

# Analysis of Information Leakage in Phenotype and Genotype Datasets

Arif Harmanci<sup>1,2</sup>, Mark Gerstein<sup>1,2,3</sup>

1 Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA  
2 Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, USA  
3 Department of Computer Science, Yale University, New Haven, CT, USA  
Corresponding author: Mark Gerstein [mark.gerstein@gersteinlab.org](mailto:mark.gerstein@gersteinlab.org)

[[Title]]  
[[The Figure/Section References]]  
[[Use past tense, all around]]  
[[At most 30 equations in the online methods?]]  
[[Methods only references?]]  
[[Number of references: 10, 20, 30, 40]]  
[[Figure Captions]]  
[[Read one more time]]  
[[PDF Figure: <http://www.jstor.org/stable/4147411>]]  
[[Check HTML rendering in the browser]]

Formatted: Font: 3 pt

Deleted: 26

Formatted: Font: 3 pt

Formatted: Font: 3 pt

Formatted: Font: 3 pt

Deleted: [[1000 => 1,000]]¶

## ABSTRACT

Privacy is receiving much attention with the increase in the breadth and depth of personalized biomedical datasets. Studies on genomic privacy are mainly focused on protection of variants. Molecular phenotype datasets can also contain substantial amount of sensitive information. Although there is no explicit genotypic information in them, subtle genotype-phenotype correlations can be used to statistically link the phenotype and genotype. The links can then be used to characterize individuals' sensitive phenotypes. Here, we first develop a formalism for the quantification of individual characterizing information leakage in a linking attack. We analyze the tradeoff between the predictability of the genotypes and the amount of leaked information that can be used for individual characterization. Then we present a general three step procedure that can be used to practically instantiate an accurate attack. We develop a particular realization of the attack for outlier cases and study different aspects of the attack.

Deleted: datasets

Deleted: a highly

## 1 INTRODUCTION

Genomics has recently emerged as one of the major foci of studies on privacy. This can be attributed to high throughput biomedical data acquisition that brings about a surge of datasets<sup>1-3</sup>. Among these, molecular phenotype datasets, like functional genomics measurements, substantially grow the list of the quasi-identifiers<sup>4</sup> which may lead to re-identification and characterization<sup>4-6</sup>. In general, statistical analysis methods are used to discover genotype-phenotype correlations<sup>7,8</sup>, which can be utilized by an

Deleted: bring

Deleted: 1,2.

Deleted: 3

Deleted: 3-5.

Deleted: , which can be utilized by an adversary for linking the entries in genotype and phenotype datasets, and revealing sensitive information.

HOWEVER

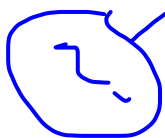
adversary for linking the entries in genotype and phenotype datasets, and revealing sensitive information. The availability of a large number of correlations increases the possibility of linking<sup>9,10</sup>.

has

Along with the initial genotype-phenotype association studies, the protection of privacy of participating individuals emerged as an important issue. Several studies addressed the problem of detecting whether an individual, with known genotype, participated in a study<sup>11</sup>. As study participants choose to remain anonymous, the detection of an individual causes privacy concern<sup>12-15</sup> by revealing their existence in the study cohort. We refer to these systematic breaches as "detection of a genome in a mixture" attacks (Supplementary Fig. 1). As the number and size of phenotype and genotype datasets increase, the detection of individuals in these datasets will be irrelevant since any individual will already have their genotype or phenotype information stored in a dataset, i.e., participation will already be known. Consequently, an adversary can then aim at pinpointing an individual among multiple, seemingly independent, genotype and phenotype datasets by linking the entries in these datasets. As personalized genomics gain more prominence, e.g. large genotype and phenotype datasets are used in medicine, the attackers will focus on gaining access to these datasets, then aim at linking different datasets that can reveal sensitive information. We will refer to these attacks as "linking attacks"<sup>4,5</sup>. One well-known example of these attacks is the attack that matched the entries in Netflix Prize Database and the Internet Movie Database (IMDb)<sup>16</sup>, and revealed sensitive information. For research purposes, Netflix released an anonymized dataset of movie ratings of thousands of viewers, which was assumed to be secure as the viewer's names were removed. However, Narayanan et al used the Internet Movie Database (IMDb), a seemingly unrelated and very large database of movie viewers, linked two databases, and revealed identities and personal information of many viewers in the Netflix database. This attack is underpinned by the fact that both Netflix and IMDb host millions of individuals and any individual who is in one dataset is very likely to be in the other dataset. As the size and number of the genotype and phenotype datasets increase, number of potentially linkable datasets will increase, which can render similar scenarios a reality in genomic privacy (Supplementary Note). Different aspects of genomic privacy, pertaining linkability of high dimensional phenotype datasets to genotypes, are yet to be explored.

2.2.1 THREAT ANALYSIS

STOLEN



## 2 RESULTS

### 2.1 Linking Attack Scenario

In the linking attacks, the attacker aims at characterizing sensitive information about a set of individuals in a genotype dataset (Fig. 1). For each individual in the genotype dataset, she aims at querying the publicly available anonymized phenotype datasets in order to characterize their sensitive phenotypes. For this, she first utilizes a public quantitative trait loci (QTL) dataset that contains phenotype-genotype correlations. She statistically predicts genotypes using the phenotypes and QTLs. Then she compares the predicted genotypes to the genotype dataset and links the entries that have good genotype concordance. Consequently, the sensitive information for the linked individuals in genotype dataset is revealed to the attacker.

Deleted: <sup>6,7</sup>

Field Code Changed

Deleted: With the initial genotype-phenotype studies, the protection of privacy of participating individuals emerged as an important issue. Several studies addressed "detection of a genome in a mixture" attacks, which can potentially detect participating individuals<sup>89</sup> (Supplementary Fig. 1). With large sets of seemingly independent datasets, the attacker can also link two or more datasets to pinpoint individuals in datasets and reveal sensitive information. One well-known example of these "linking attacks" is the attack that matched the entries in Netflix Prize Database and the Internet Movie Database (IMDb)<sup>10</sup>, and revealed sensitive information. As the size and number of the genotype and phenotype datasets increase, number of potentially linkable datasets will increase, which can make similar scenarios a reality in genomic privacy (Supplementary Note). Several studies addressed multiple scenarios for individual re-identification and genome. In addition, different formalisms are proposed for protection of sensitive information (Supplementary Note). Different aspects of genomic privacy, pertaining linkability of high dimensional phenotype datasets to genotypes, are yet to be explored. ¶ In this paper, we focus on characterizability of the individuals' sensitive information in the context of linking attacks, where the adversary exploits the genotype-phenotype correlations to link datasets and reveal sensitive information. In general, the high dimensional phenotype datasets harbor a number of phenotypes that contain sensitive information, like disease status, and other phenotypes, while not sensitive, may be used for linking to genotypes. Many public quantitative trait loci (QTL) datasets will inevitably help these linkings. Although each QTL has <sup>11-13</sup>. First we study quantification of characterizing information leakage versus risk of characterization. Then, we present a three step analysis framework for analysis of linking attacks. Then we show that the linking can be performed in a much simplified setting by just utilizing the outliers in the data. Compared to a previous publication<sup>14</sup> that relates to our study, we aim at providing quantification measures and privacy analysis frameworks and also demonstrate the accuracy of simplified linking attack under different scenarios. We finally discuss different strategies for risk management against the linking attacks. ¶

Deleted: Individual Characterization by

Deleted: Attacks

Deleted: There are three datasets in

Deleted: context of the breach by

Deleted: (Fig. 1). First dataset contains

Deleted: phenotype

Deleted: for

Deleted: . The

Deleted: can include sensitive information such as disease status in addition to several molecular

Deleted: such as gene expression levels. The second dataset contains the genotypes

Deleted: the identities for another set of

Deleted: . The third dataset contains correlations between one or more of the phenotypes in the phenotype dataset and the genotypes. Each entry contains a phenotype, a variant, and the ...

Among the QTL datasets, the abundance of eQTL datasets makes them most suitable for linking attacks. In an eQTL dataset, each entry contains a gene, a variant, and correlation coefficient, denoted by  $\rho$ , between the expression levels and genotypes. We assume that the attacker aims to build a genotype prediction model that utilizes the relation between expression levels and genotypes (Fig. 2a, Supplementary Fig. 2). As a representative dataset for reporting results and for performing mock linking attacks, we use the eQTLs and gene expression levels from the GEUVADIS project<sup>17</sup>, and the genotypes from the 1000 Genomes Project<sup>18</sup>.

## 2.2 Genotype Predictability and Information Leakage

We assume that the attacker will behave in a way that maximizes his or her chances of correctly characterizing the most number of individuals. Thus, she will try and predict the genotypes, using the phenotype measurements, for the largest set of variants that she believes she can predict correctly. The most obvious way that the attacker does this is by first sorting the genotype-phenotype pairs with respect to decreasing strength of correlation then predicting the genotypes for each variant (Supplementary Fig. 3). The attacker will encounter a tradeoff: As she goes down the list, more individuals can be characterized (more genotypes can characterize more individuals), but it also becomes more likely that she makes an error in the prediction, since the correlation decreases going down the list. This tradeoff can also be viewed as the tradeoff between precision (fraction of the linkings that are correct) and recall (fraction of individuals that are correctly linked). We will propose two measures, individual characterizing information ( $ICI$ ) and genotype predictability ( $\pi$ ), to study this tradeoff.

$ICI$  can be interpreted as the total amount of information in a set of variant genotypes that can be used to pinpoint an individual in a linking attack. This quantity depends on the joint frequency of the variant genotypes. For example, if the set contains many common genotypes, they will not be very useful for pinpointing individuals. On the other hand, rare variant genotypes would give much information for linking. Thus, the information content of a set of genotypes is inversely proportional to the joint frequency of the variant genotypes. We utilize this property to quantify  $ICI$  in terms of genotype frequencies (Online Methods, Fig. 3).

For a set of variants,  $\pi$  measures how predictable genotypes are given the gene expression levels. Since genotypes and expression levels are correlated, knowledge of the expression enables one to predict the genotype more accurately than predicting genotype with no knowledge. In order to quantify the predictability, we use an information theoretic measure for randomness left in genotypes, given gene expression levels (Online Methods, Fig. 3). Although the reported correlation coefficient is a measure of predictability, it is computed differently in different studies and there is no easy way to combine and interpret the correlation coefficients when we would like to estimate the joint predictability of multiple eQTL genotypes.

We first considered each eQTL and evaluated the genotype predictability versus the characterizing information leakage. We use the GEUVADIS dataset as a representative dataset for this computation. We computed, for each eQTL, average  $\pi$  and average  $ICI$  over all the individuals (Fig. 4a). Most of the data points are spread along the anti-diagonal: The eQTL variants with high major allele frequencies have high predictability and low  $ICI$  and vice versa for eQTL variants with lower major allele frequencies.

*THIS IS A PAPER*

**Deleted:** as the representative phenotype dataset. As explained earlier

**Deleted:** gene expression-genotype correlation (

**Deleted:** )

**Deleted:** these datasets

**Moved (insertion) [1]**

**Deleted:** In an eQTL dataset, genes and variants and a correlation coefficient, denoted by  $\rho$ , between the expression levels and genotypes are reported. The absolute value of  $\rho$  indicates the strength of association between the eQTL genotype and the eQTL expression level. The sign of  $\rho$  represents the direction of association, i.e., which homozygous genotype corresponds to higher expression levels. We assume that the attacker captures this relation with a prediction model and builds the *a posteriori* distribution of the eQTL genotypes given the eQTL expression levels (Fig.

**Moved up [1]:** 2a, Supplementary Fig. 2).

**Deleted:** As a representative dataset for reporting results, we utilize the eQTLs, and gene expression levels from the GEUVADIS project<sup>15</sup> and the genotypes from the 1000 Genomes Project<sup>16</sup> Genotypes Predictability and Information Leakage¶

**Deleted:** /

**Deleted:** he or

**Deleted:** he or

**Deleted:** he or

**Deleted:** he or

**Deleted:** he or

**Deleted:** In this section we

**Deleted:** quantify

**Deleted:** The attacker aims to correctly characterize  $n_p$  individuals in the expression dataset among  $n_g$  individuals in the genotype dataset. In order to correctly characterize an individual, given the individual's expression levels, the attacker should predict the genotypes for a set of selected eQTLs correctly such that the predicted set of genotypes are not shared by more than one individual, i.e., the predicted genotypes can be matched to the correct individual. In other words, the joint frequency of the set of predicted genotypes should be  $\frac{1}{n_g}$ . Equivalently, if the genotypes contain at least  $\log_2(n_g)$  bits of information, the individual is vulnerable to characterization of his/her phenotypes. Following this idea, we quantify the leakage of individual characterizing

**Deleted:** (correct)

**Deleted:** of the eQTL

**Deleted:** from

**Deleted:** , it is necessary to uniformly estimate predictability f

**Deleted:**  $\rho$

**Formatted:** Font color: Text 1

**Deleted:** these

**Deleted:** values

**Deleted:** We will utilize the exponential of negative conditiona

(Fig. 4b). This is expected because the genotypes of the high frequency variants can be predicted, on average, easily (most individuals will harbor one dominant genotype) and consequently does not deliver much characterizing information and vice versa for the eQTLs with smaller major frequency alleles. The eQTLs with high absolute correlation (Fig. 4b) deviate from the anti-diagonal, compared to shuffled data (Fig. 4c). These eQTLs have high ICI and high predictability,

When multiple genotypes are utilized, the information leakage is greatly increased. To study this, we computed ICI (in bits) and predictability for increasing number of eQTLs (Supplementary Note, Fig. 4d). As expected, the predictability decreases with increasing ICI leakage. Inspection of mean predictability versus mean ICI enables us to estimate the number of vulnerable individuals at different predictability levels. For example, at 20% predictability, there is approximately 8 bits of ICI leakage. At this level of leakage, the adversary can pinpoint an individual, with 20% accuracy, within a sample of  $2^8 = 256$  individuals. Thus, within any sample of 256 individuals, we expect the attacker to be correctly link  $256 \times 20\% = 51$  individuals. At 5% predictability, the leakage is 11 bits and the attacker can pinpoint an individual in a sample of  $2^{11} = 2048$  individuals. This corresponds to approximately 100 individuals getting correctly linked (5% of 2048). Auxiliary information can be easily added into ICI. For example, gender information, which can be predicted with high accuracy from many molecular phenotype datasets brings 1 bit of additional auxiliary information to ICI (Supplementary Note).

### 2.3 Framework for Instantiation of Linking Attacks

We present a three step framework for practical instantiation of linking attacks (Fig. 2b). This framework can be used to perform mock linking attacks on datasets for evaluating whether they will be effective for risk assessment purposes. We use this framework to simulate mock attacks in the following sections for assessing their accuracies. The input is the phenotype measurements for an individual, who is being queried for a match to individuals in the genotype dataset (Fig. 1). In the first step, the attacker selects the QTLs, which will be used in linking. The selection of QTLs can be based on different criteria. As discussed earlier, the genotype predictability ( $\pi$ ) is the most suitable QTL selection criterion. Although the attacker cannot practically compute predictability using only the QTL list, any function of predictability would still be useful to the attacker for selecting QTLs. For example, the most accessible criterion is selection based on the absolute strength of association,  $|\rho|$ , between the phenotypes and genotypes. The second step is genotype prediction for the selected QTLs using a prediction model. The third and final step of a linking attack is comparison of the predicted genotypes to the genotypes of the individuals in genotype dataset to identify the individual that matches best to the predicted genotypes. In this step, the attacker links the predicted genotypes to the individual in the genotype dataset (Online Methods).

### 2.4 Individual Characterization by Linking Attacks

Using the three step approach, we first evaluated the accuracy of linking when the attacker utilizes genotype prediction where the attacker builds the posterior distribution of genotypes given the expression levels (Supplementary Note) where the attacker knows the exact joint distribution of genotypes and expression. The attacker builds the posterior distribution of genotypes given expression levels from the joint distribution. Finally, she predicts each genotype by selecting the genotype with

BASELINE  
REF

Deleted: 4c), with

Deleted: , which will be targeted by an attacker

Deleted: To estimate ICI leakage with

Deleted: first sorted the

Deleted: with respect to

Formatted: Font color: Text 1

Deleted: reported correlation,  $|\rho|$ . Then for top  $n=1,2,3,\dots,20$  eQTLs, we estimated mean  $\pi$  and mean

Formatted: Font color: Text 1

Formatted: Font: Not Italic, Font color: Text 1

Deleted: over all the samples as illustrated in Fig S2a.

Deleted:  $\pi$

Formatted: Font color: Text 1

Deleted: for each  $n$  (Fig. 4d)

Deleted:

Deleted: correctly link all individuals, on average

Deleted: chance, in

Deleted: characterizable

Deleted: size is

Deleted: , which can be interpreted as a higher risk of characterizability. In addition, auxiliary

Deleted: leaking information

Deleted: bits

Deleted: .

Deleted: Analysis

Deleted: Individual Characterization

Deleted: individual characterization in the context

Moved (insertion) [2]

Deleted: 2c).

Deleted: . The aim of the attacker

Formatted: Font color: Text 1

Deleted: to link the disease state of the individual to the correct individual

Deleted: .

Formatted: Font color: Text 1

Deleted: described in

Deleted: previous section

Deleted: In the case of eQTLs, this is the reported correlation coefficient,  $|\rho|$ .

Deleted: We assume that the attacker performs maximum  $\alpha$  ...

Deleted: individual characterization

Deleted:  $\pi_v$

Deleted: Using the three step approach, we first evaluated the ...

maximum *a posteriori* probability given gene expression level (Supplementary Note, **Supplementary Fig. 4**). For several eQTL selections with changing correlation threshold, the linking accuracy is above 95% and gets close to 100% when auxiliary information is available (**Fig. 5a**).

In general, knowledge or correct reconstruction of the exact joint genotype expression distribution may not be possible because the the genotype-phenotype correlation coefficient alone is not sufficient to perfectly reconstruct the genotype distribution given the expression levels. The attacker can, however, utilize a priori knowledge about the relation between gene expression levels and genotypes and build the joint genotype-expression distributions using models with varying complexities and parameters (Online Methods, Supplementary Note, **Supplementary Fig. 5**). We focus on a highly simplified model where the attacker exploits the knowledge that the eQTL genotypes and expression levels are correlated such that the extremes of the gene expression levels (highest and smallest expression levels) are observed with extremes of the genotypes (homozygous genotypes). We use a measure, termed extremity, to quantify the outlierness of expression levels (Online Methods, Supplementary Note, **Supplementary Fig. 6, 7**). Based on the extremity of expression level and the gradient of association, the attacker first builds an estimate of the joint genotype-expression distribution, then constructs the posterior distribution of genotypes and finally chooses the genotypes with maximum *a posteriori* probability (Online Methods, Supplementary Note, **Fig. 2a, b**).

The prediction methodology assigns zero probability to heterozygous genotype, and assigns only homozygous genotypes to variants, for which the associated gene's expression level has absolute extremity higher than a threshold. We performed linking attack using this prediction method (in 2<sup>nd</sup> step of linking). In the 1<sup>st</sup> step of the attack, we used absolute correlation and extremity thresholds for eQTL selection. The linking accuracy is higher than 95% for much of the eQTL selections (**Fig 2a, Supplementary Fig. 6d**). We also observed that changing extremity threshold does not affect the linking accuracy substantially compared to changing absolute correlation threshold. We thus focus on attack scenarios where the absolute extremity threshold is set to zero. With this approach, the genotype prediction accuracy increases with increasing absolute correlation threshold, as expected (**Supplementary Fig. 6c**). We next performed linking attack with this model where we used the correlation based eQTL selection in step 1, then extremity based genotype prediction in step 2. In the step 3, we evaluated two distance measures for linking the predicted genotypes to the individuals in genotype dataset (**Online Methods, Supplementary Fig. 8**). More than 95% of the individuals (**Fig. 5c, d**) are vulnerable for most of the parameter selections. When the auxiliary information is present, the fraction of vulnerable individuals increases to 100% for most of the eQTL selections. These results show that linking attack with extremity based genotype prediction, although technically simple, can be extremely effective in characterizing individuals. We evaluated whether the attacker can estimate the reliability of the linkings so as to focus on highly reliable linkings. We observed that the measure we termed, *first distance gap*, denoted by  $d_{1,2}$  (Online Methods), serves as a good reliability estimate for each linking. We computed the positive predictive value (PPV) versus sensitivity of the linkings in the testing set with changing  $d_{1,2}$  threshold (Online Methods). Compared to random sortings, the attacker can link a large fraction (79%) of the individuals at a PPV higher than 95% (**Fig. 5d, Supplementary Fig. 9**). We also studied several biases that can affect linking accuracy. First when the eQTLs are discovered

Deleted: *a posteriori*  
Deleted: of the genotypes given expression levels  
Deleted: knowledge of only  
Deleted: via eQTLs  
Deleted: enough  
Deleted: regenerate  
Deleted: *a posteriori*  
Deleted: of genotypes  
Deleted: used

Deleted: can coarsely  
Deleted: (

Moved up [2]: 2b).  
Deleted: therefore  
Deleted: .  
Deleted: , as expected,  
Deleted: (  
Moved down [3]: 5b).  
Deleted: To perform  
Deleted: ,  
Deleted: utilized

Deleted: . First is based on comparison of the predicted genotypes to all the genotypes in genotype dataset. Second is based on comparison of the predicted genotypes to only the homozygous genotypes in the genotype dataset, which is motivated by the fact that the attacker only predicts homozygous genotypes in the genotype prediction step (Online Methods). For each measure, the attacker links the predicted genotypes to the individual whose genotypes minimize the selected distance measure (**Supplementary Fig. 8**). More than 95% of the individuals (**Fig. 5c**,

Deleted: We will focus on homozygous genotype matching based distance computation in the rest of the paper for simplicity of presentation. ¶  
To test reproducibility of linking accuracy, we generated an eQTL training set (210 individuals) and identified eQTLs. The linking of expression and genotype datasets for remaining 211 individuals (testing set) is around 95% (**Supplementary Fig. 9a**), which shows that the linking attack is still effective when testing and training sets do not match. In addition, extremity based linking requires much less information compared to previously proposed method<sup>14</sup> (Supplementary Note, **Supplementary Table 1c**). To study these further, we

Deleted: **Supplementary Fig. 9b**).

on a sample set that and the linking attack is performed on another sample set, the accuracies are still very high (Supplementary Note, **Supplementary Fig. 9a**). Moreover, attacks are accurate when there is mismatch between the tissue or population of eQTL discovery sample set and tissue or population of linking attack sample set (Supplementary Note, **Supplementary Table 1a, b**). In addition, we observed that the extremity attack is still effective when genotype sample size is very large (Supplementary Note, **Supplementary Fig. 9c, d**), which points out the applicability on large sample sizes. We also observed that the extremity attack may link close relatives to each other, which can create potential privacy concerns for the family (**Supplementary Fig. 10**).

### 3 DISCUSSION

In genomic privacy, it is necessary to consider the basic premise of sharing any type of personal information: There is always an amount of leakage in the sensitive information<sup>19</sup>. In addition, as shown by previous studies, we often cannot propose black-and-white solutions to problems in privacy which mainly roots from the multifaceted nature of privacy. We believe these make it necessary for the genomic data sharing and publishing mechanisms to incorporate statistical quantification methods to objectively quantify risk estimates before the datasets are released. The quantification methodology and the analysis frameworks presented here and in future studies can be applied for analysis of the information leakage in the datasets where the correlative relations can be exploited for performing linking attacks (Supplementary Note, **Supplementary Fig. 11**).

Our study focuses on the individual privacy breaches in the context of linking attacks, where an individual's existence in two seemingly independent databases (e.g., phenotype and the genotype) can cause a privacy concern when an attacker links statistically the databases using the a priori information about correlation of different entries in the databases. The obvious risk management strategy against these attacks is restricting access to the phenotype datasets. This approach has, however, high cost in terms of lost research opportunities. Another approach is serving encrypted data, where data analysis is performed directly on the encrypted data, for example, using homomorphic encryption<sup>20</sup>. This approach has very high compute requirements and not practical yet. One other approach is to utilize statistical techniques like k-anonymization<sup>21</sup>. These can be employed on the phenotype datasets before being published. For this, it is necessary to develop new approaches and heuristics that can effectively circumvent high computational requirements<sup>22</sup>. In addition, several other approaches have addressed scenarios where k-anonymization may fail to protect data<sup>23,24</sup>. These scenarios must be properly handled in the risk management strategies. The anonymization strategies can use the estimates of leakage and predictability from our study to determine the QTLs that cause most leakage and anonymize the phenotype data accordingly. Another approach is to serve phenotypic data from a statistical database. In this context, differential privacy has been proposed as an optimal way for privacy aware data serving from statistical databases. The data release mechanisms in a differentially private scenario can benefit from the estimates of ICI leakage in each QTL. Differentially private data serving may, however, decrease the biological utility of the data significantly<sup>25</sup>. We believe new studies should address protection and risk management strategies for serving utility maximized and privacy aware high dimensional phenotype datasets.

EARLIER

**Deleted:** In order to study the effect of genotype dataset size, we simulated the genotypes of the eQTLs (from the training set) for 100,000 individuals using genotype frequencies from 1000 Genomes Project (Supplementary Note). We then merged simulated dataset with testing dataset (of 100,211 individuals) and observed that linking of testing expression samples to large genotype dataset is around 96% accurate (**Supplementary Fig. 9c**). In addition, the attacker can correctly link 55% of individuals with more than 95% PPV (**Supplementary Fig. 9d**). We next separated the GEUVADIS samples to 5 populations and identified eQTLs for each population and performed linking using population specific eQTLs (**Supplementary Table S1a**). The eQTLs from close European populations (CEU, GBR, TSI, FIN) enable high linking accuracy (higher than 95%) for European populations. When the eQTLs are identified from an African population (YRI), the linking accuracies in European populations are smaller. To also study the mismatch between the tissues where linking and eQTL discovery is performed, we downloaded the eQTLs for 6 tissues from the GTex Project<sup>13</sup> and linked the GEUVADIS expression samples to genotypes samples. The accuracy is generally high (>75%) and is highest for blood eQTLs, as expected (**Supplementary Table 1b**). We also observed that when close relatives of individuals in the expression dataset (30 CEU trio dataset from the HAPMAP project<sup>17,18</sup>) are in the genotype dataset, the linking attack assigns much lower ranks to relatives compared to the random individuals (Supplementary Note, **Supplementary Fig. 10**), which may cause privacy concerns for the family of individuals.

**Deleted:** We showed that the attacks can be effective even when the eQTLs are identified in tissues or populations that are different than the samples being linked, and also the attacks can target not only the linked individual but their families, which increases the impact of breach. The obvious risk management strategy against these attacks is restricting access to the phenotype datasets. The statistical techniques like k-anonymization and differential privacy can also be utilized. These, however, have associated drawbacks about loss of biological utility, and high computational complexity. Moreover, some studies also demonstrated that there are still risks associated with linkability of the anonymized data<sup>20-23</sup>. We believe new studies should address protection and risk management strategies for serving utility maximized and privacy aware high dimensional phenotype datasets.

QTL

## 4 DATASETS

The normalized gene expression levels for 462 individuals and the eQTL dataset are obtained from GEUVADIS mRNA sequencing project<sup>17</sup>. The eQTL dataset contains all the significant (Identified at most 5% false discovery rate) gene-variant pairs with high genotype-expression correlation. To ensure that there are no dependencies between the variant genotypes and expression levels, we used the eQTL entries where gene and variants are unique. In other words, each variant and gene are found exactly once in the final eQTL dataset (Section S4). The genotype, gender, and population information datasets for 1092 individuals are obtained from 1000 Genomes Project<sup>18</sup>. For 421 individuals, both the genotype data and gene expression levels are available. For tissue analysis, the publicly available significant eQTLs for 6 tissues that are computed by the GTex project are downloaded from the GTex Portal.

Deleted: [[ACCESSION NUMBERS??]]

**Deleted:** The normalized gene expression levels for 462 individuals and the eQTL dataset are obtained from gEUVADIS mRNA sequencing project<sup>15</sup>. The eQTL dataset contains all the significant (Identified at most 5% false discovery rate) gene-variant pairs with high genotype-expression correlation. To ensure that there are no dependencies between the variant genotypes and expression levels, we used the eQTL entries where gene and variants are unique. In other words, each variant and gene are found exactly once in the final eQTL dataset (Section S4). The genotype, gender, and population information datasets for 1092 individuals are obtained from 1000 Genomes Project<sup>16</sup>. For 421 individuals, both the genotype data and gene expression levels are available. For tissue analysis, the publicly available significant eQTLs for 6 tissues that are computed by the GTex project are downloaded from the GTex Portal.¶  
[[Accession numbers]]¶

## 5 ACKNOWLEDGEMENTS

Authors would like to thank Akdes Serin Harmanci for constructive comments and discussions on study design and running of external tools. Authors also would like to thank the anonymous reviewers for constructive criticism that made the manuscript complete and comprehensive.

## 6 AUTHOR CONTRIBUTIONS

A.H. designed the study, gathered datasets, performed experiments, and drafted the manuscript. M.G. conceived the study, oversaw the experiments and wrote the manuscript. Both authors approved final manuscript.

Authors declare no conflict of financial interest.

## 7 REFERENCES

1. Sboner, A., Mu, X., Greenbaum, D., Auerbach, R. K. & Gerstein, M. B. The real cost of sequencing: higher than you think! *Genome Biol.* **12**, 125 (2011).
2. Rodriguez, L. L., Brooks, L. D., Greenberg, J. H. & Green, E. D. The Complexities of Genomic Identifi ability. *Science (80- )*. **339**, 275–276 (2013).
3. Erlich, Y. & Narayanan, A. Routes for breaching and protecting genetic privacy. *Nat. Rev. Genet.* **15**, 409–21 (2014).
4. Sweeney, L., Abu, A. & Winn, J. Identifying Participants in the Personal Genome Project by Name. *SSRN Electron. J.* 1–4 (2013). doi:10.2139/ssrn.2257732
5. Sweeney, L. *Uniqueness of Simple Demographics in the U.S. Population, LIDAP-WP4. Forthcom. B.*

Deleted: 4

entitled, *Identifiability Data*. (2000).

6. Golle, P. Revisiting the uniqueness of simple demographics in the US population. in *Proc. 5th ACM Work. Priv. Electron. Soc.* 77–80 (2006). doi:http://doi.acm.org/10.1145/1179601.1179615

Deleted: 5

7. Consortium, T. G. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–5 (2013).

Deleted: 6

8. Ardlie, K. G. *et al.* The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science (80-. J.)*. **348**, 648–660 (2015).

Moved (insertion) [4]

9. Pakstis, A. J. *et al.* SNPs for a universal individual identification panel. *Hum. Genet.* **127**, 315–324 (2010).

10. Wei, Y. L., Li, C. X., Jia, J., Hu, L. & Liu, Y. Forensic Identification Using a Multiplex Assay of 47 SNPs. *J. Forensic Sci.* **57**, 1448–1456 (2012).

Deleted: 7

11. Gymrek, M., McGuire, A. L., Golan, D., Halperin, E. & Erlich, Y. Identifying personal genomes by surname inference. *Science* **339**, 321–4 (2013).

Deleted: 8

12. Homer, N. *et al.* Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. *PLoS Genet.* **4**, (2008).

13. Im, H. K., Gamazon, E. R., Nicolae, D. L. & Cox, N. J. On sharing quantitative trait GWAS results in an era of multiple-omics data and the limits of genomic privacy. *Am. J. Hum. Genet.* **90**, 591–598 (2012).

Deleted: 9

14. Lunshof, J. E., Chadwick, R., Vorhaus, D. B. & Church, G. M. From genetic privacy to open consent. *Nat. Rev. Genet.* **9**, 406–411 (2008).

Deleted: 10

15. Church, G. *et al.* Public access to genome-wide data: Five views on balancing research with privacy and protection. *PLoS Genet.* **5**, (2009).

Deleted: 11. Xia, K. *et al.* SeeQTL: A searchable database for human eQTLs. *Bioinformatics* **28**, 451–452 (2012).¶  
¶  
12. . Consortium, T. G. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–5 (2013).¶  
¶  
13.

16. Narayanan, A. & Shmatikov, V. Robust de-anonymization of large sparse datasets. in *Proc. - IEEE Symp. Secur. Priv.* 111–125 (2008). doi:10.1109/SP.2008.33

Moved up [4]: Ardlie, K. G. *et al.* The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science (80-. J.)*. **348**, 648–660 (2015).¶

17. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501**, 506–11 (2013).

Deleted: 14. . Schadt, E. E., Woo, S. & Hao, K. Bayesian method to predict individual SNP genotypes from gene expression data. *Nat. Genet.* **44**, 603–608 (2012).¶  
¶  
15



18. The 1000 Genomes Project Consortium. An integrated map of genetic variation. *Nature* **135**, 0–9 (2012).

Deleted: 16

19. Narayanan, A. *et al.* *Redefining Genomic Privacy: Trust and Empowerment*. *bioRxiv* (2014). doi:10.1101/006601

Deleted: 17. Stranger, B. E. *et al.* Relative impact of nucleotide and copy number variation on gene expression phenotypes. *Science* **315**, 848–853 (2007).¶

20. Vaikuntanathan, V. *Computing Blindfolded: New Developments in Fully Homomorphic Encryption*. *2011 IEEE 52nd Annu. Symp. Found. Comput. Sci.* 5–16 (2011). doi:10.1109/FOCS.2011.98

¶  
18. The International HapMap 3 Consortium. Integrating common and rare genetic variation in diverse human populations. *Nature* **467**, 52–8 (2010).¶  
¶

21. SWEENEY, L. k-ANONYMITY: A MODEL FOR PROTECTING PRIVACY. *Int. J. Uncertainty, Fuzziness Knowledge-Based Syst.* **10**, 557–570 (2002).

22. Meyerson, A. & Williams, R. On the complexity of optimal K-anonymity. in *Proc. twentythird ACM SIGMOD-SIGACT-SIGART Symp. Princ. database Syst. Pod. 04* 223–228 (2004). doi:10.1145/1055558.1055591

23. Ninghui, L., Tiancheng, L. & Venkatasubramanian, S. t-Closeness: Privacy beyond k-anonymity and  $\ell$ -diversity. in *Proc. - Int. Conf. Data Eng.* 106–115 (2007). doi:10.1109/ICDE.2007.367856

Moved down [5]: Machanavajjhala, A., Kifer, D., Gehrke, J. & Venkatasubramanian, M. L-diversity. *ACM Trans. Knowl. Discov. Data* **1**, 3–es (2007).¶  
¶

24. Machanavajjhala, A., Kifer, D., Gehrke, J. & Venkatasubramanian, M. L-diversity. *ACM Trans. Knowl. Discov. Data* **1**, 3–es (2007).

Deleted: 21

Moved (insertion) [5]

25. Fredrikson, M., Lantz, E., Jha, S. & Lin, S. Privacy in Pharmacogenetics: An End-to-End Case Study of Personalized Warfarin Dosing. in *23rd USENIX Secur. Symp.* (2014). at <http://www.biostat.wisc.edu/~page/WarfarinUsenix2014.pdf>

Deleted: 22. Wong, R. C.-W. W., Fu, A. W.-C. C., Wang, K. & Pei, J. Minimality attack in privacy preserving data publishing. in *Proc. 33rd Int. Conf. Very large data bases* 543–554 (2007). at <http://dl.acm.org/citation.cfm?id=1325851.1325914>¶

¶  
23

26. Cover, T. M. & Thomas, J. A. *Elements of Information Theory*. *Elem. Inf. Theory* (2005). doi:10.1002/047174882X

Deleted: 24

27. Shabalin, A. A. Matrix eQTL: Ultra fast eQTL analysis via large matrix operations. *Bioinformatics* **28**, 1353–1358 (2012).

Deleted: 25. Herbert A. Sturges. The Choice of a Class Interval. *J. Am. Stat. Assoc.* **21**, 65–66 (1926).¶

¶  
26

## 8 FIGURE LEGENDS

**Figure 1:** Illustration of the linking attack. (a) Phenotype dataset contains  $q$  different phenotype measurements and the HIV Status for a list of  $n$  individuals. Genotype dataset contains the variants genotypes for  $m$  individuals. Phenotype-Genotype correlation datasets contains  $q$  phenotypes, variants, and their correlations. The attacker does genotype prediction for all the variants. The attacker then links

Formatted: Not Highlight

Deleted:

Formatted: Not Highlight

the phenotype dataset to the genotype dataset by matching the genotypes. The linking potentially reveals the HIV status for the subjects in the genotypes dataset. The IDs and HIV Status are colored to illustrate how the linking combines the entries in the two datasets. The non-shaded columns are used for linking.

**Figure 2:** Illustration of genotype-expression associations and linking attacks (a) Schematic representation of genotype and expression associations. Genotype (y-axis) and expression (x-axis) are correlated, indicated by line fit and  $\rho$ . The rectangles represent conditional distribution of expression given genotype values. (b) Illustration of extremity based genotype prediction. Expression range is divided into two equal ranges (separated by  $e_{mid}$ ). The blue rectangles represent the distribution used for prediction. Given the distribution of expression (tri-model distribution on right), the positive extremity is assigned genotype 2, and negative extremity is assigned genotype 0. (c) Three step linking process. First step is selection of phenotypes and genotypes to be used in linking. Second step is prediction of genotypes. Last step is linking of predicted genotypes to the genotype dataset.

**Figure 3:** Illustration of individual characterizing information (ICI) and correct predictability of genotypes. (a) Graphical representation of ICI formulation. ICI for a set of  $n$  variant genotypes is computed in terms of population genotype frequencies. Each genotype contributes to ICI additively with the logarithm of reciprocal of the genotype frequency (illustrated by the genotype distributions). (b) Graphical representation of  $\pi$ . Given the joint distribution of genotype and expression (shown below), the conditional distribution of genotypes given expression level  $e$  is computed. The exponential of the conditional distribution entropy is used for computing the predictability.

**Figure 4:** ICI versus  $\pi$  for each eQTL. Plots show, for each eQTL, the information leakage (x-axis) versus correct genotype predictability (y-axis). For each eQTL, the estimated ICI leakage and genotype predictability are plotted. The dots are colored with respect to the major allele frequency (a) and with respect to absolute correlation of the eQTL (b). ICI versus  $\pi$  for shuffled data (red) is compared to the real dataset (blue) in (c).

**Figure 5:** Accuracy measures for linking attacks. (a) Linking accuracy with MAP genotype predictions. Absolute correlation threshold (x-axis) versus fraction of vulnerable individuals (y-axis). The yellow arrow indicates the maximized position of linking accuracy. Red, green, and cyan plots show linking accuracy with gender, population, and gender + population as auxiliary information. (b) The genotype prediction accuracy. The genotype prediction accuracy (y-axis) of with changing absolute correlation threshold (x-axis). (c) Linking accuracy with extremity based linking with all genotypes. (d) Linking accuracy with extremity based linking with homozygous genotypes.

## [[FOLLOWING ARE NOT UPDATED, YET]]

**Table S1:** Linking accuracy of extremity based linking attack using the eQTLs are identified in different populations and different tissues. (a) The table shows the linking accuracies (for populations shown in the rows) when the eQTLs that are identified using data (indicated in each column) from different

Formatted: Not Highlight

Formatted: Not Highlight

populations. (b) The linking accuracy of individuals in GEUVADIS project when eQTLs identified from different tissues are used in linking.

**Table S2:** Linking attack accuracy comparison. The table shows linking accuracy for Schadt et al and extremity based linking attack methods. Each row corresponds (for Schadt et al Method) to a different number of data points in the training datasets that is input to Schadt et al method.

**Supplementary Figure 1:** Schematic comparison of linking attacks (Left) and detection of a genome in a mixture attacks (Right). Each box in the figure represents a dataset in the form of a matrix. Multiple boxes next to each other correspond to concatenation of matrices. Linking attacks aim at linking genotype and phenotype datasets. The phenotype datasets contain both “predicting” phenotypes and other phenotypes, some of which can be sensitive. The attacker first predict genotypes for each of the predicting phenotype. The predicted genotypes are then compared with the genotypes in the genotype dataset. After the linking, all the datasets are concatenated where the identifiers can be matched to the sensitive phenotypes. Different colors indicate how the linking merges different information. The detection of a genome in a mixture attacks start with a genotype dataset. The attacker gets access to the statistics of a GWAS or genotyping dataset (for example, regression coefficients or allele frequencies). Then the attacker generates a statistic and tests it against that of a reference population. The testing result can be converted into the study membership indicator (attended/not attended) which shows whether the tested individual was in the study cohort or not.

**Supplementary Figure 2:** Representation of the eQTLs. (a) The average ICI leakage versus the genotype predictability is shown for real (red) and shuffled (blue) eQTL dataset is shown. (b) The absolute correlation versus predictability is shown.

Figure shows the attacker’s presumed strategy for linking attack. (a) The phenotype and variant pairs are sorted with respect to decreasing absolute correlations values. For the top  $n$  pairs, joint predictability and ICI are computed. (b) Illustration of prior, joint, and posterior distributions of genotypes and expression levels. Leftmost figure shows the distribution of genotypes over the sample set, which is labelled as the prior distribution. Middle figure shows the joint distribution of genotypes and expression levels. Notice that there is a significant negative correlation between genotype values and the expression levels. Rightmost figure shows the posterior distribution of genotypes given that the gene expression level is 10. The posterior distribution has a maximum (MAP prediction) at genotype 2, which is indicated by a star.

**Supplementary Figure 3:** The distribution of ranks of the individuals in the linking step. At each gradient threshold, the box plots show, for each individual, their ranks in the genotype comparison in the 3<sup>rd</sup> step of linking attack with MAP genotype prediction. Notice that at around 0.35 correlation threshold, the assigned ranks are minimized, i.e., most of the individual are linked correctly.

**Supplementary Figure 4:** The median absolute gene expression extremity statistics over 462 individuals in GEUVADIS dataset. (a) For each individual, the extremity is computed over all the genes (23,662 genes) reported in the expression dataset. The median of the absolute value of the extremity is plotted. X-axis shows the sample index and y-axis shows the extremity. The absolute median extremity fluctuates around 0.25, which is exactly the midpoint between minimum and maximum values of absolute extremity. (b) For each individual, we count the number of genes above the extremity threshold. The plot shows the extremity threshold versus the median number of genes (over 462 individuals) above the extremity threshold. Around half of the genes (indicated by dashed yellow lines) have higher than almost 0.3 extremity on average over all the individuals. Also, around median number of 1000 genes over the samples have higher than 0.45 extremity (indicated by dashed red lines).

**Supplementary Figure 5:** Illustration of linking for  $j$ th individual. The attacker first predicts the genotypes ( $\hat{v}_{\cdot,j}$ ) which are then used to compute the distance to all the individuals in the genotype dataset. The computed distances are then sorted in decreasing. The top matching individual (in the example, individual  $a$ ) is assigned as the linked individual. The first distance gap,  $d_{1,2}$ , is computed as the difference between the second ( $d_{j,(2)}$ ) and the first ( $d_{j,(1)}$ ) distances in the sorted list.

**Supplementary Figure 6:**

**Supplementary Figure 7:** A representative example of extremity based linking. The phenotype dataset (Consisting of gene expression levels for 6 genes) is shown above. Each phenotype measurement is represented by blue (negative extreme), yellow (positive extreme), or grey (non-extreme) dots. Based on the extremity of phenotypes, the attacker performs prediction of genotypes, which are shown below in (2). She uses the eQTL dataset (with genes and SNPs) for prediction. Blue and brown triangles correspond to the correct genotype predictions. The grey crosses correspond to the incorrect or unavailable genotype predictions. The attacker compares the predicted genotypes to the genotype dataset in (3), where triangles show the genotypes, and performs linking. The attacker links the predicted genotypes to the genotype dataset. 3 individuals (Bob, Alice, and John) are highlighted. The attacker can link Bob and John by matching them to their genotypes. The correct prediction of rs7274244 (in yellow dashed rectangle) enables the attacker to distinguish between correct entries and reveal both of their disease status as positive. For Alice, the predicted genotypes are equally matching at two entries both of which match at 2 genotypes; PID-b and PID-k (with negative and positive disease status) thus the attacker cannot exactly reveal Alice's disease status.

**Supplementary Figure 8:** Illustration of risk assessment procedure for joint genotyping/phenotyping data generation. There are two paths of risk assessment to be performed. The first path evaluates the risks associated with release of the QTL datasets. The genotype and phenotype data (on the left) is first used for quantitative trait loci identification (QTL identification box). This generates the significant QTLs. These are then utilized, in addition to the list of external QTL databases, in quantification of leakage versus predictability, as presented in Section 2.2. These results are then relayed to the risk assessment procedures. The second risk assessment procedure evaluates the release of genotype and phenotype datasets. For this, the datasets are input to application of a list of linking attacks (Presented in Sections

- Formatted: Not Highlight
- Formatted: Not Highlight
- Formatted: Not Highlight
- Formatted: Not Highlight
- Formatted: Not Highlight
- Formatted: Not Highlight
- Formatted: Not Highlight
- Formatted: Not Highlight
- Deleted: green
- Formatted: Not Highlight
- Deleted: red
- Formatted: Not Highlight
- Deleted: He or she
- Formatted: Not Highlight

2.3, and 2.4, and other linking attacks in the literature) for evaluation of characterization risks. The results are then relayed to risk assessment procedures.

**Supplementary Figure 9:** Models of joint genotype-expression distribution with varying numbers of parameters for a positively correlated eQTL. (a) shows the true distribution where grey boxes represent the expression distributions given different genotypes. Red line show the gradient of correlation between genotype and expression. First simplification of the model is shown in (b). The expression distribution can be modeled with Gaussians with different means and variances with total of 6 parameters. The variances can be assumed same for different genotypes (c), where 4 parameters are required. (d) illustrates a representation of the uniform expression distribution given genotypes, where 4 parameters are required. The conditional distribution of expression is uniform (cross shaded rectangles) over the ranges  $(e_1, e_2)$ ,  $(e_2, e_3)$ , and  $(e_3, e_4)$  given genotypes 0, 1, and 2, respectively. The transparent grey rectangles shows the original distributions. (e) is a simplification of (d) where no conditional probability of expression is assigned given genotype is 1. In this model, only one parameter ( $e_{mid}$ ) is necessary. The conditional probability of expression given genotypes 0 and 2 are uniform for expression levels below  $e_{mid}$  and above  $e_{mid}$  respectively (shown with cross shaded rectangles). The original distribution is included with grey rectangles for comparison. Extremity based prediction is an instantiation of the model in (e).

## 9 ONLINE METHODS

### 9.1 Genotype, Expression, and eQTL Datasets

The eQTL, expression, and genotype datasets contain the information for linking attack (**Supplementary Fig. 2**). The eQTL dataset is composed of a list of gene-variant pairs such that the gene expression levels and variant genotypes are significantly correlated. We will denote the number of eQTL entries with  $q$ . The eQTL (gene) expression levels and eQTL (variant) genotypes are stored in  $q \times n_e$  and  $q \times n_v$  matrices  $e$  and  $v$ , respectively, where  $n_e$  and  $n_v$  denotes the number of individuals in gene expression dataset and individuals in genotype dataset. The  $k$ th row of  $e$ ,  $e_k$ , contains the gene expression values for  $k$ th eQTL entry and  $e_{k,j}$  represents the expression of the  $k$ th gene for  $j$ th individual. Similarly,  $k$ th row of  $v$ ,  $v_k$ , contains the genotypes for  $k$ th eQTL variant and  $v_{k,j}$  represents the genotype ( $v_{k,j} \in \{0,1,2\}$ ) of  $k$  variant for  $j$ th individual. The coding of the genotypes from homozygous or heterozygous genotype categories to the numeric values are done according to the correlation dataset (Online Methods). We assume that the variant genotypes and gene expression levels for the  $k$ th eQTL entry are distributed randomly over the samples in accordance with random variables (RVs) which we denote with  $V_k$  and  $E_k$ , respectively. We denote the correlation between the RVs with  $\rho(E_k, V_k)$ . In most of the eQTL studies, the value of the correlation is reported in terms of a gradient (or the regression coefficient) in addition to the significance of association (p-value) between genotypes and expression levels.

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Deleted: ,

## 9.2 Quantification of Characterizing Information and Predictability

The genotype RV  $V_k$  takes 3 different values,  $\{0,1,2\}$ , where the genotype coding is done per counting the number of alternate alleles in the genotype. Given that the genotype is  $g_{k,j}$ , we quantify the individual characterizing information in terms of *self-information*<sup>26</sup> of the event that RV takes the value  $g_{k,j}$ :

$$ICI(V_k = g_{k,j}) = I(V_k = g_{k,j}) = -\log_2(p(V_k = g_{k,j})) \quad (1)$$

where  $V_k$  is the RV that represents the  $k$ th eQTL genotype,  $p(V_k = g_{k,j})$  is the probability (frequency) of that  $V_k$  takes the value  $g_{k,j}$ , and  $ICI$  denotes the individual characterizing information. Given multiple eQTL genotypes, assuming that they are independent, the total individual characterizing information is simply summation of those:

$$\begin{aligned} ICI(\{V_1 = v_{1,j}, V_2 = v_{2,j}, \dots, V_N = v_{N,j}\}) \\ = -\sum_{k=1}^N \log_2(p(V_k = v_{k,j})). \end{aligned} \quad (2)$$

The genotype probabilities are estimated by the frequency of genotypes in the genotype dataset. We measure the predictability of eQTL genotypes using an entropy based measure. Finally, the base of logarithm that is used determines the units in which ICI is reported. When base two logarithm is used as above, the unit of ICI is bits.

Given the genotype RV,  $V_k$ , and the correlated gene expression RV,  $E_k$ ,

$$\pi(V_k | E_k = e) = \exp(-H(V_k | E_k = e)) \quad (3)$$

where  $\pi$  denotes the predictability of  $V_k$  given the gene expression level  $e$ , and  $H$  denotes the entropy of  $V_k$  given gene expression level  $e$  for  $E_k$ . The extension to multiple eQTLs is straightforward. For the  $k$ th individual, given the expression levels  $e_{k,j}$  for all the eQTLs, the total predictability is computed as

$$\begin{aligned} \pi(\{V_k\}, \{E_k = e_{k,j}\}) = \exp(-H(\{V_k\} | \{E_k = e_{k,j}\})) \\ = \exp\left(-\sum_k H(V_k | E_k = e_{k,j})\right) \end{aligned} \quad (4)$$

In addition, this measure is guaranteed to be between 0 and 1 such that 0 represents no predictability and 1 representing perfect predictability. The measure can be thought as mapping the prediction process to a uniform random guessing where the average correct prediction probability is measured by  $\pi$ .

## 9.3 Extremity Based MAP Genotype Prediction

Using an estimate of the joint distribution, the attacker can compute the *a posteriori* distribution of genotypes given gene expression levels. To quantify the extremeness of expression levels, we use a

Field Code Changed

Deleted: <sup>24</sup>

Deleted:  $\log(p(V_k = g_{k,j}))$

Deleted:  $\sum_{k=1}^N \log(p(V_k = v_{k,j}))$

Formatted Table

Formatted Table

Deleted: Estimation of

Deleted: Entropy

Deleted: We estimate the genotype entropy using the Shannon's entropy<sup>24,¶</sup>

statistic we termed *extremity*. For the gene expression levels for  $k^{\text{th}}$  eQTL,  $e_{k,j}$  extremity of the  $j^{\text{th}}$  individual's expression level,  $e_{k,j}$  is defined as

$$ext(e_{k,j}) = \frac{\text{rank of } e_{k,j} \text{ in } \{e_{k,1}, e_{k,2}, \dots, e_{k,n_e}\}}{n_e} - 0.5. \quad (5)$$

Extremity can be interpreted as a normalized rank, which is bounded between -0.5 and 0.5. The average median extremity is uniformly distributed among individuals (Supplementary Fig. 6a). In addition, around half of the genes (10,000) in each individual have higher than extremity value of 0.3. Also, around 1000 genes have higher than 0.45 absolute extremity (Supplementary Fig. 6b). In other words, each individual harbors substantial number of genes whose expressions are at the extremes within the population. These can potentially serve as quasi-identifiers. It is worth noting, however, that not all of these extreme genes are associated with eQTLs.

Following from the above discussion, the adversary builds the posterior distribution for  $k^{\text{th}}$  eQTL genotypes as

$$P(V_k = 0 | E_k = e_{k,j}) = \begin{cases} 1 & \text{if } |ext(e_{k,j})| > \delta, ext(e_{k,j}) \times \rho(E_k, V_k) < 0 \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

$$P(V_k = 2 | E_k = e_{k,j}) = \begin{cases} 1 & \text{if } |ext(e_{k,j})| > \delta, ext(e_{k,j}) \times \rho(E_k, V_k) > 0 \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

$$P(V_k = 1 | E_k = e_{k,j}) = 0. \quad (8)$$

From the *a posteriori* probabilities, when the sign of the extremity and the reported correlation are the same, the attacker assigns the genotype value 2, and otherwise, genotype value 0. Finally, the genotype value 1 is never assigned in this prediction method, i.e., the *a posteriori* probability is zero. As yet another way of interpretation, the genotype prediction can be interpreted as a rank correlation between the genotypes and expression levels and choosing the homozygous genotypes that maximize the absolute values of the rank correlation. Thus, this process can be generalized as a rank correlation based prediction. We are focusing on the extremes and heterozygous genotype is observed at medium levels of expression. The posterior distribution of genotypes in equations (4-6) can be derived from a simplified model of the genotype-expression distribution that utilizes just one parameter (Online Methods). We used the posterior genotype probabilities in extremity based prediction and assessed the genotype prediction accuracy. As expected, the accuracy of genotype predictions increases with increasing correlation threshold (Fig. 5b).

The slight decrease of genotype accuracy at correlation thresholds higher than 0.7 is caused by the fact that the accuracy (fraction of correct genotype predictions within all genotypes) is not robust at very small number of SNPs. Although we expect very high accuracy, even one wrong prediction among small number of total genotypes decreases the accuracy significantly.

**Deleted:**  $H(V_k) = -\sum_{v \in \{0,1,2\}} p(V_k = v) \times \log(p(V_k = v))$

**Formatted Table**

**Deleted:** where  $V_k$  represents the RV for  $k^{\text{th}}$  eQTL variant genotypes and  $p(V_k = v)$  represents the probability that  $V_k$  takes the value  $v$ . This probability can be also interpreted as the population frequency of the genotype  $v$  at the  $k^{\text{th}}$  eQTL's variant locus. These probabilities are estimated from the distribution of genotypes over all the samples. As the genotypes are discrete valued, the above formula can be computed in a straightforward way by the summation after the probabilities are estimated. In the formulation for conditional predictability of genotypes given expression levels, we also use the conditional specific entropies<sup>24</sup> of the genotypes given the gene expression levels. For this, we use the following formulation:

**Formatted Table**

**Deleted:**  $H(V_k | E_k = e_{k,j}) = -\sum_{v \in \{0,1,2\}} p(V_k = v | E_k = e_{k,j}) \times \log(p(V_k = v | E_k = e_{k,j}))$

**Moved (insertion) [3]**

**Deleted:** where  $p(V_k = v | E_k = e_{k,j})$  represents the conditional probability that  $V_k$  takes the value  $v$  under the condition that the RV representing gene expression level for  $k^{\text{th}}$  eQTLs ( $E_k$ ) is  $e_{k,j}$ . Since the gene expression levels are continuous, to estimate the conditional probabilities of genotypes given expression levels; we start with the joint distribution of  $E_k$  and  $V_k$ , then bin the gene expression levels. For this, we use Sturges' rule<sup>25</sup> to choose the number of bins. This rule states that the number of bins should be selected as  $n_b = \lceil \log_2(n_e) \rceil + 1 = \lceil \log_2(421) \rceil + 1 = 10$ . The binning is done for each gene by first sorting the expression levels for all the individuals, then the range of gene expression levels are divided into  $n_b = 10$  bins of equal size and each expression level is mapped to a value between in  $[0, n_b - 1]$ . The expression level of  $k^{\text{th}}$  gene in  $j^{\text{th}}$  individual,  $e_{k,j}$ , is mapped to

## 9.4 First Distance Gap Statistic Computation

Following the previous section, the attacker computes, for each individual, the distance to all the genotypes in genotype dataset, then identifies the individual with smallest distance. Let  $d_{j,(1)}$  and  $d_{j,(2)}$  denote the minimum and second minimum genotype distances (among  $d^H(\tilde{\mathbf{v}}_{\cdot,j}, \mathbf{v}_{\cdot,\mathbf{a}})$  for all  $\mathbf{a}$ ) for  $j$ th individual. We propose using the difference between these distances, termed *first distance gap statistic*, as a measure of reliability of linking. For this, the attacker computes following difference:

$$d_{1,2}(j) = d_{j,(2)} - d_{j,(1)} \quad (9)$$

Deleted: 10

First distance gap can be computed without the knowledge of the true genotypes, and is immediately accessible by the attacker with no need for auxiliary information (**Supplementary Fig. 8**). The basic motivation for this statistic comes from the observation that the first distance gap for correctly linked individuals are much higher compared to the incorrectly linked individuals.

## 9.5 eQTL Identification with Matrix eQTL

For identification of eQTLs, we used Matrix eQTL <sup>27</sup> [method](#). We first generated the testing and training sample lists by randomly picking 210 and 211 individuals, respectively, for testing and training sets. We then separated the genotype and expression matrices into training and testing sets. Matrix eQTL is run to identify the eQTLs using the training dataset. In order to decrease the run time, Matrix eQTL is run in cis-eQTL identification mode. After the eQTLs are generated, we filtered out the eQTLs whose FDR (as reported by Matrix eQTL) was larger than 5%. We finally removed the redundancy by ensuring that each gene and each SNP is used only once in the eQTL final list. To accomplish this, we selected the eQTL that is correlated with highest association with each gene. The association statistic reported by Matrix eQTL was used as the measure of strength of association between expression levels and genotypes. Similar procedure is applied when eQTLs for 30 trios are identified.

Deleted: <sup>26</sup> method.

## 9.6 Modeling of Genotype-Phenotype Distribution

In the second step of the linking attack, the genotype predictions are performed. The genotype predictions are used, as an intermediate information, as input to the third step (**Fig. 2c**), where linking is performed. The main aim of attacker is to maximize the linking accuracy (not the genotype prediction accuracy), which depends jointly on the genotype prediction accuracy and the accuracy of the genotype matching in the 3<sup>rd</sup> step. Other than the accuracy of linking, another important consideration, for risk management purposes, is the amount of auxiliary input data (like training data for prediction model) that the genotype prediction takes. The prediction methods that require high amount of auxiliary data would decrease the applicability of the linking attack as the attacker would need to gather extra information before performing the attack. On the other hand, the prediction methods that require little or no auxiliary data makes the linking attack much more realistic and prevalent. It is therefore useful, in the risk management strategies, to study complexities of genotype prediction methods and evaluate how these translate into assessing the accuracy and applicability of the linking attack. We study different simplifications of genotype prediction, and illustrate different levels of complexity for genotype prediction.



In MAP based genotype prediction and linking attack, we assume that the attacker estimates the posterior distribution of genotypes and utilizes the maximum *a posteriori* estimate of the genotype as the general prediction method. For this, attacker must first model the joint genotype-phenotype distribution and then build the posterior genotype distribution (**Supplementary Fig. 5a**). The first level level of model can be built by decomposing the conditional distribution of expression with independent variances and means (**Supplementary Fig. 5b**). Assuming that mean and variance are sufficient statistics for the conditional distributions (e.g., normally distributed), the joint distributions can be modeled when the 6 parameters (3 means and 3 variances) are trained. The training can be performed using unsupervised methods like expectation maximization or can be performed using training data. This would, however, increase the required auxiliary data and decrease the applicability of the linking attack. A simplification of the model by assuming the variances of the conditional expression distributions are same for each genotype (**Supplementary Fig. 5c**). This decreases the number of parameters to be trained to 4 (3 means and 1 variance). An equally complex model with 4 parameters can be built assuming the conditional distributions are uniform at non-overlapping ranges of expression for each genotype (**Supplementary Fig. 5d**). This model requires 4 parameters to be trained corresponding to the expression range limits. Another simplification of the genotype prediction can be performed (**Supplementary Fig. 5e**), which requires only one parameter to be trained. In this model, the prediction only assigns uniform probability for homozygous genotypes when expression levels higher or lower than  $e_{mid}$  and assigns 0 conditional probability to the heterozygous genotypes, which brings up an important point: This simplified model is exactly the distribution that is utilized in the extremity based genotype prediction. In the extremity based prediction, we estimate  $e_{mid}$  simply as the mid-point of the range of gene expression levels within the expression dataset (Supplementary Note).

## 9.7 Code Availability

All the analysis code that is used to generate results can be obtained from <http://privaseq.gersteinlab.org>

## 10 METHODS ONLY REFERENCES