

Learning Your Identity and Disease from Research Papers: Information Leaks in Genome Wide Association Study

Rui Wang, Yong Fuga Li, XiaoFeng Wang, Haixu Tang, Xiaoyong Zhou
Indiana University Bloomington
Bloomington, IN

Abstract

Genome-wide association studies (GWAS) aim at discovering the association between genetic variations, particularly single-nucleotide polymorphism (SNP), and common diseases, which is well recognized to be one of the most important and active areas in biomedical research. Also renowned is the privacy implication of such studies, which has been brought into the limelight by the recent attack proposed by Homer et al. Homer's attack demonstrates that it is possible to identify a GWAS participant from the allele frequencies of a large number of SNPs. Such a threat, unfortunately, was found in our research to be significantly understated. In this paper, we show that individuals can actually be identified from even a relatively small set of statistics, as those routinely published in GWAS papers. We present two attacks. The first one extends Homer's attack with a much more powerful test statistic, based on the correlations among different SNPs described by coefficient of determination (r^2). This attack can determine the presence of an individual from the statistics related to a couple of hundred SNPs. The second attack can lead to complete disclosure of hundreds of participants' SNPs, through analyzing the information derived from published statistics. We also found that those attacks can succeed even when the precisions of the statistics are low and part of data is missing. We evaluated our attacks on the real human genomes and concluded that such threats are completely realistic.

Categories and Subject Descriptors

K.6.5 [Security and Protection]: Unauthorized access

General Terms

Security

Keywords

Genome Wide Association Study, Single Nucleotide Polymorphism, Test Statistics, Markov Model, Integer Programming

1. INTRODUCTION

The rapid advancement in genome technology has revolutionized the field of human genetics by enabling the large-scale applications

of genome-wide association study (GWAS) [7], a study that aims at discovering the association between human genes and common diseases. To this end, GWAS investigators determined the genotypes of two groups of participants, people with a disease (cases) and similar people without (controls) in an attempt to use statistical testing to identify genetic markers, typically single-nucleotide polymorphisms (SNP), that are associated to disease susceptibility genes [46]. If the variation of a SNP is found to be significantly higher in the case group than that in the control group, it is reported as a potential marker of the disease. Of great importance to such a study is privacy of the participants, whose sensitive information, personally identifiable genetic markers in particular, should not be leaked out without explicit consent. So far, this has been enforced through an informed consent from participants [9] and an agreement from investigators to ensure proper use of data according to the consent. Unfortunately, while this process prevents explicit misuse of participants' DNA data, it turns out to be insufficient for deterring information leaks in a more implicit way. Particularly, this paper reports a surprising finding of our research: even the test statistics computed over a small set of SNPs, like those routinely published in GWAS papers, could reveal a substantial amount of genetic information about participants, and even lead to disclosure of their identities.

The inadequacy of privacy protection in current genome research has also been pointed out by other researchers. For example, Mallin et al [45] show that even after removal of explicit identifiers (e.g., name, social security number), an individual could still be identified from a genetic database by examining the genetic markers related to her phenotypes (e.g., eye, skin and hair color). More seriously, Homer et al [39, 60, 5] recently proposes a statistical attack that could determine the presence of an individual in a group (e.g. the cases) from the aggregate *allele frequencies*, i.e. fractions of individual variants (called *allele*) occurred at each SNP site, of the whole group. Unlike the threats that have been extensively studied in statistical disclosure control [52, 27], Homer's attack takes advantage of the rich background information related to human genome available in the public domain, as well as the particular statistical properties of genomic data: it compares the victim's SNP profile against the allele frequencies of two populations, a "mixture" such as the case group in a GWAS and a reference population that can be acquired from public sources such as the International HapMap project [8]; given the profile of a sufficient number of independent SNPs from the victim (at least 10,000), her affiliation in the mixture can be determined with high confidence. The impact of this finding is significant. As an example, the NIH was reported to swiftly remove all aggregate data of GWAS, including allele frequencies, from public websites [60, 5].

Homer's attack made an important step towards better under-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CCS'09, November 9–13, 2009, Chicago, Illinois, USA.

Copyright 2009 ACM 978-1-60558-352-5/09/11 ...\$10.00.

Table 1: GWAS Terminologies used in this paper.

Terminologies	Description
Polymorphism	The occurrence of two or more genetic forms (e.g. alleles of SNPs) among individuals in the same population of a species.
Single Nucleotide Polymorphism(SNP)	The smallest possible polymorphism, which involves two types of nucleotides out of four (A, T, C, G) at a single nucleotide site in the genome.
Allele	One of the two sets of DNAs in a human individual’s genome. In a population, most SNP sites have two types of alleles. The one with higher frequency is defined as the major allele (denoted by 0), the other as a minor allele (denoted by 1).
Genotype	The combination of two set of alleles. For a SNP site with two common alleles in human population, there are three possible genotypes: two homozygotes, 00 and 11, and one heterozygote 01.
Locus(plural loci)	The surrounding regions of a SNP site in the genome .
Haplotype	Haplotype, also referred to as SNP sequence, is the specific combination of alleles across multiple neighboring SNP sites in a locus. Each individual has two haplotypes inherited from the parents. At the population level, some haplotypes are more common than others.
Linkage disequilibrium(LD)	Non-random association of alleles among multiple neighboring SNP sites.

standing of the privacy risks involved in publishing personal or aggregated genomic data. However, its impacts on GWAS remain uncertain, as most data released by a GWAS are test statistics such as p -values and r square (r^2) instead of allele frequencies used in the attack. Moreover, the attack needs a large amount of genome data, including a highly-dense genomic profile ($> 10,000$ independent SNPs) from the victim and their corresponding allele frequencies from a mixture, which can be hard to come by.

In this paper, we show that in the absence of proper protection, even a moderate disclosure of those test statistics, as did most GWAS papers, could pose a privacy risk that cannot be ignored. We present two attacks on these statistics. The first one could statistically identify an individual in the case group from a small set of statistics, which in some cases are only related to a single *locus*, the surrounding region in the genome where a disease susceptibility SNP is discovered in a GWAS. Like Homer’s attack, our technique also needs a reference population, which can be obtained from the reports of the same study conducted over other populations, and a SNP profile from the victim. Unlike the prior approach, however, our attack utilizes r^2 , a measurement of the correlation between SNPs (referred to as linkage disequilibrium (LD)), which is much more powerful than allele frequencies of single SNPs. As a result, the presence of an individual can be determined from the statistics involving only a couple of hundred SNPs at some locus. The second attack utilizes integer programming to analyze the pair-wise correlations of SNPs (measured by r^2 or D'), which in some cases could recover hundreds of participants’ SNPs. We also describe a technique that reverse engineers the statistics (e.g., p -values, r^2 and D') to calculate pair-wise allele frequencies, a necessary step in both of the attacks.

We believe that our paper makes the following contributions:

- Novel identification attacks on GWAS statistics.* We developed novel techniques to recover personally identifiable information from the test statistics published by GWAS papers. These techniques are powerful, capable of restoring hundreds of SNPs and identifying an individual using a much smaller number of SNPs than the prior attack [39]. This suggests that privacy threats in genome research are much more realistic than we thought.

- Study of countermeasures.* Our research shows that a simple countermeasure like publishing only coarse-grained statistics might not work well, as the relations among various statistics can still give away a sufficient amount of information for restoring fine-grained data. The strength of such a “correlation” attack demands a well-thought-out response that is built upon analysis of the connections among these statistics.

- Implementation and evaluations.* We implemented the proposed attacks and evaluated them on real human genome data.

The attack technologies we present in this paper, though powerful, are not without limitations. Specifically, our current approaches are not designed to directly work on *genotype*, the genetic constitution of individuals, but on *haplotype*, the SNP sequences derived from genotype. These two concepts are elaborated in Table 1, together with other genomic terminologies necessary for understanding this paper. With the maturity of genotype phasing techniques [57, 55, 19, 14, 25, 28], genome research today is increasingly based upon haplotypes, which offer more accurate descriptions of LD and therefore are more useful to the research. However, many previous GWAS studies use genotypes, which are more difficult to analyze. The feasibility of attacks on genotypes is left to our future research.

The rest of the paper has been organized as follows: Section 2 introduces the background knowledge; Section 3 elaborates the attacks; Section 4 reports evaluation results; Section 5 discusses the limitations of our techniques and possible defense; Section 6 surveys the related research and Section 7 concludes the paper.

2. BACKGROUNDS

2.1 GWAS: Steps and Test Statistics

A GWAS takes multiple steps to unravel the association between genetic variation and a common disease. Researchers first genotype participants from the case and control groups to extract a set of SNP profiles on selected sites. Usually two alleles can be found at each SNP site, referred to as the major and the minor alleles, denoted by 0 and 1. After proper quality control, the allele frequencies of these SNPs (i.e. the frequencies of 1 or 0) are computed over the case and control respectively. These frequencies are then used as inputs to an association test.

Association tests are used to detect the SNPs significantly associated to the cases (thus potentially to the disease) under the study. Examples of the tests include Pearson’s chi square (χ^2) [50], logistic regression [13], Fisher’s exact test [33] or Cochran-Armitage’s test for trends [15]. These tests are performed over the allele frequencies of the case and control groups, and their outcomes are used to calculate a p -value for each SNP. The SNPs with sufficiently small p -values, for example, below 10^{-7} , are selected as putative markers. Note that SNP-disease associations (the p -values) can be evaluated not only at the single SNP genotype level, but also at the level of haplotype (allele combinations involving multiple SNPs), which was shown a stronger statistical power than the genotype-based tests. To achieve this, a class of algorithms called the phasing algorithms [2] will be first applied to infer the most likely haplotypes of the individuals (two haplotypes for each individual) in the case and control group from their genotypes.

Table 2: GWAS Statistics. Here C_{pq} represents the count of an allele pq ($p, q \in \{0, 1\}$).

2 × 2 contingency table for Disease - SNP association study of two SNPs			2 × 2 contingency table for the linkage disequilibrium analysis of two SNPs			
SNP	Disease		SNP 1	SNP 2		Sum
	Control	Case	Alleles	0	1	
0 (Other allele)	C_{00}	C_{01}	0	C_{00}	C_{01}	C_{0*}
1 (Risk allele)	C_{10}	C_{11}	1	C_{10}	C_{11}	C_{1*}
			Sum	C_{*0}	C_{*1}	

Formula for related measures of association and LD		
Statistics	Formula	Asymptotic distribution
log odds ratios	$\log \frac{C_{00}C_{11}}{C_{01}C_{10}}$	Gaussian
r^2	$\frac{(C_{00}C_{11} - C_{01}C_{10})^2}{C_{0*}C_{1*}C_{*0}C_{*1}}$	Chi-square
D'	$\frac{C_{00}C_{11} - C_{01}C_{10}}{D_{max}}$	-

Once putative markers have been detected, the study often moves on to map their associations with other SNPs in the same loci, which is referred to as *linkage disequilibrium* (LD) [51]. This typically involves measurement of pair-wise allele frequencies (the frequencies of specific two-SNP allele combinations), and calculating statistics such as sensitivity indices (D') or coefficients of correlation (r^2) over the frequencies. These statistics can help identify other SNPs also related to the disease.

One last step of GWAS is to replicate the study on other case and control groups to verify whether the association between the SNP markers and the disease, as identified in prior steps, can also be observed from those populations.

Table 2 lists the formula for calculating aforementioned test statistics. These statistics are routinely published in GWAS papers [29, 54, 53, 57, 40]. Typically, p -values of tens or sometimes hundreds of SNPs have been reported. Thousands of r^2 or D' that reflect the LD among these SNPs are often illustrated in figures, and can be acquired from authors without any restrictions. Sometimes, detailed accounts of replication studies are also made public, which discloses multiple populations with identical allele-frequency distributions, and can therefore be used in Homer’s attack and our attack elaborated in the follow-up sections.

2.2 Homer’s Attack

A statistical attack recently proposed by Homer, et al [39] is believed to threaten the privacy assurance in current GWAS. Following we describe this attack on phased genotypes, i.e., haplotypes or SNP sequences.

In Homer’s attack, the attacker is assumed to already have a high-density SNP profile of the victim, which can be extracted from a small amount of blood sample. This assumption is realistic, as the cost of genotyping is becoming increasingly affordable [4]. What the attacker wants to determine is the presence of an individual in the case group, an indicative of her contraction of a disease. To this end, the attacker measures the distances between the allele frequency of every SNP j on the profile, $Y_j \in \{0, 1\}$, and the corresponding frequencies in a reference Pop_j , and a mixture M_j respectively. These distances are used to compute the following statistic:

$$D(Y_j) = |Y_j - Pop_j| - |Y_j - M_j| \quad (1)$$

Assuming that the distributions of individual SNPs’ allele frequencies in the mixture and the reference population are identical, $D(Y_j)$ will have the same distributions across all independent SNP j . As a result, their sum, according to the central limit theorem, will converge to a normal distribution. The mean of the distribution is zero if the victim is not in the case group, and positive otherwise.

Using statistical hypothesis testing, the authors found that 25,000 SNPs of a member in a mixture built from the HapMap offered a p -value below 10^{-6} , given the null hypothesis that she does not belong to the mixture.

3. OUR ATTACKS

3.1 From Statistics to Allele Frequencies

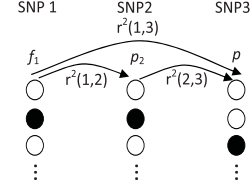


Figure 1: Recover allele frequencies.

A SNP contains a major allele, denoted by 0, and a minor allele, 1. Their individual frequencies and the allele frequencies of SNP pairs (00, 01, 10, and 11) contain a large amount of information. The former is the main ingredient of Homer’s attack, and both are needed in our attacks, as elaborated in Section 3.2 and 3.3. A GWAS typically reports the allele frequencies for the SNPs identified as putative genetic markers for a disease. For other SNPs, however, only their p -values and LD statistics, r^2 or D' , are published. Therefore, the first problem that an attacker needs to solve is how to recover the frequencies from these statistics.

This step is by no means trivial. For example, though the mapping between a p -value and the outcome of an association test is typically one-on-one, there could be multiple frequencies that correspond to that outcome. For example, using Pearson’s chi square, a p -value of 0.01 was found to associated with 80 possible SNP frequencies in a population studied in our research. The problem is further complicated by the fact that published statistics typically have low precisions, which makes determination of their input frequencies even more difficult.

On the other hand, the statistics published by a GWAS are often connected: for example, the p -values of individual SNPs are bridged by the r^2 between them. This relation can be leveraged to reveal their allele frequencies, as elaborate below.

Inference of single SNP allele frequencies. We found that the allele frequencies of single SNPs can be restored by “propagating” a genetic marker’s frequencies to other SNPs through their linkage disequilibrium, often described by r^2 or D' . Loosely speaking, such a correlation tells us how likely one SNP of an individual can be used to infer some of her other SNPs. The r^2 between two SNPs, as illustrated in Table 2, can also be calculated as follows:

$$r^2 = \frac{(C_{00}N - C_{*0}C_{0*})^2}{C_{0*}C_{1*}C_{*0}C_{*1}} \quad (2)$$

where N is the size of the population, C_{00} is the counts of the pairs of the SNPs’ major alleles, C_{0*} and C_{1*} are the counts of the first SNP’s major and minor alleles respectively, and C_{*0} and C_{*1} are the allele counts of the second SNP. In the case that the first SNP is actually a marker, whose counts (C_{0*} and C_{1*}) are known, C_{*0} and C_{*1} , which are interdependent, can be determined by the r^2 once C_{00} is known. Actually, C_{00} also relates to C_{*0} , C_{0*} and N . Using these relations as constraints, we can find solutions for C_{*0} and C_{*1} through constraint solving [1]. This can be done efficiently, as the searching space here is bounded by N^2 .

It is possible, however, that those constraints are satisfied by more than one solution. To make the solution unique, we need to

find more constraints. An immediate one is the p -value of the second SNP, which is computed over the SNP's allele counts in both case and control populations. Other constraints come from the relations among SNPs, as illustrated in Figure 1. Consider three SNPs, S_1 , S_2 and S_3 , where S_1 is a marker. Besides the aforementioned constraints that exist in pairs (S_1, S_2) and (S_1, S_3) , r^2 for (S_2, S_3) can also be leveraged. Similarly, more constraints can be found by looking into the correlations among more SNPs. This approach was demonstrated to be very effective in our research: it completely recovered single SNP frequencies from the statistics with only moderate precisions (r^2 rounded to 2 decimal places, given a population of 200 individuals). Though we discuss our technique here using r^2 , the same approach can also be applied to D' .

Recovery of pair-wise frequencies. From the allele frequencies of individual SNPs, pair-wise frequencies can be directly calculated. Specifically, solving Equation 2 with r^2 , C_{0*} , C_{1*} , C_{*0} and C_{*1} gives us C_{00} . Note that there are actually two solutions for C_{00} in Equation 2, for positive and negative r respectively. However, typically only one of them is an integer, the legitimate value for C_{00} . The counts of other pair-wise alleles, C_{01} , C_{10} and C_{11} , are found from the following linear equations:

$$\begin{aligned} C_{0*} &= C_{00} + C_{01} \\ C_{1*} &= C_{10} + C_{11} \\ C_{*0} &= C_{00} + C_{10} \\ C_{*1} &= C_{01} + C_{11} \end{aligned} \quad (3)$$

Inaccurate statistics. A practical hurdle for our attacks is that the published statistics are often of low-precision. This can be handled by changing Equation 2 to an inequality, giving r^2 a range of acceptable values. Specifically, we use the following constraint to find single-SNP allele counts:

$$L < \frac{(C_{00}N - C_{*0}C_{0*})^2}{C_{0*}C_{1*}C_{*0}C_{*1}} < U \quad (4)$$

where L and U are the lower and upper limits of r^2 respectively.

Signs. An important piece of information for our statistical attack (Section 3.2) is *sign*, which is determined by the equality $C_{00}C_{11} > C_{01}C_{10}$: it is positive if the inequality holds, and negative otherwise. It is conceivable that signs are much easier to recover than allele frequencies, which can actually be used to compute the signs. In our research, we first ran our constraint solver on $C_{00}C_{11} > C_{01}C_{10}$, together with other constraints for inferring frequencies, and then on $C_{00}C_{11} < C_{01}C_{10}$. The sign is recovered if no solution is found in one of these two cases.

3.2 A Statistical Attack

We follow the strategy proposed by Homer et al. to design our statistical attack. The goal of our attack is to determine the likelihood of a victim to be in a case group of a GWAS study based on one given SNP sequence of a victim. To achieve this goal, we first establish a *reference* group (Figure 2), consisting of the SNP sequences from a group of individuals, drawn from a reference population with the same genetic background as the case group. The International HapMap project [8] provides a large source of samples for this exercise, containing individuals from various ethnic groups, including Nigeria (Yoruba), Japan/China and US residents with ancestry from Northern and Western Europe. In Section 4.2, we will show that the statistical power of our attack relies on the selection of the reference group. However, even when the reference group does not completely mimic the genetic background of the case group, the attack still works, although with a lower power.

Once the reference group is established, we propose a hypothesis

test on the SNP sequence of the victim to determine her presence in the case group. A high confidence estimated by this test indicates that the victim's SNP sequence is significantly closer to the reported LD structure (measured by the pair-wise r^2 values) than any other sequence in the reference group, and it is unique enough for identifying the victim. Below we describe this test statistic in details.

Test statistic. Given N sequences of M neighboring SNPs in the genome, we define the *signed* allele correlations r_{ij} between two SNPs i and j ($1 \leq i < j \leq M$) as $r_{ij} = \frac{C_{11}C_{00} - C_{01}C_{10}}{\sqrt{C_{1*}C_{0*}C_{*1}C_{*0}}}$, where C_{pq} is the pair-wise allele counts, i.e. the number of haplotypes with allele p ($p \in \{0, 1\}$) at SNP i and allele q ($q \in \{0, 1\}$) at SNP j , and $C_{p*} = C_{p0} + C_{p1}$ and $C_{*q} = C_{0q} + C_{1q}$.

The signed allele correlations can be computed solely from a set of given SNP sequences, and thus can be easily computed from the reference group (denoted as r^R , Figure 2). On the other hand, although the signed allele correlations of the case individuals (denoted as r^C) usually are not reported in GWAS papers (r^2 is published instead), they can be derived from the reported data, i.e., r^2 with signs recovered from constraint solving (Section 3.1). To test on the SNP sequence of the victim, $H = h_1h_2\dots h_m$ (where $h_i \in \{0, 1\}$ is an allele of SNP i), we use an "indicator" of the allele of a SNP pair ij , Y_{ij}^{pq} , where pq is an allele pair ($p, q \in \{0, 1\}$) of the SNP pair. Specifically, if $h_i = p$ and $h_j = q$, $Y_{ij}^{pq} = 1$; otherwise, $Y_{ij}^{pq} = 0$. Note that for a specific SNP pair, only one $Y^{pq} = 1$ and the other three $Y^{pq} = 0$ among four possible allele pairs (11, 10, 01 and 00). Now we are ready to define the hypothesis test statistic, T_r as,

$$\begin{aligned} T_r &= \sum_{1 \leq i < j \leq N} T_{ij} \\ &= \sum_{1 \leq i < j \leq N} (|(Y_{ij}^{00} + Y_{ij}^{11}) - (r_{ij}^R + 1)/2| \\ &\quad - |(Y_{ij}^{00} + Y_{ij}^{11}) - (r_{ij}^C + 1)/2|) \\ &= \sum_{1 \leq i < j \leq N} (r_{ij}^C - r_{ij}^R) \cdot (Y_{ij}^{00} + Y_{ij}^{11} - Y_{ij}^{01} - Y_{ij}^{10}) \end{aligned} \quad (5)$$

The statistic is designed to make the test of an individual's presence in the case group *valid*. Intuitively, this means that the contribution of a case participant's SNP sequence to T_r needs to be *positive* so as to statistically distinguish her relation with the case from that with the reference she is not in. Formally, based on the null hypothesis that the to-be-tested sequence (of the victim) is not within the case group, we have $E(T_r) = 0$; however, if the sequence is indeed in the case, the expected contribution of this instance to a specific signed allele correlation r_{ij}^C , i.e., $\frac{\partial T_{ij}}{\partial C_{h_i h_j}}$, is non-negative for any pair of SNPs i and j and any pair of alleles 11, 10, 01 or 00 (see Appendix 1 for a proof). Therefore, the sum statistic of T_r is valid. Note that because the signed allele correlation r ranges from -1 to 1 , for the validity of the test statistic, we use $(r + 1)/2$ to map it to a value between 0 and 1. In the Appendix, we also prove that T_r is close to an optimal test statistic assuming r approximately follows a normal distribution. Notably, though the test statistic proposed here has a similar form as the one proposed by Homer et al. [39], it sums over $\binom{m}{2}$ variables (i.e. signed allele correlations) instead of m independent SNPs, and hence, as we show in Section 4.2, it results in a much more powerful attack.

Markov model estimation. However, since the signed allele correlations (r_{ij}) are not completely independent, we cannot assume the distribution of T_r under the null hypothesis is normal. With the limited size of available reference population from the HapMap project, we resorted to the Markov chain modeling and sampling techniques to simulate the data for estimating the confidence of our statistical attack. Markov models [36, 43, 47] have been used ex-

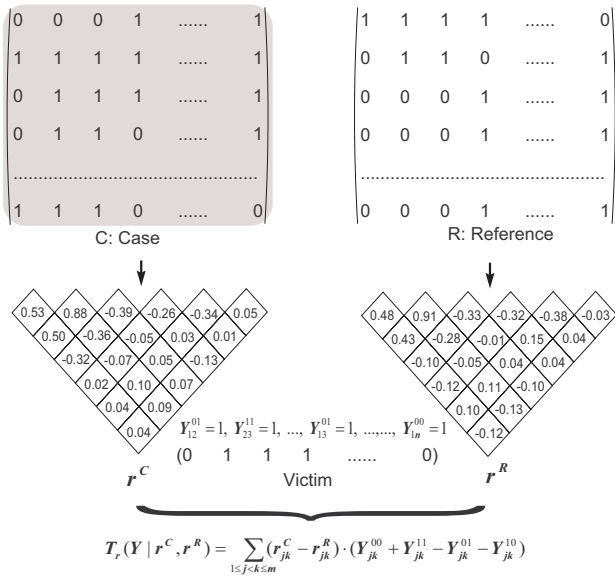


Figure 2: A schematic illustration of a statistical attack on the reported r^2 values between pairs of SNPs. Assume the attacker has obtained the haplotype of a victim and his goal is to determine how likely this victim is within a case group in a GWAS study. To achieve this goal, the attacker will compute a test statistics (T_r) that measures the distance between two set of signed allele correlations, r^C derived from the reported (r^2) values computed from the haplotypes of the individuals in the case group (left panel), and r^R computed directly from haplotypes of the reference group (right panel). T_r has the same property as the test statistic T_p proposed by Homer et al.: it is significantly larger when the victim is within the case group than otherwise. Using an empirically deduced distribution of T_r given the null hypothesis that the victim is not in the case group, the attacker can then estimate the confidence level of identifying the victim.

tensively in modeling SNP sequences. In our research, we constructed an inhomogeneous 1st order Markov model by a maximum likelihood approach from a limited number of haplotypes obtained from HapMap website (as the training dataset). The model parameters, including one set of initial probabilities and $m - 1$ set of transition probabilities, can be directly estimated based on the counts for single SNPs and the neighboring SNP pairs in the training dataset. Once the Markov model is built, it can be used to draw unlimited number of haplotypes with a similar genetic background as the ones in the training dataset. In each simulation, at least 1000 case and reference groups (each with 50-1000 haplotypes) are sampled, and the test statistics T_r is computed for each of these paired case/reference groups. Using this technology, we studied how the power of the test statistics is affected by various GWAS aspects such as the SNP sequence of the victim, the case group size, the resemblance between the reference and case group, the precision of reported r^2 and missing values of r^2 .

Reference population. A practical attack needs a reference population that preferably resembles the case group in the distribution of r . Such a population can come from various sources, for example, replication studies on different populations, or reuse of the case data for a new study. Among these, of a particular interest is the possibility that one could actually acquire the genomes of the reference population through signing an agreement. Though privacy

agreements legally bind researchers against revealing the genomic data under study, it does nothing to stop an attacker using the collected data as a reference population to compromise the privacy of individuals involved in other datasets, for which she has no obligation. Actually, in Section 4.2, we demonstrate that even a bad reference, with an r distribution deviated from the case group, can allow us to identify some individuals with high confidence.

Encoding nucleotides. A problem an attacker has to tackle is how to translate the victim's nucleotides (A, T, C, G) into alleles (0 or 1). These alleles are determined from the population studied in a GWAS: for the two nucleotide values a SNP can take, the popular one is encoded as 0, and the other as 1. The attacker, who does not have direct access to the population, needs to figure out how to encode the victim's nucleotides before launching the attack. A solution is to utilize the HapMap to identify major alleles, which is often very effective, as observed in our research.

3.3 An Integer-Programming Attack

Given pair-wise allele frequencies for a whole locus that involves tens of or even a couple of hundred SNPs, the idea of recovering individuals' SNPs becomes really enticing. We believe that this can actually be done with proper techniques. In this section, we report such an attack based upon integer programming.

Illustrated in Figure 3 is the sequences of N individuals, each containing M SNPs. These SNPs form an N by M matrix, below which are their LDs. Our attack attempts to recover the whole or part of the matrix from the LDs that are in the form of pair-wise allele frequencies. To this end, we designed a "divide-and-conquer" approach described below.

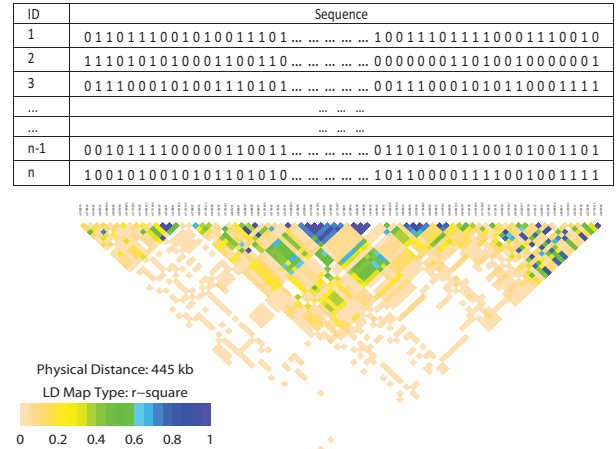


Figure 3: Recover individuals' SNPs

Human DNA information is passed from one generation to another in a way that recombinations occur much more frequently between two SNP blocks, called *haplotype block*, than within a block. As a result, a haplotype block typically contains only a small number of combinations (i.e., haplotypes) of SNP values. To leverage this property, our attack partitions a SNP sequence into blocks according to their LDs, to ensure that the SNPs on the same block have strong connections. Then, integer programming is utilized to find out the haplotypes within individual blocks that satisfy the constraints of pair-wise frequencies and single-SNP frequencies. Finally, haplotypes in different blocks are connected based upon the LDs between these blocks.

Recovering individual blocks. GWAS participants' SNP sequences within a block are actually described by the counts of different haplotypes. Therefore, the first stage of our attack focuses on determin-

Sequence			Step 1			Step 2			Step 3		
0	0	0	$\begin{cases} x_0 + x_1 = 1 \\ x_2 + x_3 = 1 \\ x_4 + x_5 = 1 \\ x_6 + x_7 = 2 \\ x_0 + x_4 = 1 \\ x_1 + x_5 = 1 \\ x_2 + x_6 = 2 \\ x_3 + x_7 = 1 \\ x_0 + x_2 = 1 \\ x_1 + x_3 = 1 \\ x_4 + x_6 = 2 \\ x_5 + x_7 = 1 \end{cases} \Rightarrow \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \end{pmatrix} = \begin{pmatrix} \bar{a} \\ 1 \\ 1 \\ -1 \\ 1 \\ -1 \\ -1 \\ 1 \end{pmatrix} x_7 + \begin{pmatrix} \bar{b} \\ 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 2 \\ 0 \end{pmatrix}$			$\begin{aligned} \min x_7 &\Rightarrow x_7 \in [0,1] \\ \max x_7 &\Rightarrow \begin{cases} x_0 \geq 0 \\ x_1 \geq 0 \\ x_2 \geq 0 \\ x_3 \geq 0 \\ x_4 \geq 0 \\ x_5 \geq 0 \\ x_6 \geq 0 \\ x_7 \geq 0 \end{cases} \Rightarrow \begin{cases} -x_7 + 1 \geq 0 \\ x_7 \geq 0 \\ x_7 \geq 0 \\ -x_7 + 1 \geq 0 \\ x_7 \geq 0 \\ -x_7 + 1 \geq 0 \\ x_7 + 2 \geq 0 \\ x_7 \geq 0 \end{cases} \end{aligned}$			$\begin{cases} \text{if } x_7 = 0, \text{solution} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 2 \\ 0 \end{pmatrix} \\ \text{if } x_7 = 1, \text{solution} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 1 \end{pmatrix} \end{cases}$		
1	0	1									
1	1	0									
0	1	1									
1	1	0									
Input											
x_0	000										
x_1	001										
x_2	010										
x_3	011										
x_4	100										
x_5	101										
x_6	110										
x_7	111										

Figure 4: Recover individual blocks

ing those counts. Let $[x_0 \cdots x_{2l-1}]$ be a vector of integer variables that represent the counts of individual haplotypes, where l is the number of the SNPs on a block. Note that individual SNP can only have two values, and thus 2^l is the total number of different haplotypes in the block. To find a solution for the vector, our approach takes the following three steps:

- Step 1: We first build a system of linear equations upon the equality constraints derived from pair-wise allele counts and the number of the participants. Figure 4 describes an example that works on a 3-SNP block. The block can have total 8 possible haplotypes whose counts are denoted by $\vec{X} = [x_0 \cdots x_7]$. Given 12 pair-wise frequencies (each pair can take 4 alleles), we can build 12 linear equations: the first equation represents the count of the allele “00” for the first two SNPs, the second is the count of “01” for the same SNP pair, and other pair-wise counts are described by the remaining equations. Such a system is then solved using Gaussian elimination, and its solution set can be represented in a parametric form that includes a set of free variables. In the example, the solutions are described by $\vec{X} = \vec{a}x_7 + \vec{b}$, where \vec{a} and \vec{b} are two constant vectors as illustrated in Figure 4, and x_7 is a free variable.

- Step 2: The ranges of the free variables are determined by integer programming. We first add in inequality constraints that require all the solutions to be non-negative. Then, for every free variable, a pair of integer programming problems are solved to minimize and maximize its value under the constraints. This gives us the acceptable range of the variable. Figure 4 displays the inequality constraints for the example, under which $\min x_7$ and $\max x_7$ reveal that the variable can take 0 or 1.

- Step 3: In the presence of multiple free variables, we need to exhaustively search their value ranges to find a combination that satisfies all the constraints. This is an exponential problem. However, its scale can be controlled by adjusting block size. Every combination found is fed into the parametric solution of the linear equation system, which gives us the haplotype counts we are looking for.

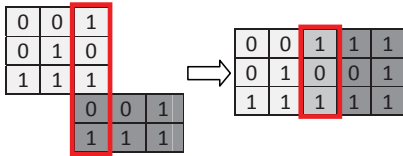


Figure 5: Connect two blocks through overlapping.

Connecting different blocks. After restoring haplotypes for individual blocks, we move on to link different blocks together based on the LDs between them. The SNPs in different blocks typically have weak correlations. However, the aggregate connections from

individual SNP pairs can be strong between two large blocks. We took a strategy in our research that first merges blocks with strong LDs into a large block, and then bridges it to other blocks. To connect two blocks, we again use a vector of variables to represent the counts for individual combinations of the haplotypes from different blocks. For example, consider two 3-SNP blocks, each with two haplotypes: (“001” and “011”) for one and (“100” and “110”) for the other; our approach generates 4 variables to represent the counts for (“001” “100”), (“001” “110”), (“011” “100”) and (“011” “110”) respectively. In general, two blocks with m and n different haplotypes respectively bring in mn variables. The solutions for those variables are computed through integer programming, as does the step for recovering individual blocks.

A problem arises when the number of haplotypes in each block is large, which makes the number of variables even larger. As a result, the time integer programming takes to find a solution can increase exponentially. In our research, we adopted a simple technique to mitigate this problem. For the part of a SNP sequence where LDs are weak, our approach cuts blocks in a way that allows two blocks to share a set of SNPs. Those SNPs let us look at the common part of the haplotypes from different blocks, and as a result, help reduce the number of variables. Consider two blocks, A and B , with an overlap that includes k haplotypes. For each haplotype i ($1 \leq i \leq k$), suppose that it is attached to m_i different haplotypes in A and n_i in B . This gives us at most $\sum_k m_i n_i$ different ways to connect the haplotypes from different blocks. Figure 5 presents an example in which two blocks with 3 and 2 haplotypes respectively have an overlap involving 2 haplotypes (“0” and “1”). We only need 3 variables to describe all possible haplotype combinations between these blocks, instead of 6.

Identification of an individual. Once we obtain a solution (a set of SNP sequences), we could compute the corresponding test statistic T_r (based on the r values for the case, and a reference sample available) and estimate the identification confidence of each of the sequences. If one sequence receives a small p -value, it is probably a correct solution, and the individual with the sequence (known to us now) is probably in the case group; on the other hand if a sequence receives a large p -value, the sequence could either be wrong or not unique enough to determine a person. Actually, a sequence with a high identification confidence can reveal a lot of information about its owner. For example, if it happens to contain the SNPs related to observable phenotypes, Malin’s attack [45] could be applied to identify the individual associated with the sequence, even if the attacker does not have the victim’s DNA profile *a priori*.

3.4 Simple Defense

Low-precision statistics. As described in Section 3.1, our approach recovers allele frequencies by propagating genetic markers’

frequencies to other SNPs through their LDs. The markers' frequencies are among the most important outcomes of a GWAS and therefore have to be released. What the defender can do here is to downgrade the precisions of LD statistics. This, unfortunately, is often insufficient for blocking the information that can be used to recover allele frequencies. The fundamental problem here is that all the statistics, p -values of individuals SNP and r^2 for SNP pairs, are correlated. Such correlations, together with the marker's frequencies and the size of a case group, can make up for the information loss caused by coarse-grained statistics. Consider the example in Figure 1: the LD between SNPs (S_1, S_2) , affects the relation between (S_2, S_3) , given the constraints of the total number of participants; the LDs of (S_1, S_2) and (S_2, S_3) further constrain that of (S_1, S_3) . Our approach leverages such relations and therefore is very robust to inaccurate statistics. An experimental study reported in Section 4 shows that given r^2 rounded to 2 decimal places (only 1 decimal place for r), we still restored more than 50% of pair-wise allele frequencies and all the signs.

On the other hand, the attacks we propose do not rely on fine-grained statistics. Our statistical attack only needs r rounded to 1 decimal place to establish one's presence in the case (Section 4.2). The integer-programming attack can still work in the absence of some frequency constraints, though this can result in multiple solutions. These properties of our attacks were evaluated through an experimental study, which is reported in Section 4.

Thresholds. Publishing less data can certainly make it more difficult for an attacker to infer sensitive information. However, it equally renders GWA papers less informative. An obvious solution is to use a threshold to remove the data deemed insignificant to the research. The question is how to set such a threshold. For example, in Figure 3, all the r^2 values below 0.01 are dropped. However, we can still figure out haplotype frequencies for individual haplotype blocks using the r^2 within blocks, and connect different blocks together by running a maximum likelihood estimator over the remaining LDs between blocks. Moreover, if the remaining data contains sufficient information for recovering signs, our statistical attack still works.

We believe that a one-size-fits-all threshold does not work for GWAS. Techniques need to be developed to assess the outcomes of individual research to compute a customized threshold that enables dissemination of the findings of a study without compromising the privacy of its participants.

Noises. Adding carefully-designed noises to genome data can impede our analysis. Particularly, one can append individual pair-wise SNP counts with decimal components, which removes the integer constraint we use to uniquely identify C_{00} from the two solutions of Equation 2 (Section 3.1). However, this can still be insufficient to defeat our attack. Again, the constraints imposed by the relations among multiple SNP pairs can still enable us to pinpoint the counts closest to satisfying the constraints. Biocomputing tools such as R [10] use maximum likelihood methods to compute pair-wise allele frequencies directly from genotypes. As a result, these frequencies are inconsistent with each other, and with population size and p -values, which mitigates the threat of our attack. However, these frequencies are inaccurate and therefore much less useful to genome research than those computed from haplotypes.

4. EVALUATIONS

This section reports an experimental study of the techniques we propose. Like the prior work [39], our study was based upon real haplotypes from the HapMap project [41] (<http://www.hapmap.org/>). More specifically, we used SNP sequences at

Table 3: Infer Frequencies and Signs

Statistics Precision		Recovered Information %		
r^2	p -value	single SNP frequency	pair-wise frequency	sign of r
0.1	0.1	12.1	1.8	6.7
0.1	0.00001	40.6	11.7	31.7
0.01	*	100	50.1	98.7
0.001	*	100	90.4	100
0.0001	*	100	95.1	100

FGFR2 locus (around SNP rs1219648) in the HapMap phase 3 release 2, which covers 200kb region from SNP rs12354864 (human reference genome b36 location 123189345) to SNP rs7900009 (human reference genome b36 location 123450068). The locus was chosen because it was recently reported in a GWAS paper [40] to be associated with the risk of sporadic postmenopausal breast cancer, and the linkage disequilibrium plot using r^2 for the 174 SNPs in the locus is directly available in the paper (as did routinely in most GWAS articles), which makes it a proper target of our attack. Note the real haplotypes used in our evaluation are *not* from the individuals involved in that GWAS [40], but instead from the HapMap database, due to privacy concerns.

In our experiment, we randomly chose a population of 200 with Africa backgrounds, half as cases and other half as controls. Association statistics of the population were first computed, including pair-wise r^2 of the case group and p -values of individual SNPs estimated from Pearson's chi-square. Our attacks were evaluated over these statistics, together with the knowledge of the sizes of the case and control groups and the allele frequencies of a SNP treated as a genetic marker (SNP rs1219648). Such information is typically reported in a GWAS paper and its supplementary materials, or can be conveniently acquired from authors of the paper. Following we describe this evaluation study.

4.1 Inferring Allele Frequencies and Signs

The first step of our attack was to infer the allele frequencies for both individual SNPs and SNP pairs from the statistics. Our approach started with the marker and propagated its frequencies to other SNPs through their r^2 . During this process, a constraint solver (cream [1]) was used to search for the unique solutions for single and pair-wise allele counts within the population under the constraints posed by the r^2 , the p -values and the population sizes, as described in Section 3.1. The solver also recovered the signs of r , if they could be uniquely determined. In our experiment, we evaluated this approach against the statistics with various precisions. The outcomes are illustrated in Table 3.

The experiment demonstrates the prowess of our technique. Given an r^2 rounded to 2 decimal places, which amounts to 1 decimal place for r , we successfully recovered all the single-SNP frequencies, half of pair-wise frequencies and almost all the signs for r , regardless of the precisions of p -values. 95% of pair-wise frequencies were identified when the precision of r^2 reached 4 decimal places. An interesting observation is that p -values could make up for the loss of precision in r^2 : an r^2 of 1 decimal place, once paired with a p -value of 5 decimal places, could be used to restore more than 30% of signs. On the other hand, GWAS papers typically offer p -values rounded to 4 to 5 decimal places (often in the supplementary materials [53, 49, 59, 37]), and r^2 even to 8 decimal places [40]. This is more than enough for our attack.

4.2 Identifying Individuals

Evaluation of the power of T_r based on simulations. Based on reference groups drawn from the first order inhomogeneous Markov chain (see section 3.2), we estimated the power of our statistical attack. If not mentioned otherwise, the simulation model is built

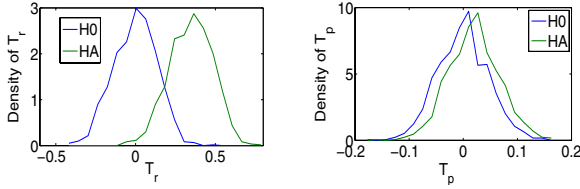


Figure 6: Comparison of estimated probability density functions of the test statistic on pairwise r values (T_r , left) vs. the test statistic based on single SNP frequencies (T_p , right). The probability density is estimated based on 1000 rounds of simulations with 1 fixed victim, 200 cases and 200 references.

from the FGFR2 loci with 230 SNP sequences of the YRI population (Yoruba in Ibadan, Nigeria) from the HapMap database. We compared our attack with the one proposed by Homer et al. [39], and found that the statistical attack based on T_r is much more powerful than the attack on single SNP profiles, which uses the statistic $T_p = \sum_j D(Y_j)$, where $D(Y_j)$ is defined in equation 1 (see Figure 6). Let H_0 be the null hypothesis (the victim is not in the case group), and let H_A be the alternative hypothesis (the victim is in the case group). Note that the power of the statistic is higher if the overlapping portion of the two distributions (under H_0 and H_A , respectively) is smaller (Figure 6). We estimated the distribution of H_0 by sampling the victim and the cases independently, and the distribution of H_A by including the victim’s SNP sequence into the sampled case group. For 200 sequences of 174 SNPs at the FGFR2 locus in each of the case and control groups, T_r can identify 80% individuals in the case group at fixed type I error of 0.05 (or 95% confidence), while T_p can only identify around 9%. Note that the lower bound of the power for any statistic equals type I error (5%); thus we can pessimistically estimate that T_r correctly identified at least 75% of case individuals while T_p only worked out 4%, indicating that T_r is about 19 times more powerful. For T_p to reach a similar power, we found that around 30 times more SNPs were required (i.e. 5000 SNPs).

Using the simulation based on the Markov model, we found that T_r approximately follows a normal distribution (Figure 6), and the *standard deviation* of the distribution under the null hypothesis is nearly constant with respect to various SNP sequence of the victim and various ethnic group the cases belong to, as long as the sample size and number of SNPs are fixed. This implies that we may not need to re-derive the null distribution for every test. In a set of simulations with 100 SNPs, 200 cases and 200 references, we found that the power of our attack varied from 30% to 80% (with average 62% and standard error 11%) depending on the SNP sequences of the cases. When we looked into the ethnic group of the cases, we found that the power for cases from Africa (population YRI, Yoruba in Ibadan, Nigeria) was higher than those from central Europe or east Asia (data not shown), indicating that some individuals (from some ethnic groups) are easier to be identified than others, presumably because they carried more sensitive SNPs than the others. Finally, we examined the power of the attack on cases with various sizes, ranging from 50 to 1600 (table 4), which revealed that even when the case was very large (e.g. 1600, larger than typical GWAS studies), there were still a significant proportion (18.1%) of cases who could be identified confidently.

We evaluated our attack against the *low precision* approach (only providing low precision values for r in the GWAS paper) and the *threshold* defense (removing r values below the threshold). We found the performance of the statistical test was very robust to both countermeasures. At a very strong defense level, e.g. threshold 0.1 for $|r|$, or the precision level of 1 decimal place, majority of the

Table 4: Dependence of the statistical power of the T_r on sample size based on the simulation on the FGFR2 locus. N is the number of SNP sequences in the case groups (the same as the number in the reference group).

N	50	100	200	400	800	1600
power (%)	99.9	85.7	67.2	40.4	36.2	18.1

Table 5: Percentage of statistical power (at 0.05 Type I error) left at various precision of input data r^C . The power is estimated based on 1000 rounds of simulated attacks. The number of cases and controls are both 200. The individual SNP sequences were randomly drawn from the inhomogeneous Markov Chain built on 230 SNP sequences of the FGFR2 locus from the HapMap phase 3 YRI (Africa) population.

Precision of r^C	0.5	0.2	0.1	0.01	0.001
% power π left	12	74	85	100	100

statistical power (79% and 85% respectively) was retained. The results are shown in Table 5.

The power of realistic attacks. The simulations described above can generate a large set of SNP sequences, which enabled us to systematically analyze the power of our statistical attack. A further step is to evaluate the performance of our attack under a realistic setting. Three types of attacks were experimented in our research (Figure 7). In each of them, three groups (C - case group, R - reference group, and T - test group) of real SNP sequences of FGFR2 locus were retrieved from different populations in HapMap database (see section 3.2 for details); then, for each SNP sequence H from *Case*, *Reference*, and *Test*, T_r was computed using H and the r values from *Case* and *Reference* (see figure 2). Note that the individuals in the test group were completely independent from the case and the reference, and therefore the mean of their T_r was zero. In the first two attacks, *Case* and *Reference* individuals were taken from a same population (YRI, Yoruba in Ibadan, Nigeria, or JPT+CHB, Japanese in Tokyo, Japan, and Han Chinese in Beijing, China). This mimics the situation where we have a good reference sample that resembles the case group very well. In the third experiments, *Case* (from YRI) and *Reference* (from ASW - African ancestry in Southwest USA) were taken from a different but related population. This mimics the situation where we do not have a good reference sample, but have to use an *average* reference available from a public data source. For each experiment, with the case and reference SNP sequences fixed, T_r was computed for each SNP sequence in the case, reference, and the test groups.

The results from these realistic attacks are promising. For the two experiments with good references, many individuals in the case group received higher T_r values than those from other groups (Figure 7A), suggesting that many SNP sequences in the case group can indeed be identified with relatively high confidence, with p -values from 0.01 to 10^{-5} (Figure 7C). While for the situation where an average reference was used, there was a drop of the discrimination power (Figure 7B vs. Figure 7A right), however there was still a significant amount of power left, making the attack still possible. Note that in many GWAS studies, data for more than one loci are provided, which makes the attack even more powerful.

To estimate the confidence (p -values) for each victim, we only need to estimate the variance of T_r under the null hypothesis. One approach is to build a Markov chain from the reference sample (which is known to the attacker), and then estimate the variance by the simulation experiment (see 3.2 for details). A second approach is to estimate the variance directly from T_r values obtained from individuals of the test sample. We found in our research that the second approach gave relatively more a conservative estimate, whereas

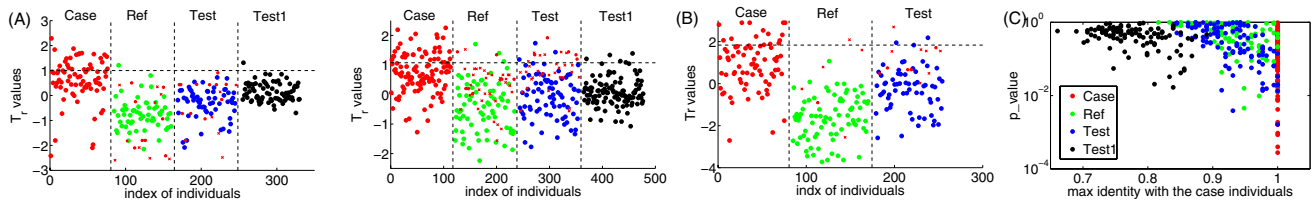


Figure 7: Evaluation of the power of the statistical attack. (A) Realistic attacks using good references. SNP sequences retrieved from Hapmap database are used as cases (*Case*, red dots), references (*Ref*, green dots) as well as the tests (*Test*, blue dots and *Test1*, black dots). Individuals in *Ref/Test1* are marked if they have 100% (*) or > 95% (x) identity with some individuals in *Cases*. In the left plot, *Case/Ref/Test* SNP sequences are taken from *YRI* (Yoruba in Ibadan, Nigeria), while *Test1* haplotypes are from a different population *JPT+CHB* (Japanese in Tokyo, Japan, and Han Chinese in Beijing, China). In the right plot, *Case/Ref/Test* are from *JPT+CHB*, with *Test1* from *YRI*. (B). Realistic attacks using average references. Legend should be interpreted the same as in A except that there is no *Test1*. Here cases are taken from *yri* population, whereas references are taken from a different (but related) population *ASW* (African ancestry in Southwest USA) (C). Estimated p -values for the attack shown in the right plot of A. The variance under the null distribution were estimated by T_r values from the groups of *Test* and *Test1*.

the first approach offered an over-optimistic estimation that might lead to a lot of false positives. Figure 7C gives an example of estimated p -values for the attack with a good reference. One interesting observation is that, many non-case individuals could also receive relatively small p -values, especially when they were not completely independent from some case individuals (Figure 7C). This is expected, as essentially, our statistical attack is to test the (genetic) independence of the victim from the case individuals.

4.3 Recovering SNPs

We implemented the Integer-Programming attack using Matlab, based upon two toolboxes, *rref* [6] for solving systems of linear equations, and *bintprog* [3] for integer programming. This implementation was run on our dataset to recover the SNP sequences for all 100 individuals from their pair-wise allele frequencies. In the attack, we partitioned the sequences into small blocks according to the LDs of the locus, as demonstrated by the r^2 . Block size varied around 10 SNPs. Haplotypes within individual blocks were first restored through solving linear-equation systems and integer programming, and then connected together using overlapping blocks in accordance to the LDs, as described in Section 3.3. This attack was run on a system with 2.80GHz Core 2 Duo processor and 3GB memory. Within 12 hours, we successfully restored the 174 SNPs for all 100 participants. This demonstrates that the threat of the integer programming attack is realistic.

5. DISCUSSION

A GWAS can either analyze individuals' genotypes or phase them into haplotypes before the analysis happens. Attack on genotypes is more challenging, as we can no longer use the integer constraint and consistent LDs to recover pair-wise allele frequencies. Nevertheless, some of the techniques we propose, particularly the way to estimate single allele frequencies and the statistic attack, can still work. Further study of such an attack is left to our future research. On the other hand, the maturity of phasing technologies makes the use of haplotypes, which carries more information, increasingly a trend [57, 55, 19, 14, 25, 28]. This makes our attack a realistic threat.

The attacks we propose could be defeated by well-planned countermeasures. A potential approach can be adding noise to published dataset. This, however, needs to be done carefully, because the noise can undermine the scientific values of a paper, making others hard to repeat the experiments reported by a GWAS. One technique we can use is to adjust some of the published LDs to the extent that the signs of r are just changed. This requires a careful selection of a set of SNP pairs, whose LDs are not significantly affected by

switch of signs. Examples of the SNPs are those with very close $C_{00}C_{11}$ and $C_{01}C_{10}$. Also important here is the assurance that other unperturbed LDs do not give an attacker sufficient statistical powers to identify an individual. Another approach is to selectively remove some data. For example, dropping some LD statistics can interrupt the path for propagating markers' frequencies, and as a result, makes allele frequencies more difficult to recover. Fundamentally, effective defense against our attacks can be achieved by enforcing differential privacy: that is, the presence of an individual in the case group changes nothing but negligible statistical features of the group. This will be investigated in our follow-up research.

6. RELATED WORK

The problem of releasing aggregate data in privacy-preserving ways has been intensively studied in the areas like privacy-preserving data analysis [30, 34], statistical disclosure control [18, 20, 32], inference control [24], privacy-preserving data mining [11, 12], and private data analysis [31, 48].

Privacy problems in GWAS, however, are related to special structures of genome data (linkage disequilibrium) and background information (reference populations), which have not been studied in prior research. Recently, research has been conducted on privacy-preserving genome computing [16, 42, 22]. Those approaches are more to do with preventing a party from accessing sensitive data, than protecting sensitive information from being inferred from the outcome of a computation. A recently proposed concept highly related to our attacks is *differential privacy* [30]. Loosely speaking, the concept ensures that removing or adding an individual's record to a database does not substantially changes the statistics calculated from the database. So far, this has been achieved through adding noise [17, 31, 20]. Such an approach, if carefully designed, could mitigate the threat of the attacks we propose.

Recovering individuals' SNP sequences is essentially the problem of contingency table release, which has been studied in statistics community for decades [21, 38, 23, 58, 26] in the context of census. Consider an individual's record in a database as a row, which consists of k binary attributes. A contingency table is a vector that describes, for each combination of k attributes (called a *setting*), the counts of the rows satisfying this setting. Typically, only the counts (called *marginals*) of different settings are published, instead of the table. What an attacker wants to do is to infer the table from the marginals.

In our integer programming attack, individual SNP frequencies and pair-wise frequencies are actually such marginals. However, existing techniques [58] cannot be directly applied to recover SNPs from them, because of the scale of our problem: up to our knowl-

edge, prior research can only handle as many as 16 attributes [26], while our attack needs to work on 174 attributes. This was achieved in our research by taking advantage of special properties of genomic data. Our “divide-and-conquer” approach made full use of the correlations among different SNPs, which can be observed from published r^2 . More specifically, we first partitioned a SNP sequence into small blocks according to their LDs, which significantly increases the chance for restoring these blocks because the LDs within the same blocks are strong; then, individual blocks were connected using their aggregate correlations. The techniques designed for this purpose, solving linear equations and use of overlapping blocks in particular, are novel, up to our knowledge.

Few researcher has looked into the privacy risks in GWAS [44, 35]. The most relevant work other than Homer et al. is the association study based on the *pooled* genotypes, i.e. the aggregate genotype profiles of a group of cases and controls (rather than for each individual of them) were used in the SNP marker discovery, e.g. by Yang et.al. [56]. These approaches, though adopting a similar experimental setting, aim to address a distinct problem as our approach, that is, to retrieve SNPs (or SNP sequences) significantly associated with the disease, which tend to be the common genetic features of the disease population rather than the features specific to the individuals in the case group, and hence cannot be used for identification of individuals.

7. CONCLUSION

GWAS is among the most active research areas in biomedical research. It is also the area fraught with privacy concerns. The recent work by Homer et al. [39] demonstrates that personal identification is plausible by analyzing a large number of allele frequencies related to GWAS. The privacy threat of this kind has been found in our research to be even more realistic than expected. In this paper, we report two new techniques that can lead to identification of the participants of a GWAS from a small set of statistics, as those routinely published in GWAS papers. One of the techniques can statistically determine the presence of an individual in the case group, based upon the LDs among as few as a couple of hundred SNPs. The other attack can even recover all participants' SNP sequences related to the statistics. We also show that these attacks work on coarse-grained statistics. Our experimental study further justifies the concerns of such threats, which were shown to be capable of cracking statistics computed from real genome data.

A further step in this important direction is study of potential attack techniques that work on genotypes. On the defense side, we believe that research in statistical disclosure control, differential privacy in particular, can offer an effective guideline to mitigate and ultimately eliminate the privacy threat to GWAS.

8. ACKNOWLEDGEMENTS

We thank anonymous reviewers for their valuable comments. This research has been supported in part by the Center of Applied Cybersecurity Research at the Indiana University.

9. REFERENCES

- [1] Cream: Class library for constraint programming in java. <http://bach.istc.kobe-u.ac.jp/cream/>, 2003.
- [2] Haplotype estimation and association. http://slack.scr.man.ac.uk/theory/association_hap.html, 2005.
- [3] bintprog: solve binary integer programming problems. <http://www.mathworks.com/access/helpdesk/help/toolbox/optim/index.html?access/helpdesk/help/toolbox/optim/ug/bintprog.html>, 2008.
- [4] Company to sequence 1,000 human genomes in 2009 for \$5,000 each. <http://www.medicalnewstoday.com/articles/124572.php>, 2008.
- [5] Nih background fact sheet on gwas policy update. http://grants.nih.gov/grants/gwas/background_fact_sheet_20080828.pdf, 2008.
- [6] rref: Reduced row echelon form. <http://www.mathworks.com/access/helpdesk/help/techdoc/index.html?/access/helpdesk/help/techdoc/ref/rref.html>, 2008.
- [7] Genome-wide association studies. <http://grants.nih.gov/grants/gwas/>, 2009.
- [8] International hapmap project. <http://www.hapmap.org/>, 2009.
- [9] Policy for sharing of data obtained in nih supported or conducted genome-wide association studies (gwas). <http://grants.nih.gov/grants/guide/notice-files/NOT-OD-07-088.html#principles>, 2009.
- [10] The r project for statistical computing. <http://www.r-project.org/>, 2009.
- [11] D. Agrawal and C. C. Aggarwal. On the design and quantification of privacy preserving data mining algorithms. In *PODS '01: Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 247–255, New York, NY, USA, 2001. ACM.
- [12] R. Agrawal and R. Srikant. Privacy-preserving data mining. *SIGMOD Rec.*, 29(2):439–450, 2000.
- [13] A. Agresti. *Categorical Data Analysis*. New York: Wiley-Interscience, 2002.
- [14] C. I. Amos et al. Genome-wide association scan of tag snps identifies a susceptibility locus for lung cancer at 15q25.1. 40(5):616–622, 2008.
- [15] P. Armitage. Tests for linear trends in proportions and frequencies. *Biometrics*, 11(4):375–386, 1955.
- [16] M. J. Atallah, F. Kerschbaum, and W. Du. Secure and private sequence comparisons. In *WPES '03: Proceedings of the 2003 ACM workshop on Privacy in the electronic society*, pages 39–44, New York, NY, USA, 2003. ACM.
- [17] B. Barak, K. Chaudhuri, C. Dwork, S. Kale, F. McSherry, and K. Talwar. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *PODS '07: Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 273–282, New York, NY, USA, 2007. ACM.
- [18] L. L. Beck. A security mechanism for statistical database. *ACM Trans. Database Syst.*, 5(3):316–3338, 1980.
- [19] T. Becker and C. Herold. Joint analysis of tightly linked snps in screening step of genome-wide association studies leads to increased power. 2009.
- [20] A. Blum, C. Dwork, F. McSherry, and K. Nissim. Practical privacy: the sulq framework. In *PODS '05: Proceedings of the twenty-fourth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 128–138, New York, NY, USA, 2005. ACM.
- [21] C. E. Bonferroni. Teoria statistica delle classi e calcolo delle probabilit . *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, 8(1):3–62, 1936.
- [22] F. Bruekers, S. Katzenbeisser, K. Kursawe, and P. Tuyls. Privacy-preserving matching of dna profiles. Technical Report Report 2008/203, ACR Cryptology ePrint Archive, 2008.
- [23] Y. Chen, P. Diaconis, S. P. Holmes, and J. S. Liu. Sequential monte carlo methods for statistical analysis of tables. *Journal of the American Statistical Association*, 100:109–120, 2003.
- [24] F. Y. Chin and G. Ozsoyoglu. Auditing and inference control in statistical databases. *IEEE Trans. Softw. Eng.*, 8(6):574–582, 1982.
- [25] T. DA et al. Genome-wide haplotype association study identifies the slc22a3-lpa2-lpa gene cluster as a risk locus for coronary artery disease. *Nature Genetics*, 41(3):283–285, February 2009.
- [26] A. Dobra and S. E. Fienberg. Bounds for cell entries in contingency tables induced by fixed marginal totals. *Statistical Journal of the United Nations ECE*, 18:363–371, 2001.

[27] J. Domingo-ferrer. Advances in inference control in statistical databases: An overview. In *In Inference Control in Statistical Databases: From Theory to Practise* (2002, pages 1–7. Springer Verlag, 2002.

[28] A. Doring et al. Slc2a9 influences uric acid concentrations with pronounced sex-specific effects. 40(4):430–436, 2008.

[29] R. H. Duerr et al. A genome-wide association study identifies il23r as an inflammatory bowel disease gene. *Science*, October 2006.

[30] C. Dwork. Differential privacy. In *in ICALP*, pages 1–12. Springer, 2006.

[31] C. Dwork, F. Mcsherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, 2006.

[32] S. E. Fienberg. Datamining and disclosure limitation for categorical statistical databases. In *Proceedings of Workshop on Privacy and Security Aspects of Data Mining, Fourth IEEE International Conference on Data Mining (ICDM 2004*, pages 1–12. Nova Science Publishing, 2004.

[33] R. Fisher. On the interpretation of chi-square from contingency tables, and the calculation of p. *Journal of the Royal Statistical Society*, 85(1):87–94, 1922.

[34] J. Gehrke. Models and methods for privacy-preserving data analysis and publishing. In *ICDE '06: Proceedings of the 22nd International Conference on Data Engineering*, page 105, Washington, DC, USA, 2006. IEEE Computer Society.

[35] H. T. Greely. The uneasy ethical and legal underpinnings of large-scale genomic biobanks. *Annu Rev Genomics Hum Genet*, 8:343–364, 2007.

[36] G. Greenspan and D. Geiger. Modeling haplotype block variation using markov chains. *Genetics*, 172(4):2583–2599, Apr 2006.

[37] H. Hakonarson et al. A genome-wide association study identifies k1a0350 as a type 1 diabetes gene. 448(7153):591–594, 2007.

[38] W. Hoeffding. Scale-invariant correlation theory. *Masstabinvariante Korrelationstheorie, Schriften des Mathematischen Instituts und des Instituts für Angewandte Mathematik der Universität*, 5:179–233, 1940.

[39] N. Homer, S. Szelinger, M. Redman, D. Duggan, W. Tembe, J. Muehling, J. V. Pearson, D. A. Stephan, S. F. Nelson, and D. W. Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLoS Genet*, 4(8):e1000167+, 2008.

[40] D. J. Hunter et al. A genome-wide association study identifies alleles in fgfr2 associated with risk of sporadic postmenopausal breast cancer. *Nature Genetics*, 39(7):870–874, May 2007.

[41] International HapMap Consortium, K. A. Frazer, et al. A second generation human haplotype map of over 3.1 million snps. *Nature*, 449(7164):851–861, Oct 2007.

[42] S. Jha, L. Kruger, and V. Shmatikov. Towards practical privacy for genomic computation. In *2008 IEEE Symposium on Security and Privacy*, 2008.

[43] Y. Kim, S. Feng, and Z. B. Zeng. Measuring and partitioning the high-order linkage disequilibrium by multiple order markov chains. *Genet Epidemiol*, 32(4):301–312, May 2008.

[44] J. E. Lunshof, R. Chadwick, D. B. Vorhaus, and G. M. Church. From genetic privacy to open consent. *Nat Rev Genet*, 9(5):406–411, May 2008.

[45] B. Malin. Re-identification of familial database records. In *AMIA Annu Symp Proc. 2006*, 2006.

[46] M. McCarthy, D. Smedley, and W. Hide. New methods for finding disease-susceptibility genes: impact and potential. *Genome Biology*, 4(10):119, 2003.

[47] A. P. Morris, J. C. Whittaker, and D. J. Balding. Little loss of information due to unknown phase for fine-scale linkage-disequilibrium mapping with single-nucleotide-polymorphism genotype data. *Am J Hum Genet*, 74(5):945–953, May 2004.

[48] K. Nissim, S. Raskhodnikova, and A. Smith. Smooth sensitivity and sampling in private data analysis. In *STOC '07: Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 75–84, New York, NY, USA, 2007. ACM.

[49] D. G. I. of Broad Institute of Harvard, L. U. MIT, N. I.

of BioMedical Research, et al. Genome-Wide Association Analysis Identifies Loci for Type 2 Diabetes and Triglyceride Levels. *Science*, 316(5829):1331–1336, 2007.

[50] R. L. Plackett. Karl pearson and the chi-squared test. *International Statistical Review / Revue Internationale de Statistique*, 51(No.1):59–72, 1983.

[51] R. Robbins. Some applications of mathematics to breeding problems iii. *Genetics*, 3(4):375–389, 1918.

[52] N. E. Schulte. Statistical disclosure control (sdc) in practice: some examples in. *Statistical Journal of the United Nations Economic Commission for Europe*, 18(4):321–328, 2001.

[53] L. Scott et al. A genome-wide association study of type 2 diabetes in finns detects multiple susceptibility variants. *Science*, April 2007.

[54] R. Sladek et al. A genome-wide association study identifies novel risk loci for type 2 diabetes. *Nature*, February 2007.

[55] J. Sun et al. Evidence for two independent prostate cancer risk-associated loci in the hnf1b gene at 17q12. 40(10):1153–1155, 2008.

[56] H. C. Yang, C. C. Pan, C. Y. Lin, and C. S. Fann. Pda: Pooled dna analyzer. *BMC Bioinformatics*, 7:233–233, 2006.

[57] M. Yeager et al. Genome-wide association study of prostate cancer identifies a second risk locus at 8q24. *Nature Genetics*, 39(5):645–649, April 2007.

[58] I. H. D. Yuguo Chen and S. Sullivant. Sequential importance sampling for multiway tables. *The Annals of Statistics*, 34(1):523–545, 2006.

[59] E. Zeggini et al. Replication of Genome-Wide Association Signals in UK Samples Reveals Risk Loci for Type 2 Diabetes. *Science*, 316(5829):1336–1341, 2007.

[60] E. A. Zerhouni and E. G. Nabel. Protecting aggregate genomic data. *Science*, 322(5898):44a+, October 2008.

APPENDIX

1. Proof of the validity of T_r

Given the definition of the signed allele correlation r ,

$$r_{ij}^C = \frac{C_{11}C_{00} - C_{01}C_{10}}{\sqrt{C_{1*}C_{0*}C_{*1}C_{*0}}} \quad (6)$$

we have

$$\begin{aligned} \frac{\partial T_{ij}}{\partial C_{00}} &= \frac{\partial r^C \cdot Y^{00}}{\partial C_{00}} \\ &= \frac{C_{11} \cdot [1 - \frac{1}{2} \cdot (1 - \frac{C_{10}C_{01}}{C_{11}C_{00}}) \cdot (\frac{C_{00}}{C_{0*}} + \frac{C_{00}}{C_{*0}})]}{\sqrt{C_{1*}C_{0*}C_{*1}C_{*0}}} \geq 0 \end{aligned} \quad (7)$$

for any pair of SNPs i and j . Similarly, we can get $\frac{\partial T_{ij}}{\partial C_{11}}$, $\frac{\partial T_{ij}}{\partial C_{10}}$ and $\frac{\partial T_{ij}}{\partial C_{01}}$ are also non-negative for any pair of SNPs i and j .

2. Proof of the optimality of the T_r statistics

Given $r = \{\hat{r}_{ij}|i, j\}$ and assuming the signed allele correlations follow the normal distribution, the optimal statistic for H_0 : victim is in Case vs. H_A : victim is not in Case is $T_r = \sum_{i,j} \{2\mu_0(\hat{r}_{ij}^C - r_{ij}^C) + \frac{1}{n} \{[\hat{r}_{ij}^C - (\mu_0 + r_{ij}^C)]^2 - [\mu_0^2 + \frac{(1-(r_{ij}^C)^2)}{n-1}]\}\}$ by applying the Neyman-Pearson lemma, where $\mu_0 \approx \frac{\partial r_{ij}^C}{\partial C_{00}}Y_{00} + \frac{\partial r_{ij}^C}{\partial C_{01}}Y_{01} + \frac{\partial r_{ij}^C}{\partial C_{10}}Y_{10} + \frac{\partial r_{ij}^C}{\partial C_{11}}Y_{11}$, n is the number of individual in the case group, r_{ij}^C is the theoretical signed allele correlations, and \hat{r}_{ij}^C is the estimated signed allele correlations measure. By omitting the second terms in the formula, which is $1/n$ times smaller than the first term, and replacing $\frac{\partial r_{ij}^C}{\partial C_{**}}$ by their signs (the real values need pair-wise frequencies to compute, which are unknown here), we have $T_r = 2\sum_{i,j} [(Y_{00} + Y_{11} - Y_{01} - Y_{10}) \cdot (\hat{r}_{ij}^C - r_{ij}^C)]$. In practice, we estimate \hat{r}_{ij}^C by using a reference group, i.e. \hat{r}_{ij}^R .