of STRs as input and outputs maximum likelihood diploid genotypes for each STR in the genome. While HipSTR infers the entire sequence of each STR allele, we focus here on differences in repeat copy number rather than sequence variation within the repeat itself. To maximize the quality of genotype calls, individuals were genotyped jointly with HipSTR's multi-sample calling mode using phased SNP genotypes and aligned reads as input (**Online Methods**). Multi-sample calling allows HipSTR to leverage information on haplotypes discovered across all samples in the dataset to estimate per-locus error parameters and output genotype likelihoods for each possible diploid genotype. Notably, our HipSTR catalog excluded most known STRs implicated in expansion disorders such as Huntington's Disease and hereditary ataxias, since even the normal allele range for these STRs is above or near the length of Illumina reads<sup>32-35</sup>. To supplement our panel, we additionally used Tredparse<sup>36</sup> to genotype a targeted set of known pathogenic STRs in our cohort (**Supplementary Table 1**). Tredparse incorporates multiple features of paired-end reads to estimate the size of repeats longer than the read length.

An average of 1.14 million STRs passed HipSTR's default filtering settings in each sample (**Figure 1A**). We obtained at least one call for 97% of all STRs in the HipSTR reference of 1.6 million STRs and for 15 of 25 STRs in the Tredparse reference with an average overall call rate of 90% (**Figure 1B**). We applied additional stringent genotype quality filters to ensure accurate calls for downstream phasing and imputation analysis. STRs overlapping segmental duplications, with call rates less than 80%, or with genotype frequencies unexpected under Hardy-Weinberg Equilibrium were removed (**Online Methods**). We further removed STRs with low heterozygosity ( $<0.095$ ) to restrict analysis to polymorphic STRs. We found that these filters increased the quality of our calls, as evidenced by the average Mendelian inheritance rate of 99.8% and 97.9% at STRs that passed and failed quality filters, respectively (**Figure 1C**). After

filtering, 453,671 and 9 STRs from the HipSTR and Tredparse panels, respectively, remained in our catalog.

We further assessed the quality of our STR genotypes by comparing patterns of variation from SSC to previous catalogs of STR variation obtained using a distinct set of samples and STR genotyping methods. For HipSTR calls, we found that per-locus heterozygosities (**Online Methods**) were highly concordant with a catalog generated from the 1000 Genomes Project<sup>37</sup> data using lobSTR<sup>38</sup>. ( $r=0.96$ ;  $p<10^{-200}$ ;  $n=386,100$ ) (**Figure 1D**). For Tredparse calls, allele frequency spectra observed in SSC matched closely to previously reported normal allele frequencies at each STR (**Figure 1E**). For STRs genotyped both by HipSTR and Tredparse, estimated repeat lengths were highly concordant (average concordance 99.4%, **Supplementary Table 1**). Overall, these results show that our catalog consists of robust STR genotypes suitable for downstream phasing and imputation analysis.

### ***A genome-wide SNP+STR haplotype reference panel***

We examined the extent of linkage disequilibrium between STRs and nearby SNPs using two metrics. The first, termed “length  $r^2$ ”, is defined as the squared Pearson correlation between STR allele length and the SNP genotype. The second, termed “allelic  $r^2$ ”, treats each STR allele as a separate bi-allelic locus and is computed similar to traditional SNP-SNP LD (**Online Methods**). Similar to previous studies<sup>24</sup>, SNP-STR LD was dramatically weaker than SNP-SNP LD by both metrics (**Supplementary Figure 2A**) with length  $r^2$  generally stronger than allelic  $r^2$ . We additionally determined the best tag SNP (**Online Methods**) for each STR, which was on average 5.5kb away (**Supplementary Figure 2B**). Nearly all STRs were in significant LD (Length  $r^2$   $p<0.05$ ) with the best tag SNP, suggesting that phasing would result in informative haplotypes.

We developed a pipeline to phase STRs onto SNP haplotypes leveraging the quad family structure (**Figure 2A**). Based on our LD analysis, we used a window size of  $\pm 50$ kb to phase each STR separately using Beagle<sup>39</sup>, which was recently demonstrated to perform well in phasing multi-allelic STRs<sup>40</sup> and can incorporate pedigree information. Resulting phased haplotypes from the parent samples were merged into a single genome-wide reference panel for downstream imputation.

We evaluated the utility of our phased panel for imputation using a “leave-one-out” analysis in the SSC samples. For each sample, we constructed a modified reference panel with that sample’s haplotypes removed and then performed genome-wide imputation. We measured concordance, length  $r^2$ , and allelic  $r^2$  between imputed vs. observed genotypes at each STR, where “observed” refers to genotypes obtained by HipSTR or Tredparse. For each of these metrics, we additionally computed the value expected under a model where genotypes are imputed randomly (**Online Methods**) for comparison. Imputed genotypes showed an average of 96.7% concordance with observed genotypes, compared to 61.0% expected under a random model (**Table 1**). As expected, concordance was strongest at the least polymorphic STRs (**Figure 2B, Supplementary Figures 3A, 4**) and allelic  $r^2$  was highest for the most common alleles (**Supplementary Figure 3B**). Length  $r^2$  was not strongly associated with heterozygosity, although the least and most heterozygous STRs tended to have lower length  $r^2$  (**Supplementary Figure 3C**). Imputation metrics were weakly negatively correlated with distance to the best tag SNP ( $r=-0.06$ ;  $p=0.06$ ,  $r=-0.04$ ;  $p=0.27$ ; and  $r=-0.06$ ,  $p=7.5 \times 10^{-5}$  for concordance, length  $r^2$ , and allelic  $r^2$ , respectively). To further evaluate imputation performance at highly polymorphic STRs, we examined the CODIS STRs used in forensics analysis (**Supplementary Table 2**). Per-locus concordances were highly correlated with imputation results recently reported by Edge, *et al*<sup>40</sup> (Pearson  $r^2=0.93$ ;  $p=6.3 \times 10^{-6}$ ;  $n=10$ ), but were on average 8.8% higher, likely as a result of our larger and more homogenous cohort. Per-locus imputation statistics for all STRs are reported in **Supplementary Tables 3 and 4**).

We next evaluated our ability to impute STR genotypes into an external dataset. For this, we focused on samples from the 1000 Genomes Project<sup>37</sup> with high quality SNP genotypes obtained from low coverage whole genome sequencing (WGS) ( $n=2,504$ ) or genotyping arrays ( $n=2,486$  for Affy 6.0, and  $n=2,318$  for Omni 2.5). We validated imputed genotypes for subsets of 1000 Genomes samples using three orthogonal technologies: Illumina WGS+HipSTR, capillary electrophoresis, and 10X Genomics+HipSTR. In each case we evaluated performance using the orthogonal data as the “truth” set.

First, we used HipSTR to genotype STRs in separate high-coverage (30x) WGS datasets available for 150 of the samples (see **URLs**) from European ( $n=50$ ), African ( $n=50$ ), and East Asian ( $n=50$ ) backgrounds. Per-locus concordance, length  $r^2$ , and allelic  $r^2$  were highly concordant between the SSC panel and 1000 Genomes samples of European origin (Pearson



$r=0.94$ ,  $0.63$ , and  $0.85$ , respectively) (**Figure 2C**; **Supplementary Figure 5**; **Table 1**). Overall imputation performance did not vary when using phased genotypes obtained from WGS vs. Omni2.5 for imputation (**Supplementary Table 5**). Concordance was noticeably weaker in African and East Asian samples, likely due to different population background compared to the SSC samples and lower LD in African populations<sup>41</sup>.

Next, we compared imputed genotypes to capillary electrophoresis data<sup>42</sup> (see **URLs**) available for a subset of samples in our panel at highly polymorphic STRs. After filtering non-European samples and STRs that could not be reliably mapped to HipSTR notation (**Online Methods**), 41 samples and 206 STRs remained for comparison. We obtained an average overall concordance of 76.9% with capillary genotypes compared with 76.4% expected based on HipSTR analysis. Per-locus concordances based on HipSTR vs. capillary genotypes were strongly correlated ( $r=0.83$ ;  $p=1.05 \times 10^{-53}$ ;  $n=206$ ) (**Figure 2D**).

Finally, we compared imputed genotypes from the highly characterized NA12878 genome to phased data available from 10X Genomics (see **URLs**), a synthetic long read technology. We constructed a phased validation panel by calling HipSTR separately on reads from each phase and combining with phased SNP genotypes (**Online Methods, Supplementary Figure 6**). We could obtain phased 10X calls for 116,764 of the STRs in our panel. We used the nearest heterozygous SNP to each STR to match phase order between our panel and the 10X data, which allowed us to directly compare imputed alleles and evaluate phase accuracy. Overall, imputed STR alleles showed 96% concordance with those obtained from 10X and per-locus genotype concordance was consistent with concordance metrics measured in SSC (**Figure 2E**). Taken together, validation of imputed STR genotypes against three separate “truth” sets demonstrates the accuracy of our original SNP+STR haplotype panel and shows that our quality metrics are reliable indicators of per-locus imputation performance across datasets.

### ***Imputation increases power to detect STR associations***

We sought to determine whether our SNP+STR haplotype panel could increase power to detect underlying STR associations over standard GWAS. First, we simulated phenotypes based on a single causal STR and examined the power of the imputed STR genotypes vs. nearby SNPs to detect associations. We focused primarily on a linear additive model relating STR dosage, defined as the average allele length, to quantitative phenotypes (**Figure 3A**), since the majority



of known functional STRs follow similar models (e.g.<sup>17,43–45</sup>). Association testing simulations were performed 100 times for each STR on chromosome 21 in our dataset (**Online Methods**). As expected, the strength of association for each variant as measured by the negative  $\log_{10}$  p-value was linearly related with its length  $r^2$  with the causal variant (**Figure 3B**). On average, imputed STR genotypes explained 17.7% more variation in STR allele length compared to the best tag SNP (mean  $r^2=0.92$  and  $0.74$  for imputed STRs vs. SNPs, respectively). The advantage from STR imputation grew as a function of the number of common STR alleles (**Supplementary Figure 7**). Imputed genotypes showed a corresponding increase in power to detect associations at a given p-value threshold (**Figure 3C**). Similar trends were observed for case-control traits (**Supplementary Figure 8**). We additionally tested the ability of imputed STR genotypes to identify associations due to non-linear models relating STR genotype to phenotype (**Supplementary Figure 9**). While both STR and SNP-based tests had limited power to detect non-linear associations, per-allele STR association tests had higher power than the best tag SNP in 60% of simulations. Importantly, testing for complex models relating repeat length to phenotype will only be possible when allele lengths are available, thus demonstrating an additional need for STR imputation over SNP-based tests to detect these associations.

We next determined whether STR imputation could identify STR associations using real phenotypes. We focused on gene expression, given the large number of reported associations between STR length and expression of nearby genes in *cis*<sup>15,16</sup> (termed eSTRs). To this end, we analyzed eSTRs from samples in the Genotype-Tissue Expression<sup>46</sup> (GTEx) dataset for which RNA-sequencing, WGS, and SNP array data were available. As a test case, we imputed STR genotypes using SNP data for chromosome 21 and tested for association with genes expressed in whole blood. For comparison, we additionally performed each association using genotypes obtained from WGS using HipSTR (**Online Methods**). A total of 2,452 STR x gene tests were performed in each case. Association p-values were similarly distributed across both analyses and showed a strong departure from the uniform distribution expected under a null hypothesis of no eSTR associations (**Figure 3D**). For all nominally significant associations ( $p<0.05$ ), effect sizes were strongly correlated when using imputed vs. HipSTR genotypes ( $r=0.99$ ;  $p=1.01\times 10^{-79}$ ,  $n=97$ ). Furthermore, effect sizes obtained from imputed data were concordant with previously reported effect sizes in a separate cohort using a different cell type (lymphoblastoid cell lines)<sup>15</sup> ( $r=0.79$ ;  $p=0.0042$ ,  $n=11$ ) (**Figure 3E**).

We identified genes for which the STR is most likely the causal variant and tested whether STR imputation had greater power to identify causal eSTRs compared to SNP-based analyses. We used ANOVA model comparison to determine genes for which the STR explained additional variation over the top SNP (**Online Methods**). We additionally applied CAVIAR<sup>47</sup> to fine-map associations using the most strongly associated STR and the top 100 associated SNPs for each gene (**Online Methods**). We identified 3 genes with ANOVA  $p < 0.05$  for which the STR was the top variant returned by CAVIAR. One example, a CG-rich STR in the promoter of *CSTB*, was previously demonstrated to act as an eSTR<sup>48</sup> and expansions of this repeat are implicated in myoclonus epilepsy<sup>49</sup>. In each case, imputed STR genotypes were more strongly associated with gene expression compared to the best tag SNP (**Figure 3F-G, Supplementary Table 6**).

### ***Phasing and imputing normal alleles at known pathogenic STRs***

Finally, to determine whether alleles at known pathogenic STRs could be accurately imputed, we examined results of our imputation pipeline at 12 STRs previously implicated in expansion disorders that were included in our panel (**Table 2**). Our analysis focused on alleles in the normal repeat range for each STR, since pathogenic repeat expansions at these STRs are unlikely to be present in the SSC cohort. Notably, accurate imputation of non-pathogenic allele ranges is still informative as (1) long normal or intermediate size alleles may result in mild symptoms in some expansion disorders<sup>50,51,52</sup> (2) longer alleles are more at risk for expansion<sup>53</sup> and (3) allele lengths below the pathogenic range could potentially be associated with more complex phenotypes<sup>51</sup>.

Similar to the CODIS markers, these STRs are highly polymorphic with 10 or more alleles per locus. In all cases, imputed genotypes were more strongly correlated with HipSTR or Tredparse genotypes compared to the best tag SNP. Where both HipSTR and Tredparse genotypes were available, concordance results were nearly identical across all STRs (**Supplementary Table 7**). Visualization of SNP-STR haplotypes at the CAG repeat implicated in dentatorubral-pallidoluysian atrophy (DRPLA)<sup>54</sup> reveals a typical complex relationship between STR allele length and local SNP haplotype (**Figure 4A**), with the same STR allele often present on multiple SNP haplotype backgrounds. Still, for most STRs there is a clear association of specific haplotypes with different allele length ranges allowing accurate imputation across a large range of allele sizes (**Figure 4B, Supplementary Figure 10**).

Resolution of SNP-STR haplotypes can be used to infer the mutation history of a specific STR locus<sup>25,26</sup>. Notably, for many STR expansion orders it has been shown that pathogenic expansion alleles originated from a founder haplotype<sup>55-58</sup> associated with a long allele. We compared SNP haplotypes at the DRPLA locus in our dataset to a previously reported founder haplotype<sup>58</sup>. In concordance with the hypothesis of a single founder haplotype, we found that SNP haplotypes with smaller Hamming distance to the known founder haplotype had longer CAG tracts ( $r=-0.79$ ;  $p<10^{-200}$ ). This finding demonstrates that while we were unable to directly impute pathogenic expansion alleles, STR imputation can accurately identify which individuals are at risk for carrying expansions or pre-pathogenic mutations and the inferred haplotypes can reveal the history by which such mutations arise.

## Discussion

Our study combines available whole genome sequencing datasets with existing bioinformatics tools to generate the first phased SNP+STR haplotype panel allowing genome-wide imputation of STRs into SNP data. Despite their exceptionally high rates of polymorphism, 92% of STRs in our panel could be imputed with at least 90% concordance, and 38% achieved greater than 99% concordance. Imputation performance varied widely across STRs, primarily due to differences in polymorphism levels across loci. Bi-allelic STRs could be imputed nearly perfectly (average concordance >99%, compared to 74% expected by chance), whereas STRs with the highest heterozygosity, including forensics markers and known pathogenic repeats, could be imputed to around 70% concordance (compared to around 35% expected by chance). We additionally show that imputation improves power to detect STR associations over standard SNP-based GWAS and could detect both known and novel associations between STR lengths and expression of nearby genes.

A widely recognized limitation of GWAS is the fact that common SNP associations still explain only a small fraction of heritability of most traits. Multiple explanations for this have been proposed, including minute effect sizes of individual variants and a potential role for high-impact rare variation<sup>59</sup>. However, studies in large cohorts reaching hundreds of thousands of samples<sup>1-3</sup>, as well as deep sequencing studies to detect rare variants<sup>60</sup>, have so far not confirmed these hypotheses. An increasingly supported idea is that complex variants not well tagged by SNPs may comprise an important component of the “missing heritability.”<sup>10,12,61</sup>

GWAS is essentially blind to contributions from highly polymorphic STRs and other repeats, despite their known importance to human disease and molecular phenotypes. Thus STR association studies will undoubtedly uncover additional heritability that is so far unaccounted for. Notably, while autism phenotypes are available for the SSC families, this cohort is too small to perform a GWAS and was specifically ascertained for families enriched for *de novo*, rather than inherited, pathogenic mutations. In future work our panel can be applied to impute STRs into larger cohorts for autism and other complex traits for which tens of thousands of SNP array datasets are available.

Our initial haplotype panel faces several important limitations. First, the majority of samples are of European origin, limiting imputation accuracy in other population groups. Second, imputation accuracy is mediocre for the most highly polymorphic STRs, some of which will ultimately have to be directly genotyped to adequately test for associations. Notably, our work relied on existing tools originally designed for SNP imputation. Further work on computational methods specifically for imputing repeats may be able to improve performance. Finally, thousands of long STRs are filtered from our panel due to the limitation imposed by short read lengths. While we have included target STRs implicated in STR expansion disorders, many long STRs are still inaccessible using current tools. New methods are now being developed for genome-wide genotyping of more complex STRs<sup>62</sup> and longer variable number tandem repeats (VNTRs)<sup>63</sup> from short reads and can be used to expand our panel in the future.

Overall, our STR imputation framework will enable an entire new class of variation to be interrogated by reanalyzing hundreds of thousands of existing datasets, with the potential to lead to novel genetic discoveries across a broad range of phenotypes.

## URLs

Simons Simplex Collection, <https://base.sfari.org/>

HipSTR, <https://github.com/tfwillems/HipSTR>

Beagle, [https://faculty.washington.edu/browning/beagle/b4\\_0.html](https://faculty.washington.edu/browning/beagle/b4_0.html)

1000 Genomes phased Affy6.0 and Omni2.5 SNP data,  
[ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/shapeit2\\_scaffolds/hd\\_chip\\_scaffolds/](ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/shapeit2_scaffolds/hd_chip_scaffolds/)

1000 Genomes Phase 3 <http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>

1000 Genomes STR data, <http://strcat.teamerlich.org/download>

Marshfield Capillary electrophoresis data, <https://payseur.genetics.wisc.edu/strpData.htm>

Marshfield marker annotations,  
[https://web.stanford.edu/group/rosenberglab/data/pemberonEtAl2009/Pemberon\\_AdditionalFile1\\_11242009.txt](https://web.stanford.edu/group/rosenberglab/data/pemberonEtAl2009/Pemberon_AdditionalFile1_11242009.txt)

NA12878 10X Genomics data,  
[https://support.10xgenomics.com/genome-exome/datasets/2.2.1/NA12878\\_WGS\\_v2](https://support.10xgenomics.com/genome-exome/datasets/2.2.1/NA12878_WGS_v2)

High-coverage Illumina sequencing for 1000 Genomes samples,  
<https://www.ebi.ac.uk/ena/data/view/PRJEB20654>

PyVCF, <https://github.com/jamescasbon/PyVCF>

Python statsmodels, <http://www.statsmodels.org/stable/index.html>

## Acknowledgements

Research reported in this publication was supported in part by the Office Of The Director, National Institutes of Health under Award Number DP5OD024577 and by a SFARI Explorer Award Number 515568. Access to SSC data was approved for this project under request id 2405.1.1. M.G. was supported in part by NIH/NIMH grant R01 MH113715. This work used the Extreme Science and Engineering Discovery Environment (XSEDE) comet resource at the San Diego Supercomputing Center through allocations ddp268 and csd568. XSEDE is supported by National Science Foundation grant number ACI-1548562. We thank Alon Goren for helpful comments on the manuscript. We additionally thank Vineet Bafna and Vikas Bansal for helpful discussions and providing access to compute resources. We are grateful to all of the families that participated in the Simons Simplex Collection as well as the principal investigators.

## Author Contributions

M.G. conceived the study, helped design and perform analyses, and drafted the initial manuscripts. S.S. generated the reference haplotype panel, performed downstream analyses, and participated in writing the manuscript. I.M. performed simulation analyses. S.F.F. performed analyses of expression data. N.M. performed analyses of pathogenic STRs. All authors have read and approved the final manuscript.

## Competing and Financial Interests

The authors have no competing financial interests to disclose.

## Figure Legends

**Figure 1: A deep catalog of STR variation in the SSC cohort. A. Number of STRs called per sample.** Dashed line represents the mean of 1.14 million STRs per sample. **B. Call rate per locus.** Dashed line represents the mean call rate of 90%. **C. Mendelian inheritance rate at filtered vs. unfiltered STRs.** The x-axis gives the posterior genotype score (Q) returned by HipSTR. The y-axis gives the average Mendelian inheritance rate for each bin across all calls on chromosome 21. STRs that were homozygous for the reference allele in all members of a family were removed. Colors represent different motif lengths. **D. Per-locus heterozygosity in SSC vs. 1000 Genomes.** Only STRs with heterozygosity  $>0.095$  in SSC are included. Color scale gives the  $\log_{10}$  number of STRs represented in each bin. **E. Allele frequencies at pathogenic STRs obtained by Tredparse vs. previously reported normal alleles.** Blue=Tredparse, Gold=Previously reported. Boxes span the interquartile range and horizontal lines give the medians. Whiskers extend to the minimum and maximum data points. The y-axis gives the number of repeat units. Tredparse alleles are based on calls in the SSC panel. Sources of previously reported allele frequencies are described in detail in **Online Methods**.

**Figure 2: Creating a reference SNP-STR haplotype panel. A. Schematic of phasing pipeline in the SSC cohort.** To create the phased panel, STR genotypes were placed onto phased SNP haplotypes using Beagle. Any missing STR genotypes were imputed. The resulting panel was then used for downstream imputation from orthogonal SNP genotypes. Blue and red denote phased and unphased variants, respectively. Positions in gray are homozygous. **B. Concordance of imputed STR genotypes vs. heterozygosity.** Blue denotes observed

per-locus values and green denotes values expected under a random model. Solid lines give median values for each bin and filled areas span the 25th to 75th percentile of values in each bin. X-axis values were binned by 0.1. Upper gray plot gives the distribution of heterozygosity values in our panel. Concordance values are based on the leave-one-out analysis in the SSC cohort. **C. Per-locus imputation concordance in SSC vs. 1000 Genomes cohorts.** Color scale gives the  $\log_{10}$  number of STRs represented in each bin. Concordance values are based on the subset of samples from the 1000 Genomes deep WGS cohort with European ancestry. **D. Per-locus imputation concordance using HipSTR vs. capillary electrophoresis genotypes.** Each dot represents one locus. The x-axis and y-axis give imputation concordance using capillary electrophoresis or HipSTR genotypes as a ground truth, respectively. Concordance was measured in separate sets of European samples for each technology. **E. Concordance of imputed vs. 10X STR genotypes in NA12878 stratified by concordance in SSC.** STRs were binned by concordance value based on the leave-one-out analysis. Concordance in NA12878 was measured across all STRs in each bin. Dots give mean values for each bin and lines denote  $\pm 1$  standard deviation. In all cases LOO refers to the leave-one-out analysis in the SSC cohort.

**Figure 3: STR imputation improves power to detect STR associations. A. Example simulated quantitative phenotype based on SSC genotypes.** A quantitative phenotype was simulated assuming a causal STR (red). Power to detect the association was compared between the causal STR, imputed STR genotypes, and all common SNPs ( $MAF > 0.05$ ) within a 50kb window of the STR (gray). **B. Strength of association ( $-\log_{10} p$ ) is linearly related with LD with the causal variant.** For SNPs, the x-axis gives the length  $r^2$  calculated using observed genotypes. For the imputed STR (blue), the x-axis gives the length  $r^2$  from leave-one-out analysis. **C. The gain in power using imputed genotypes is linearly related to the gain in  $r^2$  compared to the best tag SNP.** Gray contours give the bivariate kernel density estimate. Top and right gray area gives the distribution of points along the x- and y-axes, respectively. Power was calculated based on the number of simulations out of 100 with nominal  $p < 0.05$ . **D. Quantile-quantile plot for eSTR association tests.** Each dot represents a single STR  $\times$  gene test. The x-axis gives the expected  $\log_{10}$  p-value distribution under a null model of no eSTR associations. Red and blue dots give  $\log_{10}$  p-values for association tests using HipSTR genotypes and imputed STR genotypes, respectively. Black dashed line gives the diagonal. **E. Comparison of eSTR effect sizes using observed vs. imputed genotypes.** Each dot



represents a single STR  $\times$  gene test. The x-axis gives effect sizes obtained using imputed genotypes. Gray dots give the effect size in GTEx whole blood using HipSTR genotypes. Purple dots give effect sizes reported previously<sup>15</sup> in lymphoblastoid cell lines. **F., G. Example putative causal eSTRs identified using imputed STR genotypes.** Left, middle, and right plots give HipSTR STR dosage (red), imputed STR dosage (blue), and the best tag SNP genotype (gray) vs. normalized gene expression, respectively. STR dosage is defined as the average length difference from hg19. One dot represents one sample. P-values are obtained using linear regression of genotype vs. gene expression. STR and SNP sequence information is shown for the coding strand. Gene diagrams are not drawn to scale.

**Figure 4: SNP haplotypes distinguish allele lengths at known pathogenic STRs. A. Example SNP-STR haplotypes inferred in European samples at a polyglutamine repeat in *ATN1* implicated in DRPLA.** Each column represents a SNP from the founder haplotype reported by Veneziano, *et al.* Each row represents a single haplotype inferred in 1000 Genomes Project phase 3 European samples, with gray and black boxes denoting major and minor alleles, respectively. Haplotypes are grouped by the corresponding STR allele. The number of SNP haplotypes for each group of STR alleles is annotated to the left of each box. Alleles seen fewer than 10 times in 1000 Genomes samples were excluded from the visualization. **B. Comparison of imputed vs. observed STR genotypes in SSC samples at the *DRPLA* locus.** The x-axis gives the maximum likelihood genotype dosage returned by HipSTR and the y-axis gives the imputed dosage. Dosage is defined as the sum of the two allele lengths of each genotype relative to the hg19 reference genome. The bubble size represents the number of samples summarized by each data point. **C. Distribution of DRPLA repeat length vs. similarity to the pathogenic founder haplotype.** The founder haplotype refers to the SNP haplotype reported by Veneziano, *et al.* on which a pathogenic expansion in *ATN1* implicated in DRPLA likely originated. The x-axis gives the Hamming distance between observed haplotypes and the founder haplotype, computed as the number of positions with discordant alleles. White dots represent the median length.

## Tables

**Table 1: Imputation performance summary.** Results indicate mean across all STRs analyzed. Allelic  $r^2$  values include all common alleles (frequency at least 5%). Values in parentheses for each metric give expected values under a random imputation model based on allele frequencies in each population. “Multi-allelic” refers to STRs with three or more common alleles.

Panel (n=number of samples)	Concordance	Length $r^2$	Allelic $r^2$
SSC - LOO (n=1,916)	96.7% (61.0%)	0.906 (0.605)	0.861 (0.552)
SSC - LOO (multi-allelic)	94.3% (48.5%)	0.888 (0.334)	0.800 (0.333)
1000 Genomes - EUR (n=49)	97.0% (63.2%)	0.921 (0.678)	0.892 (0.543)
1000 Genomes - EUR (multi-allelic)	94.8% (50.0%)	0.900 (0.334)	0.828 (0.314)
1000 Genomes - AFR (n=46)	90.6% (57.9%)	0.746 (0.619)	0.706 (0.493)
1000 Genomes - AFR (multi-allelic)	85.6% (44.4%)	0.708 (0.336)	0.653 (0.310)
1000 Genomes - EAS (n=45)	93.8% (66.0%)	0.823 (0.690)	0.781 (0.557)
1000 Genomes - EAS (multi-allelic)	89.4% (53.7%)	0.780 (0.336)	0.663 (0.313)

**Table 2: Imputation performance at known pathogenic repeats.** <sup>a</sup>HD=Huntington's Disease; SCA=spinocerebellar ataxia; DRPLA=Dentatorubral-pallidolusian Atrophy; DM1=Myotonic Dystrophy Type 1; HDL=Huntingon's Disease-Like 2. LOO refers to the leave-one-out analysis in the SSC cohort. The best tag SNP for an STR is defined as the SNP within 50kb with the highest length  $r^2$ . STRs above the black bar were genotyped using Tredparse and below the bar were genotyped using HipSTR. Values in parentheses for concordance give the expectation under a random model.

Locus	Motif	Disorder <sup>a</sup>	Length $r^2$ LOO	Conc. LOO	Best tag SNP	$r^2_{\text{bestSNP}}$
4:3076604	CAG	HD	0.47	64.3% (27.5%)	rs762855	0.11
6:16327867	CAG	SCA1	0.72	85.3% (33.8%)	rs17860797	0.04
6:170870996	CAG	SCA17	0.51	80.0% (31.5%)	rs9472489	0.15
12:112036755	CAG	SCA2	0.49	96.2% (80.2%)	rs148019457	0.28
12:7045892	CAG	DRPLA	0.86	81.2% (24.9%)	rs34199021	0.69
13:70713516	CTG/CAG	SCA8	0.87	84.7% (24.0%)	rs9564660	0.39
14:92537355	CAG	SCA3	0.88	86.4% (27.5%)	rs7144492	0.27
19:46273463	CTG	DM1	0.87	86.9% (30.8%)	rs7254351	0.44
19:13318673	CAG	SCA6	0.81	92.0% (39.2%)	rs2070737	0.63
3:63898362	CAG	SCA7	0.75	92.0% (63.9%)	rs58676857	0.57
5:146258292	CAG	SCA12	0.88	93.8% (46.3%)	rs2082405	0.64
16:87637894	CAG	HDL	0.55	88.2% (46.5%)	rs2434850	0.34

## Online Methods

### **SSC Dataset**

The SSC Phase 1 dataset consists of 1,916 individuals from 479 quad families. Aligned BAM and gVCF files for whole genome sequencing data of individuals were obtained through SFARI base (see **URLs**) and processed on Amazon Web Services (AWS). SNP genotypes were called from gVCF files using the GATK version 3 joint calling pipeline<sup>64</sup>. A total of 27,185,239 variants that passed the default GATK filters and overlapped with sites reported in the 1000 Genomes Project<sup>37</sup> phase 3 data were retained for downstream analysis.

We performed principal components analysis (PCA) using SNPs from 2,504 samples from Phase 3 of the 1000 Genomes Project<sup>37</sup> and projected SSC samples onto the resulting PCs to infer sample ancestry (**Supplementary Figure 1**). We estimated that the SSC cohort consists of 1585 Europeans, 39 East Asian, 172 South Asian, 69 African samples, and 51 individuals that did not clearly belong to any single population group.

### **Genome-wide multi-sample STR genotyping**

STRs were jointly genotyped on the AWS EC2 platform in batches of 500 STRs. We streamed the corresponding region of each BAM file and of the phased SNP VCF files to a local EBS volume attached to each EC2 instance using samtools<sup>65</sup> version 1.4 and tabix<sup>66</sup> version 1.2, respectively. HipSTR<sup>31</sup> version v0.5 was called individually per locus with default parameters. Phased SNPs were provided as input to allow HipSTR to perform physical phasing when possible. Resulting VCF files from each batch were merged to create a genome-wide callset in VCF format.

HipSTR calls were filtered using the filter\_vcf.py script in the HipSTR package with suggested parameters (--min-call-qual 0.9 --max-call-flank-indel 0.15 --max-call-stutter 0.15). We used the following criteria to remove problematic STRs from the callset: (i) STRs overlapping segmental duplications (UCSC Table Browser<sup>67</sup> hg19.genomicSuperDups table) were removed from the callset using intersectBed<sup>68</sup> v2.25.0; (ii) Pentanucleotides and hexanucleotides containing homopolymer runs of at least 5 or 6 nucleotides, respectively, in the hg19 reference genome were removed as they were found to contain an excess of indels in the homopolymer regions; (iii) STRs with call rate <80%; (iv) STRs with heterozygosity <0.095, corresponding to a minor allele frequency of 5% for biallelic markers, were removed to restrict to polymorphic STRs; (v)

STRs with significantly more or fewer heterozygous genotypes compared to expectation under Hardy-Weinberg equilibrium ( $p < 0.01$ ) as described previously<sup>69,70</sup>. After filtering, 453,671 STRs remained in our panel.

### **Genotyping clinically relevant STRs**

A total of 25 clinically relevant STRs were called using Tredparse<sup>71</sup> v0.75 from the aligned BAM files obtained through SFARI base on Amazon EC2. Default profiles containing information about the genomic position, reference repeat length, and repeat motif supplied with the software were used. We filtered STRs with call rate less than 80% or for which only a single allele was identified (**Supplementary Table 1**). 9 STRs remained after filtering.

### **Computing STR heterozygosity**

For an STR with alleles  $\{1...n\}$ , let  $p_i$  be the frequency of the  $i$ th allele computed from observed genotypes. STR heterozygosity is defined as:  $H = 1 - \sum_{i=1}^n p_i^2$ . For this study all alleles with identical length are treated as the same allele. On average each length-based allele corresponded to 1.8 sequence-based alleles.

### **Comparison to 1000G catalog**

STRs for 1000 Genomes samples as described in Willems *et al.*<sup>14</sup> were downloaded from the strcat site (see **URLs**). Heterozygosity was computed using the PyVCF package (see **URLs**) for the 1000 Genomes calls and using a custom script for the SSC data to collapse alleles of identical length into a single allele. STRs passing all filters described above included in the comparison. Analysis was restricted to STRs with at least 500 calls in the 1000 Genomes dataset.

### **Comparison to normal allele frequency spectra at clinically relevant STRs**

Control distributions for **Figure 1E** were obtained from previous studies of normal alleles at known pathogenic STRs. Allele frequencies for SCA1, SCA2, SCA3, SCA6, SCA12, SCA8, SCA17, and DRPLA were obtained from Figure 1 of Majounie, *et al.*<sup>32</sup> and are based on 307 controls of Welsh origin. Frequencies for DM1 were obtained from Figure 1 of Ambrose, *et al.*<sup>33</sup> and are based on 254 controls of Chinese origin. Frequencies for HDL were obtained from Figure 1 of Figley, *et al.*<sup>34</sup> and are based on 352 controls of North American Caucasian origin.

Frequencies for SCA7 were obtained from Figure 1 of Gouw, *et al.*<sup>35</sup> and are based on 180 controls of European origin. Frequencies for HTT are based on data in the phv00173896.v1.p1 variable of dbGaP study phs000371.v1.p1 (“Genetic modifiers of Huntington’s Disease”) based on the shorter allele of 2,802 patients with Huntington’s Disease.

### ***Phasing SNPs in the SSC***

SNP genotypes were phased using SHAPEIT<sup>72</sup> version 2.r837 with 1000 Genomes Phase 3 genotypes as a reference panel and ignoring pedigree information. SHAPEIT’s duoHMM<sup>73</sup> version 0.1.7 method was used to refine phased haplotypes using pedigree structure and correcting for Mendelian errors.

### ***Phasing STRs***

Beagle<sup>39</sup> version 4.0 was used to phase each STR separately using phased SNP genotypes, pedigree information, and unphased STR genotypes as input. In order to leverage the HipSTR genotype likelihoods (GL field), Beagle requires all samples to have GL information. To accommodate this, phasing was performed in two steps. First, samples with missing data were removed and the remaining samples were phased using the “-gl” Beagle flag. Next, missing samples were added back to the VCF and all samples were jointly phased in a second Beagle round using default parameters. In this step Beagle additionally imputed any calls with missing genotypes. Genotype values (GT field) were used for the STRs genotyped using Tredparse as it does not report genotype likelihoods, and phasing and imputation of STRs was done in a single step. Phased STRs and SNPs for only the unrelated parent samples from each locus were then merged into a single genome-wide reference panel in VCF format.

### ***Imputation performance metrics***

Let  $X = \{x_1, x_2, \dots, x_n\}$  be the true STR genotypes for samples 1..n and  $Y = \{y_1, y_2, \dots, y_n\}$  be the imputed STR genotypes. Each genotype  $x_i$  is defined as  $x_i = (x_{i1}, x_{i2})$  where  $x_{i1}$  and  $x_{i2}$  give the (unordered) lengths of the two STR alleles for a diploid sample and similarly for  $Y$ . We then define the following metrics:

**Genotype concordance:** Concordance  $c_i$  was defined as: 1 if both genotypes match ( $x_{i1} = y_{i1}$  and  $x_{i2} = y_{i2}$  or  $x_{i2} = y_{i1}$  and  $x_{i1} = y_{i2}$ ); 0 if neither imputed allele matched a true allele; else 0.5 if

one but not both imputed alleles matched the true alleles. Genotype concordance for an STR is the average over all the samples  $C = \frac{1}{n} \sum_{i=1}^n c_i$ .

**Length  $r^2$ :** Define the STR genotype dosage as the sum of the lengths of the two alleles at a given site:  $d_i = x_{i1} + x_{i2}$  and  $X_d = \{d_1, d_2, \dots, d_n\}$ . Length  $r^2$  is computed as  $cov^2(X_d, Y_d) / (Var(X_d)Var(Y_d))$ .

**Allelic  $r^2$ :** For a given allele length  $a$ , define  $X_a = \{a_1, a_2, \dots, a_n\}$  where  $a_i = \sum_{j=1}^2 1_{(x_{ij}=a)}$ . Allelic  $r^2$  is computed as  $cov^2(X_a, Y_a) / (Var(X_a)Var(Y_a))$ .

**Best tag SNP:** The best tag SNP for an STR is defined as the SNP within 50kb with the highest length  $r^2$ .

For all concordance metrics, outlier genotypes containing alleles seen less than 3 times in the entire cohort were removed from the analysis.

For each STR, we additionally computed the expected value of each metric under a model where genotypes are imputed randomly based on the frequency of underlying alleles. Expected genotype concordance was calculated as  $\sum_{i,j} f_i f_j (\sum_{k,l} f_k f_l C(i,j,k,l))$ , where  $(i,j) \in \{1, \dots, n\}^2$  and  $(k,l) \in \{1, \dots, n\}^2$ ,  $n$  is the number of alleles,  $f_x$  gives the frequency of allele  $x$ , and  $C(i,j,k,l)$  gives the concordance between genotypes  $(i,j)$  and  $(k,l)$  as defined above. For example, for a bi-allelic marker with allele frequencies  $f_1$  and  $f_2$  expected genotype concordance is given by  $f_1^2(f_1^2 + (0.5)2(f_1)(f_2)) + 2f_1f_2((0.5)f_1^2 + 2f_1f_2 + (0.5)f_2^2) + f_2^2(f_2^2 + (0.5)2f_1f_2)$ . Random model values for length  $r^2$  and allelic  $r^2$  were computed by comparing randomly imputed genotypes to true genotypes at each locus.

### **Evaluating imputation performance in the 1000 Genomes data**

STRs were imputed into SNP data downloaded from the 1000 Genomes Project site from three sources (WGS, phased SNPs from Affy6.0 array; and phased SNPs from Omni2.5 array; see **URLs** and **Supplementary Table 5**) with Beagle version 4.1 using the SSC SNP-STR



haplotype panel. For comparison to WGS, STRs were jointly genotyped in 150 high-coverage WGS datasets for 150 of the 1000 Genomes Project samples (see **URLs**) using HipSTR version 0.6 followed by the filtering steps described above for the SSC cohort.

Capillary electrophoresis genotypes for 209 samples at 721 Marshfield STRs were downloaded from the Payseur Lab website (see **URLs**). PCR product sizes were converted to length differences in bp from the reference genome using product size annotations<sup>74</sup> available from the Rosenberg Lab website (see **URLs**). Prior to comparing genotypes, offsets were calculated to match HipSTR lengths to the length of Marshfield STRs as previously described<sup>14</sup>. STRs with imperfect repeat structures were removed. Capillary genotypes were rounded down to the nearest number of repeat units.

10X Genomics data for NA12878 was obtained from the NA12878 Gemline Genome v2 available on the 10X Genomics website (see **URLs**). We extracted reads belonging to phase 1 or 2 from the phased, barcoded BAM based on the HP tag into separate BAM files. HipSTR v0.6.1 was called separately on each BAM with non-default parameters `--def-stutter-model --min-reads 5 --use-unpaired` and with `--haploid-chrs` containing a list of all autosomal chromosomes to force a haploid genotyping model. Haploid STR calls were obtained for both phases at 118,353 STRs. We identified the nearest heterozygous SNP to each STR that was genotyped in both the 10X data and in our phased panel. STRs for which the nearest SNP had discordant genotypes in the two datasets were discarded leaving 116,764 STRs for analysis.

### ***Simulations for power analysis***

We analyzed parental genotypes for 5,838 STRs across chromosome 21 that passed filtering and quality control as described above. For each STR, we simulated quantitative phenotype datasets under the model:  $P = \beta G + E$ , where  $P$  is a vector of standard normalized phenotypes,  $\beta$  gives the effect size,  $E$  gives the error term drawn from a normal distribution  $N(0, 1 - \beta)$ , and  $G$  is a vector of the sum of genotype lengths for each individual scaled to have mean 0 and variance 1. For each simulated phenotype dataset, we tested the causal STR, the imputed STR genotypes, and the best tag SNP (strongest length  $r^2$ ) within 50kb of the STR for association. Association tests were performed using the Python statsmodels library OLS method (see **URLs**).

We performed additional simulations under a case control model shown in **Supplementary Figure 8**. Phenotypes (0=control, 1=case) were drawn for each sample according to the model  $\text{logit}(p_i) = \beta X_i$  where  $p_i$  is the probability that sample  $i$  is a case and  $X_i$  is the scaled genotype for individual  $i$  as described above. Association tests were performed using the Python statsmodels Logit method.

For the non-additive phenotype example (**Supplementary Figure 9**), we performed simulations under a quadratic model:  $P = \beta G^2 + E$  where  $G$  is a vector of the squared sum of allele lengths scaled by the mean allele length, and  $P$ ,  $\beta$ ,  $E$  are as described above. Two sets of association tests were performed: the first tested for association between STR length and phenotype (**Supplementary Figure 9B**) and the second set performed a separate association test for each STR allele treating the allele as a bi-allelic locus (**Supplementary Figure 9C**).

In all cases 100 separate simulations were performed and power was defined as the percent of simulations for which the nominal association p-value was less than 0.05. Figures show results for all simulations with  $\beta$  set to 0.1.

### **eSTR analysis**

Data for eSTR analysis was obtained from the Genotype-Tissue Expression (GTEx) through dbGaP under phs000424.v7.p2. This included high coverage (30x) Illumina whole genome sequencing (WGS) data from 650 unrelated samples, Omni 2.5 SNP genotypes for 450 samples, and gene-level RPKM values for whole blood in 336 samples. STRs were genotyped from WGS data using HipSTR v0.5 and subject to the same quality filtering as SSC samples. STRs were additionally imputed to Omni2.5 data with Beagle as described above. Downstream analyses were restricted to the 336 samples with available whole blood expression data. These samples consisted of 284 European, 45 African American, 3 Asian, and 3 Amerindian samples and 2 samples with no population label available.

We performed separate eSTR analyses using HipSTR and imputed genotypes. In each case, as described previously<sup>15</sup>, we performed a separate association test between gene expression and each STR within 100kb of the gene using a model  $Y = \beta X + C + \varepsilon$ , where  $X$  denotes STR genotype lengths,  $Y$  denotes expression values,  $\beta$  denotes the effect size,  $C$  denotes various

covariates, and  $\varepsilon$  is the error term. Following our previous study<sup>75</sup>, we used “STR dosage”, defined as the sum of repeat lengths of the two alleles for each sample, to define STR genotypes. All repeat lengths are reported as length difference from the hg19 reference, with 0 representing the reference allele. STR dosages were scaled to have mean 0 and variance 1. Genes with median expression of 0 were excluded and expression values for remaining genes were quantile normalized to a standard normal distribution. We included sex, population structure, and technical variation in expression as covariates. For population structure, we used the top 15 principal components resulting from perform principal components analysis on the matrix of SNP genotypes from each sample. To control for technical variation in expression, we applied PEER factor correction<sup>76,77</sup> using 83 PEER factors.

We used model comparison to determine whether the best eSTR for each gene explained variation in gene expression beyond a model consisting of the best eSNP. As described previously<sup>75</sup>, for each gene with an eSTR we determined the lead eSNP with the strongest p-value. We then compared two linear models:  $Y \sim \text{eSNP}$  (SNP-only model) vs.  $Y \sim \text{eSNP} + \text{eSTR}$  (SNP+STR model) using the `anova_lm` function in the python `statsmodels.api.stats` module. We used CAVIAR v1.0 to further fine-map eSTR signals against the top 100 eSNPs within 100kb of each gene. Pairwise-LD between the eSTR and eSNPs was estimated using the Pearson correlation between SNP dosages (0, 1, or 2) and STR dosages (sum of the two repeat allele lengths).

### **Comparison to DRPLA founder haplotypes**

The founder haplotype for the expansion allele in *ATN1* implicated in DRPLA was taken from Table 1 of Veneziano *et al.*<sup>58</sup> and consists of rs4963516, rs1007924, rs7310941, rs7303722, rs2239167, rs34199021, rs2071075, rs2071076, and rs2159887 with hg19 alleles G, A, G, T, A, A, T, C, and C respectively. Distance from the founder haplotype was calculated as the number of mismatches.

## **Data Availability**

Phased SNP-STR haplotypes for 1000 Genomes Project phase 3 samples and example

commands for imputation are available at

[https://gymreklab.github.io/2018/03/05/snpstr\\_imputation.html](https://gymreklab.github.io/2018/03/05/snpstr_imputation.html). Upon acceptance for publication STR genotypes and phased SNP-STR haplotypes for the SSC samples will be made available at <https://base.sfari.org/>.

## Code Availability

Analysis scripts and Jupyter notebooks for reproducing the figures in this study are provided in the Github repository <https://github.com/gymreklab/snpstr-imputation>.

## References

1. Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421–427 (2014).
2. Turcot, V. *et al.* Protein-altering variants associated with body mass index implicate pathways that control energy intake and expenditure in obesity. *Nat. Genet.* **50**, 26–41 (2018).
3. Scott, R. A. *et al.* An Expanded Genome-Wide Association Study of Type 2 Diabetes in Europeans. *Diabetes* **66**, 2888–2902 (2017).
4. Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies. *Nat. Rev. Genet.* **11**, 499–511 (2010).
5. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass index in ~700,000 individuals of European ancestry. *bioRxiv* 274654 (2018).  
doi:10.1101/274654
6. Sekar, A. *et al.* Schizophrenia risk from complex variation of complement component 4. *Nature* **530**, 177–183 (2016).

7. Boettger, L. M. *et al.* Recurring exon deletions in the HP (haptoglobin) gene contribute to lower blood cholesterol levels. *Nat. Genet.* **48**, 359–366 (2016).
8. Leffler, E. M. *et al.* Resistance to malaria through structural variation of red blood cell invasion receptors. *Science* **356**, (2017).
9. Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921 (2001).
10. Press, M. O., Carlson, K. D. & Queitsch, C. The overdue promise of short tandem repeat variation for heritability. *Trends Genet.* **30**, 504–512 (2014).
11. Hannan, A. J. Tandem repeats mediating genetic plasticity in health and disease. *Nat. Rev. Genet.* (2018). doi:10.1038/nrg.2017.115
12. Hannan, A. J. Tandem repeat polymorphisms: modulators of disease susceptibility and candidates for ‘missing heritability’. *Trends Genet.* **26**, 59–65 (2010).
13. Gymrek, M., Willems, T., Reich, D. & Erlich, Y. Interpreting short tandem repeat variations in humans using mutational constraint. *Nat. Genet.* **49**, 1495–1501 (2017).
14. Willems, T. *et al.* The landscape of human STR variation. *Genome Res.* **24**, 1894–1904 (2014).
15. Gymrek, M. *et al.* Abundant contribution of short tandem repeats to gene expression variation in humans. *Nat. Genet.* **48**, 22–29 (2016).
16. Quilez, J. *et al.* Polymorphic tandem repeats within gene promoters act as modifiers of gene expression and DNA methylation in humans. *Nucleic Acids Res.* **44**, 3750–3762 (2016).
17. Hefferon, T. W., Groman, J. D., Yurk, C. E. & Cutting, G. R. A variable dinucleotide repeat in the CFTR gene contributes to phenotype diversity by forming RNA secondary structures that alter splicing. *Proceedings of the National Academy of Sciences* **101**, 3504–3509

- (2004).
18. Hui, J. *et al.* Intronic CA-repeat and CA-rich elements: a new class of regulators of mammalian alternative splicing. *EMBO J.* **24**, 1988–1998 (2005).
  19. Tollervey, J. R. *et al.* Characterizing the RNA targets and position-dependent splicing regulation by TDP-43. *Nat. Neurosci.* **14**, 452–458 (2011).
  20. Mirkin, S. M. Expandable DNA repeats and human disease. *Nature* **447**, 932–940 (2007).
  21. Sutcliffe, J. S. *et al.* DNA methylation represses FMR-1 transcription in fragile X syndrome. *Hum. Mol. Genet.* **1**, 397–400 (1992).
  22. van Blitterswijk, M., DeJesus-Hernandez, M. & Rademakers, R. How do C9ORF72 repeat expansions cause amyotrophic lateral sclerosis and frontotemporal dementia: can we learn from other noncoding repeat expansion disorders? *Curr. Opin. Neurol.* **25**, 689–700 (2012).
  23. Grünewald, T. G. P. *et al.* Chimeric EWSR1-FLI1 regulates the Ewing sarcoma susceptibility gene EGR2 via a GGAA microsatellite. *Nat. Genet.* **47**, 1073–1078 (2015).
  24. Payseur, B. A., Place, M. & Weber, J. L. Linkage disequilibrium between STRPs and SNPs across the human genome. *Am. J. Hum. Genet.* **82**, 1039–1050 (2008).
  25. Tishkoff, S. A. *et al.* Global patterns of linkage disequilibrium at the CD4 locus and modern human origins. *Science* **271**, 1380–1387 (1996).
  26. Mountain, J. L. *et al.* SNPSTRs: empirically derived, rapidly typed, autosomal haplotypes for inference of population history and mutational processes. *Genome Res.* **12**, 1766–1772 (2002).
  27. Gymrek, M., Willems, T., Reich, D. & Erlich, Y. Interpreting short tandem repeat variations in humans using mutational constraint. *Nat. Genet.* **49**, 1495–1501 (2017).
  28. Sun, J. X. *et al.* A direct characterization of human mutation based on microsatellites. *Nat. Genet.* **44**, 1161–1165 (2012).

29. Lai, Y. & Sun, F. Microsatellite mutations during the polymerase chain reaction: mean field approximations and their applications. *J. Theor. Biol.* **224**, 127–137 (2003).
30. Lai, Y., Shinde, D., Arnheim, N. & Sun, F. The Mutation Process of Microsatellites During the Polymerase Chain Reaction. *J. Comput. Biol.* **10**, 143–155 (2003).
31. Willems, T. *et al.* Genome-wide profiling of heritable and de novo STR variations. *Nat. Methods* (2017). doi:10.1038/nmeth.4267
32. Majounie, E. *et al.* Case control analysis of repeat expansion size in ataxia. *Neurosci. Lett.* **429**, 28–32 (2007).
33. Ambrose, K. K. *et al.* Analysis of CTG repeat length variation in the gene in the general population and the molecular diagnosis of myotonic dystrophy type 1 in Malaysia. *BMJ Open* **7**, e010711 (2017).
34. Figley, M. D., Thomas, A. & Gitler, A. D. Evaluating noncoding nucleotide repeat expansions in amyotrophic lateral sclerosis. *Neurobiol. Aging* **35**, 936.e1–4 (2014).
35. Gouw, L. G. *et al.* Analysis of the dynamic mutation in the SCA7 gene shows marked parental effects on CAG repeat transmission. *Hum. Mol. Genet.* **7**, 525–532 (1998).
36. Tang, H. *et al.* Profiling of Short-Tandem-Repeat Disease Alleles in 12,632 Human Whole Genomes. *Am. J. Hum. Genet.* **101**, 700–715 (2017).
37. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
38. Gymrek, M., Golan, D., Rosset, S. & Erlich, Y. lobSTR: A short tandem repeat profiler for personal genomes. *Genome Res.* **22**, 1154–1162 (2012).
39. 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).



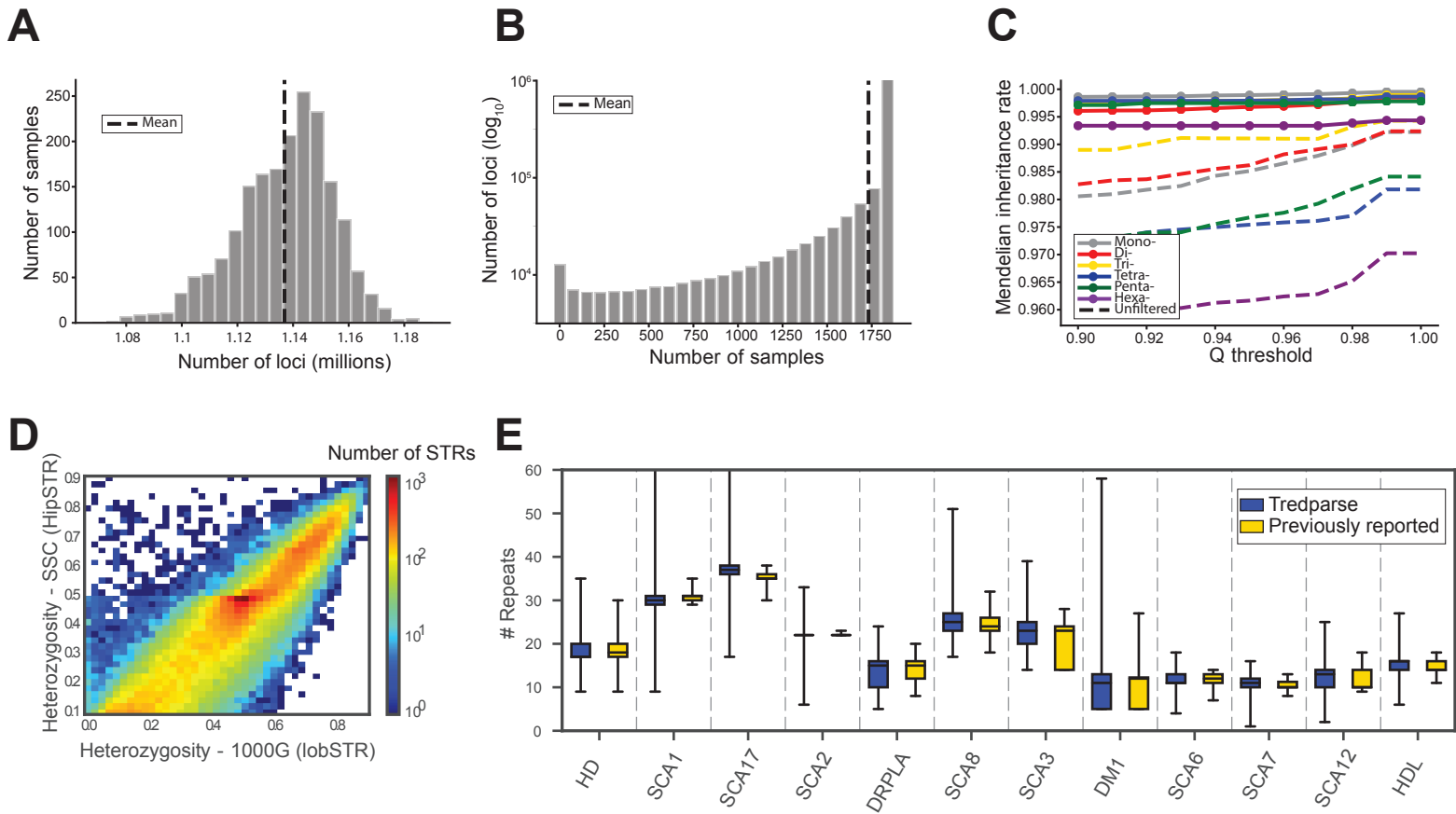
40. Edge, M. D., Algee-Hewitt, B. F. B., Pemberton, T. J., Li, J. Z. & Rosenberg, N. A. Linkage disequilibrium matches forensic genetic records to disjoint genomic marker sets. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 5671–5676 (2017).
41. International HapMap Consortium. A haplotype map of the human genome. *Nature* **437**, 1299–1320 (2005).
42. Payseur, B. A. & Jing, P. A Genomewide Comparison of Population Structure at STRPs and Nearby SNPs in Humans. *Mol. Biol. Evol.* **26**, 1369–1377 (2009).
43. Gymrek, M. *et al.* Abundant contribution of short tandem repeats to gene expression variation in humans. *Nat. Genet.* **48**, 22–29 (2016).
44. Contente, A., Dittmer, A., Koch, M. C., Roth, J. & Dobbelstein, M. A polymorphic microsatellite that mediates induction of PIG3 by p53. *Nat. Genet.* **30**, 315–320 (2002).
45. Shimajiri, S. *et al.* Shortened microsatellite d(CA)<sub>21</sub> sequence down-regulates promoter activity of matrix metalloproteinase 9 gene. *FEBS Lett.* **455**, 70–74 (1999).
46. The GTEx Consortium *et al.* The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
47. Hormozdiari, F., Kostem, E., Kang, E. Y., Pasaniuc, B. & Eskin, E. Identifying causal variants at loci with multiple signals of association. *Genetics* **198**, 497–508 (2014).
48. Borel, C. *et al.* Tandem repeat sequence variation as causative cis-eQTLs for protein-coding gene expression variation: the case of CSTB. *Hum. Mutat.* **33**, 1302–1309 (2012).
49. Lalioti, M. D. *et al.* Dodecamer repeat expansion in cystatin B gene in progressive myoclonus epilepsy. *Nature* **386**, 847–851 (1997).
50. Wheeler, A. C. *et al.* Associated features in females with an FMR1 premutation. *J. Neurodev. Disord.* **6**, 30 (2014).

51. Ha, A. D., Beck, C. A. & Jankovic, J. Intermediate CAG Repeats in Huntington's Disease: Analysis of COHORT. *Tremor Other Hyperkinet. Mov.* **2**, (2012).
52. Brenman, L. M. Spinocerebellar Ataxia Type 6 (SCA6) Phenotype in a Patient with an Intermediate Mutation Range CACNA1A Allele. *J. Neurol. Neurophysiol.* **04**, (2013).
53. Lee, D.-Y. & McMurray, C. T. Trinucleotide expansion in disease: why is there a length threshold? *Curr. Opin. Genet. Dev.* **26**, 131–140 (2014).
54. Koide, R. *et al.* Unstable expansion of CAG repeat in hereditary dentatorubral-pallidoluysian atrophy (DRPLA). *Nat. Genet.* **6**, 9–13 (1994).
55. Gan, S.-R., Ni, W., Dong, Y., Wang, N. & Wu, Z.-Y. Population genetics and new insight into range of CAG repeats of spinocerebellar ataxia type 3 in the Han Chinese population. *PLoS One* **10**, e0134405 (2015).
56. Paradisi, I., Ikonomu, V. & Arias, S. Huntington disease-like 2 (HDL2) in Venezuela: frequency and ethnic origin. *J. Hum. Genet.* **58**, 3–6 (2013).
57. Laffita-Mesa, J. M. *et al.* De novo mutations in ataxin-2 gene and ALS risk. *PLoS One* **8**, e70560 (2013).
58. Veneziano, L. *et al.* A shared haplotype for dentatorubropallidoluysian atrophy (DRPLA) in Italian families testifies of the recent introduction of the mutation. *J. Hum. Genet.* **59**, 153–157 (2014).
59. Gibson, G. Rare and common variants: twenty arguments. *Nat. Rev. Genet.* **13**, 135–145 (2012).
60. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* **536**, 41–47 (2016).
61. Hannan, A. J. Tandem repeats mediating genetic plasticity in health and disease. *Nat. Rev. Genet.* (2018). doi:10.1038/nrg.2017.115

62. Mousavi, N., Shleizer-Burko, S. & Gymrek, M. Profiling the genome-wide landscape of tandem repeat expansions. *bioRxiv* 361162 (2018). doi:10.1101/361162
63. Bakhtiari, M., Shleizer-Burko, S., Gymrek, M., Bansal, V. & Bafna, V. Targeted Genotyping of Variable Number Tandem Repeats with adVNTR. (2017). doi:10.1101/221754
64. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
65. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
66. Li, H. Tabix: fast retrieval of sequence features from generic TAB-delimited files. *Bioinformatics* **27**, 718–719 (2011).
67. Karolchik, D. The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* **32**, 493D–496 (2004).
68. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
69. Chakraborty, R., De Andrade, M., Daiger, S. P. & Budowle, B. Apparent heterozygote deficiencies observed in DNA typing data and their implications in forensic applications. *Ann. Hum. Genet.* **56**, 45–57 (1992).
70. Fisher, S. A., Lewis, C. M. & Wise, L. H. Detecting population outliers and null alleles in linkage data: application to GAW12 asthma studies. *Genet. Epidemiol.* **21 Suppl 1**, S18–23 (2001).
71. Tang, H. *et al.* Profiling of Short-Tandem-Repeat Disease Alleles in 12,632 Human Whole Genomes. *Am. J. Hum. Genet.* **101**, 700–715 (2017).
72. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for thousands of genomes. *Nat. Methods* **9**, 179–181 (2011).

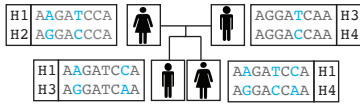
73. O'Connell, J. *et al.* A general approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genet.* **10**, e1004234 (2014).
74. Pemberton, T. J., Sandefur, C. I., Jakobsson, M. & Rosenberg, N. A. Sequence determinants of human microsatellite variability. *BMC Genomics* **10**, 1–19 (2009).
75. Gymrek, M. *et al.* Abundant contribution of short tandem repeats to gene expression variation in humans. *Nat. Genet.* **48**, 22–29 (2016).
76. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nat. Protoc.* **7**, 500–507 (2012).
77. GTEx Consortium *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).

Saini, et al. Figure 1

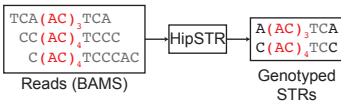


**A**

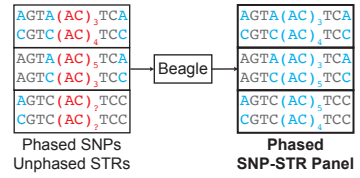
**Step 1: Family based SNP phasing**



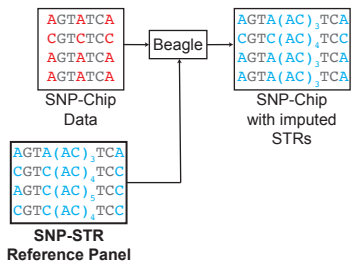
**Step 2: STR genotyping**



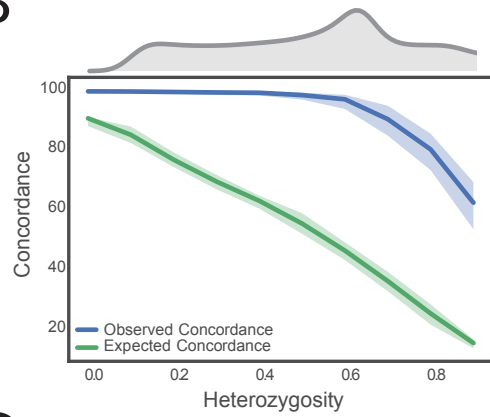
**Step 3: Joint SNP/STR phasing**



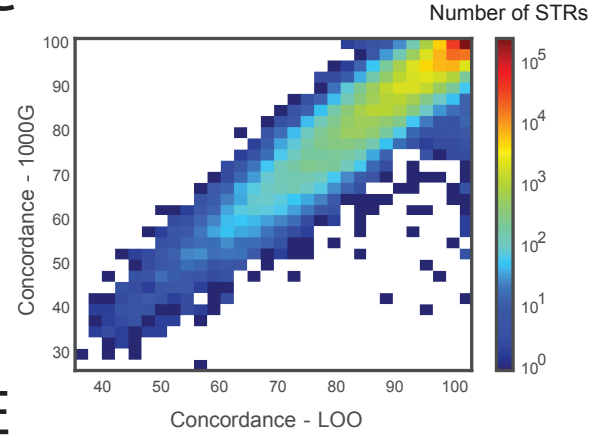
**STR imputation from SNP genotypes**



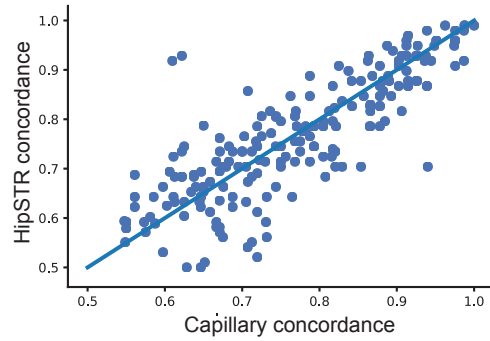
**B**



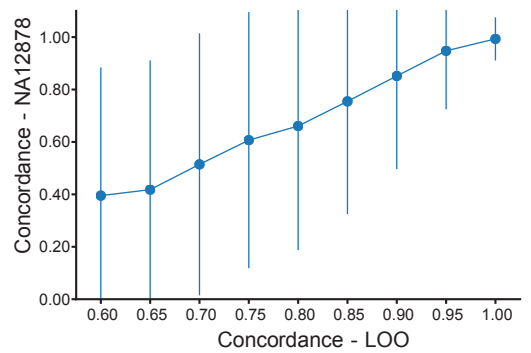
**C**



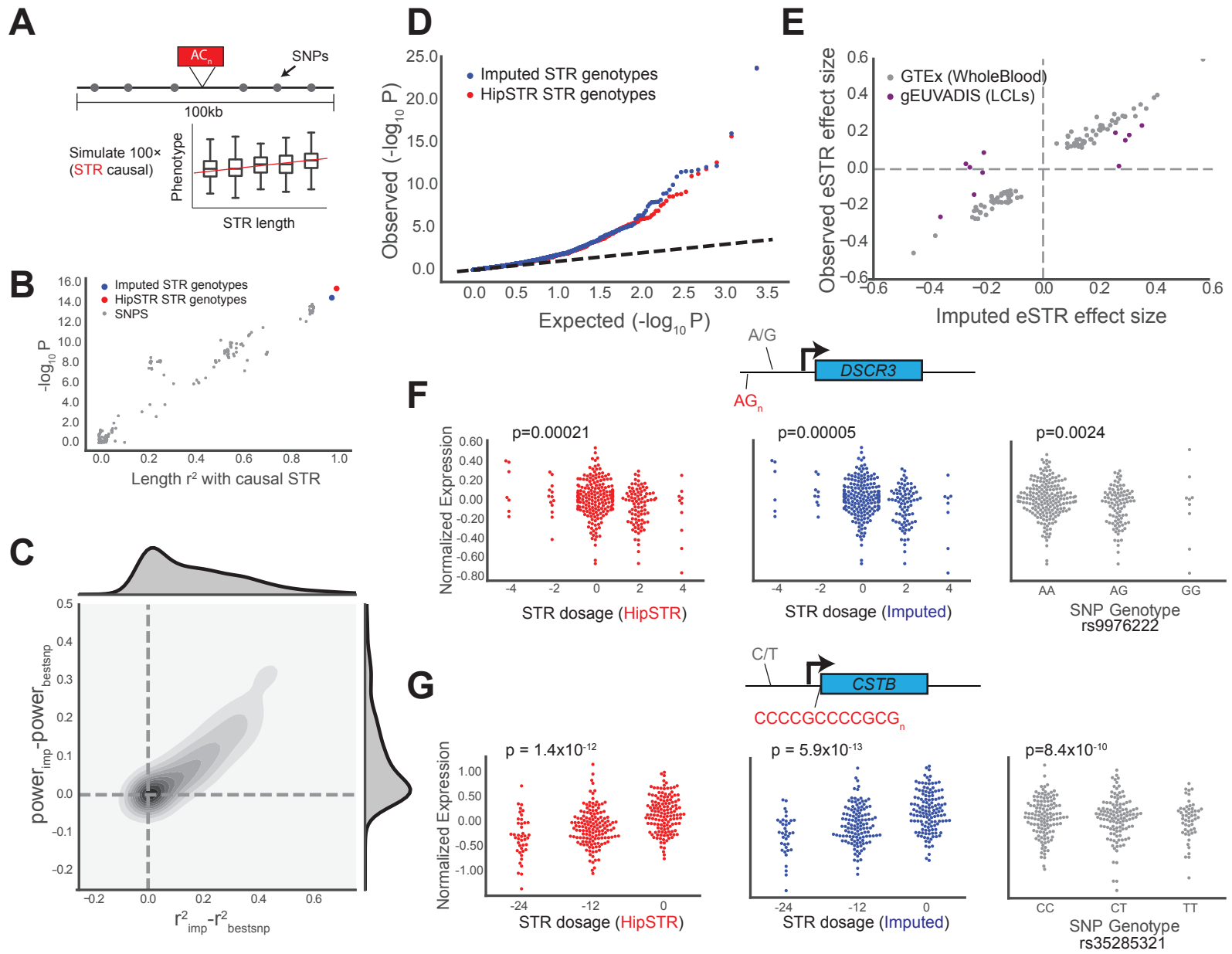
**D**



**E**



Saini, et al. Figure 3



## Saini, et al. Figure 4

