

Information theory-based measures and
privacy-preserving file formats for sensitive information
leakage from raw functional genomics data

GG et al

April 9, 2018

Abstract

Functional genomics experiments on human subjects present a privacy conundrum. On one hand, many of the conclusions we infer from these experiments are not tied to the identity of individuals but represent universal statements about disease and developmental stages. On the other hand, by virtue of the experimental procedures, the reads from them are tagged with small bits of patients' variant information, which presents privacy challenges in terms of data sharing. There is great desire to share the data as broadly as possible. Therefore measuring the amount of variant information leaked in a variety of experiments, particularly in relation to the amount of sequencing will allow us to uncover ways of reducing the information leakage, and determine an appropriate setpoint for sharing information with minimal leakage. To this end, we aimed to derive information theoretic measures for the private information leaked in experiments and develop various file format manipulations to reduce much of the leaked variants. We showed that high depth experiments such as Hi-C provide accurate genotyping that can lead to large privacy leaks. Counterintuitively, noisy and partial genotypes from low-depth experiments such as ChIP-Seq and single-cell RNA-Seq, although not useful genotypes, can be used as strong quasi-identifiers for re-identification purposes through linking attacks. We showed that these incomplete genotypes can further be used to construct an individual's complete variant set and inference of individual identifying phenotypes when combined with imputation. We then provide a proof-of-concept theoretical framework, in which the amount of leaked information can be estimated from the depth and breadth of the coverage as well as the sequencing bias of the functional genomics experiments. In order to solve the dilemma between data sharing and privacy leakage, we propose a file formatting system that enables the sharing data while protecting individuals sensitive information and preserving the utility of the data. The proposed file format can achieve different levels of privacy and utility balance. At the highest level of privacy, our file format masks all the variant information leaked from reads, which can be used to calculate signal profiles with 99% recovery of the original profiles and 100% recovery of the original gene expression levels.

1 Introduction

With the decreasing cost of DNA sequencing technologies, the number and the size of available genomic data have exponentially increased and become available to a wider group of audiences such as hospitals, research institutions and individuals [1]. Availability of genetic information gives rise to privacy concerns; for instance, genetic predisposition to diseases may bias insurance companies or create unlawful discrimination by employers [4]. In turn, privacy of individuals has become an important aspect of biomedical data science [2, 3].

Early genomic privacy studies focused on the identification of individuals in a mixture by using phenotype-genotype association [5, 6]. These studies showed that private information of an individual, such as participation in a drug-abuse study, can be revealed [5, 6]. With the increase of large-scale genomics projects such as the Personal Genome Project [7] or recreational/direct-to-consumer genomic databases, researchers showed that multiple datasets can be linked together to infer sensitive information such as participant's surnames [8] or addresses [9]. Such cross-referencing relies on quasi-identifiers, which are pieces of information that are not unique identifiers by themselves but are well correlated with unique identifiers or can be unique identifiers when combined with other quasi-identifiers [10].

Functional genomics experiments provide a wealth of information on genomic activities related to developmental stages or diseases that are essential for personalized medicine. These studies use large-scale high-throughput assays to quantify transcription (RNA-Seq) [11], epigenetic regulation (ChIP-Seq) [12] or the three-dimensional (3D) organization of genome (Hi-C) [13] in a genome-wide fashion under different conditions (e.g., samples from patients and healthy individuals). Inferring biological information from functional genomics experiments is a multi-step procedure, in which progressive summarization of the data from raw sequencing reads to the gene quantifications, transcription factor (TF) binding peaks or chromatin interaction matrices is per-

formed. Although activities of the functional genome are not necessarily tied to an individual's genotype, reads from these experiments are derived from the biosamples that belong to individuals; hence, they are tagged with individuals' variants. Public sharing of such raw data raises privacy concerns. In order to share high-utility data while preserving individuals' sensitive information, it is essential to determine a "set point", after which trade-off between the utility of the data and the privacy risk is balanced. A hurdle in determining the set point is the lack of systematic quantification of private information leakage from functional genomics data. Figure 1 summarizes the processing steps of RNA-Seq experiments as an example of how summarization decreases the risk of privacy while greatly decreasing the amount of sharing and the utility of the functional genomics data. In detail, functional genomics data analysis starts with the generation of DNA/RNA sequencing reads that are stored in a special file formats called FASTQ [14]. These files are large in size ranging from 5 GB up to 60 GB depending on the purpose of the experiment. They are then mapped to human reference genome and stored as compressed binary file types called binary alignment map (BAM) and/or compressive alignment map (CRAM) that are derived from the sequence alignment map (SAM) files in text format [15]. File formats such as CRAM have been developed to remedy the ever increasing amount of data; CRAM provides up to a ten-fold decrease when information loss is tolerated [16]. Further summarization of the mapped reads (such as signal profiles or gene expression quantification) still allows researchers to make accurate biological conclusions, while providing ~20-fold further data reduction. Although overall aggregation and averaging reduces biological information, private information leakage also decreases (Figure 1).

In particular, read alignment files (SAM / BAM / CRAM) are of great interest due to the large amount of biological data they provide, as they constitute the most important input of the majority of genome annotation pipelines. However, these files contain sequence information of the individuals that may leak sensitive data. Depending on the depth of the functional genomics experiment, raw reads can be used to identify private single nucleotide polymorphisms (SNPs), small inser-

tions and deletions (indels), and structural variants. However, current policies related to the public sharing of the BAM files are somewhat ad-hoc. For example, for the genome of the HeLa cell line, the raw reads from Hi-C experiments require special access. By contrast, reads from ChIP-Seq and RNA-Seq experiments are publicly available [17]. That is, reads from the experiments that do not require substantial depth are sometimes considered to be safe to share without privacy concerns, owing to partial and biased sequencing. However, it is not clear that these reads are leakage free. Although private information leakage from summary-level functional genomics data have been quantified previously [18, ?, ?] the lack of a systematic quantification of private data leakage from BAM files makes it difficult for biomedical data sharing policymakers to protect individuals' sensitive information in a consistent fashion. The CRAM format provides the option for the users to convert BAM files into lossy compression, in which quality scores of the alignments are manipulated. This, in turn, can be used to decrease private information leakage [16]. However, privacy leaks still occur due to the containment of mismatched information of the reads with respect to reference genome [16]. The mapped read format (MRF) was introduced as a conceptual format to remedy privacy concerns; in this keeping the sequence of reads is optional [21]. This does not only reduces the size of the data, but also makes it hard to genotype the individuals from the information in these files. However, private information leakage is not entirely removed from MRF files, as one can still infer deletions from the information in these files. Moreover, current quantification pipelines used for gene expression analysis as well as the peak calling softwares were not designed to take MRF files as inputs.

On the flip side of the coin is the utility of the mapped reads (BAM files) and challenges related to dealing with private data. Access to private data requires use agreements that have expiration dates and a tremendous amount of bureaucracy connected to them. Moreover, any secondary data product becomes private and cannot be distributed. Problems associated with the distribution of secondary data products from private biomedical data is exacerbated due to large file sizes. For

example, genome annotations that are derived from private functional genomics data require the establishment of their own databases. However, because such annotations are derived from private data, establishment and distribution of these databases require extra levels of privacy-related bureaucracy. Another example of the challenges associated with private data is that big consortia such as the Encyclopedia of DNA elements (ENCODE) [22], the Cancer Genome Atlas (TCGA) [23] or the Genotype-Tissue Expression project (GTEx) [24] are funded to enable a collaborative working environment through dedicated phone calls and meetings. In turn, participants have to go through required access procedures with their institutions. Otherwise, communication based on private data is prohibited according to data use agreements. Moreover, when multiple institutions have required access to the same data, they still cannot exchange files with each other. These challenges create a bottleneck and hinder the progress of important biomedical findings. Open data helps the advancement of biomedical data science not only by easing access to the data, but also by helping with speedy assessment of tools and methods, and in turn, reproducibility. Funding agencies and research organizations are increasingly supporting new means of data sharing and new requirements for making data publicly available while preserving participants' privacy [25]. In an attempt to consider both sides of the coin, we aimed to determine how much information is enough information to identify individuals and how we can protect the sensitive information with minimal loss of utility in a public data sharing mode. To this end, we derived novel information theory-based measures and applied these measures to quantify the amount of leaked information in various functional genomic assays from ENCODE [22] and other sources [?] at varying coverage. Based on our findings, we developed new file formats that allow the public sharing of read alignments of functional genomics experiments, while protecting the sensitive information and minimizing the amount of private data that requires special access and storage. Our file format manipulation system achieves different levels of privacy versus utility balance with an adjustable parameter.

In this study, we used an individual (NA12878) as a case example and their 1000 genomes genotypes as the gold standard [26]. We sampled reads from the sequencing data of functional genomics experiments at increasing coverage, and detected SNVs and indels using Genome Analysis Toolkit (GATK) best practices recommendations [27, 28]. We propose a new metric for quantifying the amount of information that can be obtained from sequencing data with respect to the gold standard. We next present a simple and practical instantiation of a linking attack with the assumption of adversaries accessing an increasing amount of the sequencing data. We show that individuals are vulnerable to identifications even at small coverage of sequencing data. We further show that with summation of reads from functional genomics experiments and imputation through linkage disequilibrium, the leaked number of variants can reach the total number of variants in an individual’s genome. We then provide a theoretical framework where the amount of leaked information can be estimated from depth and breadth of the coverage as well as the bias of the experiments. Finally, we focus on ways to publicly share alignment data without compromising an individual’s sensitive information. We propose privacy-enhancing file formats that hide variant information, are compressed, and have a minimal amount of utility loss.

2 Results

2.1 Information Theory to quantify private information in an individual’s genome

An individual’s genome can be represented as a set of variants. Each variant is composed of the chromosome to which it belongs, location on that chromosome, the alternative allele, and the corresponding genotype. Let $S = \{s_1, s_2, \dots, s_i, \dots, s_N\}$ be the set of variants. Then each variant can be represented as $s_i = \{v_i, g_i\}$, where v_i consists of the location and alternative allele information and g_i denotes the genotype of the variant as 1 for a heterozygous variant and 2 for a homozygous

variant. Note that we calculate the information with respect to reference genome, therefore $g_i = 0$, where the nucleotide in an individual's genome equal to the reference nucleotide is not considered.

We can then calculate the naive self-information of S in bits as

$$h(S) = - \sum_{i=1}^{i=N} \log_2(p(s_i)). \quad (1)$$

In eq. 1, N is the total number of variants in an individual's genome, $p(s_i) = n_i/n_T$ is the genotype frequency, in which n_i is the number of individuals with variant $s_i = \{v_i, g_i\}$ and n_T is the total number of individuals in the panel. Note that we denote $h(S)$ as "naive" information because it is an estimate of the real information in a situation, in which the population to which the individual belongs is unknown and the number of individuals are finite. Eq.1 holds true only if variants are independent of each other, which is not the case due to the correlation between variants in linkage disequilibrium (LD). In theory, the population to which the individual belongs to can easily be predicted by using a few variants. However, from an adversary's perspective, this will add one more layer of calculation (i.e., computational and time costs) to the identification attack. Eq.1 is also an estimate of the information when we consider all the individuals in the world (i.e., $\lim_{n_i \rightarrow \infty} h(S)$).

To understand whether naive information is a good estimate, we first calculated the information with the consideration of LD scores taken from the European population of the HapMap project [29]. LD scores are pairwise correlations between variants, which we consider as the prior information on the existence of a variant given other variants in the same LD block exist in a genome. Then, the information with LD consideration is calculated as

$$h^{LD}(S) = - \sum_{i=1}^{i=N} (1 - mLD(s_i, s_j)) h(s_i) \quad (2)$$

$mLD(s_i, s_j)$ is the maximum LD correlation of variant s_i with other variants such that $mLD(s_i, s_j) =$

$$\max_{i \neq j, j \in (1, \dots, N)} LD(s_i, s_j), \text{ where } mLD(s_i, s_j) \neq mLD(s_j, s_i).$$

Figure 2a shows a negligible difference between the naive information and information with LD consideration for the NA12878 genome. To understand the lack of difference better, we calculated the self-information of each variant in an LD block with and without LD consideration. We found that highly informative variants do not exhibit any difference due to the low LD correlations (Figure 2b). We further show that the number of variants with differences between information with or without LD consideration is small compared to the number of variants low LD correlations on average (SI Figure X). This also shows that information ($h(S)$) is driven by the rare variants.

We then estimated the information when the population size is infinite [30]. We sampled fractions on the order of 10%, 20%,..., 100% individuals from the 1000 Genomes Phase I panel (total of 2504 individuals) and calculated the information using the sampled distribution of genotypes. We repeated this calculation 100 times and calculated the mean information for each sampled fraction. The relationship between the inverse of the sample fraction and the information fit best to a power function with two terms ($y = ax^b + c$, $R = 0.99$). The y -intercept of the curve is the extrapolation of information when the population size approaches infinity ($1/\infty = 0$, Figure 2c). We again found a negligible difference between the naive information and the information when the population size is infinite (Figure 2a). We also calculated the information by starting from a single individual and adding individuals one by one to the population (SI Figure 1). These individuals were simulated using the genotype frequencies in the 1000 genomes panel and the LD information from the HapMap project (see SI methods). Both the information calculation and the KL -divergence between different-sized populations show that as the size of the population increases, the difference in the information decreases and eventually becomes negligible (SI Figure 1)

The calculations above show that naive information can be an accurate approximate to the private information content of an individual's genome when the individual's population is unknown and the population size is bound by the number of individuals in 1000 Genomes panel due to the

relationship of information at $n \rightarrow \infty \geq$ naive information \geq information with LD (Figure 2a). That is, an adversary with no prior knowledge of the population of the sample and limited number of individuals in a known genotype panel can accurately approximate the private information.

2.2 Information Theory to quantify private information leakage in functional genomics data

We next aimed to understand the relationship between the leaked information and the coverage to make a fair comparison. We sampled c amount of total nucleotides from the 24 different functional genomic experiments and from whole genome sequencing (WGS) and whole exome sequencing (WES) data of sample NA12878 (see SI Table 1). We used GATK to call SNVs and indels with the parameters and filtering suggested in the GATK best practices [27, 28]. We used the genotypes in the 1000 Genomes panel for NA1278 as the gold standard. We used “naive” point-wise mutual information (pmi) as a measure to quantify the association between the gold standard and the called variants. If $S^G = \{s_1^*, \dots, s_i^*, \dots, s_M^*\}$ is the set of variants from the gold standard and $S^F(c) = \{s_1, \dots, s_i, \dots, s_M\}$ is the set of variants called from the c total sequencing coverage of a functional genomics experiment, then the set $A = S^G \cap S^F(c)$ contains the variants that are called and are in the gold standard set. If $A = \{a_1, \dots, a_i, \dots, a_T\}$, then

$$pmi(S^G; S^F(k)) = - \sum_{i=1}^{i=T} \log_2(p(a_i)) \quad (3)$$

We then added more coverage to the sampled coverage and repeated the calculation. We repeated this procedure until we deplete all the reads of a functional genomics experiment. The overall process is depicted in Figure 2e.

2.3 Private information leakage in 24 functional genomics experiments with different coverage

We calculated pmi values for 24 functional genomics experiments with different coverage. The experiments involved whole genome approaches such as Hi-C, transcriptome-wide assays such as RNA-Seq, and targeted assays such as ChIP-Seq of histone modifications and transcription factor binding. In addition, we calculated the pmi for WGS, WES, and SNP-ChIP for comparison (Figure 3).

As expected, the Hi-C data contained almost as much information as the WGS data and more information than the SNP ChIP array data. The WGS data contained more information than the Hi-C data at the beginning of the sampling process. As we sampled nucleotides between 1.1 and 10 billion bps, the information content of the Hi-C data surpassed the WGS data (Figure 3a). We speculate that this is due to a higher quality of genotyping of the genomics regions that are in spatial proximity, as Hi-C has a bias of sequencing more reads from those regions. As expected, we could not infer as much information from the ChIP-Seq reads (Figure 3b). Surprisingly, many of the ChIP-Seq assays such as the ones targeting CTCF and RNAPII contained a large amount of information at low coverage. Furthermore, comparison between WES and different RNA-Seq experiments showed that none of the RNA-Seq experiments contained as much information as the WES data; this is due to the fact that RNA-Seq captures reads only from expressed genes in a given cell type (Figure 3c). An unexpected observation was that more information could be inferred from polyA RNA-Seq data at low coverage compared to WES and total RNA-Seq data. To make a fair comparison between each of these assays, we calculated the mean pmi per base pair depicted in Figure 3d. To do so, we normalized the pmi values by the amount of coverage (c). We then averaged each by the number of times (n) we performed sampling on that experiment ($\frac{\sum pmi(S^F(c); S^G)}{n}$). The Hi-C and ChIP-Seq experiments targeting the transcription factor HDGF provided more genotyping information per base pair compared to the WGS data. The RNA-Seq

experiments provided the least genotyping information per base pair (Figure 3d).

2.4 Genotyping accuracy

In light of our finding that genotyping can be performed using low-depth, biased functional genomics experiments, we next assessed the accuracy of genotyping by calculating the false discovery rate at different coverage. This approach also measures how much noise each assay captures. We defined the false discovery rate as the ratio between the information obtained from the incorrectly called variants ($h(S^F | S^G)$) and the information obtained from all the called variants ($h(S^F)$), namely

$$FDR(S^F(c)) = h(S^F(c) | S^G) / h(S^F(c)) \quad (4)$$

Figure 4a shows that the false discovery rate for Hi-C data was lower compared to WGS data at lower coverage. We attribute this finding to the deeper sequencing of the genomics regions in close spatial proximity. Hence, sampling more reads from regions at low coverage is more likely compared to uniform sampling of reads from WGS. ChIP-Seq data had a comparable false discovery rate to WGS and Hi-C data given the shallow sequencing depth. ChIP-Seq targeting CTCF had the lowest false discovery rate (Figure 4b). We further found that the polyA RNA-Seq experiment had the lowest false discovery rate compared to WES and total RNA-Seq. This could be attributed to the deeper sequencing of regions containing highly expressed genes and deeper sampling from these regions. In general, assays targeting the transcriptome such as WES and RNA-Seq produced noisier genotypes compared to WGS and Hi-C experiments; single-cell RNA-Seq was the noisiest among all the assays, as expected (Figure 4c).

2.5 Linking attack scenario

Linking attacks aim to re-identify an individual by cross-referencing datasets (Figure 5a). For example, in a hypothetical scenario an attacker aims to query an individual's HIV status from

his/her phenotype data. This phenotype data is released with the individuals genotype information with an anonymized identifier for each individual. We assume that the adversary obtains access to this dataset by either lawful or unlawful means. Now lets assume that the attacker has access to a biosample. This could be partial or complete mapped reads from functional genomics experiments or a saliva sample taken from a used glass. The idea is to genotype the biosample and find the matching genotype in the HIV status database. However, individuals share many common variants with each other. The number of shared variants between individuals is large within a racial population and even larger within a family. The question becomes how well an adversary should sequence an individuals genome to be able to perform successful linking. Specifically, the adversary is interested in investigating whether noisy and partial reads from functional genomics experiments can be used as quasi-identifiers and how accurate the genotyping needs to be in order to link individuals to databases.

For this, the attacker calls variants directly from the reads of anonymized functional genomic experiments. Then he/she compares the called noisy and incomplete genotypes to the genotype data panel and finds the entry with the highest pmi. This reveals the sensitive information for the linked individual to the attacker. We then consider a scenario in which the attacker has access to partial or increasing amount of reads to find out when the data crosses the set point and becomes private.

Based on the pmi values of each experiment at different coverage, we defined a metric for linking accuracy called gap_{query} . Let assume S_j^{DB} is the set of variants that belongs to the j^{th} individual in the genotype panel and $S_{query}^F(c)$ is the set of variants that was called from the functional genomics experiments of the query individual at c total sequencing coverage. We first calculate the pointwise mutual information between every individual in the panel and the query as $pmi(S^F(c); S_j^{DB})$. We

then ranked all the pmi values in a decreasing order such that;

$$pmi(S^F(k); S_i^{DB})^{(1)} > pmi(S^F(k); S_j^{DB})^{(2)} > \dots > pmi(S^F(k); S_m^{DB})^{(N)}$$

In our linking attack scenario, the assumption is that individual with the highest pmi with the query ($pmi(S^F(k); S_i^{DB})^{(1)}$) is the query if $gap_{query} > 1$. If $gap_{query} < 2$, then we assume that the query can be the first five ranked individuals, and can be further identified with auxiliary data such as gender or ethnicity.

As our query individual (NA12878) was in the panel, we could measure the accuracy of this prediction by further extending the definition of gap_{query} . We calculate the gap_{query} for three possibilities: (1) First ranked individual is NA12878, (2) first ranked individual is not NA12878, but NA12878 is in the first five ranked individuals, and (3) none of the top five matching individuals are NA12878. In the possibility (1), the attacker makes a correct prediction. The strength of this prediction is the gap_{query} , which is measured as the fold change difference between the pmi of best matching individual (correct prediction) and the second best matching individual. In the possibility (2), the strength of this prediction (gap_{query}) is measured as the fold change difference between the pmi of the real individual, that is ranked somewhere between 2nd to 5th and the pmi of the best matching individual, that is the misprediction. In the possibility (3), the attacker makes a false prediction that the query cannot be retrieved from the panel, there gap_{query} becomes 0. We can formulate this as;

$$\begin{aligned} gap_{query} &= \frac{pmi(S^F(k); S_i^{DB})^{(t)}}{pmi(S^F(k); S_j^{DB})^{(2)}}, \text{ if } S_i^{DB} = \text{query and } t = 1 \\ gap_{query} &= \frac{pmi(S^F(k); S_i^{DB})^{(t)}}{pmi(S^F(k); S_j^{DB})^{(1)}}, \text{ if } S_i^{DB} = \text{query and } t \in 2, 3, 4, 5 \\ gap_{query} &= 0, \text{ otherwise} \end{aligned} \quad (5)$$

We then defined that if gap_{query} is 0, then the individual cannot be identified as there are other individuals in the panel that have the matching genotypes. If $0 < gap_{query} \leq 1$, then the individual might be vulnerable with auxiliary data such as gender or ethnicity, because he/she is in the top five matching individuals. If $1 < gap_{query} \leq 2$, then the individual is vulnerable as we can identify him/her with a one- to two-fold difference between him/her and the second best match. Lastly, if $gap_{query} > 2$, then the individual is extremely vulnerable with more than a two-fold difference between him/her and the second best match. A detailed flowchart of the linking attack is shown in Figure 5a.

We found that NA12878 was extremely vulnerable even at the lowest sampled coverage for Hi-C and RNA-Seq data (Figure 5b). Interestingly, between ~ 1.1 and 10 billion base pairs, the Hi-C data exhibited higher linking accuracy than the WGS data, consistent with the previous observation of pmi shown in Figure 3a. The total coverage of ChIP-Seq data compared to Hi-C and RNA-Seq data was quite low (SI Table I). However, the linking accuracy of ChIP-Seq was as good as Hi-C and WGS (Figure 5b), showing extreme vulnerability of individuals with respect to a release of a small amount of data. More strikingly, the attacker can link NA12878 by using the reads of single-cell RNA-Seq data, which cover a small portion of the genome in a single cell (Figure 5d). We then added the variants of NA12878s parents to the 1,000 Genomes genotype panel and repeated the linking attack. We found that although NA12878 was still extremely vulnerable to re-identification with the presence of her parents in the database, the second-best matching individuals were her parents (SI Figure 2). This shows that, using the metric gap , an adversary can also identify individuals related to the target individual.

2.6 An individuals genome can be accurately approximated from publicly available data by imputation

To determine whether an attacker can correctly assemble an individuals variants by only using the reads from ChIP-Seq and RNA-Seq experiments, we imputed variants by using IMPUTE2 [31, 32, 33] and the variants called from ChIP-Seq and RNA-Seq experiments. We then collected all the called and imputed variants in a set. Although imputed variants did not contribute to the information due to high correlation with the called variants (SI Methods and SI Figure 3), the total number of captured variants increased significantly (SI Methods and SI Figure 3), total number of captured variants increases significantly (Figure 6a). By using shallow sequencing data of ChIP-Seq and RNA-Seq, we were able to call and impute almost as many variants as the gold standard.

We then tested if we could infer potentially sensitive phenotypes from these variants. Figure 6b shows a small set of example variants associated with physical traits such as eye color, hair color, or freckles. Many of these variants are in the called set of Hi-C, ChIP-Seq, and RNA-Seq data. The number of variants associated with traits further increased with imputation as expected.

2.7 Toy model for estimating the amount of leaked data without variant calling

Genotyping from DNA sequences is the process of comparing the DNA sequence of an individual to that of the reference human genome. To be able to successfully genotype, one needs substantial depth of sequencing reads for each base pair. According to the Lander-Waterman statistics for DNA sequencing, when random chunks of DNA are sequenced repeatedly, the depth per base pair follows the Poisson distribution with a mean that can be estimated from the read length, number of reads, and the length of the genome [34]. As functional genomics experiments aim to find highly expressed genes, TF binding enrichment, or 3D interactions of the genome, it is expected that the

sequencing depth per base pair does not follow Poisson statistics. Thus, genotyping using reads from functional genomics experiments is biased towards variants that are in the functional regions of the cell types/lines of interest.

To this end, we we hypothesized that genotyping from sequencing-based functional genomics data depends on the average depth per base pair (\bar{d}), and the total fraction of the genome that is represented at least by one read (i.e., the breadth, $b = \sum_{i=1}^N \delta(d_i)$, such that $\delta(d_i) = 1$ if $d_i > 0$, 0 otherwise and N is the total number of nucleotides in the genome), and a parameter β that estimates the sequencing bias (i.e., how much the distribution of depth per basepair deviates from the Poisson distribution, Fig. 6c). The bias parameter β is composed of two terms: (1) the negative bias β^- and (2) the positive bias β^+ . The negative bias estimates if there is an increase in the number of low depth basepairs relative to the mean with respect to the expected Poisson distribution; the positive bias estimates the increase in the number of high-depth basepairs (see SI for more details).

To quantify the genotyping accuracy from the functional genomics data, we used “naive” normalized pmi (npmi). This approach takes into account the information from the correctly identified genotypes ($pmi(S^F; S^G)$), the information missed that is in the gold standard ($h(S^G | S^F)$) and the information from the incorrectly identified genotypes (i.e FDR, $h(S^F | S^G)$) and normalizes it with the joint information of called variants and gold standard variants as;

$$npmi(S^F; S^G) = \frac{pmi(S^F; S^G)}{h(S^F, S^G)} = \frac{pmi(S^F; S^G)}{h(S^G | S^F) + pmi(S^F; S^G) + h(S^F | S^G)} \quad (6)$$

To be able to get a fit for the relationship of $npmi(S^F; S^G) = f(\bar{d}_F, b_F, \beta_F)$, we used Gaussian Process Regression (GPR) [?] to fit 40 training data points and achieved a root mean square error (RMSE) of 0.60 with the values ranging between [0,100] (Fig. 6d). We used five separate data points as a test set and achieved an RMSE of 0.47 was acheieved (Fig. 6d), see SI for more details). We performed regression learning is performed using a ten-fold cross-validation to protect against

overfitting. This toy model represents a conceptual theoretical framework limited to the small sample space available. It shows that the amount of leaked data from functional genomics experiments can be estimated without the need of performing time-consuming genotyping calculation.

2.8 Unique combination of common variants contribute significantly to the information leakage and linking accuracy

We next analyzed whether a linking attack can be prevented by removing rare variants from the datasets as their contribution to the information is the highest. We first speculated that the removal of the variants that are unique to NA12878 might be enough to prevent linking. A total of 11,472 variants along with their genotypes were observed only in NA12878, which we refer as “unique variants” (Fig. 6a). Note that we used the terminology variant not only for the location and minor allele of the SNV but also the genotypes. Therefore unique variants in this context are more than the de novo variants as [location, minor allele, genotype] as a vector is unique. After the removal of these unique variants from the NA12878 variant set, we calculated the $gap_{NA12878}$. Surprisingly, the linking accuracy was affected minimally compared to using the all of the NA12878 variants (Fig. 6b). We then created another set (“double variants”, Fig. 6a), that included the variants observed in NA12878s genome as well as one more individual in the 1,000 Genomes genotype panel (total of 16,305 genotypes). We again found that the individual was extremely vulnerable to linking attacks ($gap_{NA12878} > 2$, Fig. 6b). We then relaxed our cut-off further to remove the variants that are observed in NA12878’s genome as well as at most 1.5% of the population (“rare variants”, total of 124,093 genotypes, Fig. 6a). This also did not affect the overall linking ($gap_{NA12878} > 2$, Fig. 6b).

These rare genotypes were observed in 64 or less individuals including NA12878. A practical solution to the re-identification problem using functional genomics data would be masking or removing such rare genotypes from the reads. However, as iteratively shown here, although rare

variants are extremely informative and sufficient enough to achieve re-identification through linking attacks, their removal is not sufficient to prevent re-identification. That is, not only the rare genotypes but also the unique combination of common genotypes are identifiers of the genetic make-up of an individual. To further support this calculation, we added the genotypes of the parents of NA12878 to the panel and found that we could still link NA12878 to the correct genotypes successfully with an extreme vulnerability ($gap_{NA12878} > 2$, SI Fig. 2).

We then analyzed the contribution of small indels to the naive information and whether accurate linking was possible when we removed all the single nucleotide mutations from the data and kept the indels. Fig. ??c shows the information contribution of the indels. Although naive pmi from indels were much smaller compared to single nucleotide mutations, a high linking accuracy could be achieved by using only indels even at small coverage (Fig. 6d). This linking attack is done using one of the noisy data set we have (total RNA-Seq) to make linking more difficult.

2.9 Privacy-preserving file formats for alignments from functional genomics experiments and relation to k -anonymity

Sharing raw alignments from functional genomics experiments is extremely important in developing analysis methods and discovering novel mechanisms about the human genome. Ideally, one would share the maximal amount of information with minimal utility loss while largely maintaining an individuals privacy. As a privacy metric, we aimed to prevent leakage of any variant as well as any quasi-identifier that can lead to identification of the position of variants in the genome. We introduced a user-identified privacy-utility balance that can be adjusted according to the patients consents and institutions policies. By using the concept of k -anonymity [?], we applied a privacy-preserving transformation to the alignment files such that calling variants from transformed files is largely prevented while quantifications related to the functional genome is possible with minimal error (Fig. 8a).

A release of data possesses the k -anonymity property if the information for each person contained in the release cannot be distinguished from at least $k - 1$ individuals whose information also appear in the release. Although this concept was developed for the release of datasets with individuals, we can think of a raw alignment file (BAM) as a dataset, where information for each read is contained. Let's assume a BAM file is a dataset D , where each entry is a read. The desire is to release dataset D in a form (say D^*) such that it does not leak variants from the reads, but in the mean time any calculation f based on D and D^* retrieves almost the same result. There are two general methods to achieve k -anonymity for some value of k : suppression and generalization. If every column in D is an attribute (such as read length, cigar, sequence, or quality value), then replacing an attribute with an asterisk(*) is suppression and changing an attribute with a more general value is generalization. For example, in our file format transformation, we replaced sequence and sequence quality attributes with asterisk (suppression), and transformed the cigar of the read from partially mapped to fully mapped (generalization) to achive 3-anonymity with respect to attributes sequence, sequence quality and cigar (see SI Methods for details). Now let's say the privacy-preserving transformation is done through a function $P_{Q,r}$ such that $P_{Q,r}(D) = D^*$. Q is the operation such as "removal of small indels", "removal of mismatches", "removal of large indels" or "removal of all variants". r is the amount of reads to be manipulated given the operation Q . A calculation f can be signal depth profile calculation, TF binding peak detection or gene expression quantification (Fig. 8a). Then, we can reconstruct the eq. 7 for each unit i as

$$\frac{f(D)}{f(D^*)} = e^{\varepsilon_i}, \quad (7)$$

where a unit i can be a single basepair, an exon or a gene depending on the function f . In turn, ε_i can be calculated as the log fold change between the results derived from two datasets. This is also a quantity commonly used to compute differential gene expression [37] or ChIP-Seq binding enrichment over controls [?], and can be used as analogous in this context, where log-fold change

is the differential signal depth or expression when the manipulated data is used as an input.

Note that $|\varepsilon_i|$ is a measure of error of the new dataset D^* . We then calculated the distribution of $|\varepsilon_i|$ values over every unit and found the mean $|\varepsilon|$ per unit as the overall error. The level of privacy is controlled by the function $P_{Q,r}$, where Q determines the type of entries and r determines the number of entries of the given operation Q that are manipulated. For any particular operation, the obvious threshold could be the size of the indels, Minor Allele Frequency (MAF), or the depth of a particular unit. These thresholds can be converted into fraction of the reads affected. For example, if Q is the removal of indels and r is the reads that contain indels with $MAF < 0.5$, then only reads that have indels with $MAF < 0.5$ will be manipulated in the transformed D^* .

We constructed the privatized file format pBAM from data D^* as follows. The reads from the BAM files were categorized as perfectly mapped reads and reads with mismatches, insertions, deletions, soft- and hard-clipping. $P_{Q,r}$ replaces the sequence of all of the reads with asterisk and manipulates the cigars, alignment scores and the MD tags of the reads that are defined in Q and r . The details of how the new file format deals with reads are reported in SI Methods with a figure (SI Fig. 4). pBAM files can also be created from BAM files that are obtained by mapping sequences to the transcriptome coordinates, which is essential for gene quantification. Our transformation function $P_{Q,r}$ is general and can be applied to any alignment file types such as SAM, CRAM and MRF to create a privatized new file format. These files will be concordant to use with tools such as samtools, cramtools and mrftools.

We calculated the signal depths of each basepairs in the genome using an NA12878 RNA-Seq BAM file using STAR [?]. We then converted the BAM file into pBAMs with different qs and calculated the signal depth of each basepair. Fig. 8b shows the number of basepairs with $\varepsilon_i > 0$ with respect to the number of base pairs with no change between BAM and pBAM. We did the same calculation by averaging signal over exons as well (Fig. 8b). Furthermore, we created pBAM

files for the BAM files that are mapped to the reference transcriptome and compared the gene quantification with the gene expression levels calculated from original BAM files. We used RSEM for gene quantification and STAR for transcriptome alignment [?, ?]. We found no difference between the gene expression levels calculated using original BAM files and pBAM files (see Fig 8b and SI Methods for how we treated transcriptome alignments). Overall, when we removed all the variant leak from the BAM files, we found 0.18% difference at the basepair resolution, 0.27% difference at the exon resolution, and 0% difference at the gene level. When we removed leak associated with the mismatches, we did not see any difference as when the cigars with mismatches are manipulated, the correct mapping locations can be recovered without leakage (see SI Methods). When we removed leak associated with indels, we found 0.0016% difference at the base pair resolution, 0.0011% at the exon resolution, and 0% difference at the gene level. When we removed leak associated with split reads, we found 0.17% difference at the basepair resolution, 0.26% at the exon resolution, 0% difference at the gene level. [GG2MG: should I move these sentences to discussion?] Figure 8c shows the change in ϵ with respect to increasing r for different operations Q. When the mismatches are manipulated, the resulting signal profiles are not affected. Hence, the manipulated dataset will retrieve the same results as the original dataset regardless of the number of reads (r). However, manipulating indels and split reads will result in changes in the utility of the new file formats. This is particularly useful as for example the ENCODE consortium adopts processing pipelines, in which only highly mapped reads are used and split reads are discarded.

The pBAM file format contains necessary information to be used in functional genomics pipelines such as gene expression quantification and TF binding peak calling. The difference between the results of the ENCODE Chip-Seq TF binding peak calling pipeline (MACS2 [?]) is even more negligible when BAM and pBAM were used as input (SI Fig. 4). We then created a .diff file format that contains the original information that was manipulated in the pBAM file. With the motivation of keeping the size of private file formats relatively small, we report only differences

between BAM and pBAM in the .diff file by avoiding printing any sequence information of the reads that can be found in the reference human genome (see SI Methods). The .diff files are private files that require special permission for access. A user is able to retrieve the original BAM file when they have access to the .diff file by using our collection of scripts called ptools that can convert pBAM + .diff + reference genome into the original BAM file (Fig 8c).

2.9.1 Implementation

Conversion of BAM files to pBAM and pBAM+diff files back to BAM files are implemented as a series of scripts in bash scripting language and Python. Diff files are encoded in a compressed format to save space. For convenience, pBAM files are saved as BAM files with manipulated content and saved with a p.bam extension. That is, any pipeline that uses BAM as an input can take p.bam as an input as well. Running times and associated file sizes for alignments from RNA-Seq experiments and ChIP-Seq experiments are documented in Table 1. Our file format manipulation has been adopted by the ENCODE Consortium Data Coordination Center. Codes for the calculation of information leakage, scripts for file manipulation as well as examples of BAM, pBAM, and diff files can be found at privaseq3.gersteinlab.org.

Table 1: pTools performance. Note that when RNA-Seq reads mapped to the transcriptome, STAR and RSEM get rid of the split reads and reads with indels, only reads with mismatches are reported.

Experiment	BAM size	Q	ϵ	pBAM size	.diff size	BAM to pBAM runtime	pBAM+diff+hg to BAM runtime
RNA-Seq genome		all reads					
RNA-Seq genome		reads with mismatches					
RNA-Seq genome		reads with indels					
RNA-Seq genome		split reads					
RNA-Seq transcriptome		all reads					
RNA-Seq transcriptome		reads with mismatches					

3 Discussion

Functional genomics experiments using large-scale, high-throughput, sequencing-based assays provide a large amount of biological data. Although these experiments aim to answer questions related to genomic activities such as gene expression, TF binding, or the 3D organization of the genome, public sharing of sequencing data from these experiments can lead to recovery of genotype information and, in turn, raise privacy concerns. The systematic quantification of private information content of the functional genomics BAM files and open access to such data without compromising individuals identity have not been well studied. Current policies regarding public sharing of functional genomics BAM files are ad-hoc. The experiments that require a high depth of sequencing such as Hi-C and sometimes RNA-Seq are considered to be private, whereas relatively low-depth BAM files such as those from CHIP-Seq are often shared publicly. In this study, we derived information theory-based measures to systematically quantify the sensitive information leakage in the BAM files of functional genomics experiments in low- and high-depth experiments.

Instantiation of linking attacks by genotyping of partial or complete functional genomics data showed that even at low coverage of low-depth experiments such as CHIP-Seq, linking individuals to the databases can be done without error. When we compared the linking accuracy to the false discovery rate, we found that it is easier to link individuals to the databases than genotyping them accurately using functional genomics experiments. The implication is that noisy quasi-identifiers (i.e., low-quality SNP calling) can be used to link the data to the high-quality genotypes. For example, according to our calculations, reads from single-cell RNA-Seq data carry the largest amount of noise. This is likely due to the bias towards expressed genes in such small amounts of cells, mapping issues of splice sites, false positives from RNA editing sites, and amplification bias. However, the noisy genotypes called from a small amount of cells, even when the number of reads is only a million, are quasi-identifiers that result in very high linking accuracy. This is worrisome in terms of biomedical data sharing as the number of individuals in genotype databases

is increasing exponentially with the decreasing cost of sequencing. Furthermore, rich information about an individual's identity and his/her sensitive phenotypes can also be inferred by combining the reads from low-depth functional genomics experiments and through genotype imputation.

Another implication of the false discovery rate of genotyping in privacy is the relationship between the accuracy of the genotypes and the amount of information gained from the genotypes. For example, if the query individual is not in the genotype panel, any genotypes of the query that are not in the panel will be de novo variants and will greatly contribute to the information gained. However, these de novo variants can be rich in artifacts and sequencing errors. Conversely, any common genotype of the query will be highly accurate while poor in information. Consequently, from an adversary's perspective, the most valuable genotypes will be the rare genotypes in the panel to make accurate inferences about the query's identity and sensitive phenotypes, despite the fact that most information is gained from the de novo variants by definition. One way to correct for this is to count any de novo variant genotyped as a false discovery, which changes the false discovery rate values in Fig. 5 greatly for different functional genomics experiments and is presented in SI Fig. x.

In this manuscript, we also discuss the concept of a set point in determining the data production steps, where sensitive information leakage and utility of the data are balanced (Fig. 1). Setting a set point is possible by systematic genotyping and quantification of information. Although it is obvious that any DNA read contains variants, it is not trivial to understand the amount and the quality of sequencing to perform accurate genotyping. Moreover, we showed that genotyping accuracy of a functional genomics sample and the ability to link individuals to the databases using the same sample are not necessarily correlated. It is easier to link individuals to the databases and infer their complete variant sets than genotyping a sample with accuracy and minimal false discovery. For example, a complete set of variants from HeLa's genome may not be obtained by genotyping HeLa BAM files from functional genomic experiments. However, using only a small

number of reads from the same BAM files, accurate linking attacks are plausible. That is, noisy and incomplete genotyping from partial sequencing experiments can serve as strong quasi-identifiers, which is not straightforward to predict at first. Nevertheless, policies governing the public sharing of HeLa genome versus HeLa functional genomics reads is ad-hoc and contradictory. Therefore, it is essential to quantify the information in samples and determine the set point accurately. Importantly, functional genomics experiments advance our understanding of health and disease by revealing functions of the genome under different conditions. The quantification, analysis, and interpretation of functional genomics data is still an evolving field; hence, extensive public sharing of functional genomics data will accelerate collaborative research and reproducibility by removing the complexities associated with data accession procedures.

The increasing incentive to share data for the advancement of biomedical research and the corresponding increasing privacy concerns have led researchers to look for more complex solutions to overcome the bottleneck between data-sharing and privacy preserving means. Solutions such as differential privacy have been proposed [?, ?, ?]. Studies have shown that retrieving summary information from private statistical databases without revealing some amount of an individuals information is impossible [?]. We further studied if the concept of differential privacy can be utilized to create leakage-free raw functional genomics data (see SI Methods). Furthermore, an entire database can be inferred by using a small number of queries. Differential privacy ensures a high level of privacy such that an adversary retrieves a similar result with or without the addition of the individuals data to the database by adding perturbations or noise to the queries [?]. We further studied if the concept of differential privacy can be utilized to create leakage-free raw functional genomics data (see SI Methods). Although such a concept is useful for sharing summary statistics of functional genomics data from multiple individuals, it is conceptually hard to apply to the raw mapped read sharing from functional genomics experiments taken from a single individual. Although further research will be fruitful on how to extract useful information from genomics

data that are noisy and perturbed, we envision there will be more applications of privacy concepts like differential privacy in genomics data sharing such as releasing population-based genotype-phenotype data [36, ?].

Furthermore, an entire database can be inferred by using a small number of queries. Differential privacy ensures a high level of privacy such that an adversary retrieves a similar result with or without the addition of the individuals data to the database by adding perturbations or noise to the queries [?]. We further studied if the concept of differential privacy can be utilized to create leakage-free raw functional genomics data (see SI Methods). Although such a concept is useful for sharing summary statistics of functional genomics data from multiple individuals, it is conceptually hard to apply to the raw mapped read sharing from functional genomics experiments taken from a single individual. Although further research will be fruitful on how to extract useful information from genomics data that are noisy and perturbed, we envision there will be more applications of privacy concepts like differential privacy in genomics data sharing such as releasing population-based genotype-phenotype data [20]. There is also leakage from gene expression quantifications, which was shown to be connected with variants through the eQTLs [19]. Quantification of the leakage in all levels of data-processing steps of an RNA-Seq experiment is tabulated in Table 2 and in SI Fig. x. We also anticipate more leakages to be discovered as new functional genomics experiments are developed. Combined with the increasing attention to genomic privacy, we expect future studies will lead to novel privacy-preserving solutions in an open data-sharing mode.

Table 2: Quantification of leakage in different sources. We quantified the number of variants that are leaked from each source. We then calculated the number of accessible variants to an average RNA-Seq experiment.

Leakage Source	Leaking Variants	Number of Variants	Average Leakage per Variant (bits)	Maximum Leakage per Variant (bits)	Number of Accessible Variants Per Individual	Total Leakage (bits)
Raw Reads	Exonic Variants	2,772,064	0.10±0.28	9.88±2.12	221,293	22,129.30
RNA-Seq Signal Profiles	Exonic Deletions	51,408	0.28±0.45	7.97±2.42	1,067	298.76
Gene Expression	eQTLs	3,175	1.19±0.36	4.00±1.92	158	188.02

4 Figure Legends

[GG2MG: We haven't gone through these yet]

Figure 1: **Schematic of data types from functional genomics experiments.** (a) The flow for RNA-Seq data processing from mapped reads to the gene quantifications. (b) Different layers of produced data from functional genomics experiments. Red line denotes the set point, where privacy and utility trade-off balanced.

Figure 2: **Comparison of naive information measure with information with LD consideration and sample size correction.** (a) The process of sampling reads from functional genomics experiments for the calculation of pointwise mutual information between 1000 genomes gold standard variants for NA12878 in different coverage. (b) The maximum LD score for each variant are plotted against information per variant. Highly informative variants do not exhibit difference when information is calculated using naive approach vs. with LD consideration. (c) Naive information vs. inverse fraction of the data sampled from the 1000 genomes population. y-intercept is extrapolated from the fitted curve and denotes the information when the population size is infinite. Error bars are calculated using 100× bootstrapping. (d) Difference between the naive information, information with LD consideration and extrapolated information when population size is infinite.

Figure 3: **The pointwise mutual information calculated for 24 different functional genomics assays and WGS, WES and SNP ChIP data using NA12878 1000 genomes variants as gold standard.** (a) The pmi values for WGS and three different primary Hi-C experiments plotted at different coverages. The information contents of the gold standard and SNP ChIP are added for comparison. (b) The pmi values for 20 different ChIP-Seq experiments targeting histone modifications and transcription factor binding plotted at different coverage. (c) The pmi values for WES, total RNA-Seq, polyA RNA-Seq and single-cell RNA-Seq from two different cells plotted at different coverage. (d) The pmi values per basepair plotted using the mean of all the ratios between the pmi and the corresponding coverage.

Figure 4: **False discovery rate of functional genomics experiments at different coverages** (a) FDR comparison for Hi-C and WGS data at different sampled coverage. (b) FDR comparison for different ChIP-Seq experiments at different coverage. (c) FDR comparison for WES and different RNA-Seq experiment at different coverage.

Figure 5: **Illustration of a linking attack and the accuracy of linking.** (a) The publicly available anonymized reads from functional genomics experiments contains a set of variants and HIV status for the sample that the functional genomics experiment was performed at increasing coverage. The panel of genotypes contains the variants and associated genotypes for m individuals. The attacker links the inferred variants and genotypes to the panel of genotypes by using the best matched point-wise mutual information. The linking potentially reveals the HIV status for the linked individual. (b) Comparison of *gap* for NA12878 at different coverage for Hi-C and Total/PolyA RNA-Seq reads. WGS and SNP-ChIP are also added for comparison. (c) Comparison of *gap* for NA12878 at different coverage for 20 different ChIP-Seq experiments. (d) Comparison of *gap* for NA12878 at different coverage for single-cell RNA-Seq experiments.

Figure 6: **Individual's genome can be approximated and sensitive phenotypes can be inferred from publicly available data by imputation and a theoretical framework for prediction of amount of leaked data** (a) Number SNVs called from WGS data and all of the ChIP-Seq and RNA-Seq data together with and without imputation. (b) Variants associated with physical traits and if they present in the called variants from different functional genomics experiments before and after imputation. (c) Features of the theoretical framework - write more. (d) Accuracy of fitted model on training set- write more (e) Accuracy of fitted model on test set - write more

Figure 7: **Removal of rare variants and linking** (a) Information of the variant before and after addition of NA12878 to the population. We iteratively removed variants from the set as (I) only the variants that is only NA12878 specific, (II) the variants that have an information of 11 or higher bits after removal of NA12878 from the population, (III) the variants that have an information of 6 or higher bits after removal of NA12878 (b) Linking accuracy for every iteration of removal of NA12878 variants from the gold standard set. (c) Information of all the variants that are called from total RNA-Seq reads vs. the information of the indels that are called from total RNA-Seq reads. (d) Linking accuracy when we consider all the variants that are called from total RNA-Seq reads vs. the linking accuracy when we consider only indels called from Total RNA-Seq reads.

Figure 8: **Privacy-preserving file formats for mapped reads - ready after finalizing the figure** (a) The schematic of the difference between the signal calculated from original BAM and transformed pBAM files and the concept of ϵ for the error. (b) (c) (d)

References

- [1] Sboner A, Mu X, Greenbaum D, Auerbach RK, Gerstein MB. The real cost of sequencing: higher than you think! *Genome Biology*, 2011;12(8):125.
- [2] Joly Y, Dyke SOM, Knoppers BM, Pastinen T. Are Data Sharing and Privacy Protection Mutually Exclusive? *Cell*, 2016;167(5):1150-1154.
- [3] Erlich Y, Narayanan A. Routes for breaching and protecting genetic privacy. *Nat. Rev. Genet.*, 2014;15(6):409-421.
- [4] Joly Y, Feze IN, Song L, Knoppers BM. Comparative Approaches to Genetic Discrimination: Chasing Shadows? *Trends Genet*, 2017;33(5):299-302.
- [5] Homer N, Szelling S, Redman M, Duggan D, Tembe W, Muehling J, Pearson JV, Stephan DA, Nelson SF, Craig DW. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. *PLoS Genet.*, 2008;4(8):e1000167.
- [6] Im HK, Gamazon ER, Nicolae DL, Cox NJ. On sharing quantitative trait GWAS results in an era of multiple-omics data and the limits of genomic privacy. *Am. J. Hum. Genet.*, 2012;90(4):591-598.
- [7] Church GM. "The Personal Genome Project". *Molecular Systems Biology*, 2005;1(1):E1E3.
- [8] Gymrek M, McGuire AL, Golan D, Halperin E, Erlich Y. Identifying personal genomes by surname inference. *Science*, 2013;339(6117):321-324.
- [9] Sweeney L. k-anonymity: a model for protecting privacy. *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 2002;10(5):557-570.

- [10] Sweeney L. Simple demographics often identify people uniquely. *Carnegie Mellon University, unpublished*, 2000.
- [11] Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat. Rev. Genet.*, 2009;10(1):57-63.
- [12] Pepke S, Wold B, Mortazavi A. Computation for ChIP-seq and RNA-seq studies. *Nat. Rev. Genet.*, 2009;6:S22S32.
- [13] Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragozy T, Telling A, Amit I, Lajoie BR, Sabo PJ, Dorschner MO, Sandstrom R, Bernstein B, Bender MA, Groudine M, Gnirke A, Stamatoyannopoulos J, Mirny LA, Lander ES, Dekker J. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science*, 2009;326(5950):289-293.
- [14] Cock PJA, Fields CJ, Goto N, Heuer ML, Rice PM. The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Research*, 2009;38(6):1767-1771.
- [15] Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 2009;25(16):2078-2079.
- [16] Hsi-Yang Fritz M, Leinonen R, Cochrane G, Birney E. Efficient storage of high throughput DNA sequencing data using reference-based compression. *Genome Res.*, 2011;21(5):734-740.
- [17] Beskow LM. Lessons from HeLa Cells: The Ethics and Policy of Biospecimens. *Annu Rev Genomics Hum Genet.*, 2016;17:395-417

- [18] Schadt EE, Woo S, Hao K. Bayesian method to predict individual SNP genotypes from gene expression data. *Science*, 2012;44(5):603-608.
- [19] Harmanci A, Gerstein M. Quantification of private information leakage from phenotype-genotype data: linking attacks. *Nature Methods*, 2016;13(3):251-256.
- [20] Harmanci A, Gerstein M. Analysis of Sensitive Information Leakage in Functional Genomics Signal Profiles through Genomic Deletions. *Nature Communications*, 2017
- [21] Habegger L, Sboner A, Gianoulis TA, Rozowsky J, Agarwal A, Snyder M, Gerstein M. RSEQtools: a modular framework to analyze RNA-Seq data using compact, anonymized data summaries. *Bioinformatics*, 2011;27(2):281-283.
- [22] ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 2012;489(7414):57-74.
- [23] Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB, Shaw KR, Ozenberger BA, Ellrott K, Shmulevich I, Sander C, Stuart JM. The Cancer Genome Atlas Pan-Cancer analysis project. *Nature Genetics*, 2013;45(10):1113-1120.
- [24] The GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nature Genetics*, 2013;45(6):580-585.
- [25] National Institute of Health data sharing policy. <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-17-110.html>
- [26] The 1000 Genomes Project Consortium. A map of human genome variation from population-scale sequencing. *Nature*, 2010;467(7319):1061-1073.
- [27] DePristo M, Banks E, Poplin R, Garimella K, Maguire J, Hartl C, Philippakis A, del Angel G, Rivas MA, Hanna M, McKenna A, Fennell T, Kernysky A, Sivachenko A, Cibulskis K,

- Gabriel S, Altshuler D, Daly M. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics*, 2011;43(5):491-498.
- [28] Van der Auwera GA, Carneiro M, Hartl C, Poplin R, del Angel G, Levy-Moonshine A, Jordan T, Shakir K, Roazen D, Thibault J, Banks E, Garimella K, Altshuler D, Gabriel S, DePristo M. From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best Practices Pipeline. *Current Protocols in Bioinformatics*, 2013;43:11.10.1-33.
- [29] International HapMap Consortium. The International HapMap Project. *Nature*, 2003;426(6968):789-796.
- [30] Strong SP, Koberle R, de Ruyter van Steveninck RR, Bialek W. Entropy and Information in Neural Spike Trains. *Phys. Rev. Lett.*, 1998;80:197.
- [31] Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genetics*, 2009;80:5(6):e1000529.
- [32] Howie BN, Marchini J, Stephens M. Genotype imputation with thousands of genomes. *G3: Genes, Genomics, Genetics*, 2011;1(6):457-470.
- [33] Howie BN, Fuchsberger C, Stephens M, Marchini J, Abecasis GR. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nature Genetics*, 2012;44(8):955-959.
- [34] Lander ES, Waterman MS. Genomic mapping by fingerprinting random clones: a mathematical analysis. *Genomics*, 1988;2(3):231-239.
- [35] Rasmussen CE, Williams CKI. Gaussian Processes for Machine Learning. *MIT Press*, 2006;ISBN 0-262-18253-X.

- [36] Dwork C. Differential Privacy: A Survey of Results. *Springer Berlin Heidelberg*, 2008;Theory and Applications of Models of Computation. Lecture Notes in Computer Science. pp. 1-19
- [37] Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, Szczesniak M, Gaffney D, Elo LL, Zhang X, Mortazavi A. A survey of best practices for RNA-seq data analysis. *Genome Biol.*, 2016;17:13-14.
- S. Fienberg, A. Slavkovic, and C. Uhler. Privacy preserving GWAS data sharing. In ICDM, pages 628635, 2011.
- A. Johnson and V. Shmatikov. Privacy-preserving data exploration in genome-wide association studies. In KDD, pages 10791087, 2013.
- F. Yu, S. E. Fienberg, A. B. Slavkovi, and C. Uhler. Scalable privacy-preserving data sharing methodology for genome-wide association studies. *Journal of Biomedical Informatics*, 2014.
- Irit Dinur and Kobbi Nissim. 2003. Revealing information while preserving privacy. In Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems (PODS '03). ACM, New York, NY, USA, 202-210.