

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

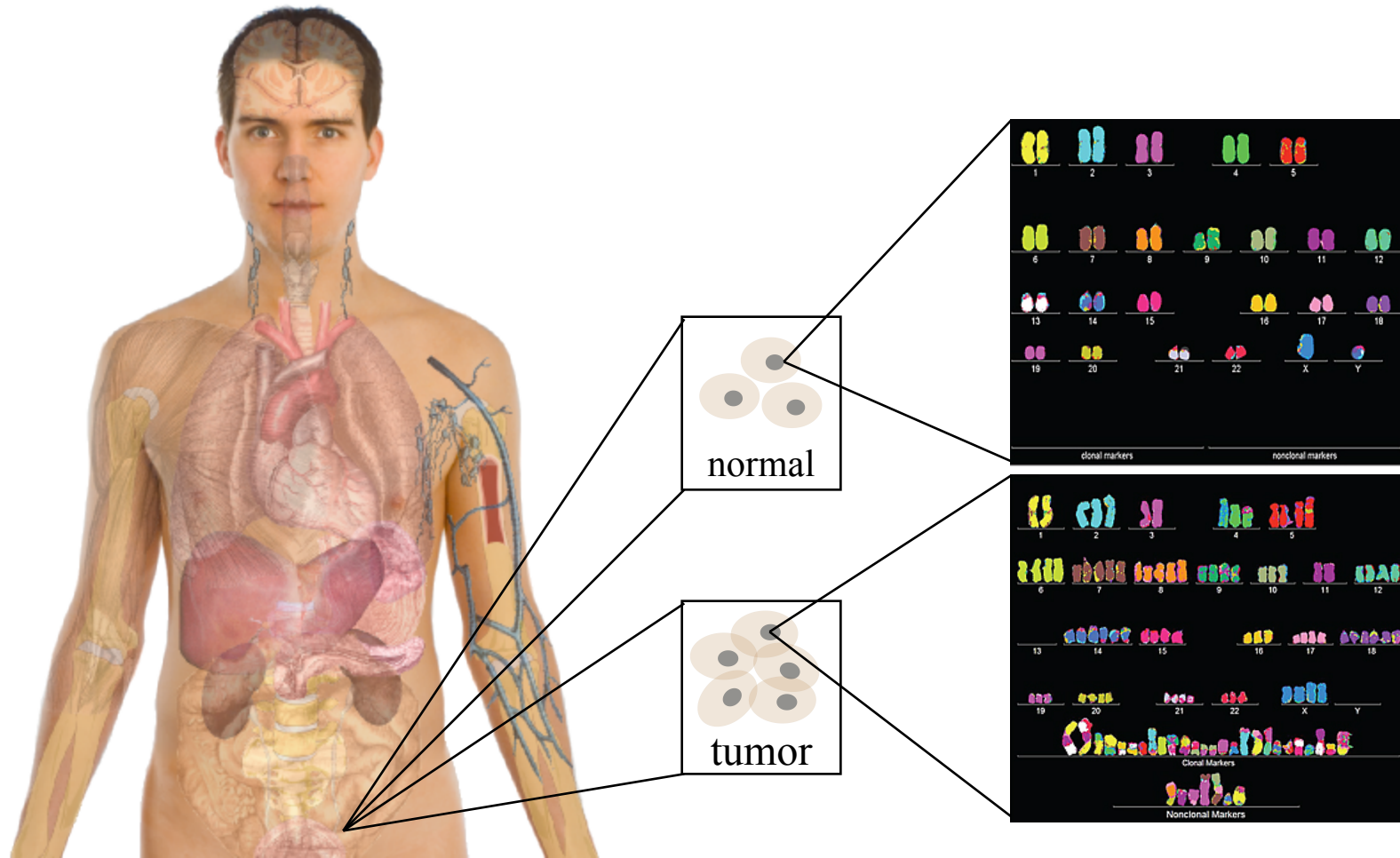
Slides freely downloadable from
Lectures.GersteinLab.org &
“tweetable” (via [@markgerstein](https://twitter.com/markgerstein)).
See last slide for more info.



Personal Genomics & Transcriptomics as a Gateway into Biology

Personal genomes (& Transcriptomes) soon will become a commonplace part of medical research & eventually treatment (esp. for cancer).

They will provide a primary connection for biological science to the general public.



Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- Details on Relevant Hacks: Genomic, Computer Security, & Netflix

• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack using extreme expression levels
- Quantifying accuracy of prediction, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

The Conundrum of Genomic Privacy: Is it a Problem?

Yes

Genetic Exceptionalism :

The Genome is very fundamental data, potentially very revealing about one's identity & characteristics

Identification Risk: Find that someone participated in a study [eg Craig, Erlich]

Characterization Risk: Finding that you have a particular trait from studying your identified genome [eg Watson ApoE status]

No

Shifting societal foci

No one really cares about your genes

You might not care



[Klitzman & Sweeney ('11), J Genet Couns 20:98]; Greenbaum & Gerstein ('09), New Sci. (Sep 23)]

Genomics has similar "Big Data" Dilemma in the Rest of Society

- Sharing & "peer-production" is central to success of many new ventures, with the same risks as in genomics
 - **EG web search**: Large-scale mining essential



- We confront privacy risks every day we access the internet
- (...or is the genome more exceptional & fundamental?)

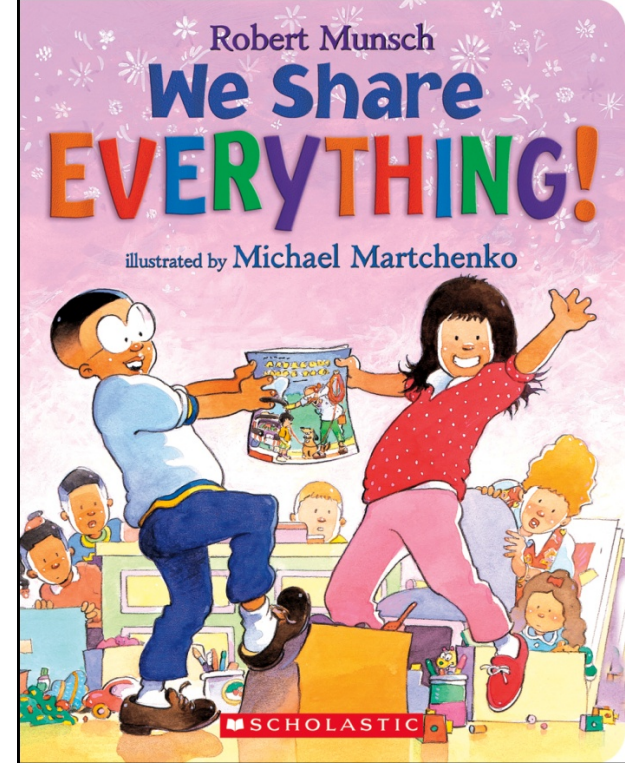
Tricky Privacy Considerations in Personal Genomics

- **Personal Genomic info. essentially meaningless currently but will it be in 20 yrs? 50 yrs?**
 - Genomic sequence very revealing about one's children. Is true consent possible?
 - Once put on the web it can't be taken back
- **Ethically challenged** history of genetics
 - Ownership of the data & what consent means (Hela)
 - Could your genetic data give rise to a product line?
- **Culture Clash:** Genomics historically has been a proponent of “open data” but not clear personal genomics fits this.
 - Clinical Medline has a very different culture.



The Other Side of the Coin: Why we should share

- Sharing helps **speed research**
 - Large-scale mining of this information is important for medical research
 - Privacy is cumbersome, particularly for big data
- Sharing is important for **reproducible research**
- Sharing is useful for **education**
 - More fun to study a known person's genome
 - Eg Zimmer's Game of Genomes in STAT



[Yale Law Roundtable ('10). *Comp. in Sci. & Eng.* 12:8; D Greenbaum & M Gerstein ('09). *Am. J. Bioethics*; D Greenbaum & M Gerstein ('10). *SF Chronicle*, May 2, Page E-4; Greenbaum et al. *PLOS CB* ('11)]

CARL ZIMMER'S
GAME OF GENOMES
SEASON 1





The Dilemma

[Economist, 15 Aug '15]

- The individual (harmed?) v the collective (benefits)
 - But do sick patients care about their privacy?
- How to balance risks v rewards - Quantification
 - What is acceptable risk? What is acceptable data leakage?
Can we quantify leakage?
 - Ex: photos of eye color
 - Cost Benefit Analysis: how helpful is identifiable data in genomic research v. potential harm from a breach?

Current Social & Technical Solutions

• **Closed Data** Approach

- Consents
- “Protected” distribution via dbGAP
- Local computes on secure computer

• Issues with Closed Data

- Non-uniformity of consents & paperwork
 - Different international norms, leading to confusion
- Encryption & computer security creates burdensome requirements on data sharing & large scale analysis
- Many schemes get “hacked”

• **Open Data**

- Genomic “test pilots” (ala PGP)?
 - Sports stars & celebrities?
- Some public data & data donation is helpful but is this a realistic solution for an unbiased sample of ~1M

Strawman Hybrid **Social** & **Tech** Proposed Solution?

- Fundamentally, researchers have to keep genetic secrets.
 - **Need for an (international) legal framework**
 - Genetic Licensure & training for individuals (similar to medical license, drivers license)
- Technology to make things easier
 - Cloud computing & enclaves (eg solution of Genomics England)
- Technological barriers shouldn't create a social incentive for “hacking”
- **Quantifying Leakage & allowing a small amounts of it**
- **Careful separation & coupling of private & public data**
 - **Lightweight, freely accessible secondary datasets coupled to underlying variants**
 - Selection of stub & "test pilot" datasets for benchmarking
 - Develop programs on public stubs on your laptop, then move the program to the cloud for private production run

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

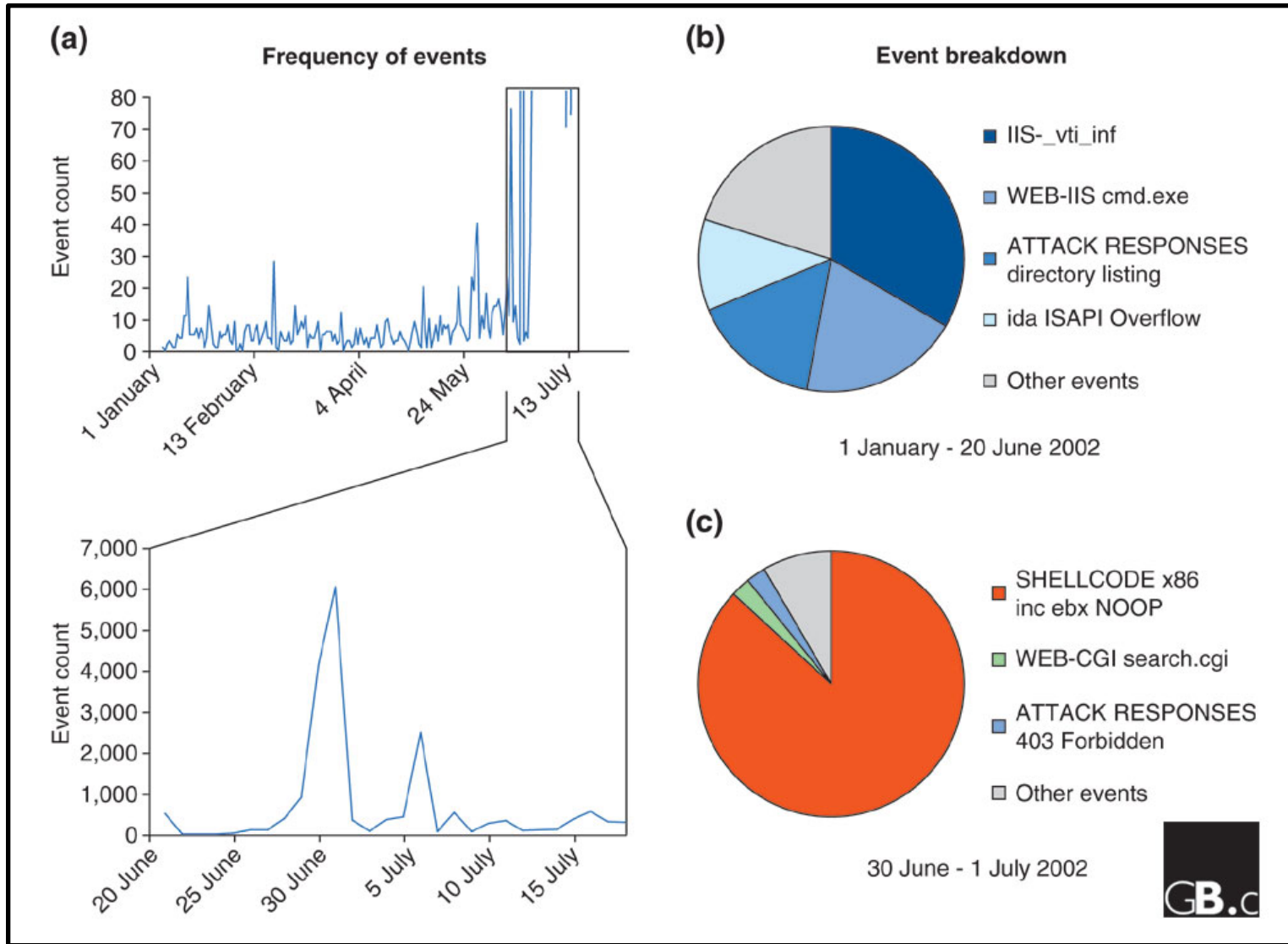
• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

Difficulty in Securing Computers & Data



[Smith et al ('05), Genome Bio]

Genomic Privacy Hacks, Mostly Focusing on Identification

- Early genomic studies were based on small cohorts
 - Individuals give consent to participate but request anonymity
 - HAPMAP, PGP, 1000 Genomes...
 - Focus on hiding the participation of individuals
 - Attacks aimed at detecting whether an individual with known genotypes participated a study
 - “Detection of genomes in a mixture” (Homer et al 2008, Im et al 2012)
- As more people are genotyped, more individuals are in large private genomic databases
 - Detection of an individual is irrelevant, as their participation is already known
 - Current EX: “An individual’s genomic/phenotypic data is most certainly stored in their hospital”
 - Future: >1M people’s health information is part of a NIH/PMI or NHS databases
- Identification attacks now focus on pinpointing individuals by cross-referencing large seemingly independent datasets
 - Illustrates that a leaked/hacker/stolen dataset, even when anonymized, can leak information
 - Sweeney et al 2013, Gymrek et al 2013

Gymrek et al, “Identifying Personal Genomes by Surname Inference” (2013)

Homer et al, “Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays.” (2008)

Im et al, “On Sharing Quantitative Trait GWAS Results in an Era of Multiple-omics Data and the Limits of Genomic Privacy” (2012)

Sweeney et al, “Identifying Participants in the Personal Genome Project by Name” (2013)

What is a linking attack? Case of Netflix Prize

Robust De-anonymization of Large Datasets

(How to Break Anonymity of the Netflix Prize Dataset)

Arvind Narayanan and Vitaly Shmatikov

The University of Texas at Austin

The Netflix logo, consisting of the word "NETFLIX" in white, bold, sans-serif capital letters with a black drop shadow, set against a red rectangular background.The IMDb logo, consisting of the letters "IMDb" in a bold, black, sans-serif font, set against a yellow rounded rectangular background.

1. Very large datasets
2. A lot of users have a Netflix and an IMDb account
3. A user rates similar scores to a movie in Netflix and IMDb
4. A user rates a particular movie around the same date in Netflix and IMDb

What is a linking attack? Case of Netflix Prize



Movie ratings database



Anonymized Netflix Prize Training Dataset
made available to contestants

100 million ratings
500,000 users
200 movie ratings/user
5,000 users/movie rating

User (ID)	Movie (ID)	Date of Rating	Rating [1,2,3,4,5]
NTFLX-0	NTFLX-19	10/12/2008	1
NTFLX-1	NTFLX-116	4/23/2009	3
NTFLX-2	NTFLX-92	5/27/2010	2
NTFLX-1	NTFLX-666	6/6/2016	5
...
...

Linking Attacks: Case of Netflix Prize



Names available for many users!

User (ID)	Movie (ID)	Date of Grade	Grade [1,2,3,4,5]
NTFLX-0	NTFLX-19	10/12/2008	1
NTFLX-1	NTFLX-116	4/23/2009	3
NTFLX-2	NTFLX-92	5/27/2010	2
NTFLX-1	NTFLX-666	6/6/2016	5
...
...

User (ID)	Movie (ID)	Date of Grade	Grade [0-10]
IMDB-0	IMDB-173	4/20/2009	5
IMDB-1	IMDB-18	10/18/2008	0
IMDB-2	IMDB-341	5/27/2010	-
...
...
...

- Many users are shared
- The grades of same users are correlated
- A user grades one movie around the same date in two databases

Linking Attacks: Case of Netflix Prize



User (ID)	Movie (ID)	Date of Grade	Grade [1,2,3,4,5]
NTFLX-0	NTFLX-19	10/12/2008	1
NTFLX-1	NTFLX-116	4/23/2009	3
NTFLX-2	NTFLX-92	5/27/2010	2
NTFLX-1	NTFLX-666	6/6/2016	5
...
...

User (ID)	Movie (ID)	Date of Grade	Grade [0-10]
IMDB-0	IMDB-173	4/20/2009	5
IMDB-1	IMDB-18	10/18/2008	0
IMDB-2	IMDB-341	5/27/2010	-
...
...
...

- Many users are shared
- The grades of same users are correlated
- A user grades one movie around the same date in two databases
- IMDB users are public
- NetFLIX and IMdB moves are public

Linking Attacks: Case of Netflix Prize



User (ID)	Movie (ID)	Date of Grade	Grade [1,2,3,4,5]
NTFLX-0	NTFLX-19	10/12/2008	1
NTFLX-1	NTFLX-116	4/23/2009	3
NTFLX-2	NTFLX-92	5/27/2010	2
NTFLX-1	NTFLX-666	6/6/2016	5
...
...

User (ID)	Movie (ID)	Date of Grade	Grade [0-10]
IMDB-0	IMDB-173	4/20/2009	5
IMDB-1	IMDB-18	10/18/2008	0
IMDB-2	IMDB-341	5/27/2010	-
...
...
...

- Many users are shared
- The grades of same users are correlated
- A user grades one movie around the same date in two databases

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

Large-scale RNA

- Recent advent of much large scale RNA-seq (& other functional genomics data) following on DNA sequencing
- Often this is of human subjects & produced by large consortia (eg TCGA, PCAWG, GTEx) and needs to be protected
- Useful to build tools & approaches that interact with these data

The Human Genome Project



ENCODE Pilot



ENCODE Production



Comparative ENCODE



Epigenome Roadmap

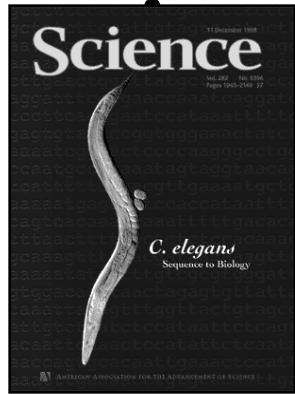


2000

2005

2010

2015



Worm Genome



modENCODE



1000 Genomes Pilot



1000 Genomes Phase 3



GTEx

2-sided nature of functional genomics data: Analysis can be very General/Public or Individual/Private



- General quantifications related to overall aspects of a condition & are not tied to an individual's genotype - ie what genes go up in cancer
 - However, data is derived from an individual & tagged with an individual's genotype
- Other calculations aim to use genotype & specific aspects of the quantification to derive general relations related to sequence variation & gene expression
- Some calculations and data derive finding very specific to the variants in a particular individual

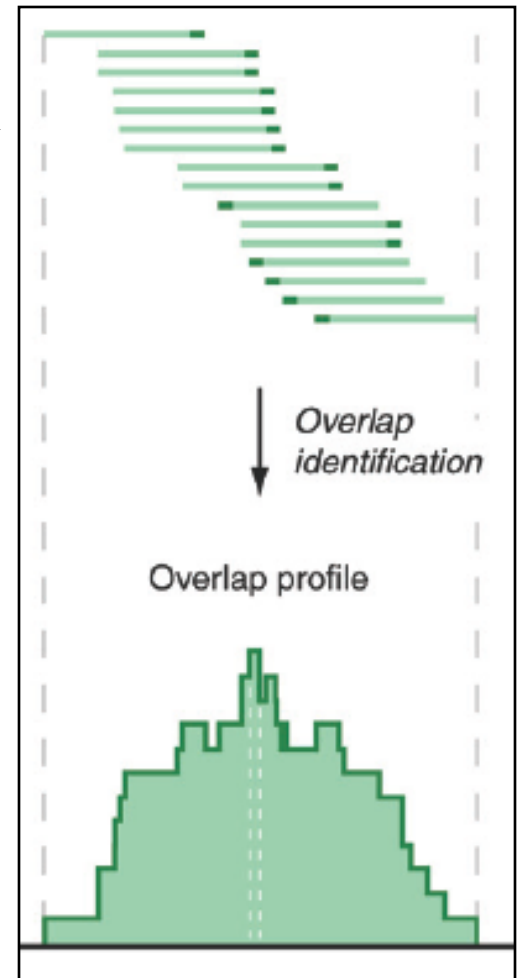
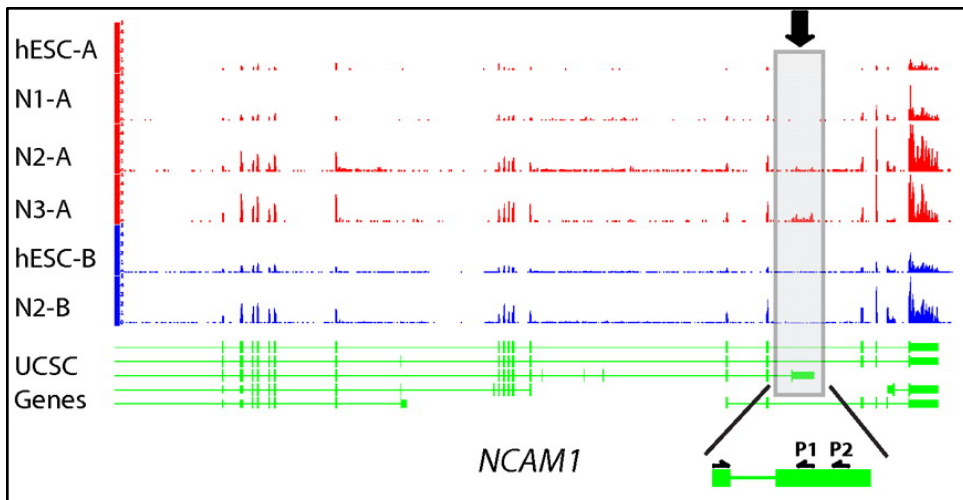
RNA-seq

RNA-seq uses next-generation sequencing technologies to reveal RNA presence and quantity within a biological sample.

```
ATACAAGCAAGTATAAGTTCGTATGCCGTCTT
GGAGGCTGGAGTTGGGGACGTATGCGGCATAG
TACCGATCGAGTCGACTGTAAACGTAGGCATA
ATTCTGACTGGTGTGCATGCTGATGTACTTAA
```

Reads (fasta)

- Quality scores (fastq)
- Mapping (BAM)
- Contain variant information in transcribed regions

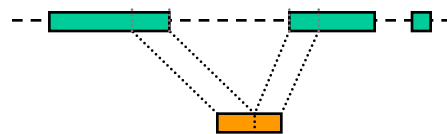
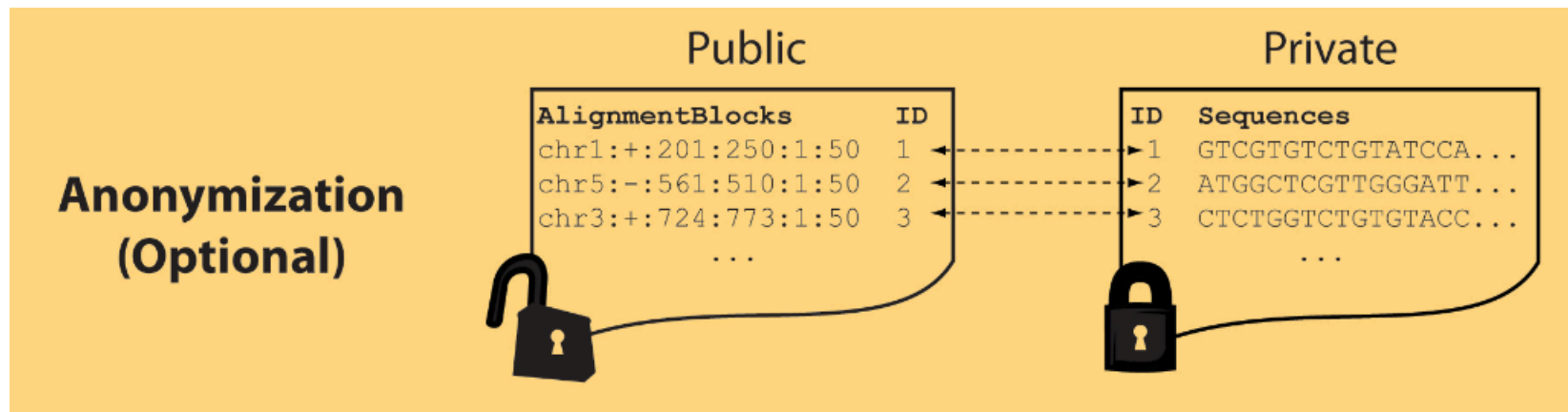


Reads => Signal

Quantitative information from RNA-seq signal: average signals at exon level (RPKM)

Light-weight formats

- Some lightweight format clearly separate public & private info., aiding exchange
- Files become much smaller
- Distinction between formats to compute on and those to archive with – become sharper with big data



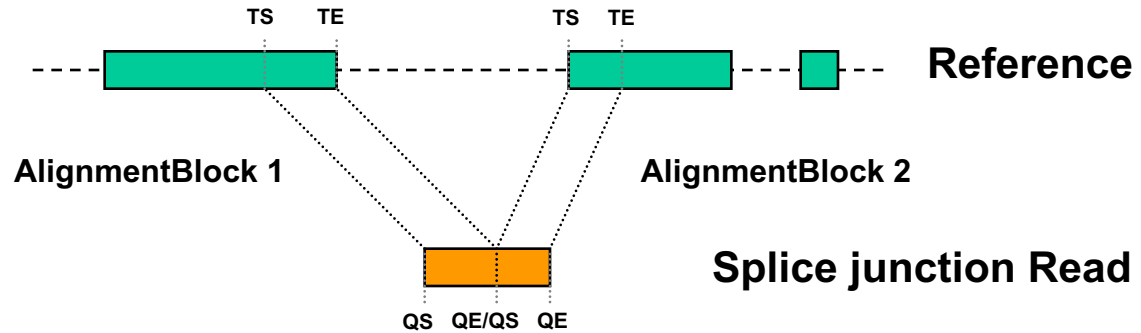
Mapping coordinates without variants (MRF)

Reads (linked via ID, 10X larger than mapping coord.)

MRF Examples

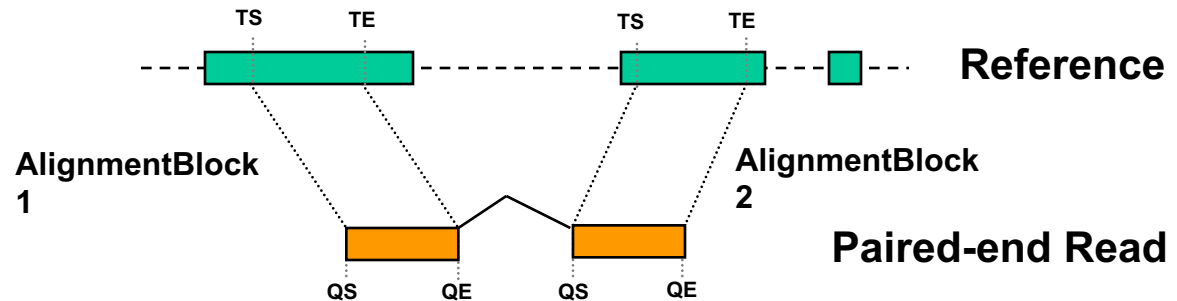
Reference based compression
(ie CRAM)
is similar but it stores actual variant beyond just position of alignment block

chr2:+:601:630:1:30,chr2:+:921:940:31:50



Legend: TS = TargetStart, TE = TargetEnd, QS = QueryStart, QE = QueryEnd

chr9:+:431:480:1:50|chr9:+:945:994:1:50



Legend: TS = TargetStart, TE = TargetEnd, QS = QueryStart, QE = QueryEnd

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

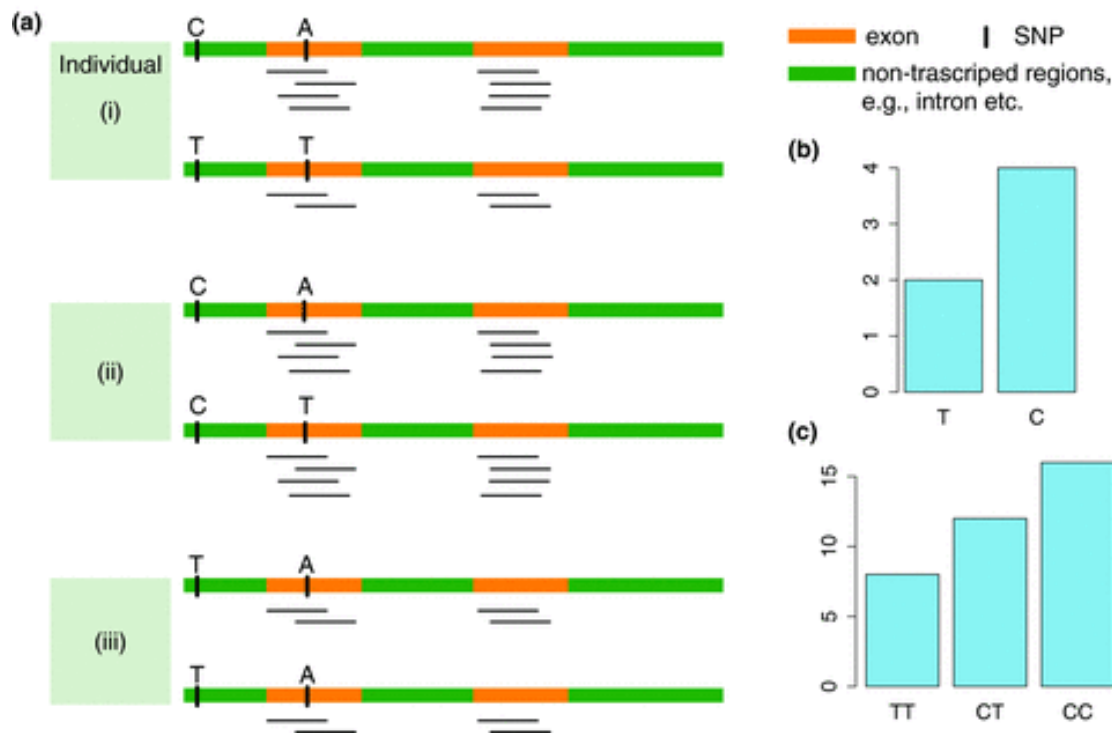
- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

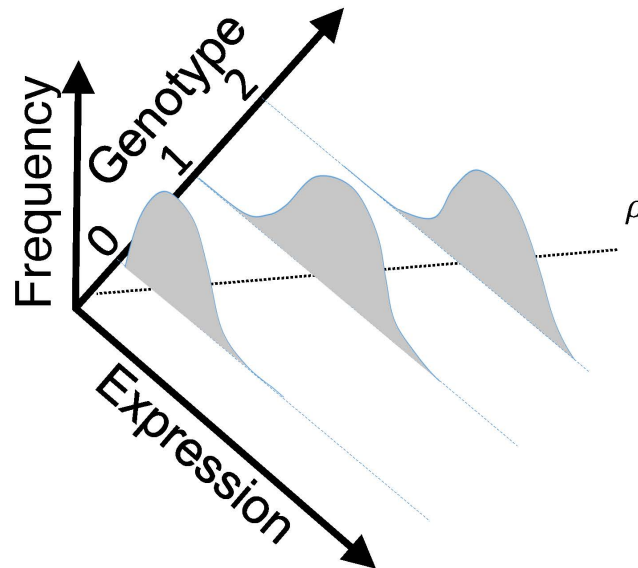
- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants



eQTL Mapping Using RNA-Seq Data

- eQTLs are genomic loci that contribute to variation in mRNA expression levels
- eQTLs provide insights on transcription regulation, and the molecular basis of phenotypic outcomes
- eQTL mapping can be done with RNA-Seq data

[*Biometrics* 68(1) 1–11]



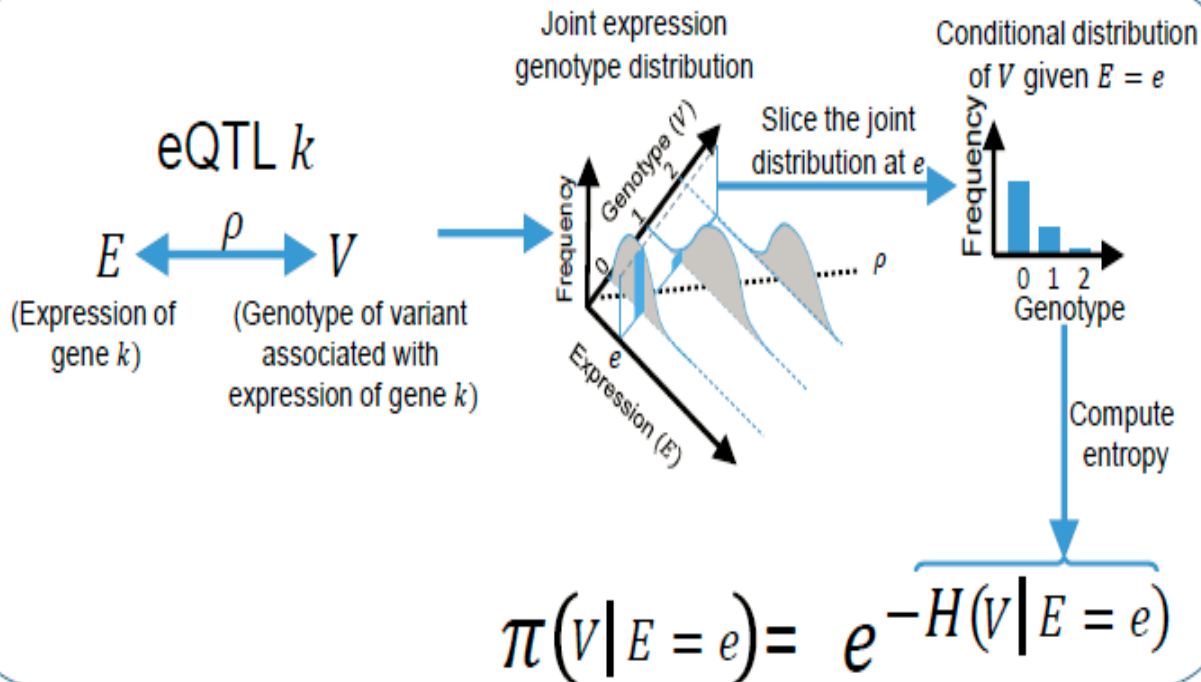
Information Content and Predictability

$$ICI \left(\begin{array}{l} \text{Individual has variant} \\ \text{genotypes } g_1, g_2, \dots, g_n \\ \text{for variants } V_1, V_2, \dots, V_n \end{array} \right) = \log \left(\frac{1}{\text{Frequency of } V_1 \text{ genotype}} \right) + \log \left(\frac{1}{\text{Frequency of } V_2 \text{ genotype}} \right) + \dots + \log \left(\frac{1}{\text{Frequency of } V_n \text{ genotype}} \right)$$

$g_1 = 2$ $g_2 = 1$ $g_n = 2$

V_1 genotype frequencies V_2 genotype frequencies V_n genotype frequencies

- Higher frequency: Lower ICI
- Lower frequency: Higher ICI
- Additive for multiple variants



- Higher cond. entropy: Lower predictability
- Lower cond. entropy: Higher predictability
- Additive for multiple eQTLs

Representative Expression, Genotype, eQTL Datasets

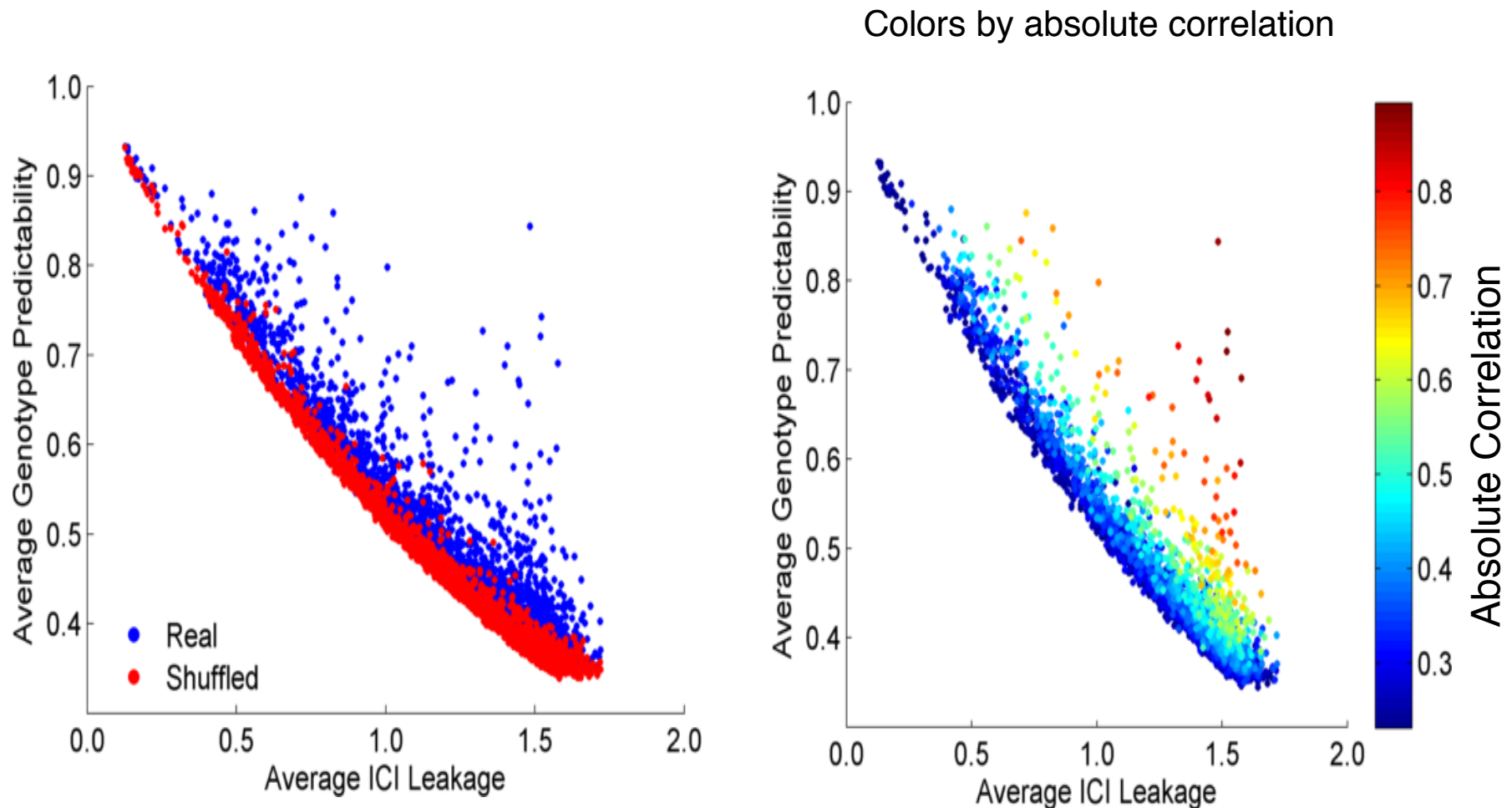
- Genotypes are available from the 1000 Genomes Project
- mRNA sequencing for 462 individuals
 - Publicly available quantification for protein coding genes
- Approximately 3,000 cis-eQTL (FDR<0.05)

1000 Genomes

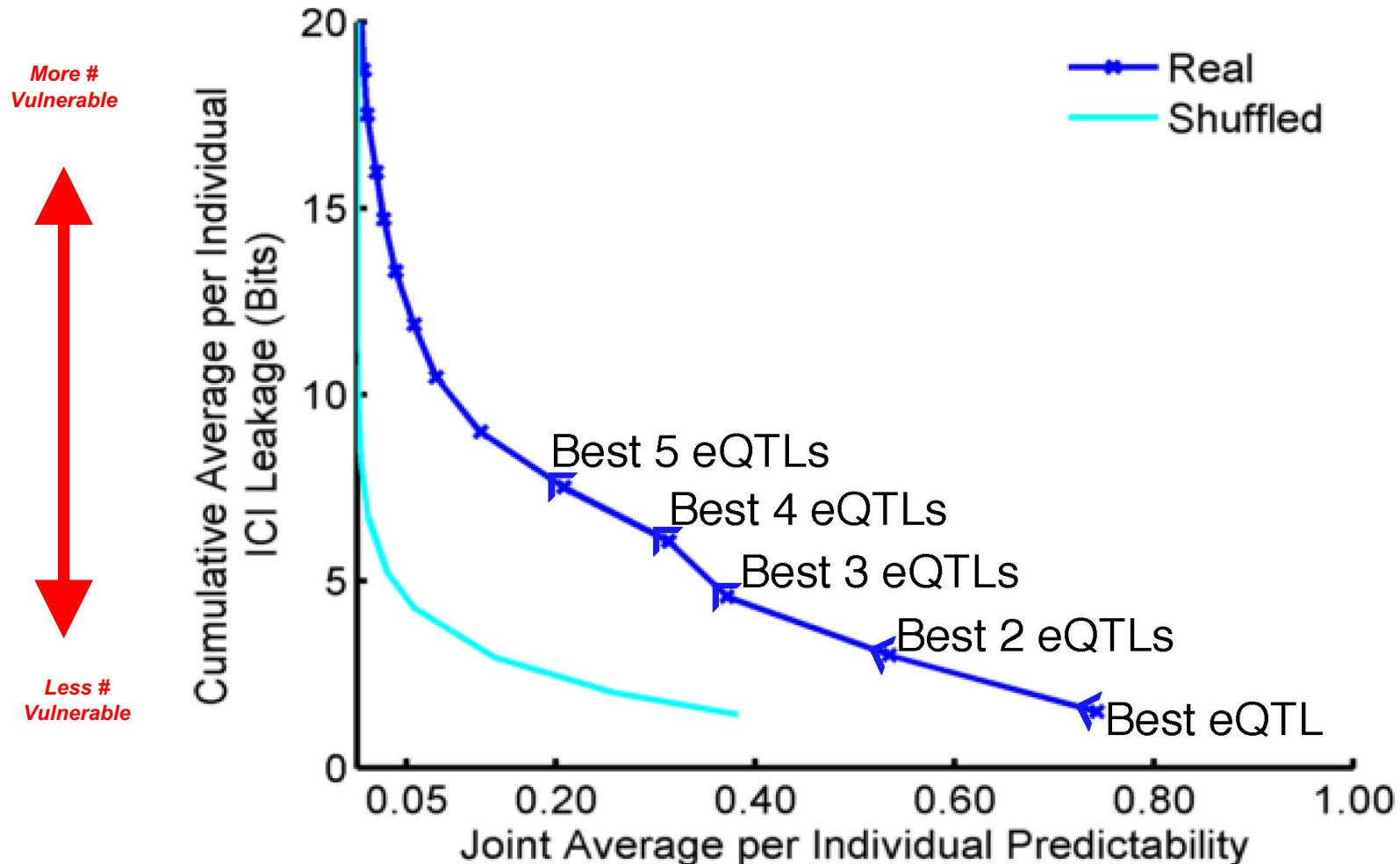
A Deep Catalog of Human Genetic Variation



Per eQTL and ICI Cumulative Leakage versus Genotype Predictability



Cumulative Leakage versus Joint Predictability



Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

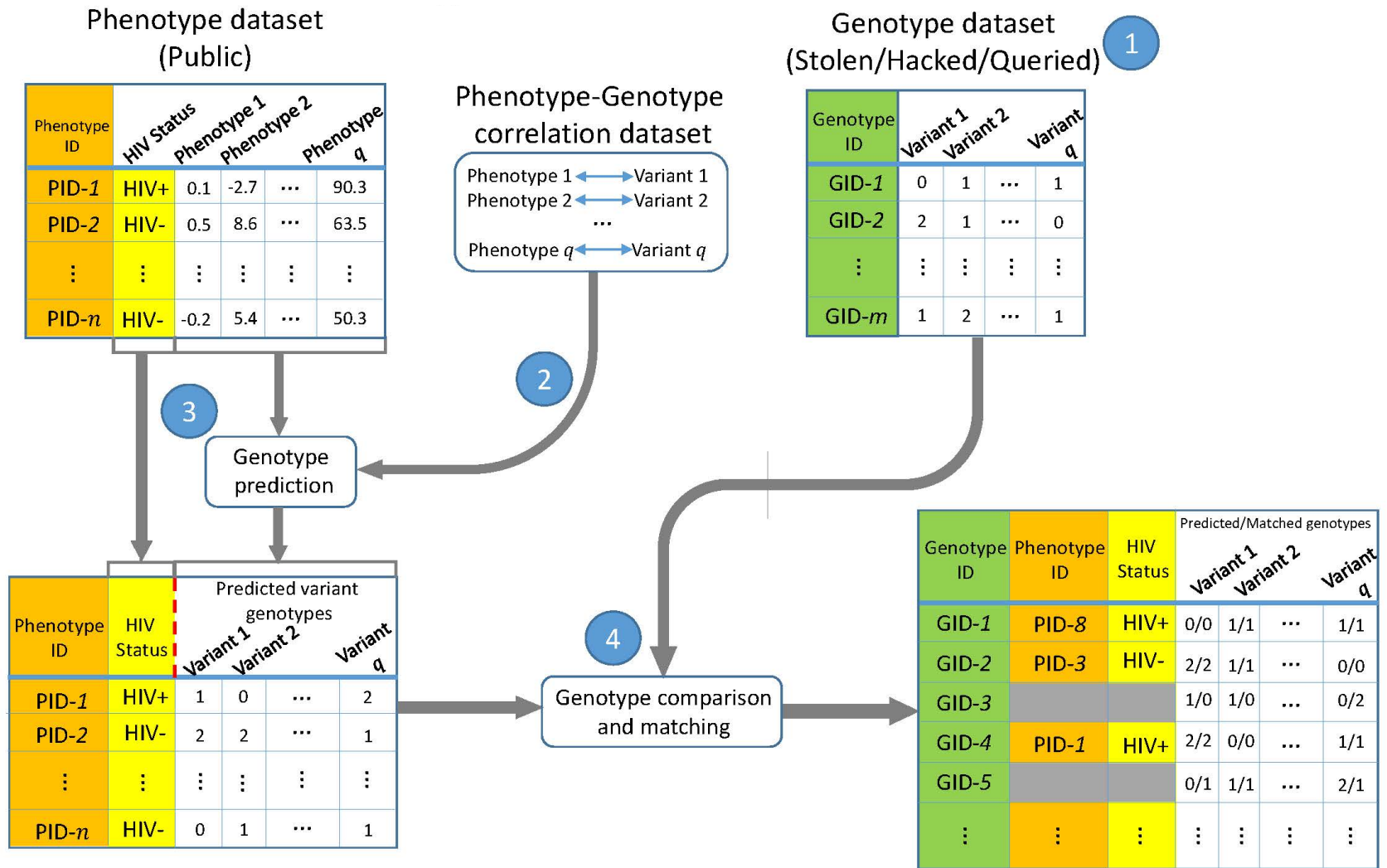
• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

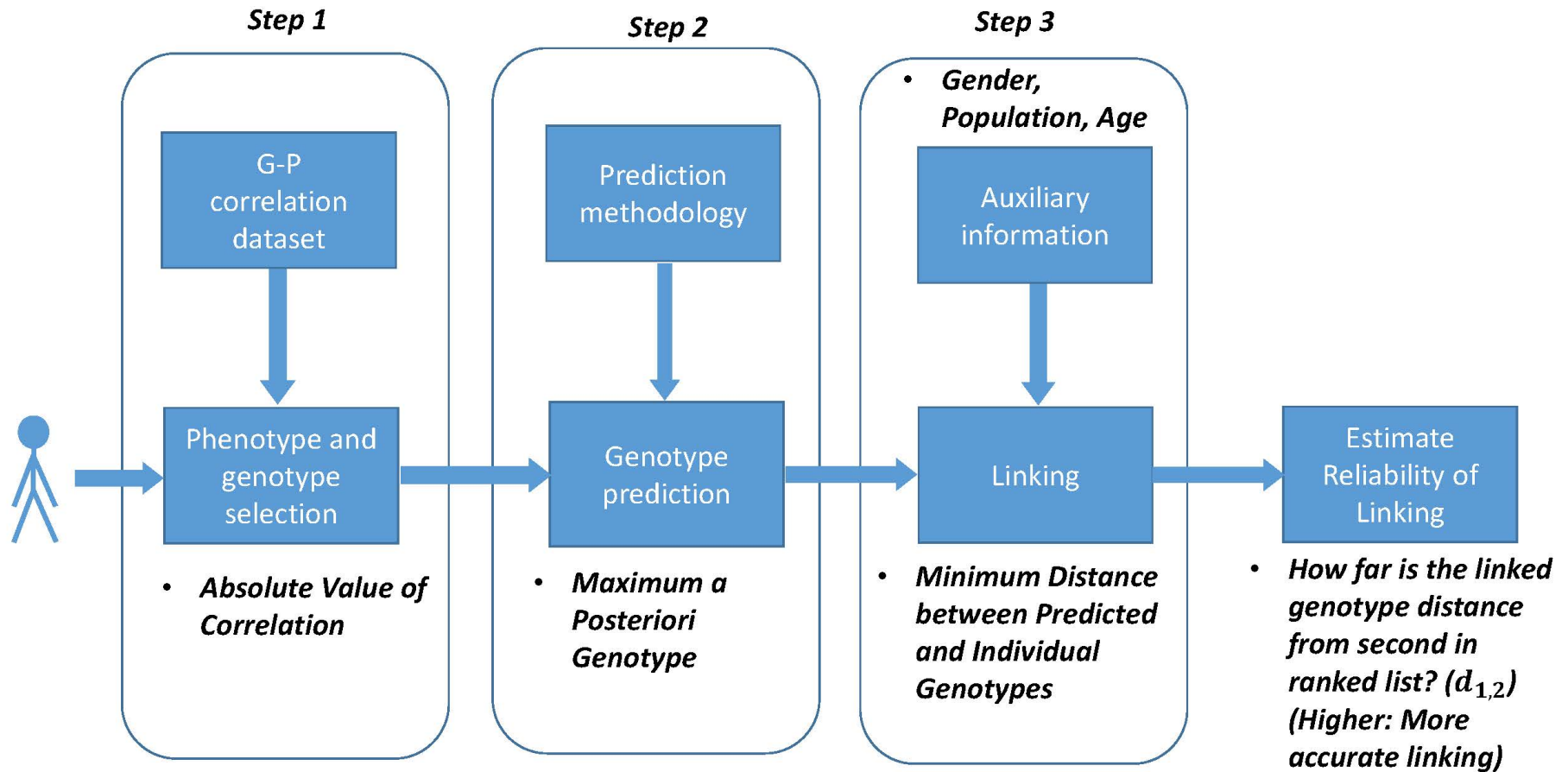
• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

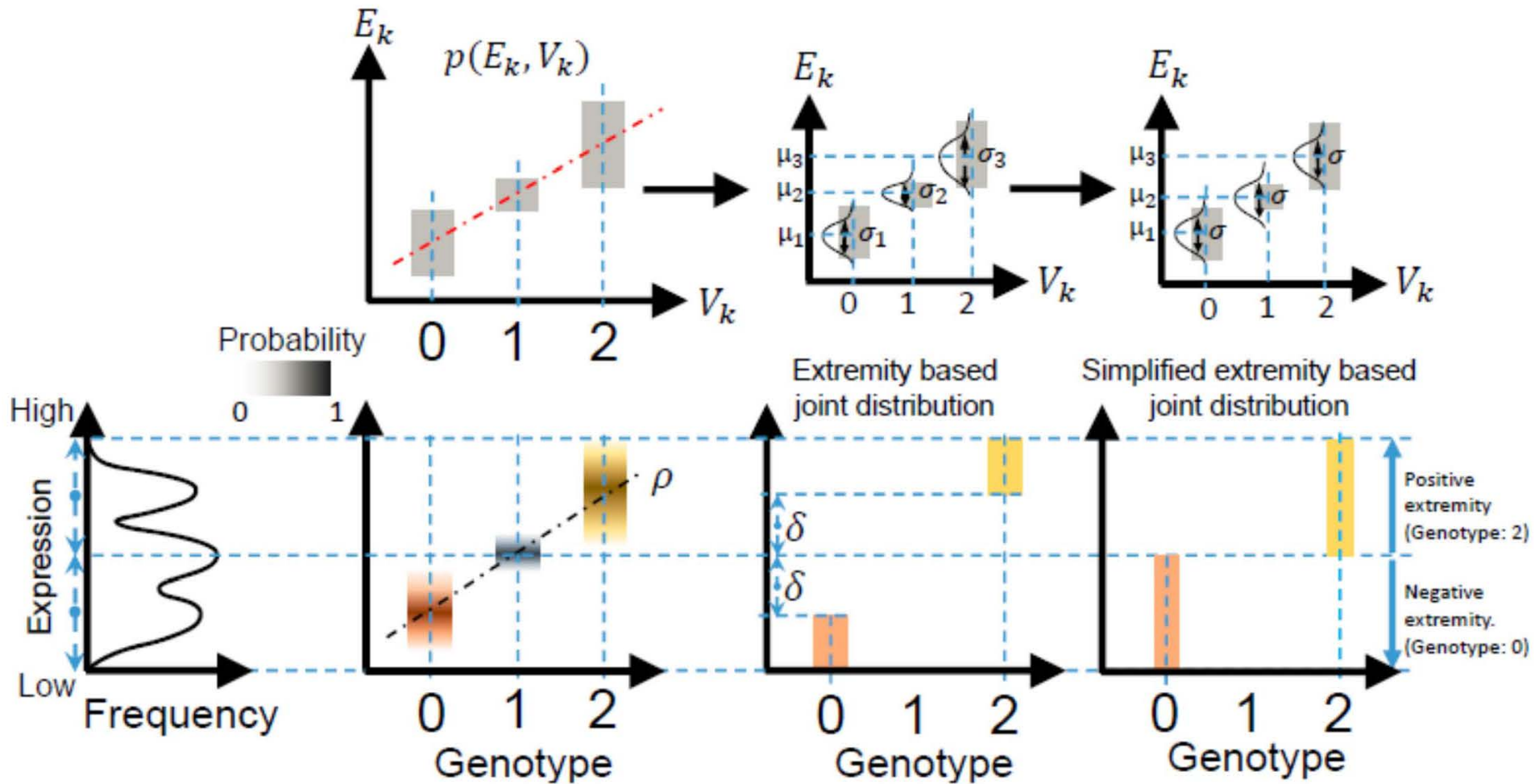
Linking Attack Scenario



Steps in Instantiation of a (Mock) Linking Attack



Levels of Expression-Genotype Model Simplifications



Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

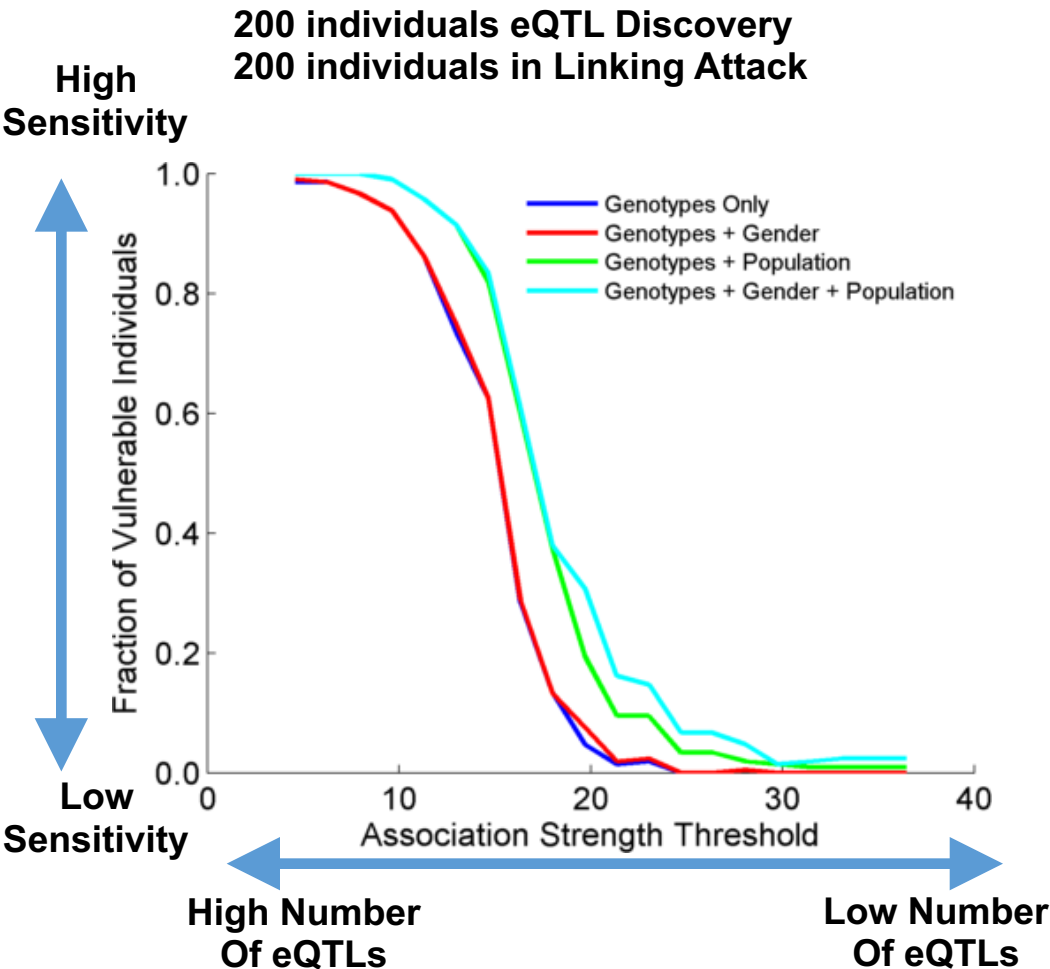
• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

Linking Attack with Extremity based Genotype Prediction

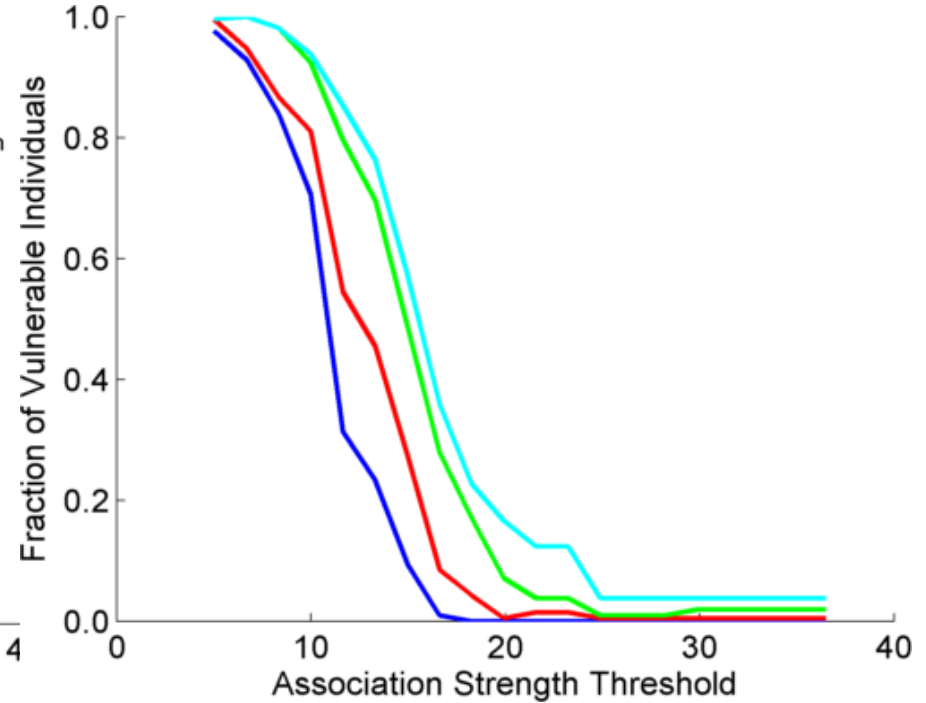
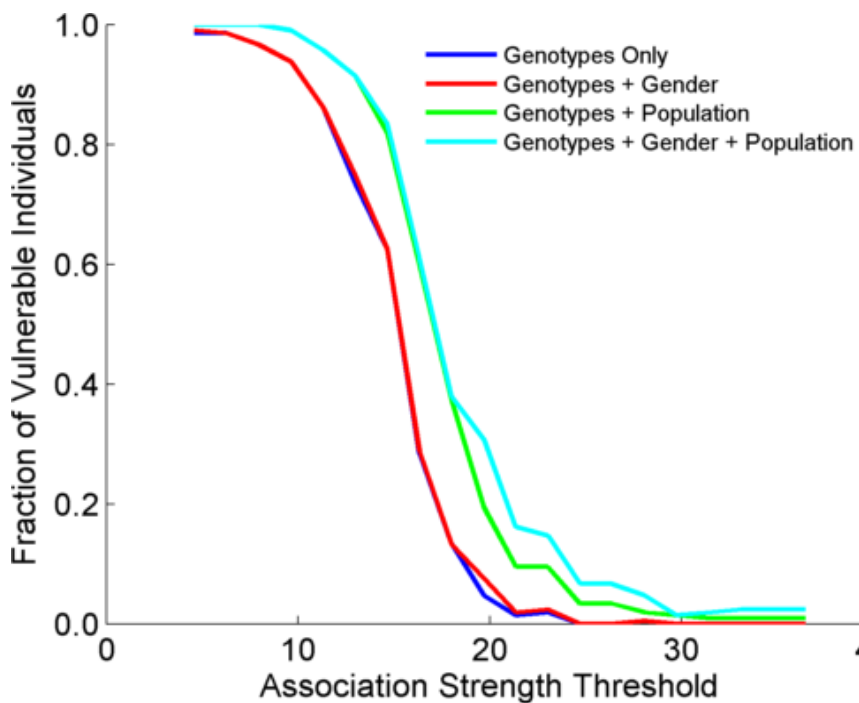


- X-axis: The threshold of association for selecting the eQTLs
 - Higher threshold: Smaller number of eQTLs
- Y-axis: Fraction of correctly linked individuals
 - Measures the **Sensitivity of the attack**

Linking Attack with Extremity based Genotype Prediction

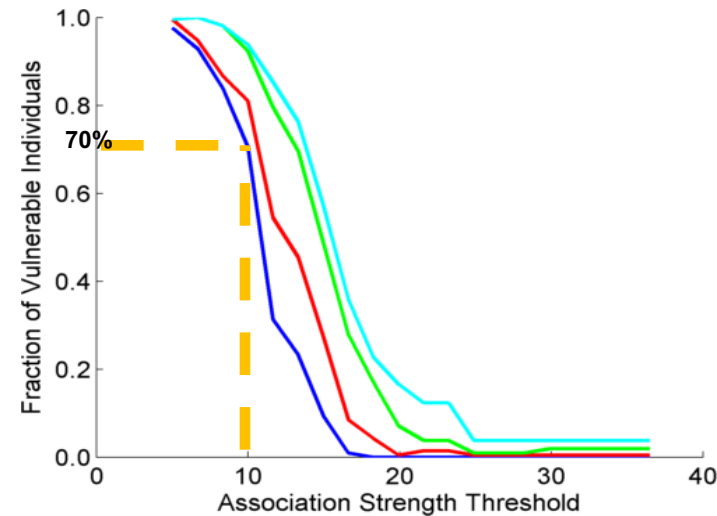
200 individuals eQTL Discovery
200 individuals in Linking Attack

200 individuals eQTL Discovery
100,200 individuals in Linking Attack

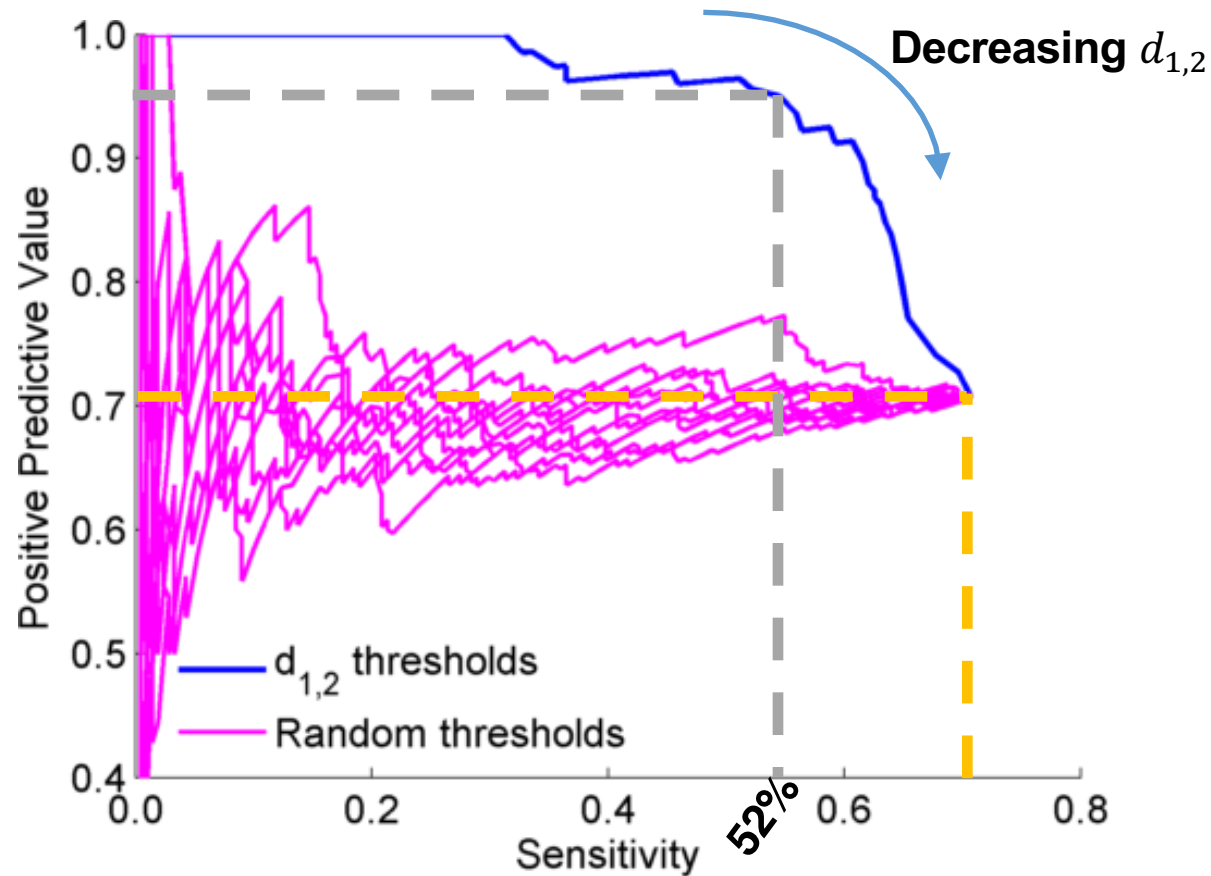


Which 70%?

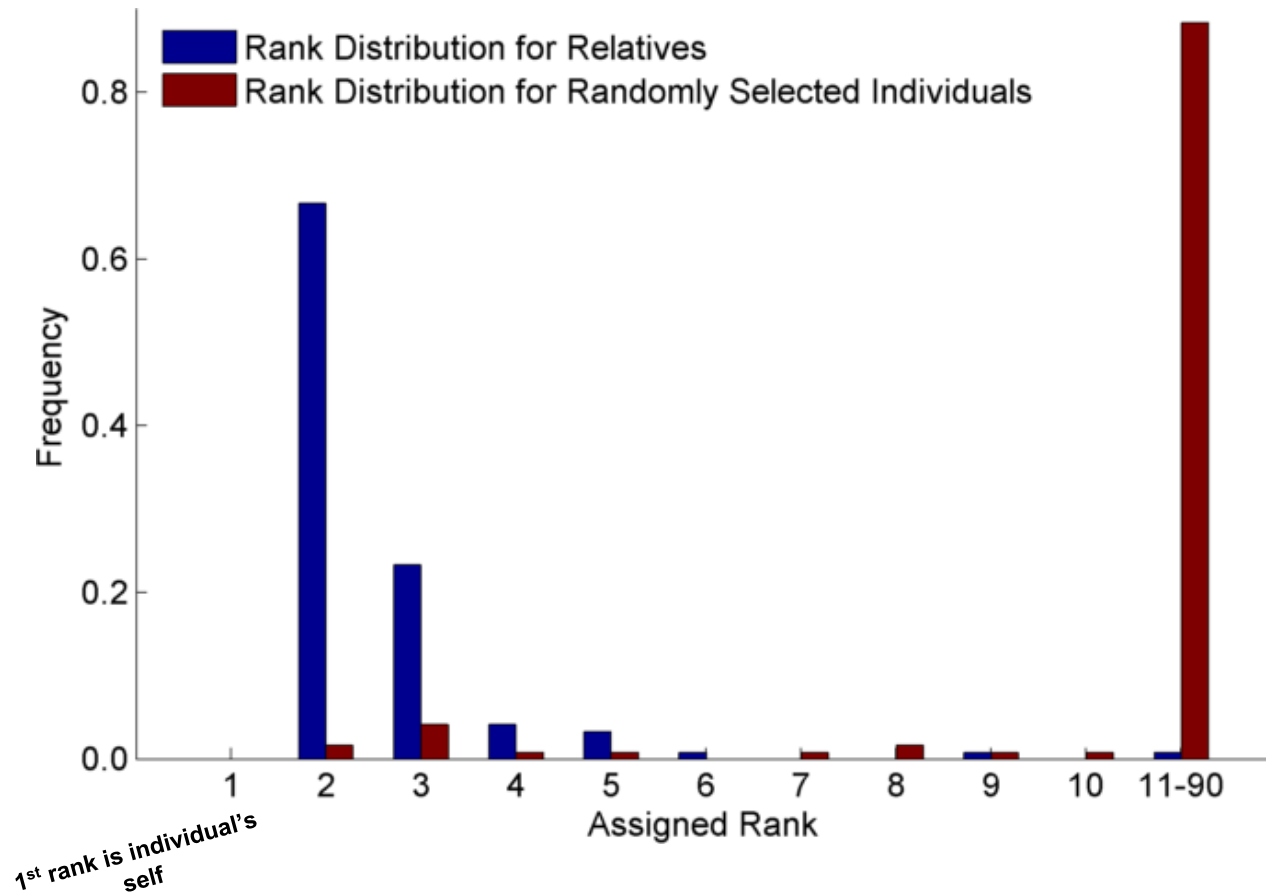
- Attacker arbitrarily selects eQTLs with association strength above 10
- 70% of the individuals are linked correctly
- But which 70%?
- Is there a way to differentiate between linkings to distinguish their reliability?
- First Distance Gap:
 - Difference between the genotype distance of second best matching and best matching individuals
 - $d_{1,2} = d_{second} - d_{first}$



Sensitivity vs PPV for Linkings selected per *first distance gap*, $d_{1,2}$



Relatives are also vulnerable (30 CEU Trios)

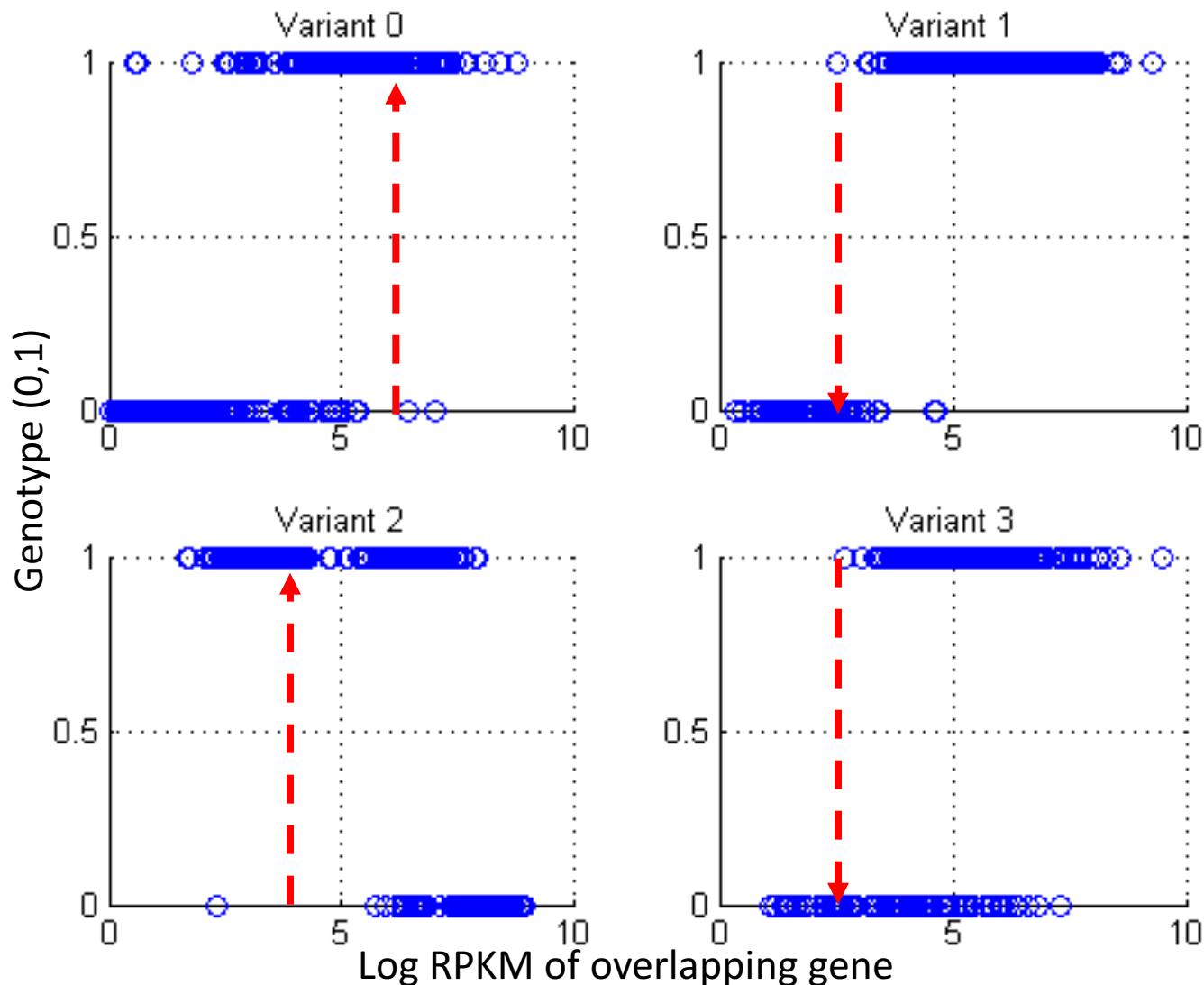


Small Data Leakage from just Gene Expression Data: 4 eQTL-SNP genotypes

Example: Vulnerable sample variants, expressions

- Variant 0 (1, 6)
- Variant 1 (0, 2)
- Variant 2 (1, 3)
- Variant 3 (0, 2)

Expression levels are outliers and are predictive of the genotype!



Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

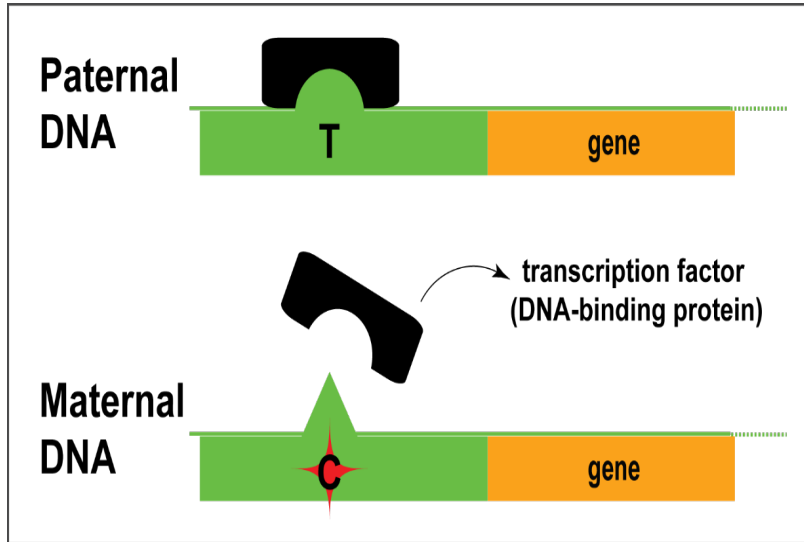
• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

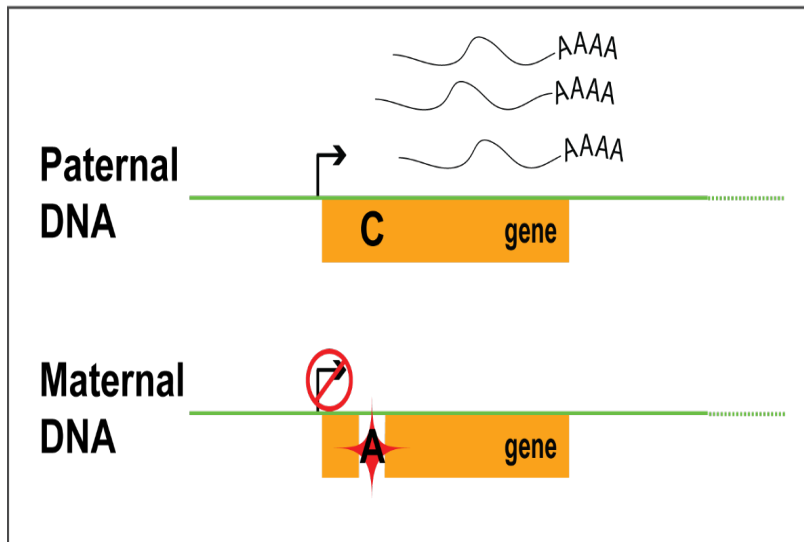
• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

Allele-specific binding and expression



Genomic variants affecting allele-specific behavior e.g. allele-specific binding (ASB)



e.g. allele-specific expression (ASE)

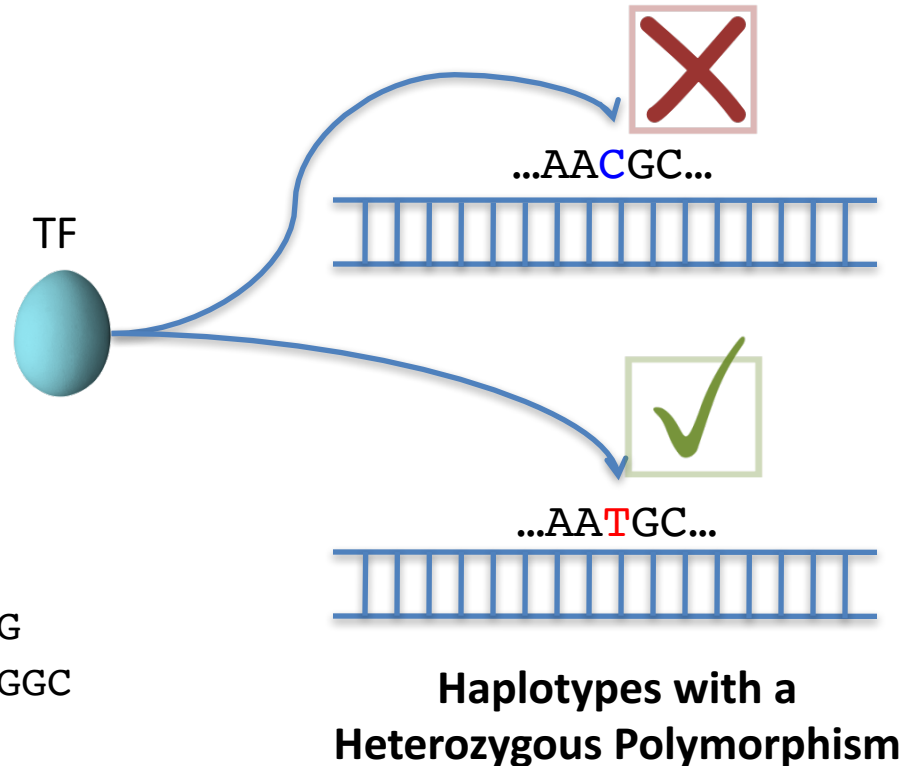
Inferring Allele Specific Binding/Expression using Sequence Reads

RNA/ChIP-Seq Reads

ACTTTGATAGCGTCAAT**T**G
 CTTTGATAGCGTCAAT**T**GC
 CTTTGATAGCGTCAAC**C**GC
 TTGACAGCGTCAAT**T**GCAC
 TGATAGCGTCAAT**T**GCACG
 ATAGCGTCAAT**T**GCACGTC
 TAGCGTCAAT**T**GCACGTCG
 CGTCAAC**C**GCACGTCGGGA
 GTCAAT**T**GCACGTCGAGAG
 CAAT**T**GCACGTCGGGAGTT
 AAT**T**GCACGTCGGGAGTTG
TGCACGTTGGGAGTTGGC

10 x **T**

2 x **C**



Interplay of the annotation and individual sequence variants

Many Technical Issues in Determining ASE/ASB: Reference Bias (naïve alignment against reference)

ASE/ASB Example:

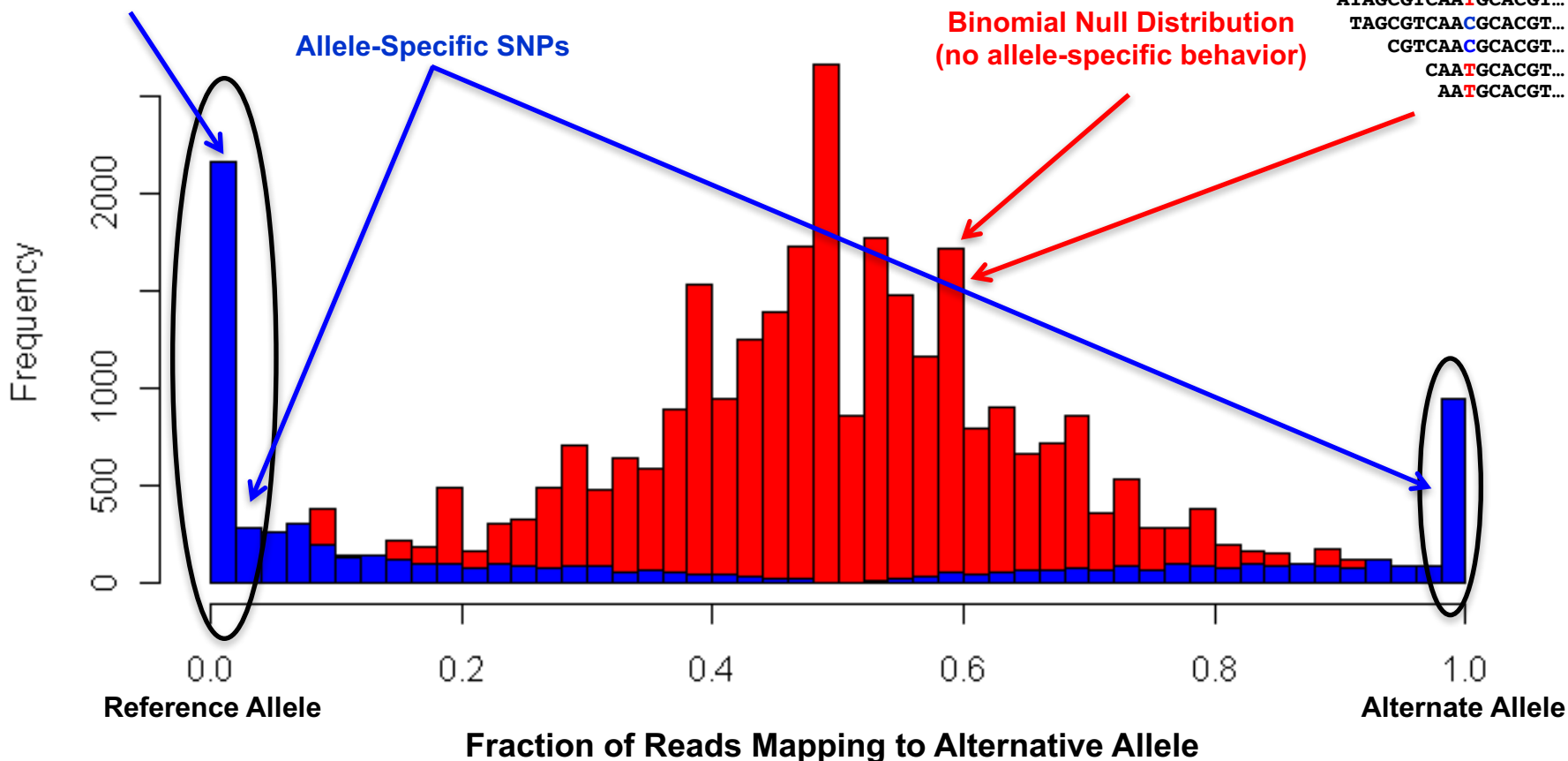
```

...GTCAATGCAC
...GTCAATGCACG
...GTCAATGCACGTC
...GTCAATGCACGTCG
...GTCAACGCACGTCGGGA
GTCAATGCACGTCGAGAG
CAATGCACGTCGGGAGTT
AATGCACGTCGGGAGTTG
    
```

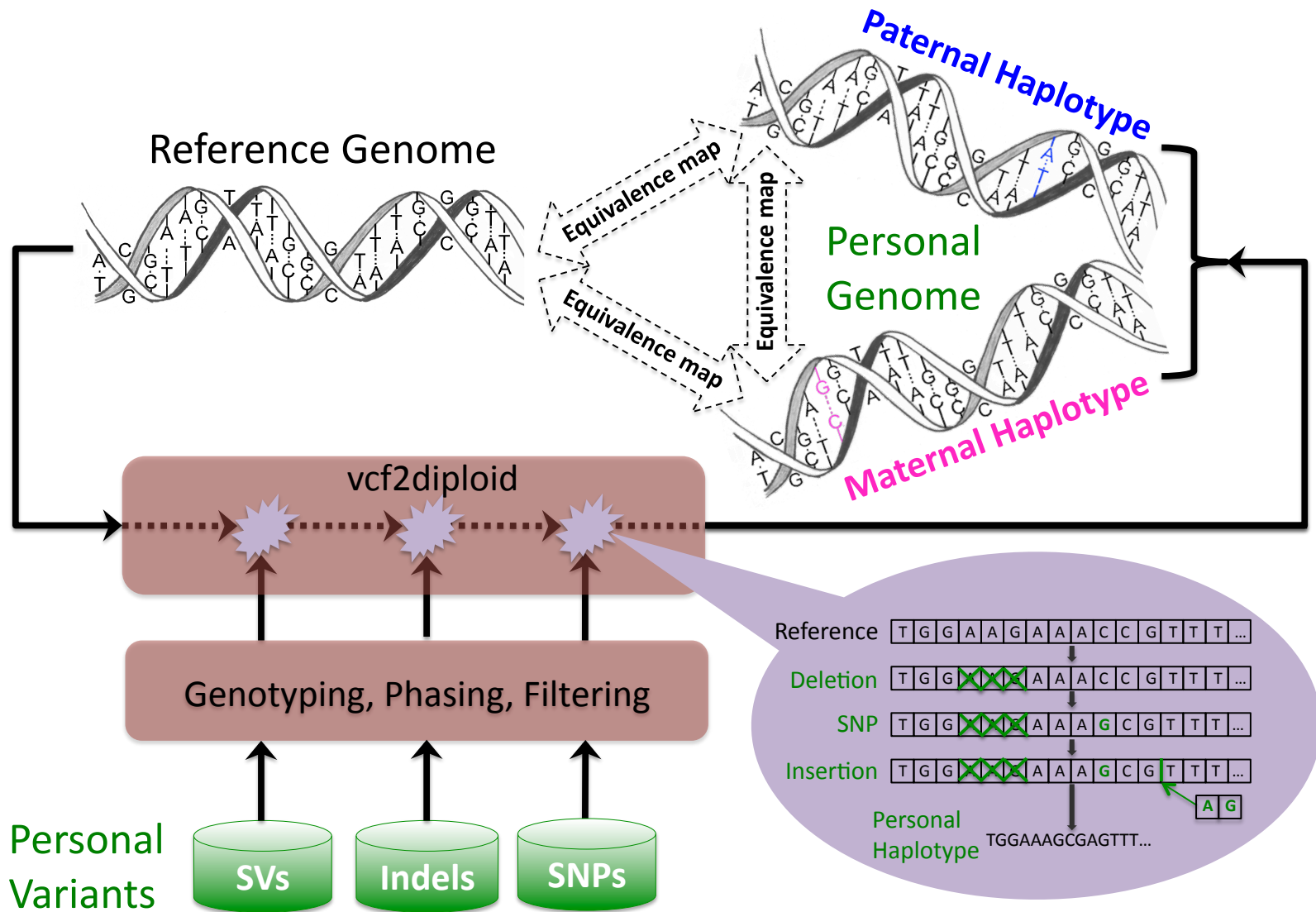
Null Example:

```

ACTTTGATAGCGTCAATG
CTTTGATAGCGTCAACGC
TTGACAGCGTCAATGCAC
ATAGCGTCAATGCACGT...
TAGCGTCAACGCACGT...
CGTCAACGCACGT...
CAATGCACGT...
AATGCACGT...
    
```



Construction of a Personal Diploid Genome & Transcriptome

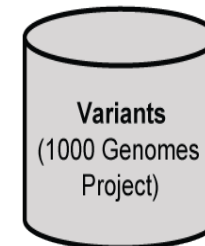


[alleleseq.gersteinlab.org]

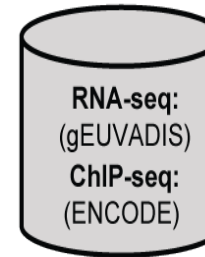
[Rozowsky et al., MSB ('11)]

AlleleDB: Building 382 personal genomes to detect allele-specific variants on a large-scale

1. Build personal genomes

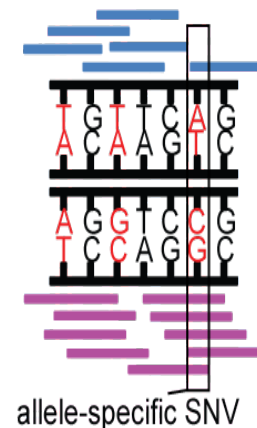


2. Align ChIP-seq & RNA-seq reads

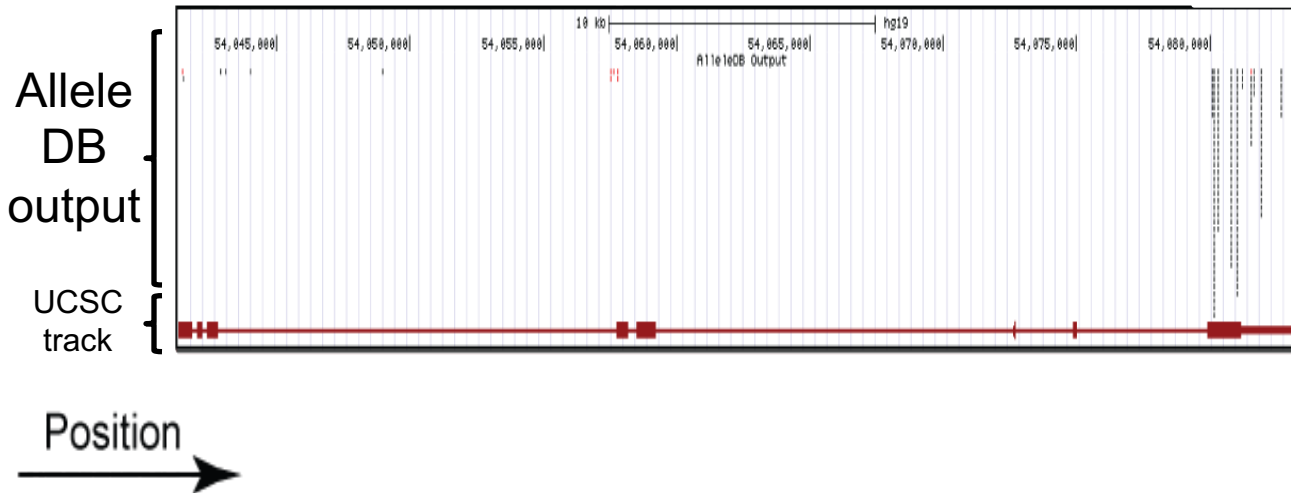


1. Detect allele-specific variants via a series of filters and tests

**Many Technical Issues:
Reference bias, Ambiguous
mapping bias, Over-dispersed
(non binomial null)**

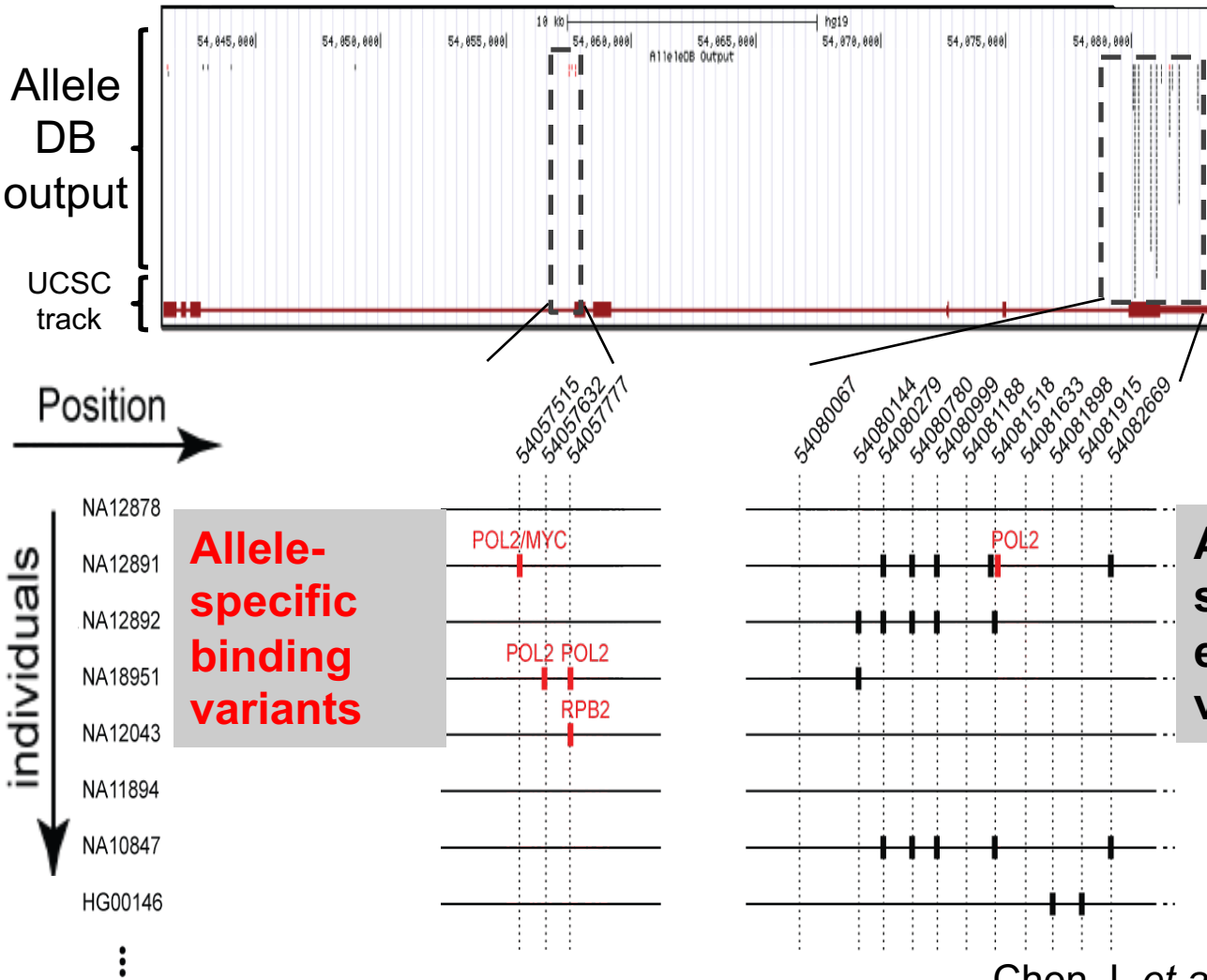


AlleleDB: Annotating rare & common allele-specific variants over a population



- Interfaces with UCSC genome browser
- Showing ZNF331 gene structure

AlleleDB: Annotating rare & common allele-specific variants over a population



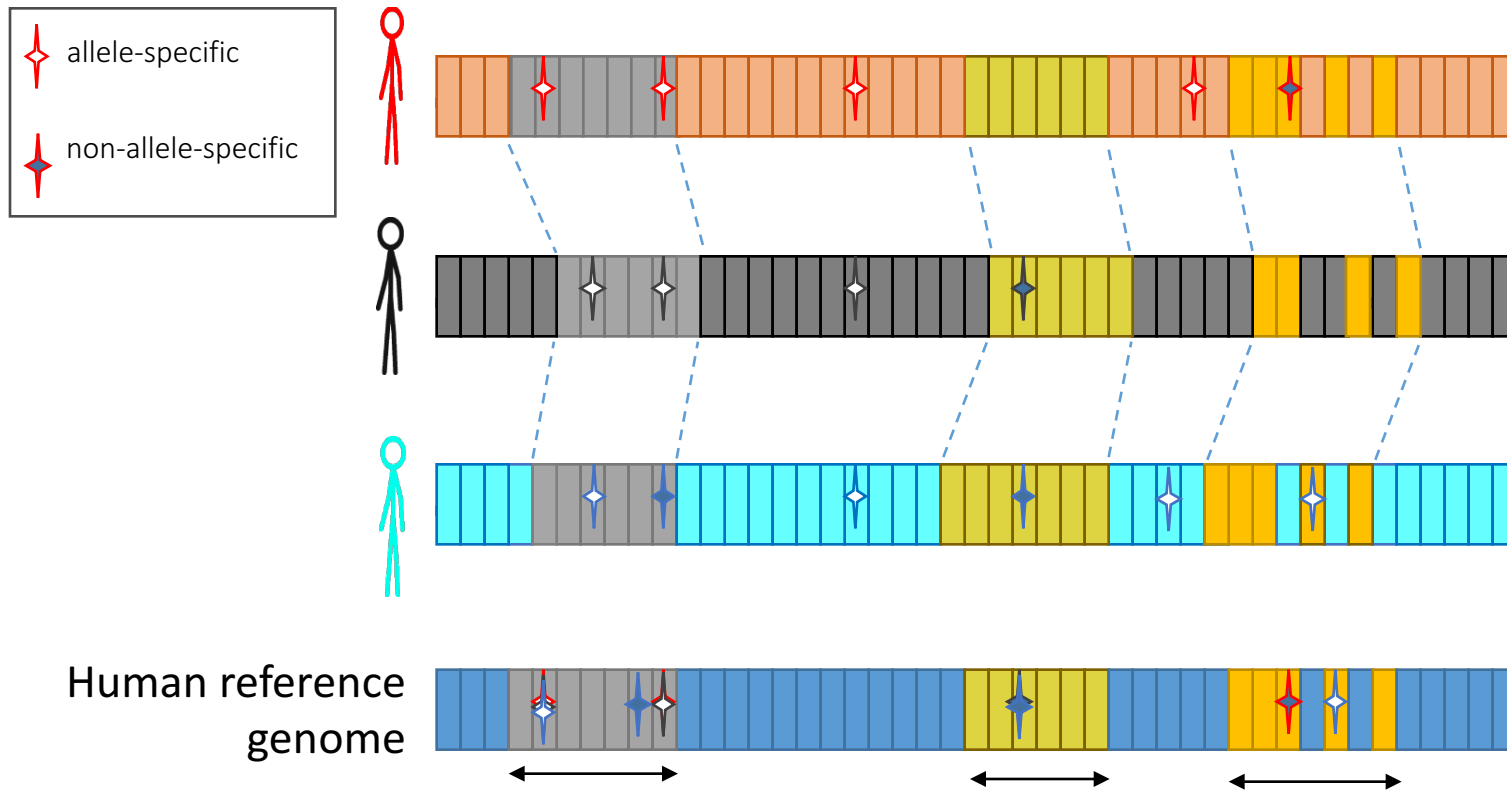
- Interfaces with UCSC genome browser
- Showing ZNF331 gene structure

Chen J. et al. (*Nature Commun.*, '16)

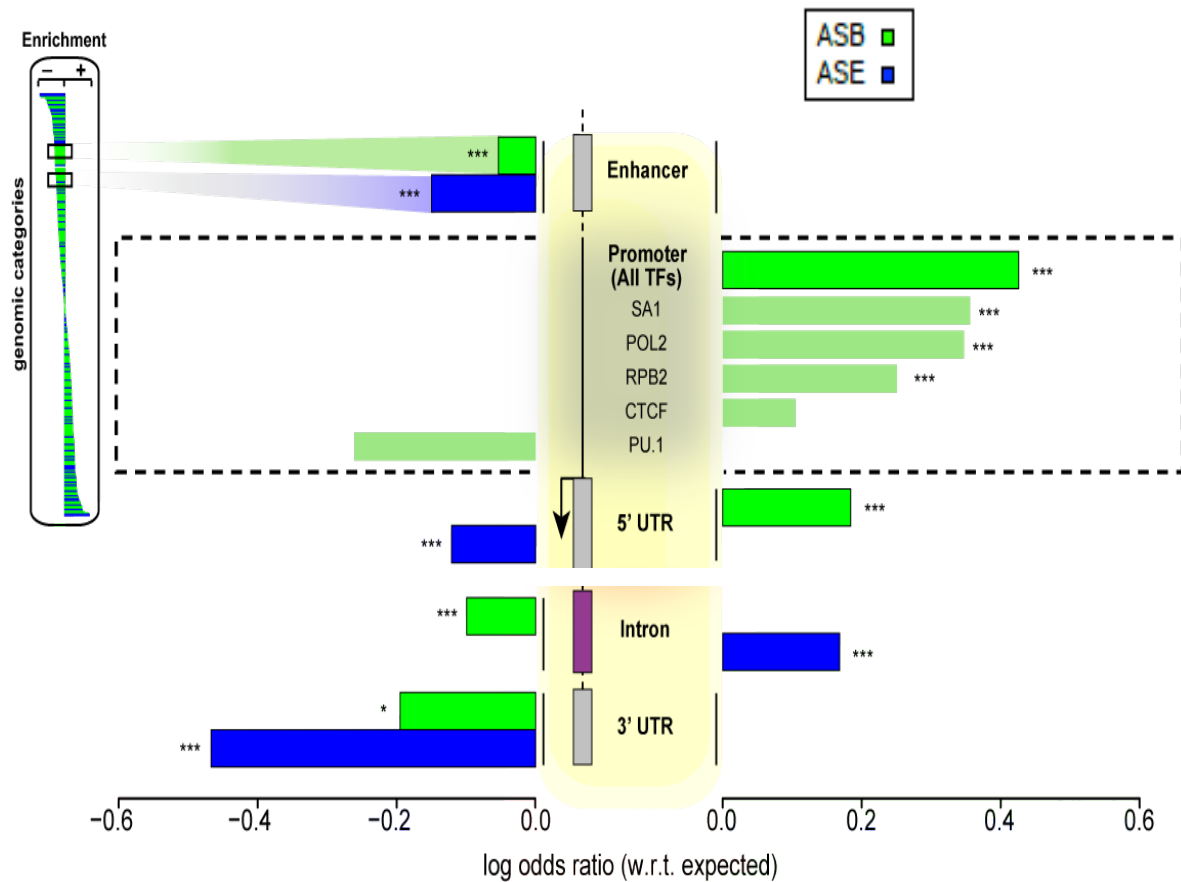
Collecting ASE/ASB variants into allele-specific genomic regions

Does a particular genomic element have a higher tendency to be allele-specific?

Fisher's exact test, for the **enrichment** of allele-specific variants in the element (with respect to non-allele-specific variants that could potentially be called as allelic)



Groups of elements that are enriched or depleted in allelic activity



Chen J. *et al.* (*Nature Commun.*, '16)

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- **Details on Relevant Hacks:** Genomic, Computer Security, & Netflix

• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + **eQTLs using ICI & predictability**
- Instantiating a **practical linking attack** using extreme expression levels
- **Quantifying accuracy of prediction**, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

Genomic Privacy & Individualized RNA-seq: Incompatible or Feasible?

• The General Dilemma of Genomic Privacy

- Fundamental, inherited info that's very private v need for large-scale mining for med. research
- Issues w/ current social & tech approaches: inconsistencies & burdensome security
- Strawman Hybrid Soc-Tech Proposal (Cloud Enclaves. Quantifying Leaks & Closely Coupled priv.-public datasets)
- Details on Relevant Hacks: Genomic, Computer Security, & Netflix

• RNA-seq: How to Publicly Share Some of it

- Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- Removing SNVs in reads w/ MRF
- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack using extreme expression levels
- Quantifying accuracy of prediction, via gap between best & 2nd best match

• Allelic Expression & Binding Activity

- Difference in molecular activity between specific alleles
- RNA-seq calculations that fundamentally involve individual variants

[papers.gersteinlab.org/subject/](http://papers.gersteinlab.org/subject/privacy) **privacy - D Greenbaum**

PrivaSeq.[gersteinlab.org](http://papers.gersteinlab.org) - **A Harmanci**

RSEQtools.[gersteinlab.org](http://papers.gersteinlab.org) [MRF]

L Habegger, A Sboner,

TA Gianoulis, J Rozowsky, A
Agarwal, M Snyder

AlleleDB.[gersteinlab.org](http://papers.gersteinlab.org)

J Chen, J Rozowsky,

T Galeev, A Harmanci,
R Kitchen, J Bedford,
A Abyzov, Y Kong, L Regan

AlleleSeq.[gersteinlab.org](http://papers.gersteinlab.org)

J Rozowsky,

A Abyzov,

J Wang, P Alves, D Raha,
A Harmanci, J Leng,
R Bjornson,
Y Kong,
N Kitabayashi,
N Bhardwaj,
M Rubin,
M Snyder

Acknowledgements

Hiring Postdocs. See gersteinlab.org/jobs !

Default Theme

- Default Outline Level 1
 - Level 2

