

# Genomic Privacy & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information



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



## Activity Patterns

- RNA Seq. gives rise to activity patterns of genes & regions in the genome

# RNA-Seq Overview

```
ATACAAGCAAGTATAAGTTCGTATGCCGTCTT
GGAGGCTGGAGTTGGGGACGTATGCGGCATAG
TACCGATCGAGTCGACTGTAAACGTAGGCATA
ATTCTGACTGGTGTCATGCTGATGTACTTAAA
```

Fastq sequence files  
~5-10 GB

Index-building + Alignment to reference genome

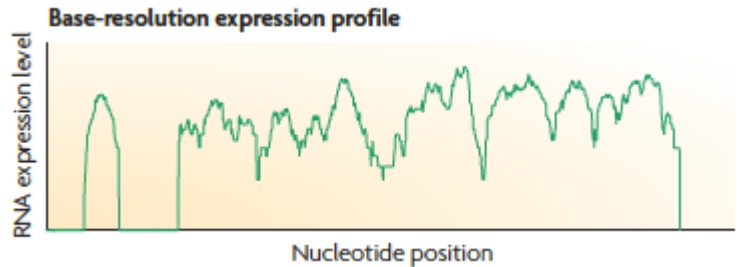
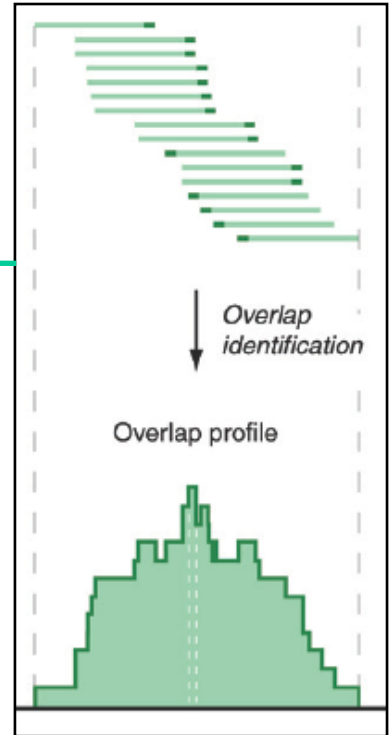
BAM files  
~1-2-fold reduction

Conversion to signal track by overlapping reads

BigWig files  
~25-fold reduction

Mapping to genes

Gene/Transcript expression matrix  
~20-fold reduction



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

Reads => Signal

Successive steps of  
Data Reduction

[NAT. REV. 10: 57; PLOS CB 4:e1000158; PNAS 4:107: 5254]



## 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 – ie gene activity as a function of:
  - Developmental stage: basic patterns and clusters of co-active genes across an organisms development
  - Evolutionary relationships: behavior preserved across a wide range of organisms
  - Tissue- and cell-type
  - Disease phenotypes: what genes go up in cancer?
- **Above are not tied to an individual's genotype. However, data is derived from an individual & tagged with an individual's genotype**
- Note, a few calculations aim to use explicitly genotype to derive general relations related to sequence variation & gene expression (eg allelic activity)

# Importance of Leakage Quantification for Genomic Privacy

- The **overall dilemma of genomic privacy**
  - From sharing information, the individual is potentially harmed but society benefits in terms of medical research
  - How to balance risks v rewards?
- **Need to quantify leakage**
  - Cost Benefit Analysis: how helpful is identifiable data in genomic research v. potential harm from a breach
  - What is acceptable risk – ie what is acceptable data leakage?
- Also, **need 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 & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information

- **Intro on RNA-seq & the General Dilemma of Genomic Privacy**

- RNA-seq Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- The need to quantify leaks

- **Quantifying RNA-seq Leakage ...from Reads**

- Almost as much as WGS
- But can remove SNVs in reads w/ MRF

- **...from eQTLs**

- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack

- **...from Indels/SVs**

- Another source of leakage in RNA-seq data & how to use these for a related linking attack

# Genomic Privacy & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information

- **Intro on RNA-seq & the General Dilemma of Genomic Privacy**

- RNA-seq Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- The need to quantify leaks

- **Quantifying RNA-seq Leakage ...from Reads**

- Almost as much as WGS
- But can remove SNVs in reads w/ MRF

- **...from eQTLs**

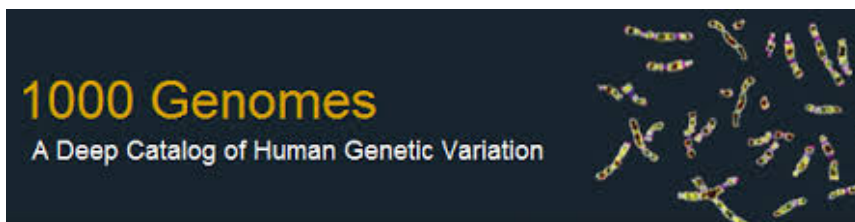
- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack

- **...from Indels/SVs**

- Another source of leakage in RNA-seq data & how to use these for a related linking attack

## Representative Expression, Genotype, eQTL Datasets on Open Datasets

- Publically available genotypes (not controlled access) are available from the 1000 Genomes Project
- mRNA sequencing for 462 individuals from gEUVADIS and ENCODE
  - Publicly available quantification for protein coding genes
- Approximately 3,000 cis-eQTL (FDR<0.05)





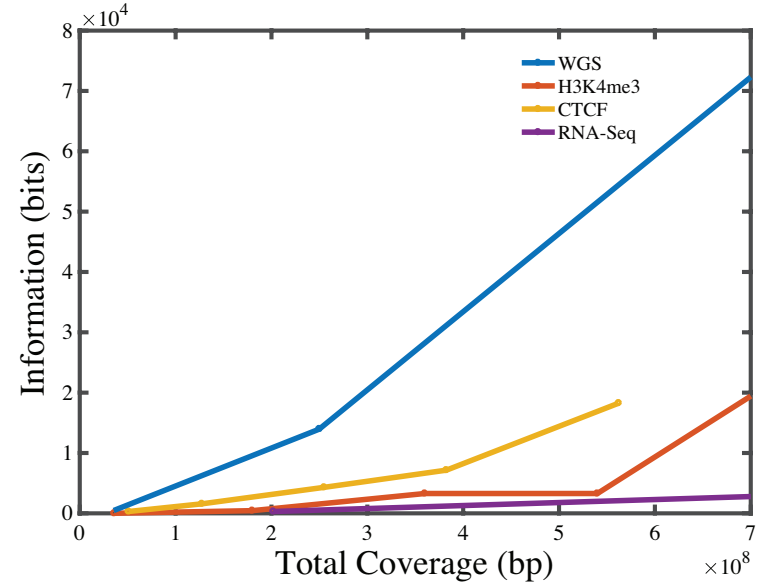
- Functional genomics data comes with a great deal of sequencing
  - NA12878 as case study - 1000 genomes variants are used as gold standard
- How much information, for example, do RNA-Seq reads (or ChIP-Seq) reads contain? Does that information enough to identify individuals?

**Variants from RNA-Seq reads**

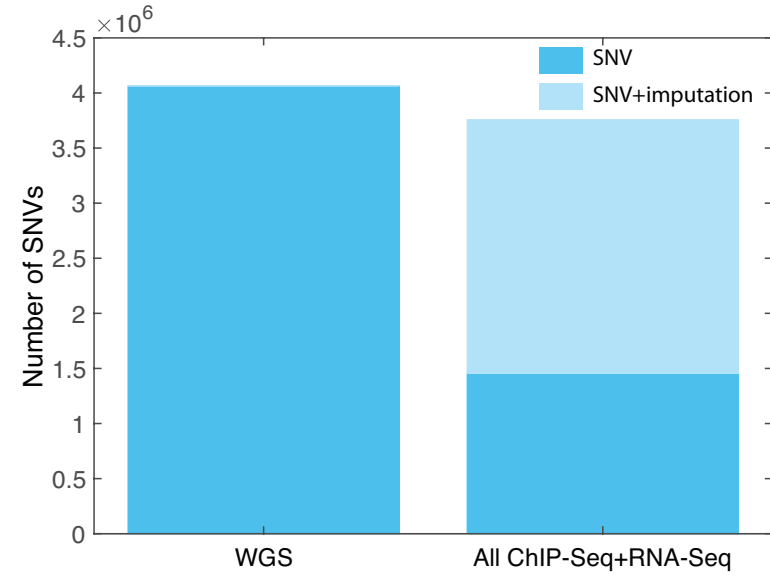
R1	start1	end1	ATAAATGAGGATTTAGAGGTGGTGACC
	reference	genome	ATAAATGAGAATTTTGAGGTGGTGACC
R2	start2	end2	T--ATTTTCTCTCATACCACCTCAACG
	reference	genome	TTTATTTTCT---ATACCACCTCAACG
R3	start3	end3	TTTATTTTCTATACCACCTCAA



# Variants directly in the reads

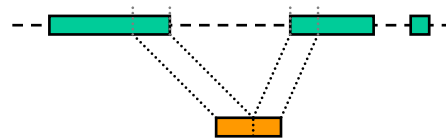
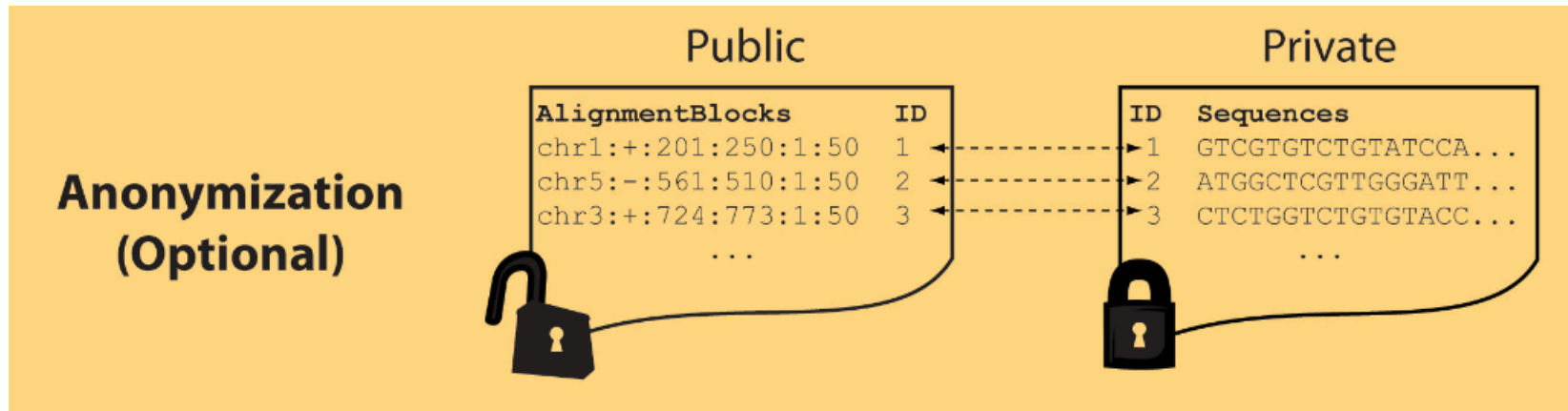


- It might seem like we don't infer much information from single ChIP-Seq and RNA-Seq experiments compared to WGS
  - However putting 10 different ChIP-Seq experiments and RNA-Seq together with imputation provides a great deal of information about the individual



# Light-weight formats to Hide Most of the Read Data (Signal Tracks)

- 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.)**

# Genomic Privacy & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information

- **Intro on RNA-seq & the General Dilemma of Genomic Privacy**

- RNA-seq Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- The need to quantify leaks

- **Quantifying RNA-seq Leakage ...from Reads**

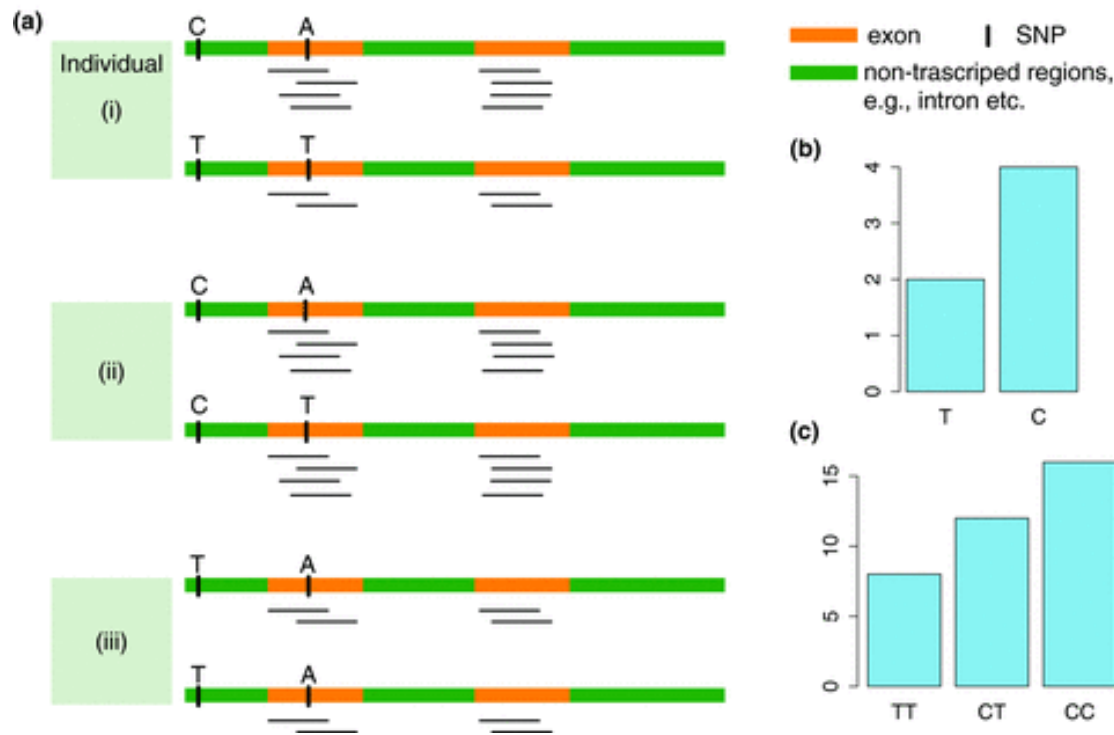
- Almost as much as WGS
- But can remove SNVs in reads w/ MRF

- **...from eQTLs**

- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack

- **...from Indels/SVs**

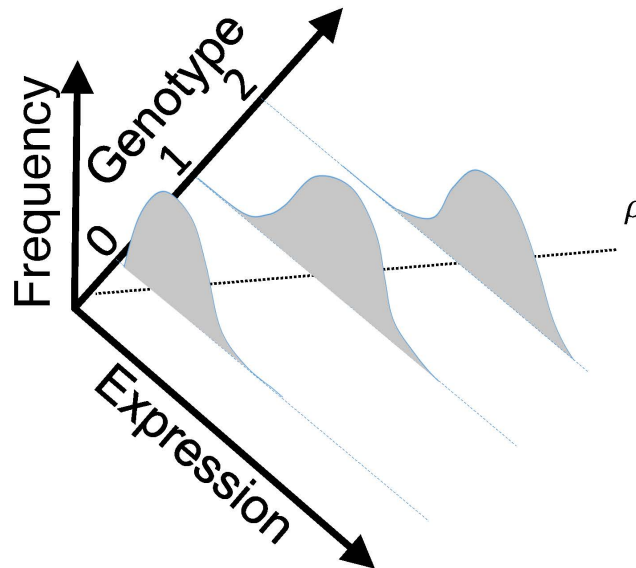
- Another source of leakage in RNA-seq data & how to use these for a related linking attack



# 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

[*Biometrics* 68(1) 1–11]



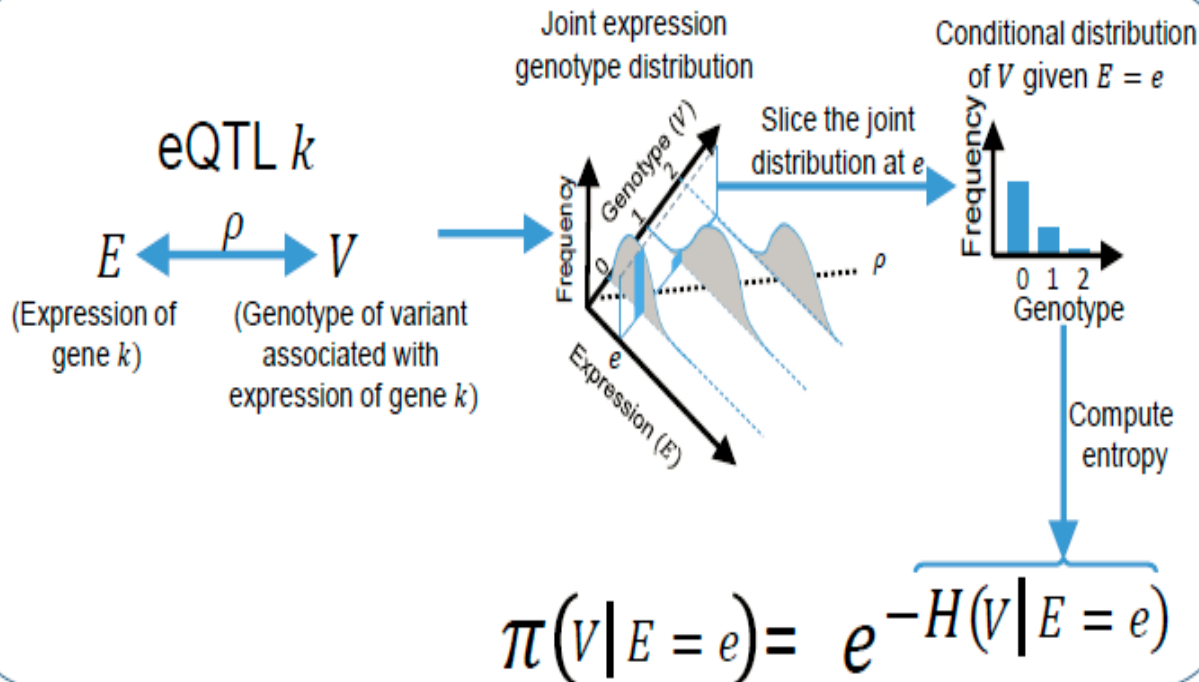
# Information Content and Predictability

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

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

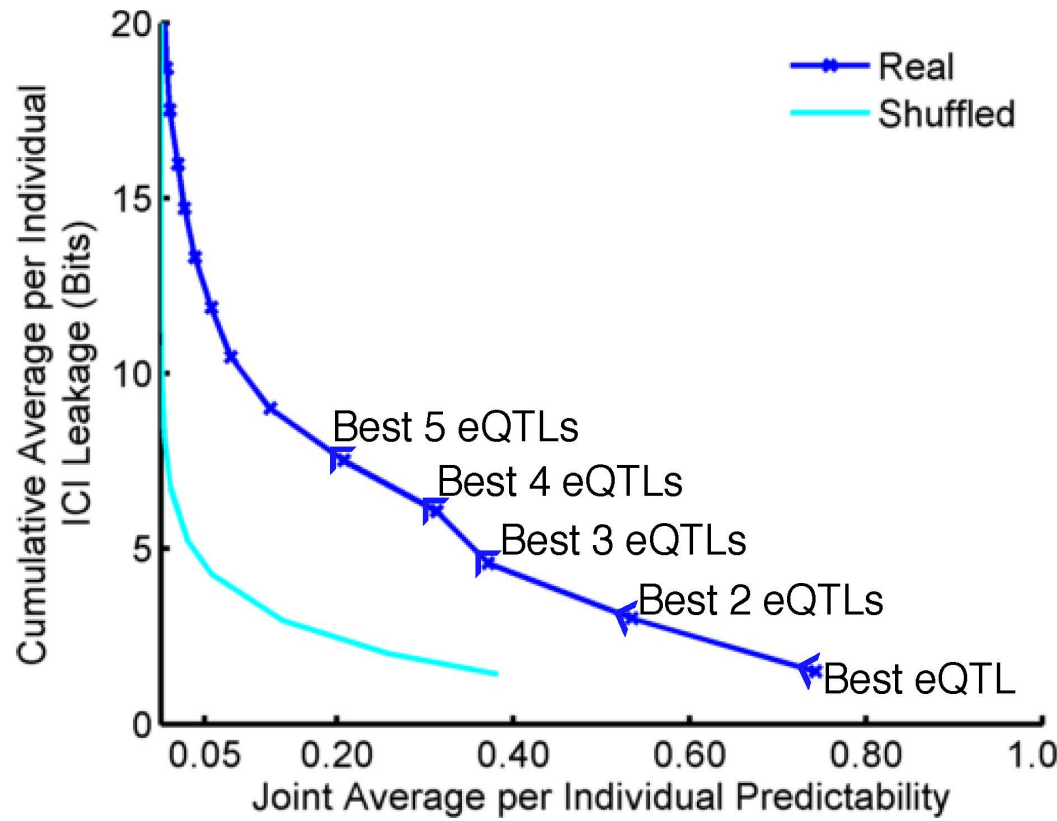
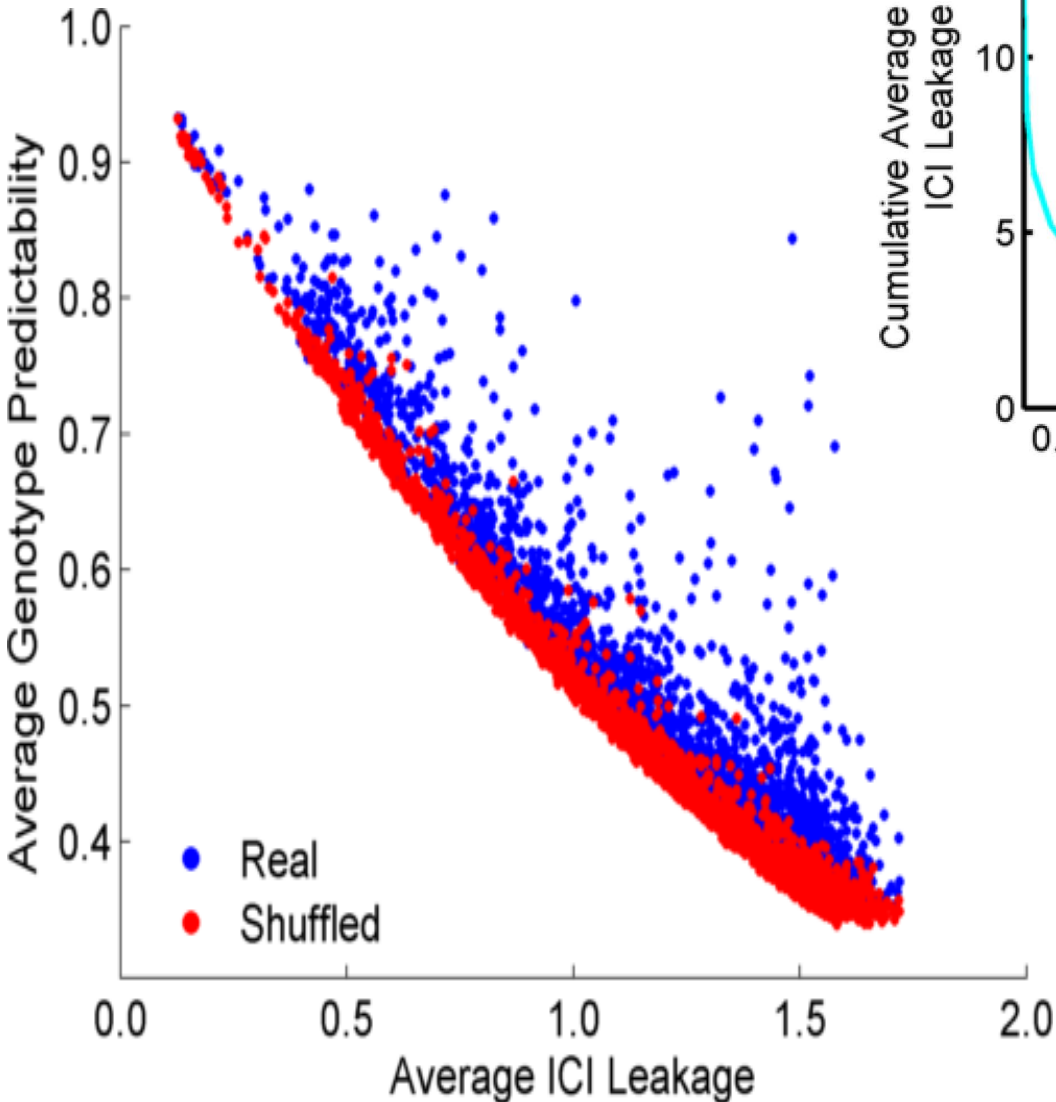
$V_1$  genotype frequencies                       $V_2$  genotype frequencies                       $V_n$  genotype frequencies

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



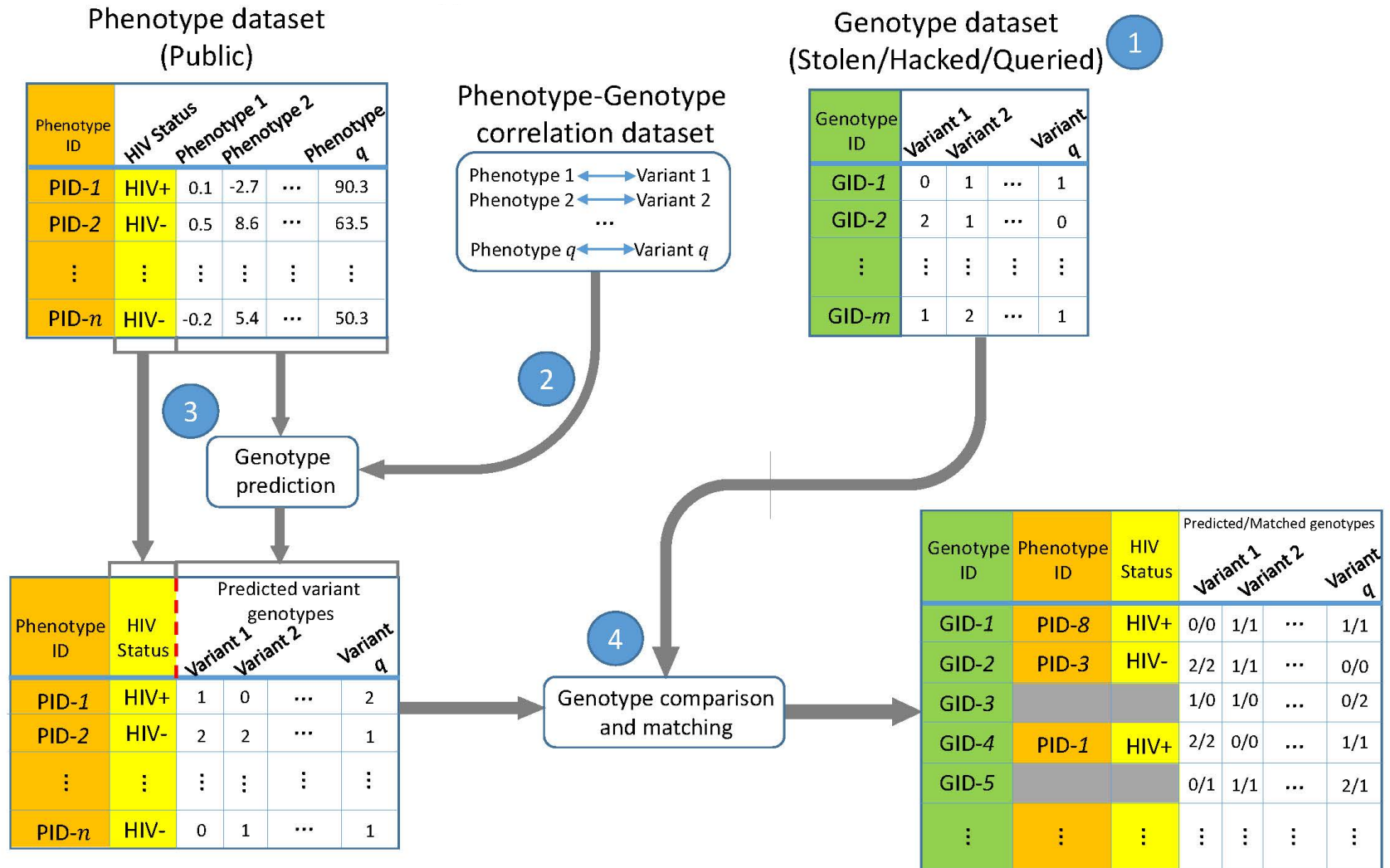
- Condition specific entropy
- Higher cond. entropy: Lower predictability
- Lower cond. entropy: Higher predictability
- Additive for multiple eQTLs



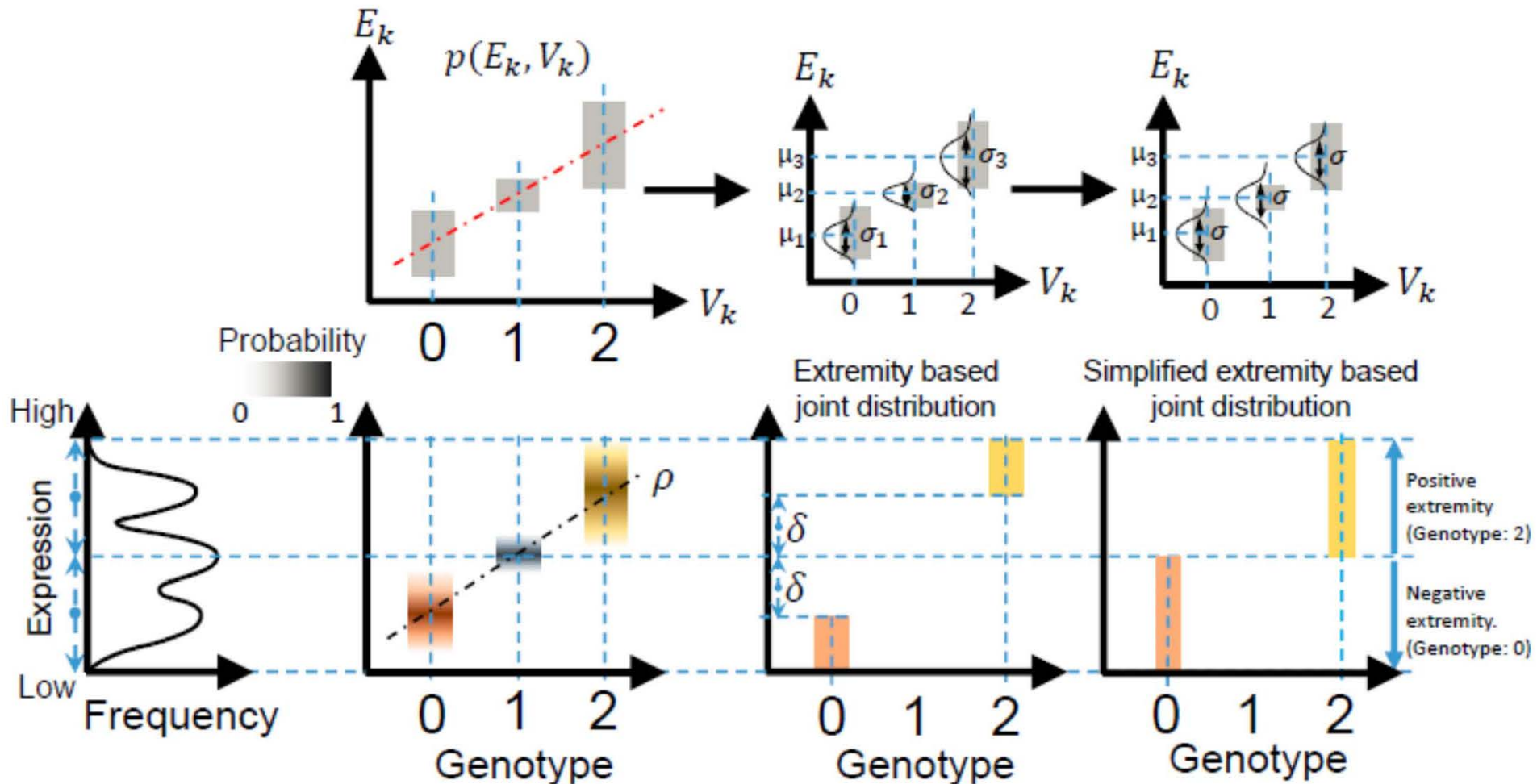


# ICI Leakage versus Genotype Predictability

# Linking Attack Scenario

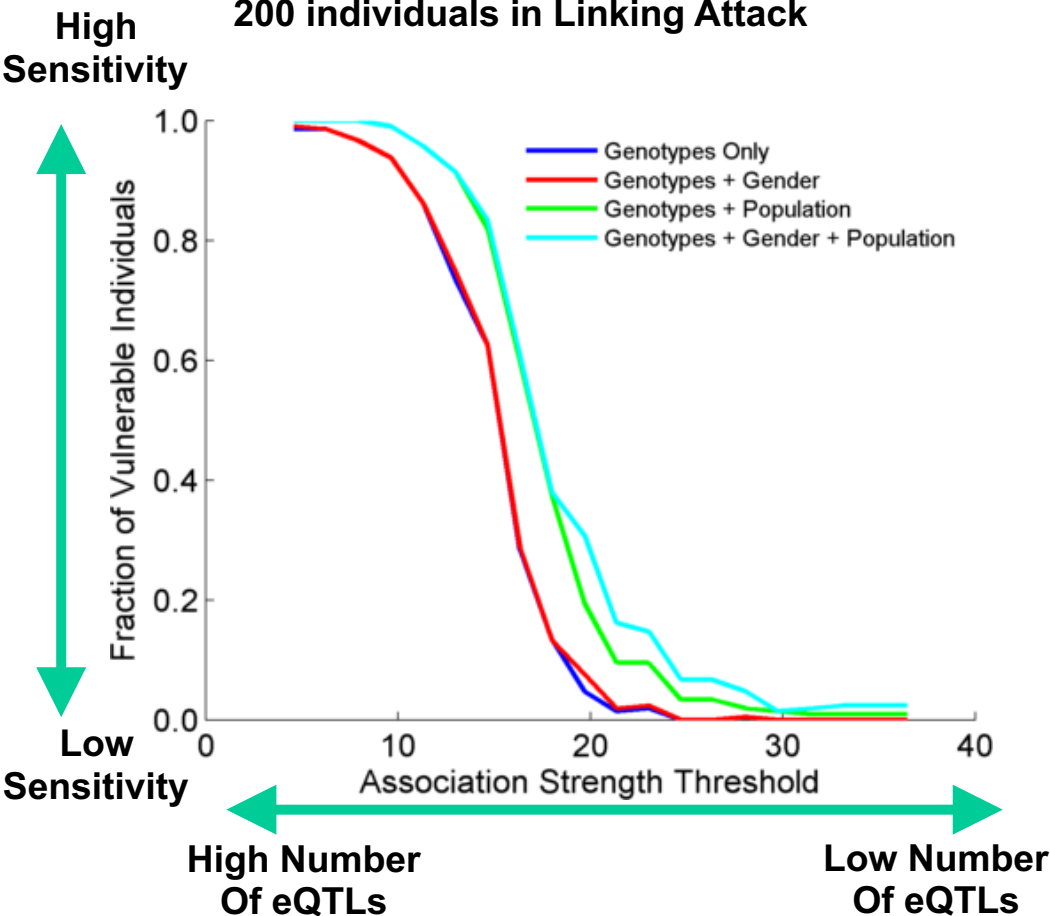


# Levels of Expression-Genotype Model Simplifications for Genotype Prediction



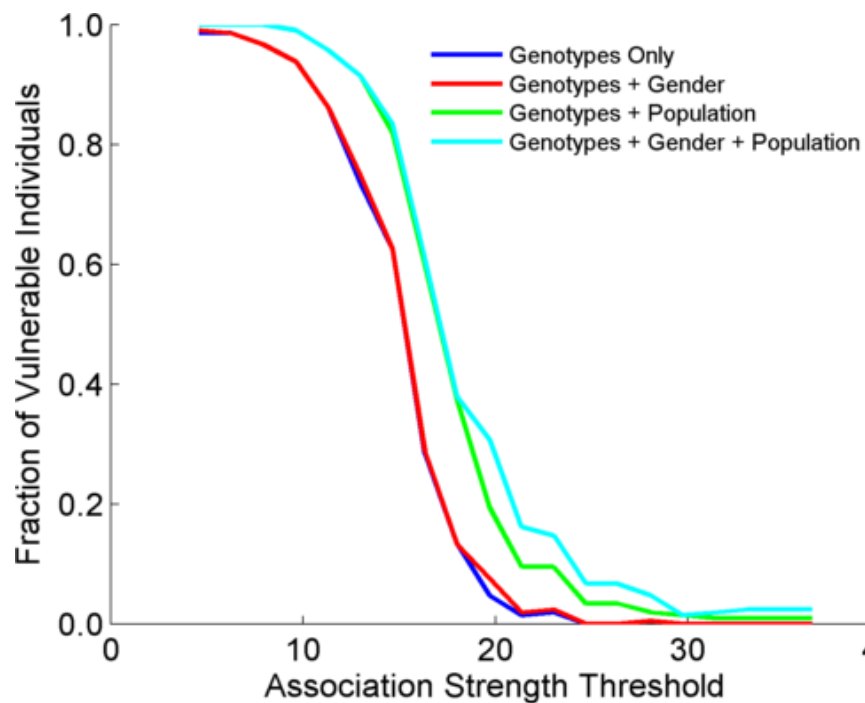
# Success in Linking Attack with Extremity based Genotype Prediction

200 individuals eQTL Discovery  
200 individuals in Linking Attack

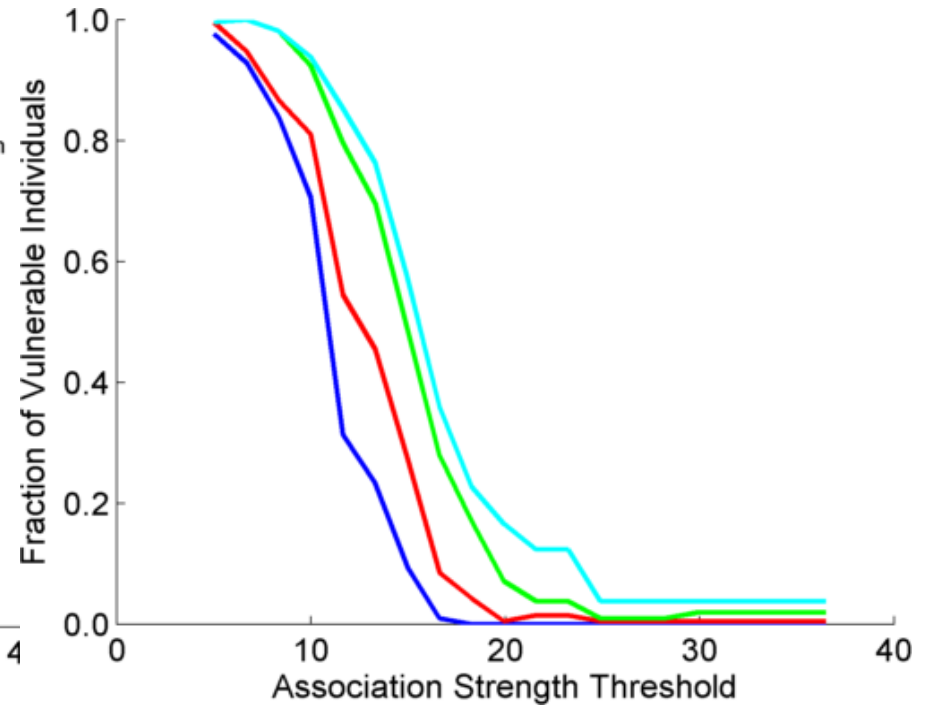


# Success in Linking Attack with Extremity based Genotype Prediction

200 individuals eQTL Discovery  
200 individuals in Linking Attack



200 individuals eQTL Discovery  
100,200 individuals in Linking Attack





# Genomic Privacy & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information

- **Intro on RNA-seq & the General Dilemma of Genomic Privacy**

- RNA-seq Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- The need to quantify leaks

- **Quantifying RNA-seq Leakage ...from Reads**

- Almost as much as WGS
- But can remove SNVs in reads w/ MRF

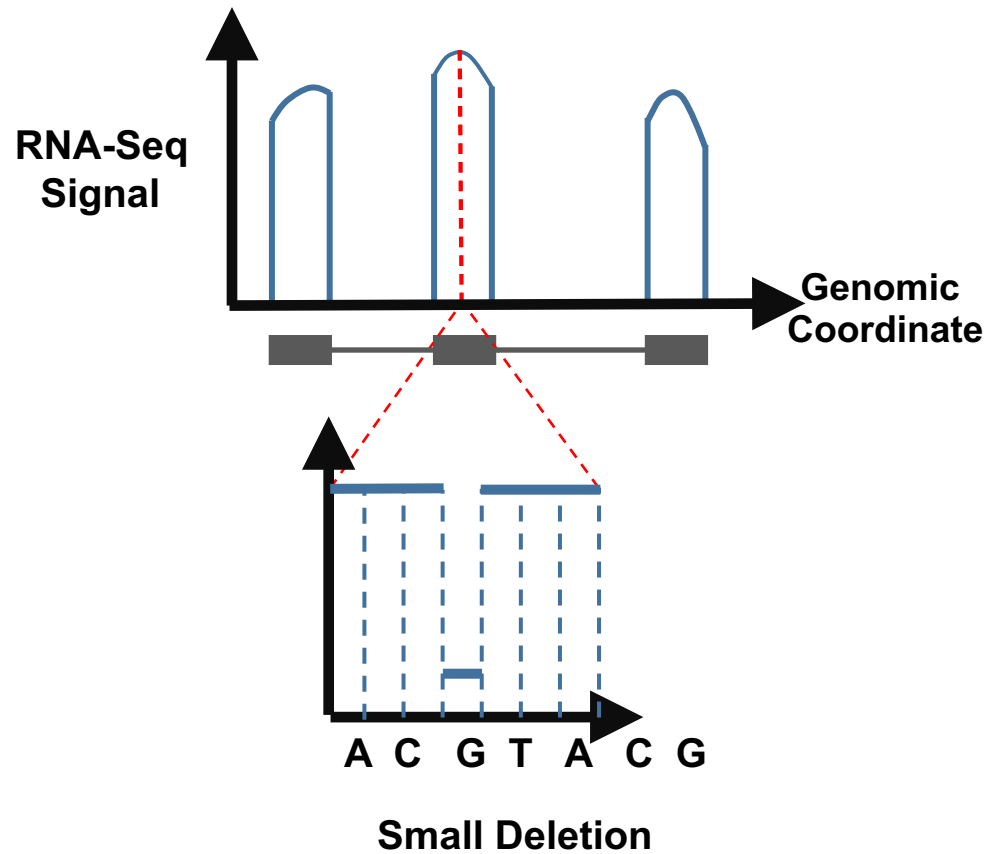
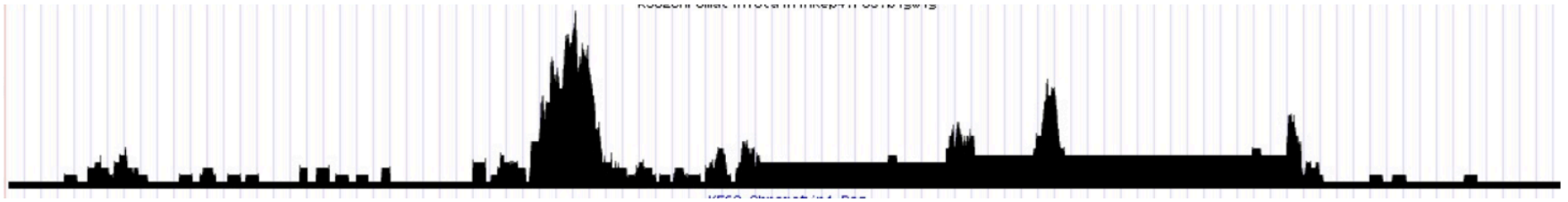
- **...from eQTLs**

- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack

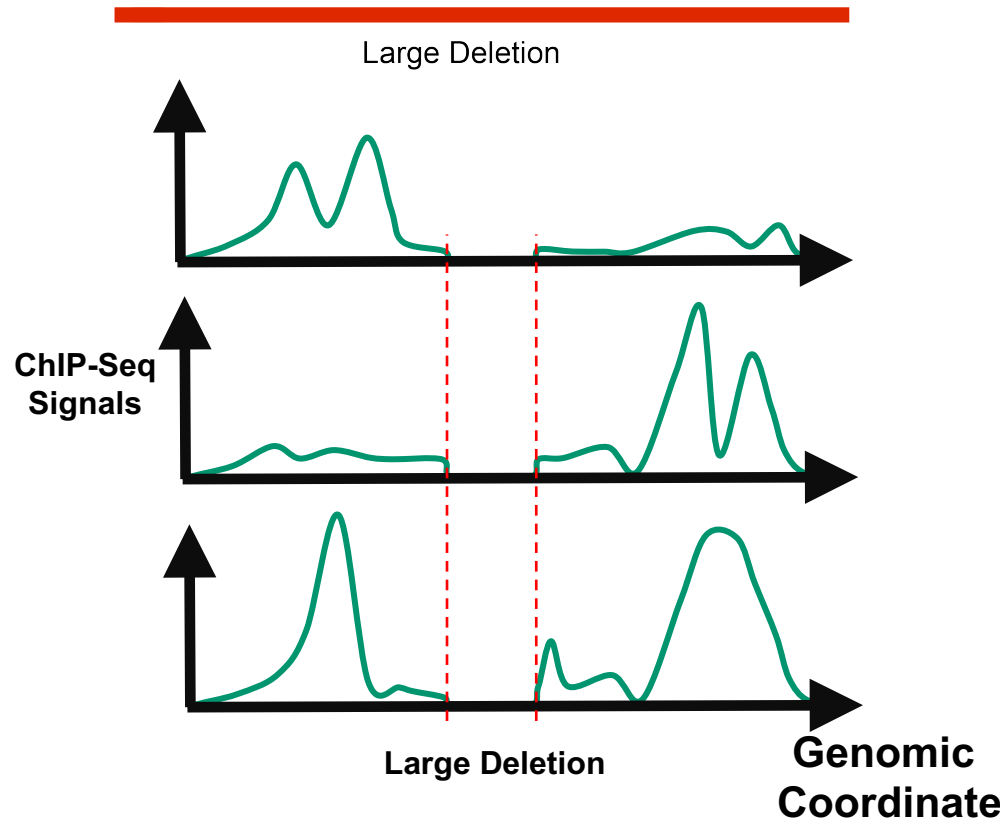
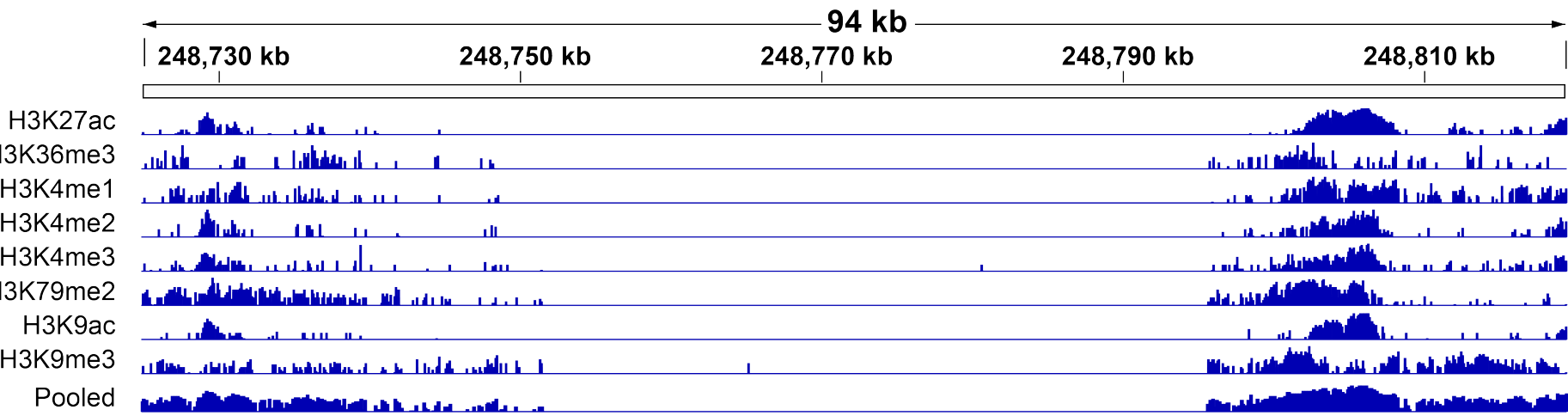
- **...from Indels/SVs**

- Another source of leakage in RNA-seq data & how to use these for a related linking attack

# Small Deletions and RNA-seq Signal

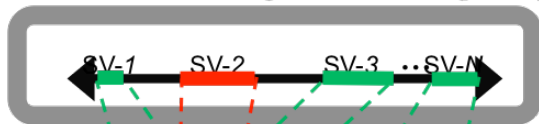


# Structural Variants are Generally Detectable from Functional Genomics Data

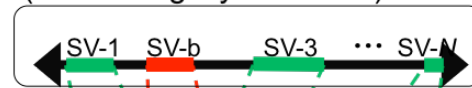


[Harmanci et al. (submitted)]

### SV Panel for Signal Profiles ( $p \downarrow S$ )



### Structural Variants Panel (Stolen/Legally Obtained) ( $p \downarrow G$ )



Anonymized Sample ID	SV-1	SV-2	SV-3	...	SV-N	HIV Status
SIND-1	0	0	2	...	2	+
SIND-2	2	0	X	...	0	-
⋮	⋮	⋮	⋮	⋮	⋮	⋮
SIND-n	0	X	X	...	0	+

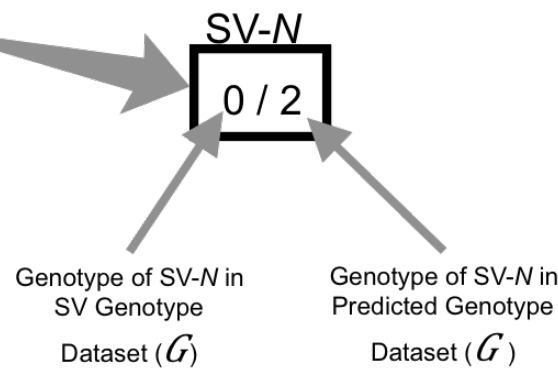
Predicted SV Genotype Dataset ( $G$ )

Patient Name	SV-1	SV-b	SV-3	...	SV-N
GIND-1	0	1	2	...	0
GIND-2	2	0	2	...	1
GIND-3	0	1	1	...	0
⋮	⋮	⋮	⋮	⋮	⋮
GIND-K	1	2	2	...	2

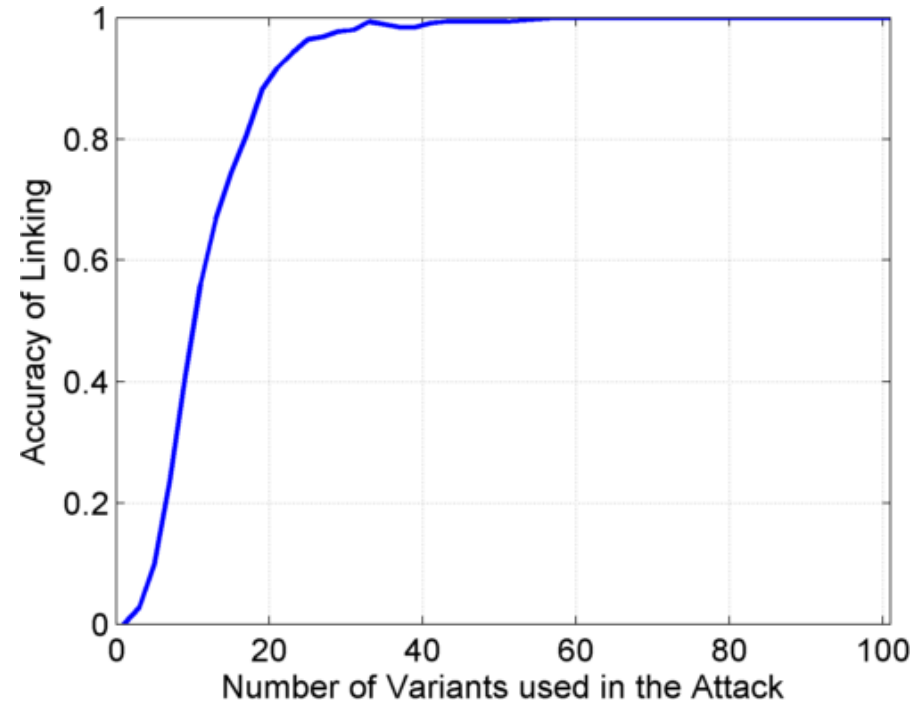
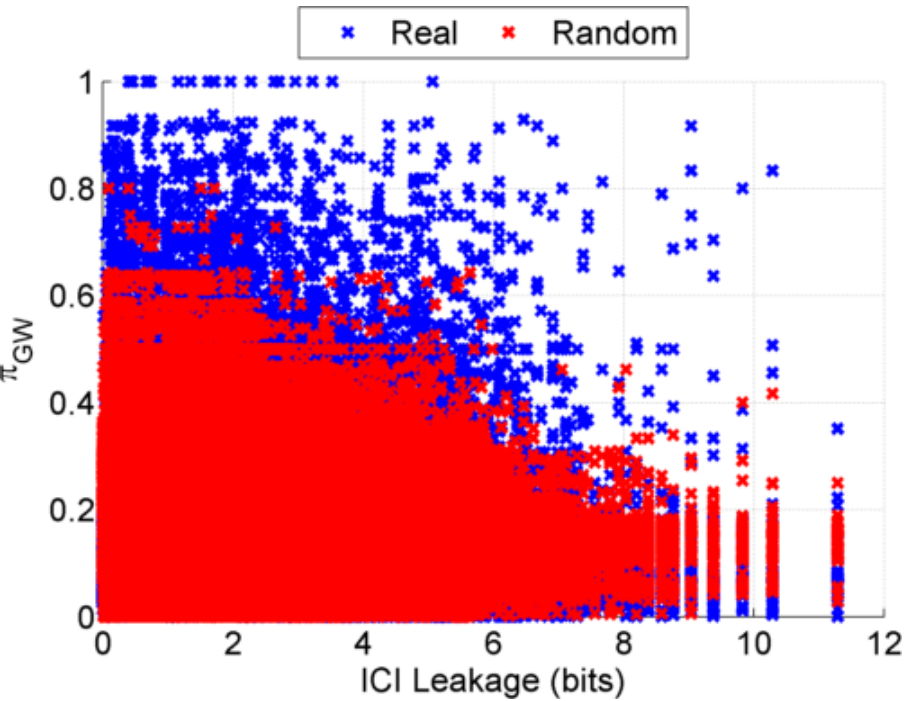
Structural Variants Genotype Dataset (Stolen/Legally Obtained) ( $G$ )

Comparison of SV Panels and Genotype Matching

Anonymized Sample ID	Patient Name	HIV Status	Genotype in $G$ / Genotype in $\hat{G}$			
			SV-1	SV-3	...	SV-N
SIND-1	GIND-2	+	0/0	1/0	...	0/2
SIND-2	GIND-1	-	0/2	1/0	...	0/0
⋮	⋮	⋮	⋮	⋮	⋮	⋮
SIND-n	GIND-3	+	0/0	1/X	...	0/0



# Predictability vs. Information Leakage & Accuracy of Linking





# Genomic Privacy & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information

- **Intro on RNA-seq & the General Dilemma of Genomic Privacy**

- RNA-seq Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
- The need to quantify leaks

- **Quantifying RNA-seq Leakage ...from Reads**

- Almost as much as WGS
- But can remove SNVs in reads w/ MRF

- **...from eQTLs**

- Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
- Instantiating a practical linking attack

- **...from Indels/SVs**

- Another source of leakage in RNA-seq data & how to use these for a related linking attack

# Genomic Privacy & the 2-sided Nature of RNA-seq: Quantifying the Leakage of Individual Information

- **Intro on RNA-seq & the General Dilemma of Genomic Privacy**
  - RNA-seq Presents a tricky privacy issue since much of the sequencing is for general, non-individual specific results yet it's tagged with individual information
  - The need to quantify leaks
- **Quantifying RNA-seq Leakage ...from Reads**
  - Almost as much as WGS
  - But can remove SNVs in reads w/ MRF
- **...from eQTLs**
  - Quantifying & removing variant info from expression levels + eQTLs using ICI & predictability
  - Instantiating a practical linking attack
- **...from Indels/SVs**
  - Another source of leakage in RNA-seq data & how to use these for a related linking attack

## Acknowledgements

**NIH BD2K  
Program**



**PrivaSeq**.gersteinlab.org –  
A **Harmanci**, G **Gürsoy**,  
F Navarro, S Wagner, X Kong

**Hiring Postdocs. See  
JOBS.gersteinlab.org !**

**Extra**



# Info about content in this slide pack

- General PERMISSIONS
  - This Presentation is copyright Mark Gerstein, Yale University, 2017.
  - Please read permissions statement at [www.gersteinlab.org/misc/permissions.html](http://www.gersteinlab.org/misc/permissions.html) .
  - Feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).
  - Paper references in the talk were mostly from Papers.GersteinLab.org.
- PHOTOS & IMAGES. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see <http://streams.gerstein.info> .
  - In particular, many of the images have particular EXIF tags, such as kwpotppt , that can be easily queried from flickr, viz: <http://www.flickr.com/photos/mbgmbg/tags/kwpotppt>