# Detecting sample swaps in diverse NGS data types using linkage disequilibrium

Nauman Javed[1], Yossi Farjoun[2], Tim Fennell[2], Charles Epstein[2], Bradley E. Bernstein[1,2], Noam Shoresh[1,2,†]

[1]Department of Pathology and Center for Cancer Research, Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts 02114, USA.

[2]Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA.

†Corresponding author. Email: nshoresh@broadinstitute.org

1  **As the number of genomics datasets grows rapidly, sample mislabeling has become a high stakes**
2  **issue. We present CrosscheckFingerprints (Crosscheck), a tool for quantifying sample-relatedness and**
3  **detecting incorrectly paired sequencing datasets from different donors. Crosscheck outperforms**
4  **similar methods and is effective even when data are sparse or from different assays.** Application of
5  **Crosscheck to 8851 ENCODE ChIP-, RNA-, and DNase-seq datasets enabled us to identify and correct**
6  **dozens of mislabeled samples and ambiguous metadata annotations, representing ~1% of ENCODE**
7  **datasets.**

8  Biomedical research is rapidly embracing large-scale analysis of next-generation sequencing (NGS)
9  datasets, often by integrating data generated by consortia or many individual research labs. Parallelized
10 NGS analysis of tissues from many different patients is also commonplace in clinical genomics pipelines.
11 In these settings, sample or data mislabeling, where datasets are incorrectly associated with a donor,
12 can lead to erroneous conclusions, misdirect future research, and affect treatment decisions[1-3] (Fig 1a).
13 Verifying the relatedness of samples that nominally share a donor is therefore a crucial quality-control
14 step in any NGS pipeline.

15 Several methods utilize genetic information from NGS datasets as an endogenous barcode to
16 verify sample relatedness[4-10]. The common logic behind these tools is that each genome harbors a
17 unique set of single nucleotide polymorphisms (SNPs) which are shared between datasets originating
18 from the same donor. A limitation of these methods is their requirement that sequencing reads from
19 both inputs overlap the exact genomic position of informative SNPs. When insufficient reads satisfy this
20 condition—for example when the input datasets are shallow or target different genomic regions (i.e
21 different transcription factors), the power to evaluate sample relatedness is compromised. Many NGS-
22 based studies now integrate multiple types of assays[11-15] and utilize shallow sequencing to reduce cost at
23 the expense of read-depth. This is commonly encountered in highly multiplexed experiments,
24 sequencing spike-ins, and large cohort sequencing efforts in population and cancer genomics (i.e. 1000
25 Genomes, structural variant calling). We therefore set out to develop a method for quantifying sample-
26 relatedness that was both robust to shallow sequencing depth and that could be systematically applied
27 to modern large-scale projects incorporating multiple data types.

28 Linkage disequilibrium (LD) is the non-random association of alleles at different loci within a
29 given population[16]. This association implies that comparing datasets across SNPs in high LD—termed LD-
30 blocks—would provide more statistical power to compare datasets than using single SNPs alone.
31 Because of LD, two non-overlapping reads from different datasets may support (or provide evidence
32 against) a common genetic background, as long as they overlap SNPs in the same LD block (Fig 1b). For
33 each input dataset, Crosscheck uses reads overlapping SNPs within each LD block to calculate a block
34 allele fraction and compute diploid genotype likelihoods, which are then compared (Methods). The

35 relative likelihood of a shared or distinct genetic background at each block is reported as a log-odds ratio
36 (LOD score). These scores are combined across all blocks to report a genome-wide LOD score. This
37 calculation relies on two approximations: that linkage between SNPs in an LD block is perfect, and that
38 SNPs in distinct blocks are independent. A positive LOD score indicates a higher likelihood that the two
39 datasets share a donor, while a negative LOD score suggests that the datasets are from distinct donors.
40 The Crosscheck calculation assumes that the two datasets are a priori equally likely to be from the same
41 donor as they are from different ones. It is possible to incorporate a different prior expectation for a
42 mismatch by shifting the LOD scores (Methods). Though the magnitude of the LOD score reflects
43 genotyping confidence, simplifying assumptions prevent direct interpretation of the LOD score as a true
44 likelihood ratio (Methods). Crosscheck is implemented as part of Picard-Tools (https://github.com/
45 broadinstitute/picard), and is routinely used for quality control by the Broad Institute's Genomics
46 Platform, using a small set of LD blocks optimized for use with whole-exome-sequencing data.

47 We reasoned that applying Crosscheck across a large, genome-wide set of LD-blocks (haplotype
48 map) would allow us to compare the genotype of diverse datasets and would be robust to low coverage
49 and sequencing errors. We constructed a map consisting of nearly 60,000 common (minor allele
50 frequency $\geq 10\%$) bi-allelic SNPs from the 1000 Genomes[11] project, the majority of which lie in LD-
51 blocks of two or more SNPs in order to maximize the probability of informative read overlap (Fig 1c,
52 Methods). SNPs within each block are highly correlated ($r^2 > 0.85$), while SNPs between blocks are
53 approximately independent ($r^2 < 0.10$). Increasing or decreasing the thresholds for within-block and
54 between-blocks correlations by 0.05 had no effect on the method's performance on a testing data set
55 (described in the next paragraph). Finally, in order to reduce bias from donor ancestry, we required that
56 LD blocks have similar allele frequencies across different human sub-populations. The pipeline for
57 creating haplotype maps exists as a standalone tool
58 (https://github.com/naumanjaved/fingerprint_maps) and can be customized to create LD blocks in
59 specific genomic areas (i.e. coding regions) and for either hg19 or GRCh38.

60 To pilot our method, we calculated LOD scores between donor-matched and donor-mismatched
61 pairs of public datasets from the ENCODE[12] database, which hosts data from thousands of diverse NGS
62 experiments (Methods). Classification performance was measured in terms of the false flag rate (FFR),
63 the fraction of donor-matched pairs incorrectly flagged as donor-mismatches, and the false match rate
64 (FMR), the fraction of donor-mismatched pairs incorrectly identified as donor-matches. Our testing set
65 comprised of all pairwise comparisons between 279 RNA-, DNase-, and ChIP-seq (targeting histones,
66 CTCF, or POL2) datasets with verified donor annotations (supplementary table S1), and all donor-
67 mismatched comparisons between 98 ChIP-seq experiments targeting transcription factors and
68 chromatin modifiers (supplementary table S2). This resulted in a final testing set of 34,336 donor-
69 mismatches, and 9,767 donor-matches. Regardless of the input assay or enrichment target, Crosscheck
70 correctly classified almost all dataset pairs with 0% FMR and 0.01% FFR, and showed a clear separation
71 between donor-mismatches (negative LOD) and donor-matches (positive LOD) (Fig. 1d). Our method
72 therefore confidently detects donor-matched and donor-mismatched dataset pairs.

73 We next quantified how using LD blocks improves classification performance. We generated two
74 equally sized subsets of our full haplotype map—one comprised solely of unlinked SNPs and the other
75 containing only LD blocks with two or more SNPs, and used these to classify the same testing dataset
76 pairs. To simulate sparse datasets generated by spike-ins and multiplexed sequencing, we conducted
77 each comparison at a range of sequencing depths, expressed as the percentage of reads subsampled

78     from the original datasets (Methods, Supplementary Fig 1a). Using LD blocks significantly decreased

79     FMR and FFR, particularly at lower read depths and for cross-assay/target comparisons (Fig 1e,

80     Supplementary Fig 1b). For example, at 5% sub-sampling ($\leq \sim 10^7$ reads), using LD blocks decreased the

81     FMR and FFR by nearly 10% relative to using single SNPs for cross-assay comparisons.

82        As mentioned above, there are other tools that quantify genetic sample relatedness. For

83     comparison purposes, we considered only methods that could be applied to the general use case that

84     Crosscheck is designed to address, namely comparing any two NGS datasets, and that can be deployed

85     at scale, so that calculating tens-to-hundreds of thousands of comparisons is tractable. Two of the

86     methods we examined, HYSIS[6] and BAM-matcher[7], did not satisfy these criteria. Two other tools,

87     Conpair[8] and BAMixChecker[9], provided inconclusive results for a high percentage of the testing-set

88     comparisons (Methods). NGSCheckmate[10](NGSC) is a model-based method that compares datasets by

89     correlating allele fractions across a panel of reference SNPs, and was the only other method that could

90     be directly compared to Crosscheck on the testing dataset. At high and intermediate read-depths, both

91     methods show similar performance. At lower read depths ($\leq 15\%$ subsampling) however, Crosscheck

92     outperforms NGSC, as indicated by a consistently lower FMR and FFR (Fig. 1f). Crosscheck is particularly

93     effective at classifying cross assay dataset pairs, where it shows a 2-3% lower FMR and FFR than NGSC at

94     5% subsampling. In these use cases, Crosscheck performs better than NGSC due to its use of LD and the

95     large number of SNPs in the haplotype map. Using LD blocks allows comparison of non-overlapping

96     reads, while using a large set of SNPs increases the chance that input datasets will contain genetically

97     informative reads. An illustrative example is a specific comparison between two ChIP-seq datasets, one

98     targeting H3K27me3 and the other H3K27ac. At 5% subsampling, these datasets cover 8% and 2% of the

99     genome respectively, and overlap at only 0.02%, which is expected from these mutually exclusive

100     histone modifications. Given this small set of potentially informative reads, NGSCheckmate wrongly

101     concludes that the datasets are derived from the same donor, while CrossCheck is still able to make the

102     correct call (Supplementary Fig. 1e).We have also tested CrossCheck, NGSC, BAMixChecker and Conpair

103     on sample pairs from 7 donors that are genetically related. We found that CrossCheck can identify all

104     pairs of samples from related individuals as donor mismatches, and is superior in this context to the

105     other tools (Supplementary Fig. 2).

106        Finally, we used the distribution of LOD scores from incorrectly classified pairs to define an

107     inconclusive LOD score range of -5 < LOD < 5, in which a dataset pair cannot be confidently classified

108     (Methods, Supplementary Fig. 1c). Outside of this range, any pair with LOD ≥ 5 is denoted a donor-

109     match, and those with LOD ≤ -5 are flagged as donor-mismatches. The inconclusive range highlights the

110     interpretability of Crosscheck's LOD score relative to NGSC's binary outputs (match or mismatch), since

111     clear donor-mismatches can be prioritized and investigated separately from inconclusive comparisons.

112     We conclude that using Crosscheck with a full haplotype map enables more accurate detection of

113     donor-mismatched pairs in diverse and shallow collections of data. To illustrate the utility of our method

114     on a consortium-scale dataset, we next analyzed the remaining datasets in ENCODE. We used our

115     method to verify the donor-annotation for all human hg19 aligned DNase-, RNA-, and ChIP-seq datasets

116     in the ENCODE database whose annotated donor was represented by at least 4 datasets – a total of

117     8,851 datasets (Fig 2a). To scale our analysis to a database of this size, we compared each dataset to a

118     set of three representative datasets from its annotated donor, and flagged any dataset with LOD < 5 for

119     further review (Methods). To exclude the possibility that the representative set for each donor

120     contained a donor-mismatch, we required that all pairwise comparisons between representative

121 datasets yield an LOD score ≥ 5. This strategy scales linearly with the size of the database, and in our
122 case results in a 1000-fold reduction in computation relative to performing all pairwise comparisons.

123 Our strategy confirmed the annotated donor for 97% of datasets. The remaining 3% (256
124 datasets) were flagged as potential donor-mismatches (LOD ≤ -5), and only ~0.1% yielded inconclusive
125 results (-5 < LOD < 5) (Fig 2b). We next compared each flagged mismatch to the representative datasets
126 for each of the ENCODE donors in order to nominate a true donor identity. We also compared each
127 flagged mismatch to all other flagged mismatches in order to identify genetically consistent clusters and
128 uncover patterns of mislabeling.

129 This analysis uncovered 3 major categories of mislabeling (as well as a small fraction, 0.4%, of
130 datasets that exhibited a pattern consistent with cross-sample contamination, as described in Methods
131 and Supplementary Fig. 3). The first is a straightforward error where cells from one donor are mistakenly
132 labeled as deriving from a different donor. The likelihood of such a mistake increases when working with
133 several cell lines that are each used in a large number of experiments. For example, out of 4 flagged
134 datasets labeled as K562, two were shown to actually derive from GM12878 cells while the other two
135 derived from HEK293 cells. This type of mislabeling may also occur for primary cells or tissues when
136 many biological samples from multiple donors are obtained from the same source, as in the case of 300
137 embryonic tissue samples processed by ENCODE from a single lab.

138 The second class of mislabeling occurs when biological samples of the same cell type from
139 multiple donors are incorrectly labeled as deriving from a single donor. This is the case with some of the
140 commercially available primary cell lines that have been deeply interrogated by the consortium over
141 more than a decade, and for which cells have been procured multiple times. For example, HUVEC cells
142 are annotated as being derived from two different donors in the ENCODE metadata. However, our
143 analysis indicates that HUVEC samples actually derive from at-least 5 distinct donors (Fig 2c). This mis-
144 annotation went undetected by ENCODE's previous quality control pipelines because all samples were
145 of the same cell type and so exhibited similar epigenetic profiles.

146 The HUVEC example also highlights the third type of labeling inaccuracy, in which a single donor
147 is accessioned multiple times by dozens of different labs over several years. This results in slight
148 variations in donor name or description, leading to genetically identical samples being incorrectly
149 attributed to distinct donors. For example, some samples deriving from putative donor A are attributed
150 to HUVEC donor 1, while other samples from donor A are attributed to the distinct HUVEC donor 2.

151 Overall, our analysis of the ENCODE dataset suggested that substantive mislabeling error
152 occurred at a rate of ~1%. For these datasets, true donor identities were confirmed using ENCODE's
153 extensive metadata records and all mislabeled datasets were corrected (Methods).

154 In conclusion, we present a robust and easy-to-use method for quantifying sample relatedness
155 which outperforms similar methods. Combined with our method for database analysis and haplotype
156 map, CrosscheckFingerprints can be readily applied for detecting sample mislabeling in large, diverse
157 databases without any optimization. We suggest it as a critical component of any NGS quality control
158 pipeline.

159 **Methods**
160 **LOD Derivation**

161 Here, a basic overview of the fingerprinting LOD score derivation is provided. A more detailed derivation
162 is available at the Picard repository at:

163 https://github.com/broadinstitute/picard/raw/master/docs/fingerprinting/main.pdf

164 Consider a LD block/locus containing a single bi-allelic SNP with major allele $A$ and minor allele $B$, and
165 two sequencing datasets $x$ and $y$. Let $\theta$ and $\varphi$ denote the diploid haplotype of datasets $x$ and $y$
166 respectively at this locus. $\theta$ and $\varphi$ can each take one of three possible haplotypes: AA, AB, or BB. Let $s$
167 be a Bernoulli random variable where $s = 1$ denotes a sample swap (indicating that $x$ and $y$ arose from
168 two independent individuals) with posterior probability $p(s = 1 | x, y)$, and $s = 0$ denotes a shared
169 genetic origin (the samples came from the same individual). Using Bayes' rule and the prior probability
170 of no-swap, the posterior odds ratio of a no-swap vs. swap is given by:

$$\frac{p(s = 0 | x, y)}{p(s = 1 | x, y)} = \frac{p(x, y | s = 0)\, p(s = 0)}{p(x, y | s = 1)\, p(s = 1)} \tag{1}$$

171 We assume that in the case of a swap, the distinct individuals are independently sampled from the
172 population and that samples from the same individual have the same genotype, allowing us to write
173 $p(\theta, \varphi | s) = p(\theta)\, p(\varphi)$ for $s = 1$, and $p(\theta, \varphi | s) = p(\theta)$ if $\theta = \varphi$. Given that $x$ is conditionally
174 independent of $\varphi$ and $y$ given $\theta$, and $y$ is conditionally independent of $\theta$ given $\varphi$, we can also write
175 $p(x, y | \theta, \varphi) = p(x | \theta)\, p(y | \varphi)$.

176 With these two expressions, we derive that:

$$p(x, y | s) = \sum_{\theta, \varphi} p(x, y | \theta, \varphi, s)\, p(\theta, \varphi | s)$$

$$= \begin{cases} \sum_{\theta} p(x|\theta)\, p(\theta) \sum_{\varphi} p(y|\varphi)\, p(\varphi) & \text{if } s = 1 \\[2mm] \sum_{\theta = \varphi} p(x|\theta)\, p(y|\varphi)\, p(\theta) & \text{if } s = 0 \end{cases} \tag{2}$$

177 Substituting the results of (2) into (1), we rewrite the posterior odds of no-swap as:

$$\frac{\sum_{\theta = \varphi} p(x|\theta)\, p(y|\varphi)\, p(\theta)}{\sum_{\theta} p(x|\theta)\, p(\theta) \sum_{\varphi} p(y|\varphi)\, p(\varphi)} \cdot \frac{p(s = 0)}{p(s = 1)} \tag{3}$$

178 Next, we consider evidence over multiple blocks $i$ with correspondingly indexed $\theta_i$, $\varphi_i$, $x_i$, and $y_i$. **We**
179 **assume that the haplotypes at distinct blocks are independent**, and that reads at one block give no
180 information about another. In practice, this assumption is enforced by guaranteeing that a single read
181 cannot be used to provide genotype evidence at more than one locus. We calculate: $p(x | \theta) =$
182 $\prod_i p(x_i | \theta_i)$ and $p(y | \varphi) = \prod_i p(y_i | \varphi_i)$, and substitute into (3) to get:

5

$$\prod_i \left( \frac{\sum_{\theta_i = \varphi_i} p(x_i \mid \theta_i) \, p(y_i \mid \varphi_i) \, p(\theta_i)}{\sum_{\theta_i} p(x_i \mid \theta_i) \, p(\theta_i) \sum_{\varphi_i} p(y_i \mid \varphi_i) \, p(\varphi_i)} \right) \cdot \frac{p(s = 0)}{p(s = 1)} \tag{4}$$

183 Finally, since the odds ratio of no-swap to swap may vary by several orders of magnitude depending on
184 the input files, we compute the base 10 logarithm in order to facilitate comparison and interpretation:

$$\begin{aligned} LOD &= \log\left( \frac{odds_{same\ individual}}{odds_{different\ individual}} \right) \\ &= \sum_i \log\left( \frac{\sum_{\theta_i = \varphi_i} p(x_i \mid \theta_i) p(y_i \mid \varphi_i) p(\theta_i)}{\sum_{\theta_i} p(x_i \mid \theta_i) p(\theta_i) \sum_{\varphi_i} p(y_i \mid \varphi_i) p(\varphi_i)} \cdot \frac{p(s = 0)}{p(s = 1)} \right) \end{aligned} \tag{5}$$

185 **The program assumes a conservative prior of $\frac{p(s=0)}{p(s=1)} = 1$ by default.** A different prior would result in a
186 shift of the LOD score by a constant, and users may adjust the LOD score by such a constant as needed
187 on a case-by-case basis. A positive LOD (log-odds ratio) is interpreted as evidence for the two datasets $x$
188 and $y$ arising from the same individual, while a negative LOD is evidence of a sample-swap, i.e. the two
189 datasets arose from different individuals. Scores close to zero are inconclusive, and tend to result from
190 low coverage, or poor overlap between the two datasets, at the observed sites.

191 To see the expected maximal contribution of a single locus, we assume that the likelihoods in (5) are
192 vanishingly small when the data doesn't match the genotype. Thus, the LOD for a single locus reduces to
193 $-\log p(\theta)$. The expected LOD contribution needs to be marginalized over the different possible
194 genotypes, leading to a $-\sum_\theta p(\theta) \log p(\theta)$, which obtains a maximal value of $1.5 \log_{10} 2 \approx 0.45$ at an
195 allele frequency of 0.5(leading to $p(\theta = AA) = 0.25$, $p(\theta = AB) = 0.5$, and $p(\theta = BB) = 0.25$). This
196 means that when creating the haplotype map, it is most informative to choose variants with an allele
197 frequency close to 0.5.

198 There is no theoretical lower limit to the contribution of a single locus. This is because, in theory,
199 overwhelming evidence (hundreds of genetically-consistent, high-quality reads) of different genotypes
200 for two datasets at even a single locus is sufficient to rule out that the samples are derived from the
201 same donor. However, as noted below in the section on the limitations of LOD calculation, there are
202 multiple factors that this formulation does not account for. Our approach ultimately relies on
203 cumulative evidence, albeit noisy, from a large number of loci, rather than looking for the small number
204 of high-confidence cases. It is for this reason that in the implementation of equation (5) in the code, we
205 have included an explicit lower cap on the possible contribution of any single LD block. The selection of
206 the specific value at which to cap the negative contribution was guided by the following argument: We
207 consider a single specific locus, and assume a conservative prior, $(s = 0)/p(s = 1) = 1$. In addition, we
208 assume that at that locus one dataset is only compatible with a single genotype, namely $p(y \mid \theta)$ is
209 nonzero for only one value of $\theta$. In this case the contribution to the likelihood ratio for that locus
210 reduces to:

211
$$\frac{p(x \mid \theta) p(y \mid \theta) p(\theta)}{\left( \sum_{\theta_i} p(x \mid \theta_i) p(\theta_i) \right) p(y \mid \theta) p(\theta)} \gtrsim p(x \mid \theta)$$

212 If both samples are in fact from the same donor, and the discrepancy between x and $\theta$ is due to a
213 sequencing error, $10^{-3}$ is a reasonable ballpark estimate of $p(x \mid \theta)$[17]. With this, the actual score
214 calculated by CrossCheck is:

$$LOD' = \sum_i max\left( log\left( \frac{\sum_{\theta_i=\varphi_i} p(x_i \mid \theta_i)p(y_i \mid \varphi_i)p(\theta_i)}{\sum_{\theta_i} p(x_i \mid \theta_i)p(\theta_i) \sum_{\varphi_i} p(y_i \mid \varphi_i)p(\varphi_i)} \cdot \frac{p(s=0)}{p(s=1)} \right), \sigma \right) \quad (6)$$

215 Where $\sigma = -3$ by default, and is a parameter that can be set by the user.

**Calculation of data likelihoods $p(x \mid \theta)$ from sequencing reads**

217 The program assumes that sequencing data arrives in the form of reads from a single individual (i.e. not
218 contaminated), from a diploid location in the genome, and with no reference bias. Only non-secondary,
219 non-duplicate reads with mapping quality greater than 20 are used to calculate likelihoods. In addition,
220 bases must have a quality score of at least 20 and must agree with either the reference or pre-
221 determined alternate base to support observations at haplotype blocks. Since the algorithm assumes
222 that read evidence is independent, the reads should have been duplicate-marked prior to fingerprinting.
223 The algorithm doesn't use SNPs from the same read-pair twice, since this would violate the assumption
224 of independence.

225 Consider a dataset $x$ for which we observe $n$ total sequencing reads, denoted by $r_k$, at a locus
226 containing a single bi-allelic SNP with major allele $A$ and minor allele $B$. The possible block haplotypes
227 are then $\theta \in \{AA, AB, BB\}$. For each read $r_k$ which overlaps the SNP, let $o_k \in \{A, B\}$ denote the observed
228 SNP allele and let $e_k \in (0,1)$ denote the probability of error of each observation(the quality score). We
229 seek to compute the likelihood of the data (the sequencing reads $r_k$) given the haplotypes. The
230 likelihood of a single base observation $p(o_i, e_i \mid \theta)$ is expressed by:

$$p(o_k, e_k \mid \theta) = \begin{cases} I_B(o_k)e_k + I_A(o_k)(1-e_k) & \theta = AA \\ 0.5 & \theta = AB \\ I_A(o_k)e_k + I_B(o_k)(1-e_k) & \theta = BB \end{cases} \quad (6)$$

231 where $I$ is an indicator function such that $I_A(o) = \begin{cases} 1 \ if \ o = A \\ 0 \ if \ o = B \end{cases}$ and $I_B(o) = \begin{cases} 1 \ if \ o = B \\ 0 \ if \ o = A \end{cases}$ and the
232 assumption is that an error will cause a switch in the observed allele from A to B.

233 The likelihood model for all reads $r$ can then be written as:

$$p(r \mid \theta) = p(o, e \mid \theta) = \prod_{k=0}^{n} p(o_k, e_k \mid \theta) \quad (7)$$

**Incorporation of Linkage Information**

235 The calculations above assume an LD block containing a single SNP for ease of computation, but the
236 framework is easily extended to account for LD blocks containing multiple SNPs, which increases power
237 of comparison. Each LD block used for genotyping contains an "anchor" SNP which is in high linkage with
238 all other SNPs within the block, and independent of all other anchor SNPs in other blocks. Given that all

7

239 SNPs in a block are tightly linked(enforced with a strict $r^2$ correlation cutoff), **we make the simplifying**
240 **assumption that the genotype at any SNP within an LD block is perfectly correlated with the genotype**
241 **of the anchor SNP, and that all SNPs within a block have the same allele frequency, equal to that of**
242 **the anchor SNP**. Then, reads overlapping any SNP within a block can be used to infer a total block
243 haplotype, which is represented by the possible diploid genotypes of the anchor SNP. For example,
244 consider an anchor SNPs $S_1$ with major allele $A$ and minor allele $B$, and a linked SNP $S_2$ with major allele
245 $C$ and minor allele $D$. Then any observation of allele $C$ at SNP $S_2$ is taken as evidence of allele $A$ at $S_1$,
246 and any observations of allele $D$ at $S_2$ is taken as evidence of allele $B$ at $S_2$. Using this strategy, evidence
247 across all SNPs within a block can be used to infer a total block haplotype, which can be represented by
248 the 3 possible diploid genotypes of the anchor SNP. That is, for an anchor SNP with major allele $A$ and
249 minor allele $B$, the possible block haplotypes are $AA, AB$, and $BB$, with prior probabilities dependent on
250 the allele frequencies of $A$ and $B$.

**Limitations of LOD calculation**

252 Though the magnitude of the LOD score reflects greater genotyping confidence, it cannot be directly
253 interpreted as a likelihood ratio (e.g. an LOD of 200 does not correspond to a $10^{200}$ probability of a
254 shared vs. different genetic origin), as the model does not fully account for sequencing noise, data
255 quality, contamination, and relatedness. In addition, we did not model the incomplete dependence
256 between haplotype blocks, nor the incomplete dependence of SNPs within blocks.

257 Our framework also assumes that the only two sources of a base are the observed allele or a sequencing
258 error. This assumption can lead to incorrect results in the cases where a sample has particularly noisy
259 data due to pre-sequencing events (such as PCR or FFPE processing), non-conforming LD blocks, or high
260 contamination. These samples could be genotyped as heterozygous due to the noisy region or the non-
261 confirming LD block structure. Including these error modes into the model would increase robustness
262 and accuracy.

**Implementation Details**

264 Crosscheck is implemented as part of the Picard-Tools suite, a set of Java command line tools for
265 manipulating high-throughput sequencing data. It accepts VCF/BAM/SAM formatted inputs and can
266 perform comparisons at the level of samples, libraries, read-groups, or files. Crosscheck is provided
267 alongside a utility called ExtractFingerprints which for an input bam, outputs a VCF containing the
268 genotypes and genotype likelihoods across all LD blocks within the supplied haplotype map. This VCF
269 can be used to store fingerprints for downstream analyses or for use with Crosscheck. More information
270 is available at https://github.com/broadinstitute/picard

**Runtime and Memory requirements**

272 For BAM mode, running Crosscheck requires approximately 2.5 gb RAM for a single input pair of BAMs.
273 Runtime is dependent on the size of the input file. Based on our benchmarking experiments, runtimes
274 are < 10 minutes for DNAse-seq, < 30s for ChIP-seq, and are on average about 2 hours for RNA-seq
275 datasets. For VCF mode, Crosscheck requires approximately 2.5 gb of ram for a single pair of inputs, with
276 runtimes < 30s using the standard hg19 haplotype map. CrosscheckFingerprints is multi-threading
277 enabled in order to speed up comparisons and fingerprint generation when multiple input pairs are
278 provided. All comparisons were conducted on Intel(R) Xeon(R) CPU E5-2680 v2 @ 2.80GHz processors.

**Map construction overview**

Maps are constructed from 1000 Genomes[11] phase 3(1000GP3) single-nucleotide polymorphisms(SNPs) which are bi-allelic, phased, and have a minor allele frequency(MAF) $\geq 10\%$. This MAF threshold is introduced since the expected maximal LOD contribution is obtained at an allele frequency of 0.50 (intuitively, rare variants are unlikely to be present in either of two samples being compared from different individuals).  Additionally, SNPs must not differ in their MAF by more than 10% between the 5 ancestral sub-populations(AFR, SAS, EAS, EUR, AMR) present in 1000GP3. This is to correct for potential sub-population bias due to differing linkage and MAF frequency of SNPs across different populations. Using PLINK2[18], we pruned SNPs meeting these criteria in order to create an independent set of "anchor'' SNPs, between which no pairwise $r^2$ correlation exceeded a threshold of 0.10.  A window size of 10 kilobases(kb) and a slide of 5 SNPs was used for pruning. By creating this set of independent SNPs, we ensure that individual haplotype blocks are independent from each other. Next, we greedily added SNPs to the blocks represented by the anchor SNPs. Adding was done in order of LDScore[19] of the anchor SNPs, with the highest LDScoring anchor SNP first( LDScore is the sum for the $r^2$ correlations of each SNP with all other SNPs within a 1 centimorgan window on either side).  Recombination maps containing mappings between genomic coordinates and recombination rates for both the hg19 and GRCh38 assemblies were obtained from http://bochet.gcc.biostat.washington.edu/beagle/ genetic_maps/ and http://mathgen.stats.ox.ac.uk/impute/1000GP_Phase3/.  We only added SNPs if their correlation with the anchor SNP has $r^2 \geq 0.85$ and they were located within a genomic window of 10,000 kb.  In this way, we prioritize the creation of larger, more genetically informative blocks that span several kb regions. The haplotype maps used for the ENCODE database analysis and benchmarking, along with the python code used to generate them, are available at: https://github.com/naumanjaved/ fingerprint_maps.

**Constructing maps only containing LD blocks or single SNPs**

The map containing only single SNP blocks was constructed by aggregating all SNPs in the full haplotype map not in strong linkage($r^2 \geq 0.85$) to other SNPs, resulting in 20792 SNPs. To construct the map containing only blocks with size $\geq 2$ used to quantify the benefits of accounting for linkage, we sub-sampled the full haplotype map. Starting with the largest blocks by number of SNPs, blocks were successively added to this map until the total number of SNPs approximately reached the number of SNPs in the map containing only independent SNPs (20801).

**Testing set construction**

*279 ChIP-seq, RNA-seq, and DNase-seq datasets with ground-truth annotation*

To create a testing set of files to evaluate our method's performance and benchmark it against other tools, we downloaded 279 hg19 bams from RNA-seq, DNase-seq, and ChIP-seq (targeting histone modifications, CTCF, or POL2) from the ENCODE Tissue Expression (ENTEX) project. The ENTEX project contains datasets from experiments on samples derived from four different tissue donors, each of which has whole genome sequencing (WGS) data available. The WGS data for each donor can be used to verify the nominal donor of each dataset comprising the testing set. For each dataset, the corresponding hg38 alignments were compared to the hg38 WGS alignments for its nominal donor. Only datasets which yielded a positive LOD score > 5 using  CrosscheckFingerprints (with the full hg38 version of haplotype

9

319    map) and a "match" result from NGSCheckMate were included in the testing set. The final testing set of
320    files and accompanying metadata are included in supplementary table S1.

321    *98 transcription factor and chromatin modifier (CM) ChIP-seq datasets without ground-truth annotation*

322    To test Crosscheck and other methods on transcription factor and chromatin modifier datasets, we
323    downloaded 98 hg19 ChIP-seq datasets from the ENCODE project. For these datasets, there was no
324    ground-truth donor sequencing data available for the nominal donor as there was for the ENTEX
325    datasets.  In this case, the false-mismatch rate (incorrect genotyping call for a donor-matched pair)
326    cannot be assessed, since there is a non-negligible probability that one of the two datasets with the
327    same nominal donor annotation is incorrectly annotated. However, the false-match rate can still be
328    assessed, since we estimate that the probability that two datasets with different donor annotations may
329    actually share the same true donor is very low. Therefore, we only characterized the ability of
330    NGSCheckmate and Crosscheckfingerprint's to accurate classify donor-unmatched pairs for this testing
331    set. In the context of detecting sample swaps, this performance measure is also more relevant than the
332    accurate detection of donor-matched datasets. All datasets and accompanying metadata is available in
333    supplementary table S2.

334    **BAM pre-processing and down-sampling for benchmarking experiments**

335    Datasets were sorted using Samtools[20] and processed using Picard's MarkDuplicates tool with default
336    settings to remove duplicates. We noted that collapsing duplicates was especially important for RNA-seq
337    datasets since PCR bias can alter allele fractions and lead to incorrect sample classification.
338    Downsampling was conducted on the duplicate marked, sorted files using the command *samtools view –*
339    *s seed.F* with a seed value of 5.

340    **Benchmarking with NGSC and Crosscheck**

341    To speed up analysis of a large number of bams with NGSCheckmate, we followed the author
342    recommendations[10] and created VCFs for each input file using the default provided SNP panel from the
343    NGSCheckMate github and the command *samtools mpileup-I -uf hg19.fasta -l*
344    *SNP_GRCh37_hg19_woChr.bed sample.bam | bcftools call -c - > ./sample.vcf.* NGSC was then run in
345    batch mode using default settings with the hg19 reference SNP panel. For Crosscheck, we first used
346    Picard's ExtractFingerprint utility with default settings and the standard hg19 haplotype map to pre-
347    compute VCFs for each input bam. Comparisons were then conducted using Crosscheck's batch mode
348    with default settings and the standard hg19 map.

349    **Evaluation of other methods that assess genetic similarity between samples**

350    We considered the following methods:

351    • **HYSIS** is intended for tumor-normal concordance verification with a priori knowledge of
352        homozygous germline mutations in the normal tissue[6]. Without considerable modifications, HYSIS is
353        therefore not suitable to handle the general use case that Crosscheck is intended for.
354    • **Bam-matcher** is a tool intended for verifying genotype concordance for whole-genome sequencing,
355        whole exome sequencing, and RNA-sequencing data[7]. Bam-matcher calls programs such as GATK[21]
356        to call variants for each input BAM. Though the resulting variants can be cached to speed up future
357        comparisons, we did not find a way to easily call and store variants for each input bam in the testing

10

358    set, and without that, performing the hundreds of thousands of benchmarking comparisons
359    becomes unfeasible.

360    • We did apply the tools **Conpair** and **BAMixChecker** to the testing set. Conpair was run with default
361    settings using the standard hg19 SNP panel and the *–min-cov* parameter set to 1. Pileups were pre-
362    generated using GATK 4.1.7.0 with the recommended settings[8]. BAMixChecker was run with
363    standard settings for hg19[9] and using GATK 4.1.6.0 for variant calling. Conpair outputs a genotype
364    concordance percentage, which should be <50% for different donor and above 80% for same donor
365    datasets. Any genotype concordance between 50 and 80% is considered inconclusive.
366    BAMixChecker outputs a concordance score between 0 and 1 with no explicit inconclusive range.
367    However, we found that BAMixChecker outputs a concordance score of exactly 0 when there is no
368    overlap between the SNP reference panel that the program uses and the input dataset. Therefore,
369    we labeled any result from BAMixChecker with a concordance score of 0 as an inconclusive
370    genotype call. We found that both methods were unable to yield a conclusive result for more than
371    25% of the comparisons even when the full datasets are considered, and the inconclusive rates
372    became even higher at the lower subsampling rates (Supplementary Fig. 1d). We reasoned that this
373    was likely due to poor overlap between the input datasets and the predefined reference panel of
374    SNPs that both methods use.

375

376    **Familial dataset acquisition and processing**

377    Paired fastqs for RNA-seq data from CEPH/Utah Pedigree 1463 were downloaded from the Gene
378    Expression Omnibus[22] (accession GSE56961). Datasets for the following accessions were downloaded:
379    SRR8505344, SRR8505340, SRR8505343, SRR1258219, SRR1258220, SRR1258218, and SRR8505347.
380    Fastqs were aligned to the hg38 reference using STAR[23] 2.6.0c with default parameters. Before analysis,
381    bams were sorted using samtools and duplicate marked/collapsed using Picard's MarkDuplicates. All
382    comparisons were conducted using the default settings and SNP panels for each method.

383    **ENCODE data acquisition**

384    ENCODE metadata was downloaded from https://www.encodeproject.org/. Metadata was filtered to
385    yield accessions for hg19 ChIP-, RNA-, and DNase-seq ENCODE bams from donors with at-least four
386    datasets. These bams were downloaded from a Broad google bucket and processed(see below)with a
387    custom Workflow Description Language[24](WDL) script. All dataset accessions and associated metadata
388    are available in supplementary table S3.

389    **ENCODE data processing**

390    Files were first sorted using *samtools sort*, and filtered using BEDTools[25] in order to only keep reads
391    overlapping SNPs in the haplotype map. This facilitated efficient storage of files, resulting in
392    approximate 10-fold reduction in file size. Finally, duplicates were marked and removed for each file
393    using Picard's MarkDuplicates function with default settings. All comparisons were conducted using the
394    version of CrosscheckFingerprints available in commit #078b0ba of Picard.

395    **ENCODE genotyping strategy**

396    *Construction of reference set*

397    To detect mislabeled samples, each dataset is compared against a reference set of 3 samples that
398    provide a high quality representation of the "true'' genotype for each ENCODE tissue donor. To

399   construct this reference set of samples, a self-LOD score is calculated for each sample by "comparing"
400   each file to itself. This score correlates with the dataset's overlap with the haplotype map, and the
401   highest self-LOD samples are those containing the most genetic information relevant for genotyping.  To
402   ensure that the reference set of samples for each tissue donor does not contain any swapped samples,
403   all reference samples are compared against one another to ensure self-consistency, which is defined as
404   an LOD score greater than 5 for all three pairwise comparisons between the three samples. In the case
405   of one swapped sample in this reference set, two negative LOD scores and one positive LOD score will
406   be obtained.  In this case, the next highest self-LOD scoring bam replaces the putative swap, and
407   representative concordance is re-assessed. This is repeated until a concordant set is found. More
408   complex patterns of swaps in the representative set are assessed on a case-by-case basis. Finally, all
409   reference samples across all nominal donors are compared against one another in order to identify
410   larger cross-donor swaps and preclude the possibility that all reference samples for a nominal donor are
411   actually swaps from the same true donor.

412   *Comparisons of samples with reference set*

413   Each sample not in the reference set is compared against the top 3 representative samples for its
414   nominal donor. Samples yielding an LOD ≤ -5 against any of the top 3 representatives are flagged as
415   swaps for review, while those yielding an LOD score between -5 and 5 are flagged as inconclusively
416   genotyped.

417   **Contamination tests**

418   Varying numbers of randomly sampled reads from two unrelated ENCODE ChIP-seq datasets,
419   ENCFF005HON ENCFF007DFB, were mixed together to create simulated contaminated datasets. Each
420   mixed sample consisted of ~ 5 million reads and contained varying proportions of the original datasets
421   (at intervals of 10%). Mixed samples were then compared to ENCFF007NTA and ENCFF029GAR, which
422   are ChIP-seq datasets from the same donor as ENCFF005HON. Comparisons were conducted on VCF files
423   generated using Picard's ExtractFingerprint utility using Crosscheck's VCF mode with default settings.


424   **Acknowledgements**

428   **Contributions**

429   N.J. constructed the haplotype map. Y.F. and T.F. designed and wrote CrosscheckFingerprints. N.J., N.S.,
430   and Y.F. designed the experimental setup. N.J. and N.S. performed the analyses. C.E. and N.S. verified
431   genotyping findings. N.J., Y.F., B.E.B., and N.S. wrote the paper.

432   **Competing Interests**

433   The authors declare no competing financial interests.


434

**Data availability**

All data used for benchmarking and ENCODE analysis are available online at https://encodedcc.org/. Specific accessions and relevant metadata for each of the benchmarking experiments are available in tables S1 and S2. Accession IDs and metadata for all datasets from ENCODE analysis are available in table S3. Haplotype maps used for benchmarking and ENCODE analysis are available at https://github.com/naumanjaved/fingerprint_maps). RNA-seq data from CEPH/Utah Pedigree 1463 were downloaded from the Gene Expression Omnibus (https://www.ncbi.nlm.nih.gov/geo/, accession GSE56961).


**Code availability**

Crosscheck code and documentation is available at https://github.com/broadinstitute/picard. Fingerprint map generation code, along with pre-compiled maps and documentation are available at https://github.com/naumanjaved/fingerprint_maps.

**References**

1        Horbach, S. P. J. M. & Halffman, W. The ghosts of HeLa: How cell line misidentification contaminates the scientific literature. PLOS ONE 12, e0186281, doi:10.1371/journal.pone.0186281 (2017).

2        Lorsch, J. R., Collins, F. S. & Lippincott-Schwartz, J. Fixing problems with cell lines. Science (New York, N.y.) 346, 1452-1453, doi:10.1126/science.1259110 (2014).

3        Biankin, A. V., Piantadosi, S. & Hollingsworth, S. J. Patient-centric trials for therapeutic development in precision oncology. Nature 526, 361-370, doi:10.1038/nature15819 (2015).

4        Danecek, P. et al. The variant call format and VCFtools. Bioinformatics 27, 2156-2158, doi:10.1093/bioinformatics/btr330 (2011).

5        Pengelly, R. J. et al. A SNP profiling panel for sample tracking in whole-exome sequencing studies. Genome Medicine 5, 89, doi:10.1186/gm492 (2013).

6        Schröder, J., Corbin, V. & Papenfuss, A. T. HYSYS: have you swapped your samples? Bioinformatics 33, 596-598, doi:10.1093/bioinformatics/btw685 (2017).

7        Wang, P. P. S., Parker, W. T., Branford, S. & Schreiber, A. W. BAM-matcher: a tool for rapid NGS sample matching. Bioinformatics 32, 2699-2701, doi:10.1093/bioinformatics/btw239 (2016).

8        Bergmann, E. A., Chen, B.-J., Arora, K., Vacic, V. & Zody, M. C. Conpair: concordance and contamination estimator for matched tumor–normal pairs. Bioinformatics 32, 3196-3198, doi:10.1093/bioinformatics/btw389 (2016).

9        Chun, H. & Kim, S. BAMixChecker: an automated checkup tool for matched sample pairs in NGS cohort. Bioinformatics 35, 4806-4808, doi:10.1093/bioinformatics/btz479 (2019).

10        Lee, S. et al. NGSCheckMate: software for validating sample identity in next-generation sequencing studies within and across data types. Nucleic Acids Research 45, e103, doi:10.1093/nar/gkx193 (2017).

11        A global reference for human genetic variation. Nature 526, 68-74, doi:10.1038/nature15393 (2015).

12        An Integrated Encyclopedia of DNA Elements in the Human Genome. Nature 489, 57-74, doi:10.1038/nature11247 (2012).

13        Regev, A. et al. The Human Cell Atlas. eLife 6, e27041, doi:10.7554/eLife.27041 (2017).

14        Network, C. G. A. R. et al. The Cancer Genome Atlas Pan-Cancer analysis project. Nature Genetics 45, 1113-1120, doi:10.1038/ng.2764 (2013).

15        Lonsdale, J. et al. The Genotype-Tissue Expression (GTEx) project. Nature Genetics 45, 580-585, doi:10.1038/ng.2653 (2013).

16        Slatkin, M. Linkage disequilibrium--understanding the evolutionary past and mapping the medical future. Nature Reviews. Genetics 9, 477-485, doi:10.1038/nrg2361 (2008).

17      Pfeiffer, F. et al. Systematic evaluation of error rates and causes in short samples in next-generation sequencing. Scientific Reports 8, 1-14, doi:10.1038/s41598-018-29325-6 (2018).

18      Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. GigaScience 4, doi:10.1186/s13742-015-0047-8 (2015).

19      Bulik-Sullivan, B. K. et al. LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. Nature Genetics 47, 291-295, doi:10.1038/ng.3211 (2015).

20      Li, H. et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics (Oxford, England) 25, 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).

21      Poplin, R. et al. Scaling accurate genetic variant discovery to tens of thousands of samples. bioRxiv, 201178, doi:10.1101/201178 (2018).

22      Barrett, T. et al. NCBI GEO: archive for functional genomics data sets—update. Nucleic Acids Research 41, D991-D995, doi:10.1093/nar/gks1193 (2013).

23      Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. Bioinformatics 29, 15-21, doi:10.1093/bioinformatics/bts635 (2013).

24      Voss, K., Gentry, J. & Auwera, G. V. d. Full-stack genomics pipelining with GATK4 + WDL + Cromwell(not peer reviewed). ISCB Comm J 6, 1381, doi:doi.org/10.7490/f1000research.1114634.1 (2017).

25      Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. Bioinformatics (Oxford, England) 26, 841-842, doi:10.1093/bioinformatics/btq033 (2010).

# Figure 1

**Figure 1: Incorporating Linkage Information allows robust comparison of sequencing datasets**

a) Sample swaps and mis-annotations, where a sample is incorrectly attributed to the wrong donor, are a high stakes issue for large consortium projects and clinical science.

b) Our method compares reads from two datasets across a genome-wide set of linkage disequilibrium LD blocks (haplotype map). The SNPs in each block are highly correlated with each other and have low correlation with SNPs in other blocks. Reads overlapping any of the SNPs in a given block inform the relatedness of the datasets, even when reads from the two datasets do not overlap one another.

c) Haplotype maps contain many large LD blocks. LD blocks are created using common, ancestry independent SNPs from 1000 Genomes. Most SNPs lie within blocks of size > 2, which boosts the chances of reads to be informative.

d) Distribution of LOD scores for 34336 donor-mismatched (red) and 9767 donor-matched pairs (green) of public ChIP-, RNA-, and DNase-seq datasets from the ENCODE project.

e) LD-based method can correctly determine sample relatedness even at low sequencing coverage. Pairwise comparisons of reference dataset pairs at different sub-sampling percentages using two equally sized SNP panels – one panel contained only independent single SNPs, while the other contained only LD blocks. Donor-mismatched dataset pairs are colored red while donor-matched dataset pairs are green.

f) Comparison of NGSC and Crosscheck's classification of 34336 donor-mismatched and 9767 donor-matched dataset pairs. Performance was measured in terms of the false flag rate (FFR), the fraction of donor-matched pairs incorrectly flagged as donor-mismatches, and the false match rate (FMR), the fraction of donor-mismatched pairs incorrectly identified as donor-matches. Comparisons are classified as *same-assay* if the two datasets are from the same assay type, and have the same target epitope in the case of ChIP-seq datasets. All other comparisons are classified as *cross-assay.*
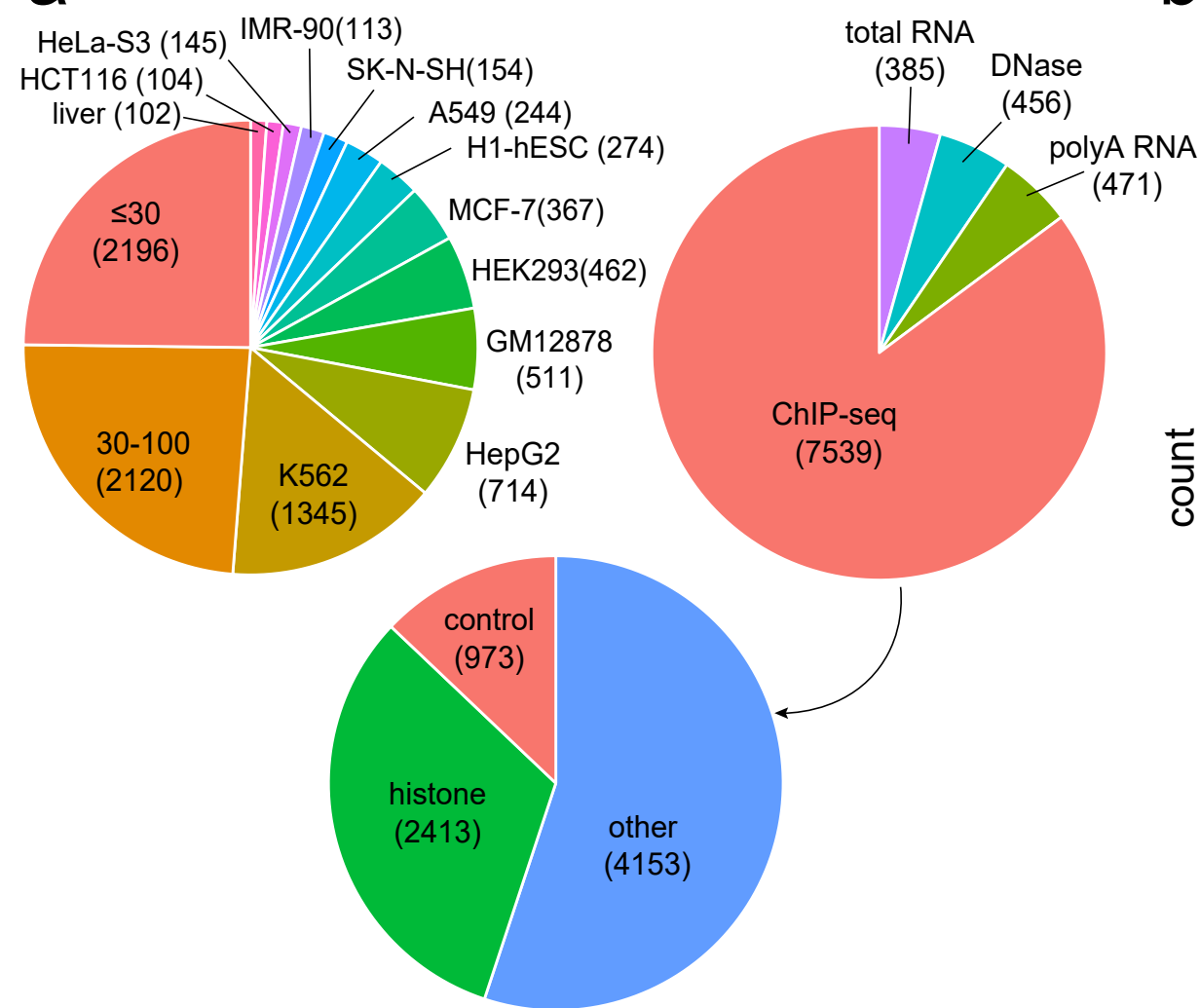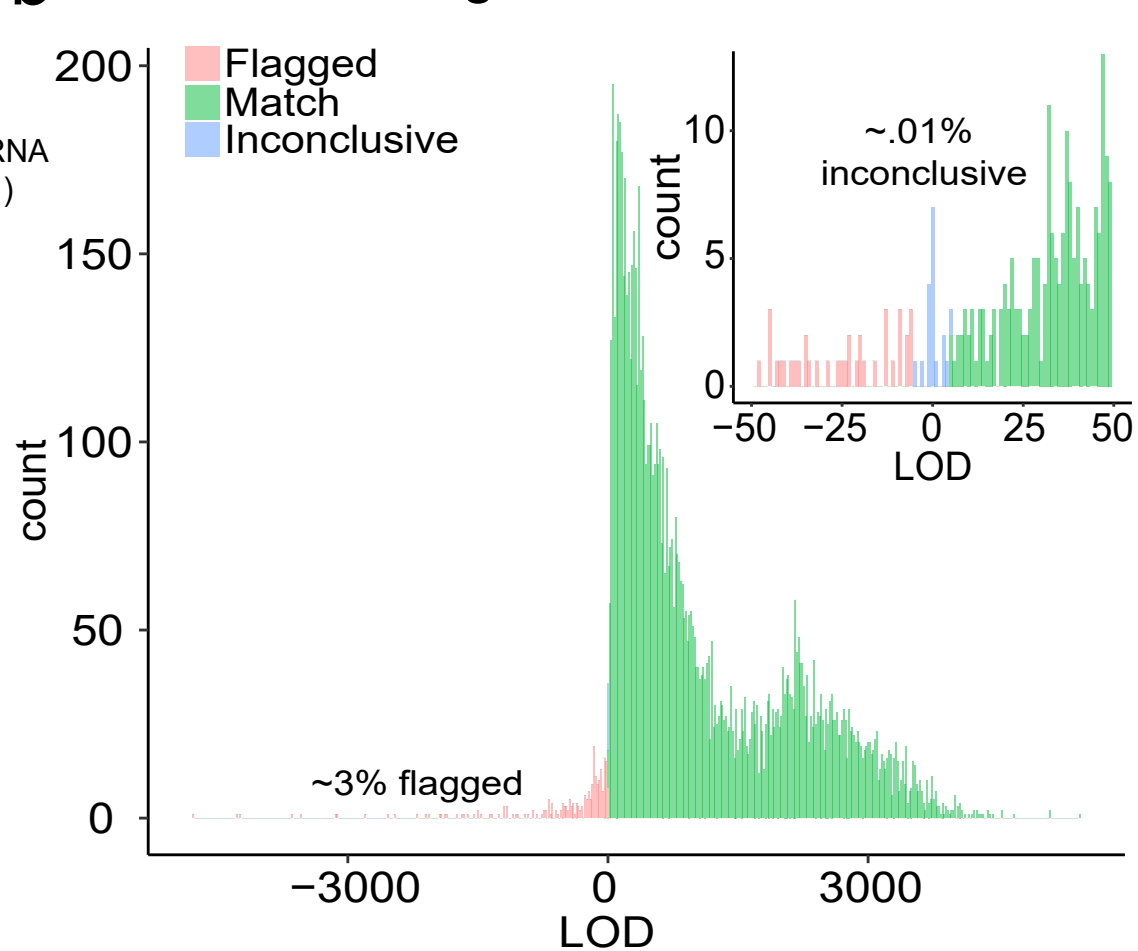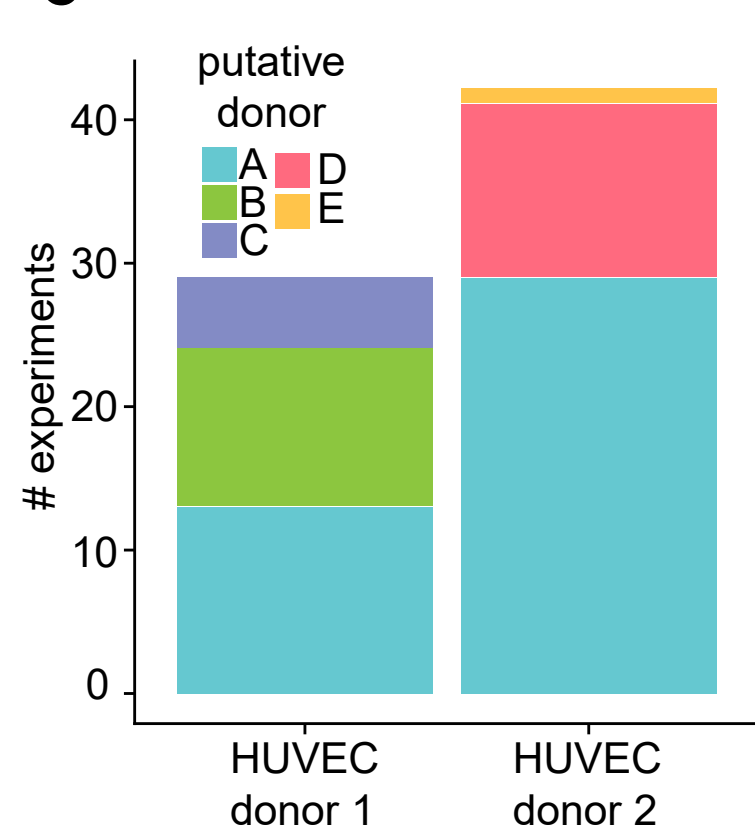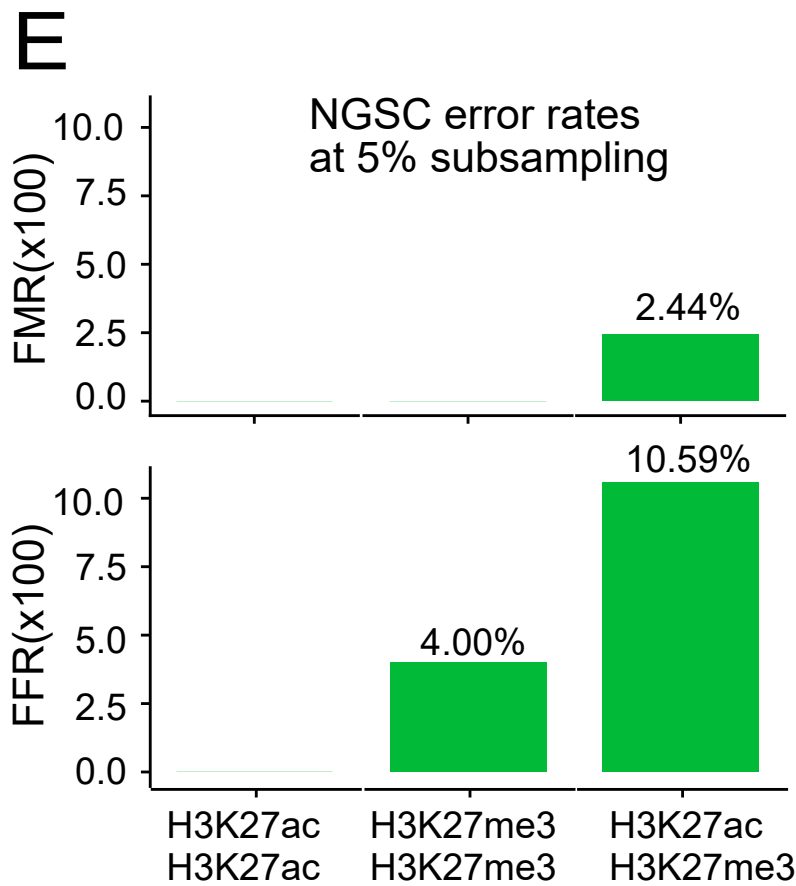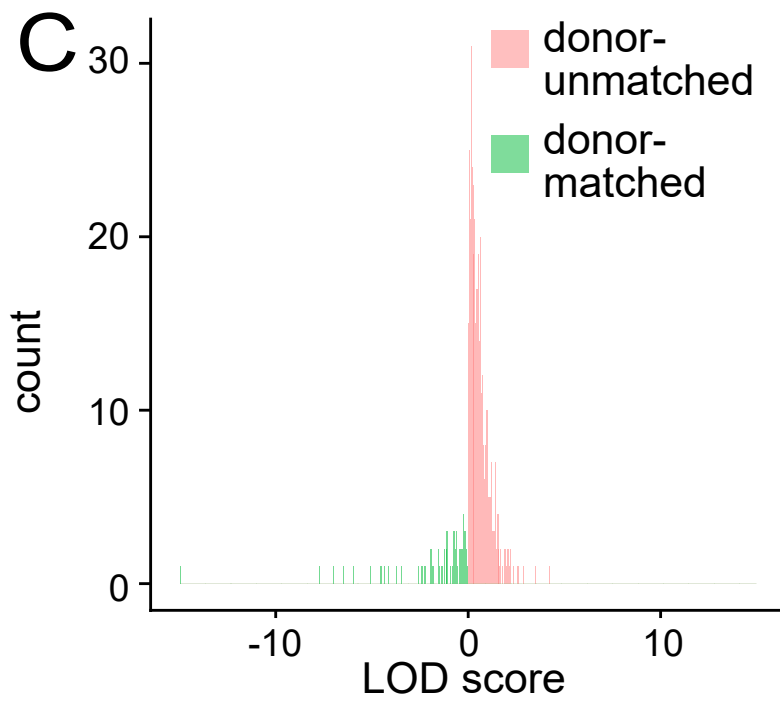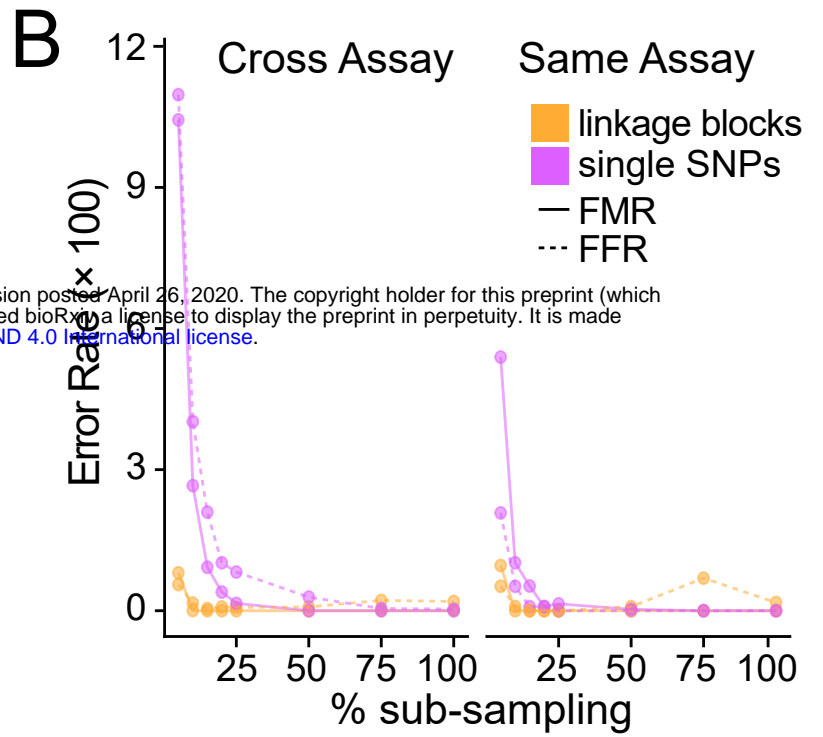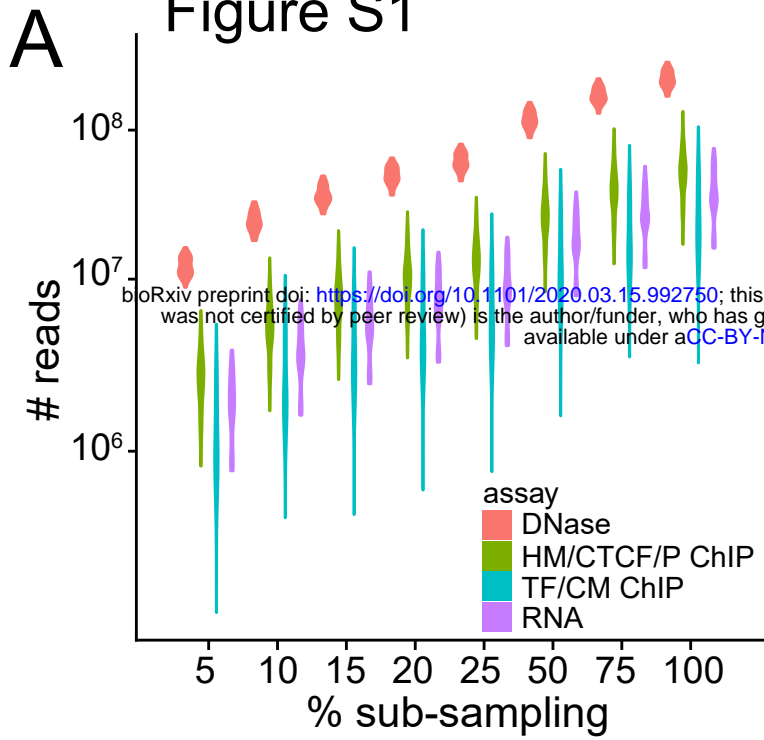
**Figure 2**

**Figure 2: Overview of ENCODE database swap detection**

a) Overview of 8851 genotyped datasets from ENCODE, partitioned by cell type (top left), assay type (top right), and by target for ChIP-seq (bottom). Cell types that had less than 100 datasets derived from them were pooled – so all the datasets from them are grouped into one of two categories. All hg19 aligned reads from total RNA-, polyA RNA-, ChIP-, and DNase-seq experiments performed on samples belonging to donors with at-least four datasets in total were included in the analysis. All ChIP-seq targets, including histone modifications(HM), transcription factors (TF), chromatin modifiers (CM), CTCF, and control experiments were included.

b) Distribution of LOD scores from ENCODE genotyping. Each dataset was compared to three representative datasets from its nominal donor. Any dataset scoring negatively against any of the three representatives was flagged for further review. A comparison resulting in an LOD score between -5 and 5 was deemed inconclusive (insufficient evidence to indicate shared or distinct genetic origin).

c) Each flagged sample was compared to all other samples from its nominal donor, as well as the representatives for all other donors in our database to nominate true donor identity and identify genetically consistent sub-clusters. Comparisons of flagged samples between two HUVEC donors reveals 5 genetically distinct clusters.
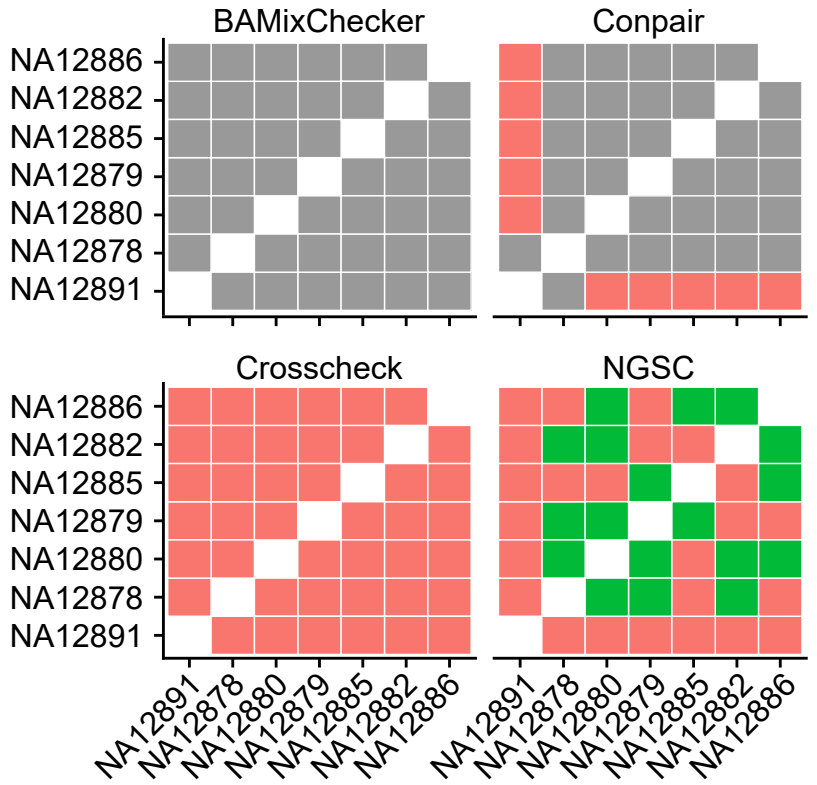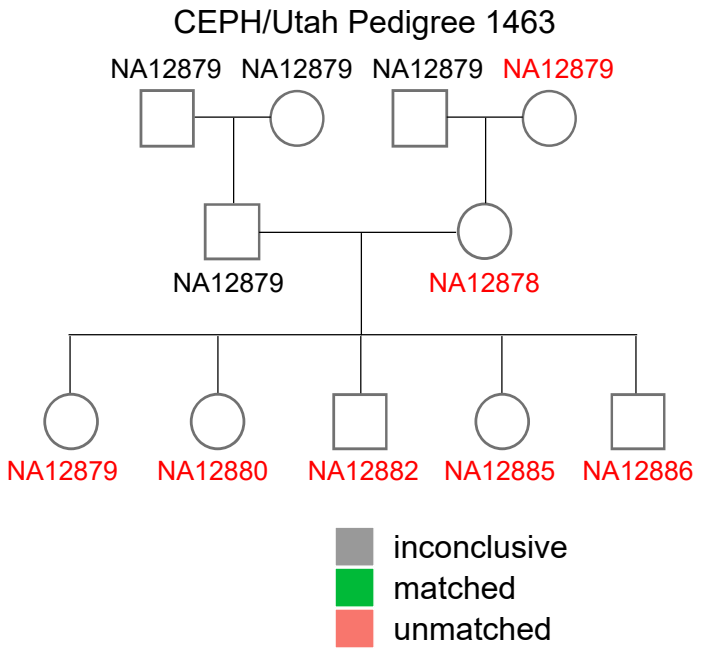
**Supplementary Figure 1**

**(A)** Distribution of number of reads in sub-sampled datasets used for benchmarking, broken down by assay type. ChIP datasets were divided into two classes – those which targeted transcription factor (TF) and chromatin modifier (CM), and those which targeted broad histone modifications (HM), POL2/POL2RA (P), or CTCF.

**(B)** Comparison of percentage false match (FM) and false flag (FF) rates for 9767 same-donor and 29573 different donor pairwise comparisons using CrosscheckFingerprints with either linkage blocks, or single SNPs only. Across different (left) and same (right) assay comparisons, incorporation of linkage information (orange line) decreases the FF and FM percentage, particularly at sub-sampling percentages. Comparisons are classified as *same-assay* if the two datasets are from the same assay type, and have the same target epitope in the case of ChIP-seq datasets. All other comparisons are classified as *cross-assay*.

**(C)** Distribution of LOD scores from false flags and false matches from benchmarking experiments. The distribution of the majority (99%) of LOD scores from these misclassifications is used to create an "inconclusive" range of LOD scores, in which donor-match or mismatch cannot be confidently called.

**(D)** Percent inconclusive genotype concordance calls for 9767 same-donor and 29573 different donor pairwise comparisons using Conpair and BAMixChecker. "Inconclusive" is defined as pairwise comparisons resulting in genotype concordances between 50 and 80% for Conpair, and a score of 0 for BAMixChecker.

**(E)** FMR and FFR for NGSC at 5% subsampling for pairwise comparisons between ChIP-seq datasets targeting the non-overlapping histone modifications H3K27ac and H3K27me3. NGSC performs worse for comparisons between H3K27ac and H3K27me3 datasets (n=41 donor-matched, n=85 donor-mismatched) than for comparisons between two H3K27ac (n=24 donor-matched, n=67) or two H3K27me3 datasets (n=11 donor-matched, n=25 donor-mismatched). In contrast, Crosscheck classifies all pairs correctly.
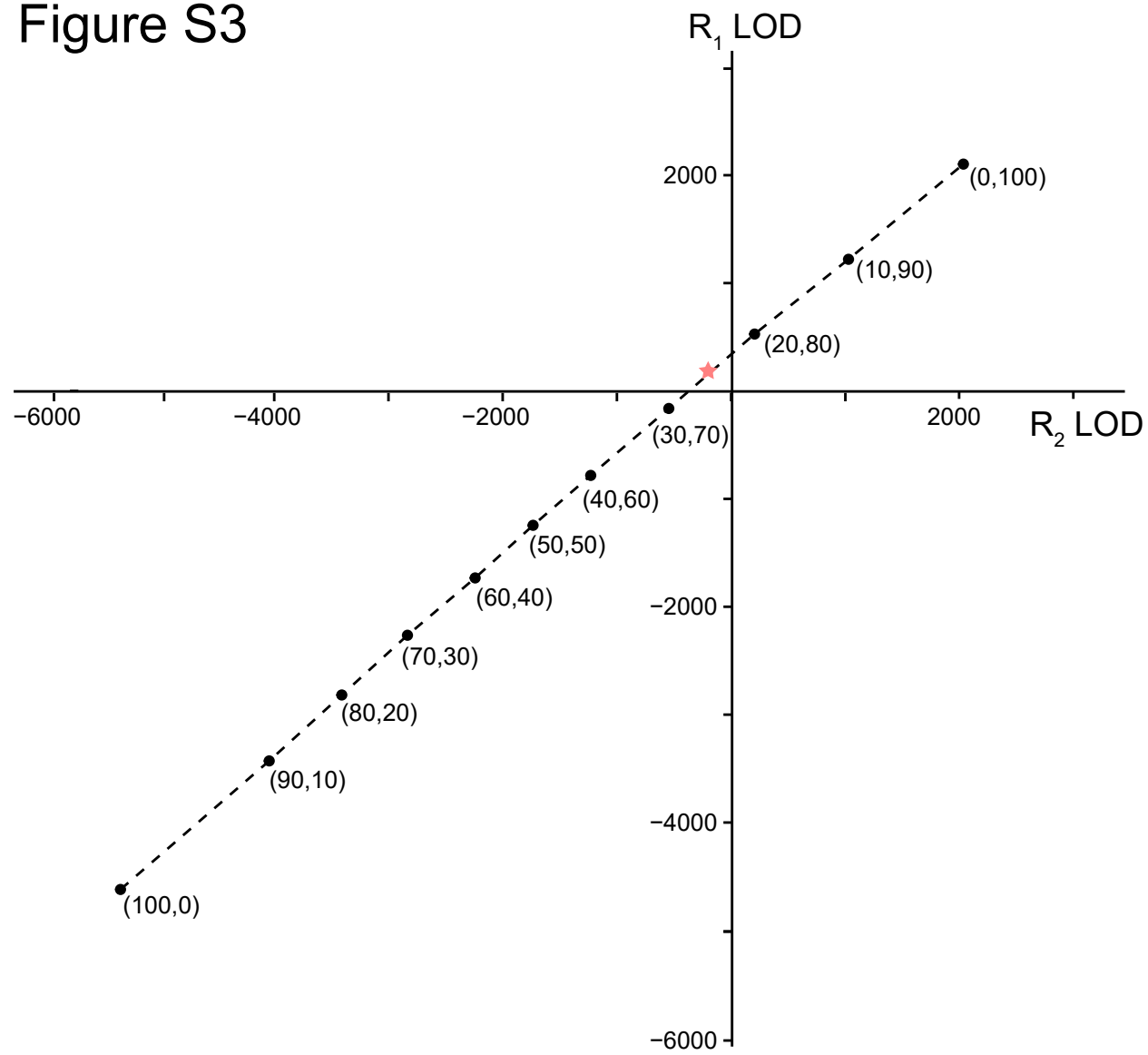
Figure S2

CEPH/Utah Pedigree 1463

**Supplementary Figure 2**

Performance of NGSC, Crosscheck, BAMixChecker, and Conpair when classifying 21 pairwise comparisons between RNA-seq datasets from 7 related individuals (indicated in red) from CEPH/Utah pedigree 1463. "Inconclusive" is defined as pairwise comparisons resulting in genotype concordance between 50 and 80% for Conpair, a score of 0 for BAMixChecker, and an LOD score between -5 and 5 for Crosscheck. NGSC incorrectly classifies 43% of pairs, while Conpair and BAMixChecker are inconclusive for 76 and 100% of pairs respectively. In contrast, Crosscheck correctly classifies all dataset pairs as mismatches.

Figure S3

**Supplementary Figure 3**

Demonstration of Crosscheck's performance for contaminated datasets. Simulated contaminated datasets were created by combining various proportions of two ENCODE ChIP-seq datasets derived from two different donors: ENCFF005HON and ENCFF007DFB. Proportions of reads deriving from ENCFF005HON and ENCFF007DFB respectively are indicated in parentheses for each mixture. Each mixture was compared to two datasets derived from the same donor as ENCFF005HON, ENCFF007NTA ($R_1$) and ENCFF029GAR ($R_2$). The star indicates a region where a contaminated sample can score as a donor match against one dataset ($R_1$), but score as a donor mismatch against a different dataset from the same donor ($R_2$).