# Quantifying Patient Privacy in Public Disease Portals

Sean Simmons

Stanley Center, Broad Institute

1/10/2017

# Acknowledgments

**Portal Team:**
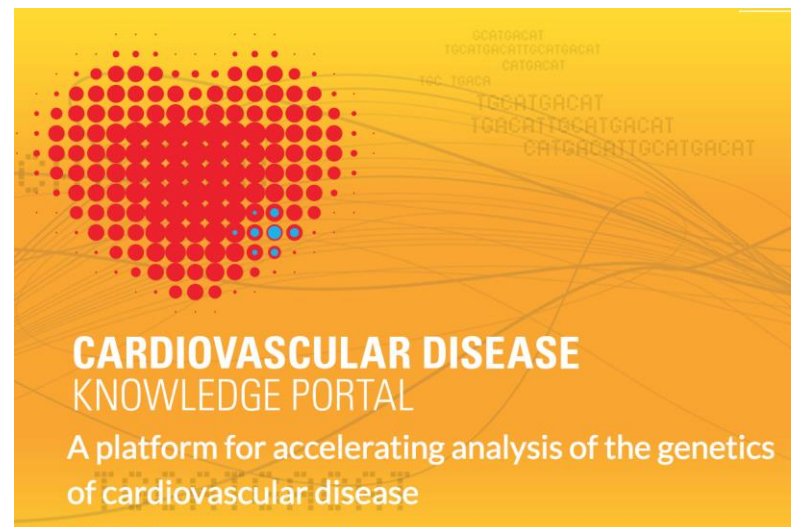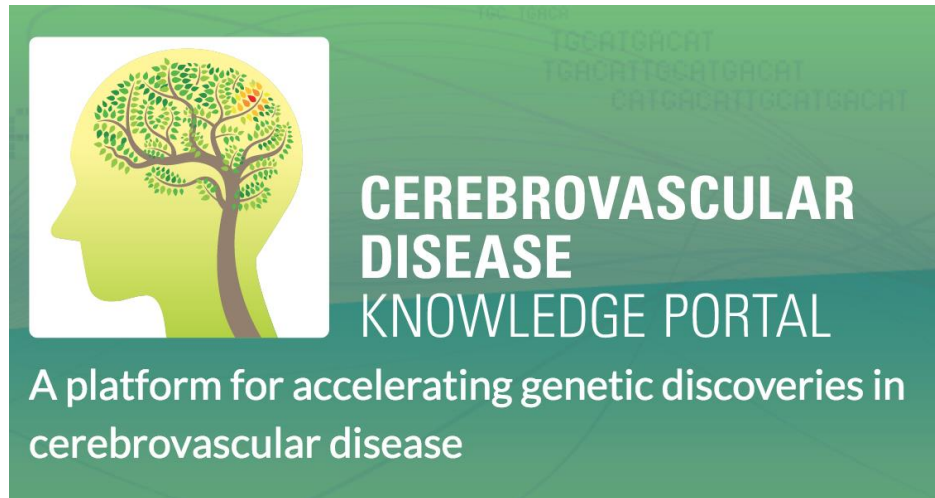- Jason Flannick
- Noel Burtt
- The rest of the team!

**Sahinalp Lab:**
- Cenk Sahinalp (see his talk later!)
- Many others!

**Berger Lab:**
- Bonnie Berger
- Many others!

# Knowledge Portals: A new way to share GWAS!

# Knowledge Portals: Sharing GWAS results

## rs13266634 associations at a glance

All type 2 diabetes associations are shown below. Click the "expand associations for all traits" button to see additional associations of at least nominal significance. Associations with related phenotypes are grouped under color-coded headers. Warning: the data sets shown share samples and thus *cannot* be combined via meta-analysis. For assistance with pooling studies for analysis, please contact us.

Color key: p < 5e-8 ❔ | p < 5e-4 ❔ | p < 0.05 ❔    Direction of effect: up ⬆ down ⬇ ❔    Dataset: sample size ❔ | frequency ❔ | count ❔

### Type 2 diabetes

| DIAGRAM Transethnic meta-analysis | GoT2D exome chip analysis | 70KforT2D GWAS | GoT2D WGS + replication | 19K exome sequence analysis | AGEN GWAS |
|---|---|---|---|---|---|
| $p = 2.7e-20$ genome-wide significant | $p = 2.73e-18$ genome-wide significant | $p = 3.01e-15$ genome-wide significant | $p = 0.00000856$ locus-wide significant | $p = 1.95e-7$ locus-wide significant | $p = 0.61$ |
| ⬇ OR = 0.877 | ⬇ OR = 0.878 | ⬇ OR = 0.871 | ⬇ OR = 0.885 | ⬇ OR = 0.881 | ⬇ OR = 0.902 |
| 110452 \| n/a \| n/a | 75670 \| 33.0% \| 49942 | 70127 \| 30.5% \| 42735 | 44414 \| 33.0% \| 29287 | 18844 \| 28.3% \| 10667 | 18817 \| n/a \| n/a |

| BioMe AMP T2D GWAS | GWAS SIGMA | METSIM GWAS | SIGMA exome chip analysis | EXTEND GWAS | CAMP GWAS |
|---|---|---|---|---|---|
| $p = 0.181$ | $p = 0.00233$ nominally significant | $p = 0.151$ | $p = 0.000147$ locus-wide significant | $p = 0.02$ nominally significant | $p = 0.0632$ |
| ⬇ OR = 0.941 | ⬇ OR = 0.895 | ⬇ OR = 0.936 | ⬇ OR = 0.874 | | ⬇ OR = 0.983 |
| 9173 \| 19.2% \| 3526 | 8891 \| 25.7% \| 4577 | 8791 \| 39.5% \| 6944 | 8214 \| 25.8% \| 4237 | 7159 \| n/a \| n/a | 3628 \| 27.8% \| 2016 |

# Knowledge Portals: Interactive analysis!

## Genetic Association Interactive Tool

The Genetic Association Interactive Tool allows you to compute custom association statistics by specifying the phenotype to test for association, a subset of samples to analyze based on specific phenotypic criteria, and a set of covariates to control for in the analysis. In order to protect patient privacy, GAIT will only allow visualization or analysis of data from more than 100 individuals.

GAIT guide

## Choose a phenotype and partitioning strategy

**Dataset**

CAMP GWAS

**Phenotype**

Type 2 diabetes

**Samples:** ☐ Filter cases and controls separately

Step 1: Select a subset of samples based on phenotypic criteria

Step 2: Control for covariates

**Launch analysis**

# Step 1: Choose who to include

# Step 2: Choose what to correct for

## Step 2: Control for covariates

Select principal components and/or phenotypes to be used as covariates in your association analysis. Principal components 1-4 are selected by default to minimize the influence of ancestry, though additional principal components may be selected to control for finer grained substructure within a population. Selecting phenotypes as covariates allows you to estimate the effects of the response phenotype independently of the covariate phenotypes.
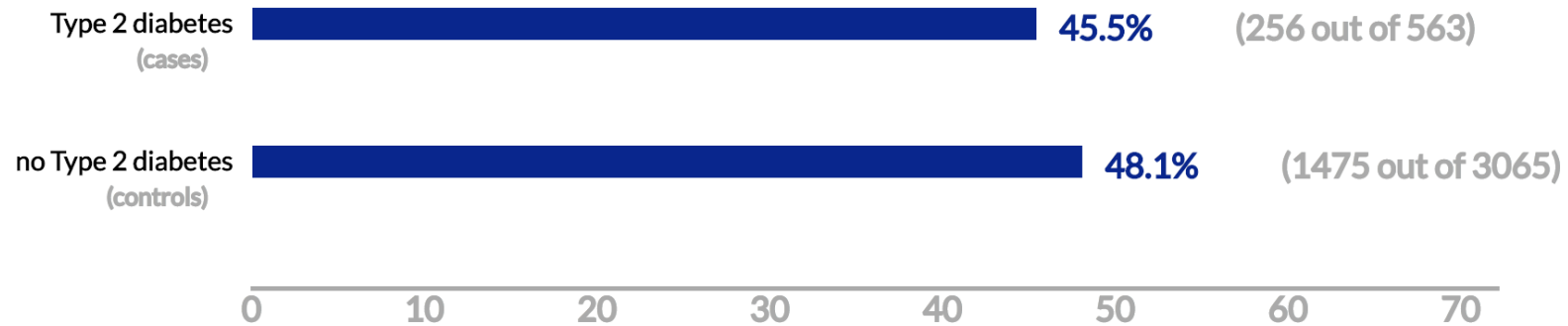
- ☑ PC-1
- ☑ PC-2
- ☑ PC-3
- ☑ PC-4
- ☐ PC-5
- ☐ PC-6
- ☐ PC-7
- ☐ PC-8
- ☐ PC-9
- ☐ PC-10

- ☑ Age
- ☑ Sex
- ☐ BMI

# Step 3: Get results!

Launch analysis

pValue = 0.0778
odds ratio = 0.873
95% CI: (0.751 to 1.02)

**Type 2 diabetes**
(cases) — **45.5%** (256 out of 563)

**no Type 2 diabetes**
(controls) — **48.1%** (1475 out of 3065)

0    10    20    30    40    50    60    70

# We know a lot about GWAS privacy…

# Difference with standard GWAS portals

- Smaller datasets analysis/ subsampling

- Repeated queries per SNP

- Ability to remove/ include confounders

# Difference with standard GWAS portals

- Smaller datasets analysis/ subsampling

- Repeated queries per SNP

- Ability to remove/ include confounders

All of these change the privacy landscape!

# Quantifying privacy

- Use a model based approach

- Looked at different types of private data leakage

- No formal guarantees (yet!)

# Quantifying privacy: Our models

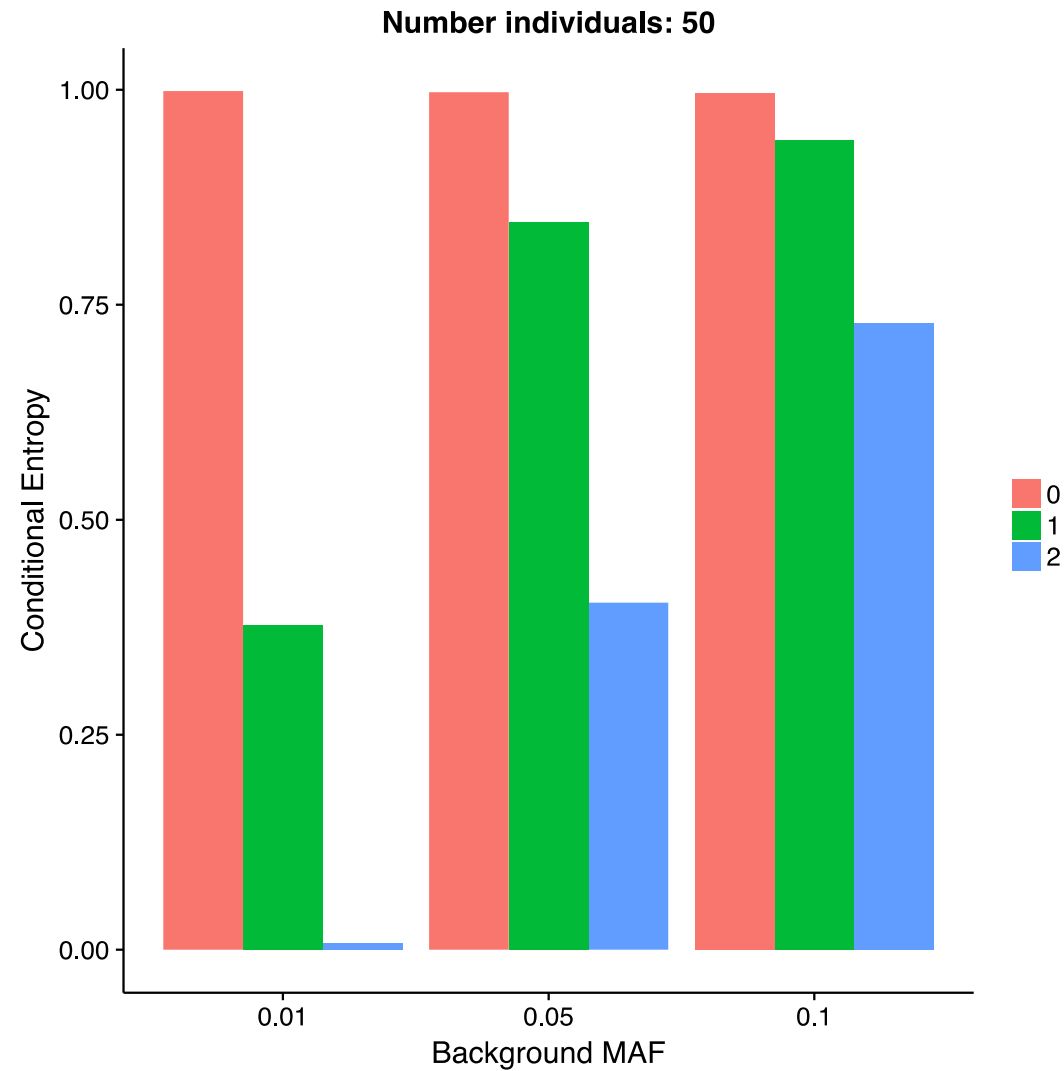| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Quantifying privacy: Adversarial uncertainty

- Adversaries knowledge modelled probabilistically

- Prior knowledge specifies probability distribution used.

- Look at entropy (among other measures) of posterior probability to determine privacy loss.

# Quantifying privacy: Conditional entropy!

$$Entropy_{y|X} = -\sum \log(\Pr(y|X))\Pr(y,X)$$

Where:

- y value of secret trait

- Pr probaility distribution representing adversaries belief

- X information being released

# Quantifying privacy: Example!

$$Entropy_{genotype|MAF} =$$

$$-\sum_{MAF,i} \log(\Pr(genotype = i|MAF))\Pr(genotype = i, MAF)$$

# Quantifying privacy: Our models

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Quantifying privacy: First Analysis

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Leaking disease participation

- Apply to private disease status information

- Assume genotype known and MAF of background population known (similar to Homer et al.)

- Privacy leakage from minor allele frequencies!

# Decrease in entropy of disease status

# Quantifying privacy: First Analysis

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Quantifying privacy: Leaking genotype

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# The concern:



Launch analysis

pValue = 0.0778
odds ratio = 0.873
95% CI: (0.751 to 1.02)

Type 2 diabetes (cases) — 45.5% (256 out of 563)

no Type 2 diabetes (controls) — 48.1% (1475 out of 3065)

0  10  20  30  40  50  60  70

ATGCAT
ATGGAT

# Entropy of the genotype

# Added concerns with interactive GWAS

- So far have focused on privacy concerns without taking advantage of the interactive aspect.

- Does interactive analysis add more?

# Added concerns with interactive GWAS

- So far have focused on privacy concerns without taking advantage of the interactive aspect.

- Does interactive analysis add more?

**Short answer? Yes**

# Quantifying privacy: Leaking genotype

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Quantifying privacy: Leaking health data

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Interactive histograms and phenotype data

# Interactive histograms and phenotype data

- Interactive histograms allow users to understand how certain traits (age, smoking status, etc) differ in dataset

- Also allow for reconstruction of phenotype information!

# Interactive histograms and phenotype data
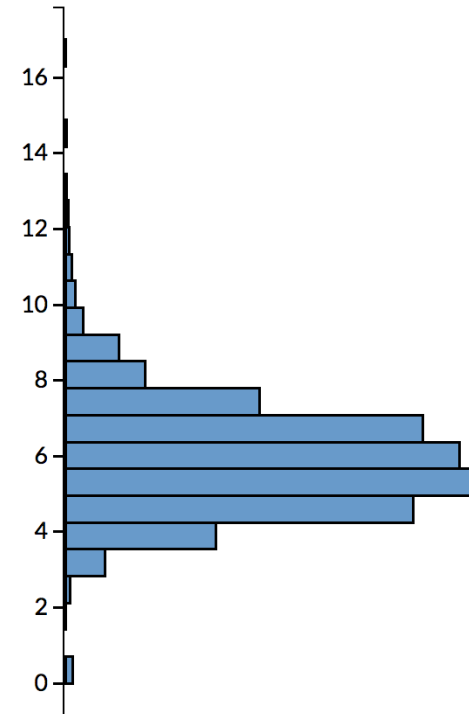
# Interactive histograms and phenotype data

Cholesterol

Diastolic blood pressure

Systolic blood pressure

16 —

**Unclear if can link to public data**

Weight (kg)

Hip circ. (cm)

2 —

Waist circ. (cm)

0 —

Hypertension          All Selected (67) ▾

Smoker                All Selected (44) ▾

Menopause             All Selected (9) ▾

Hormone treatment     All Selected (3) ▾

Hormone treatment menopause   All Selected (3) ▾

# Quantifying privacy: Leaking health data

| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Quantifying privacy: Leaking genotype data

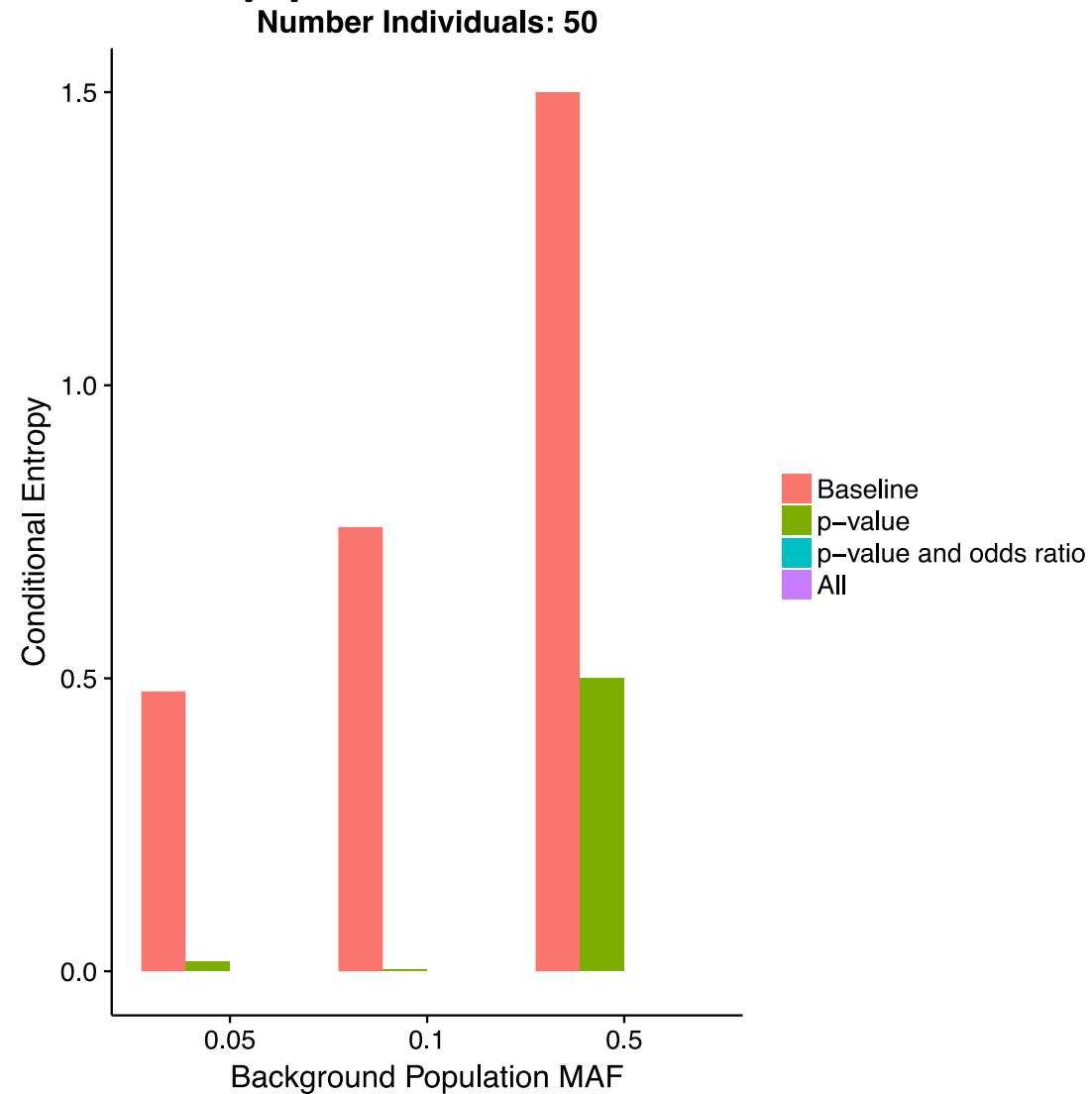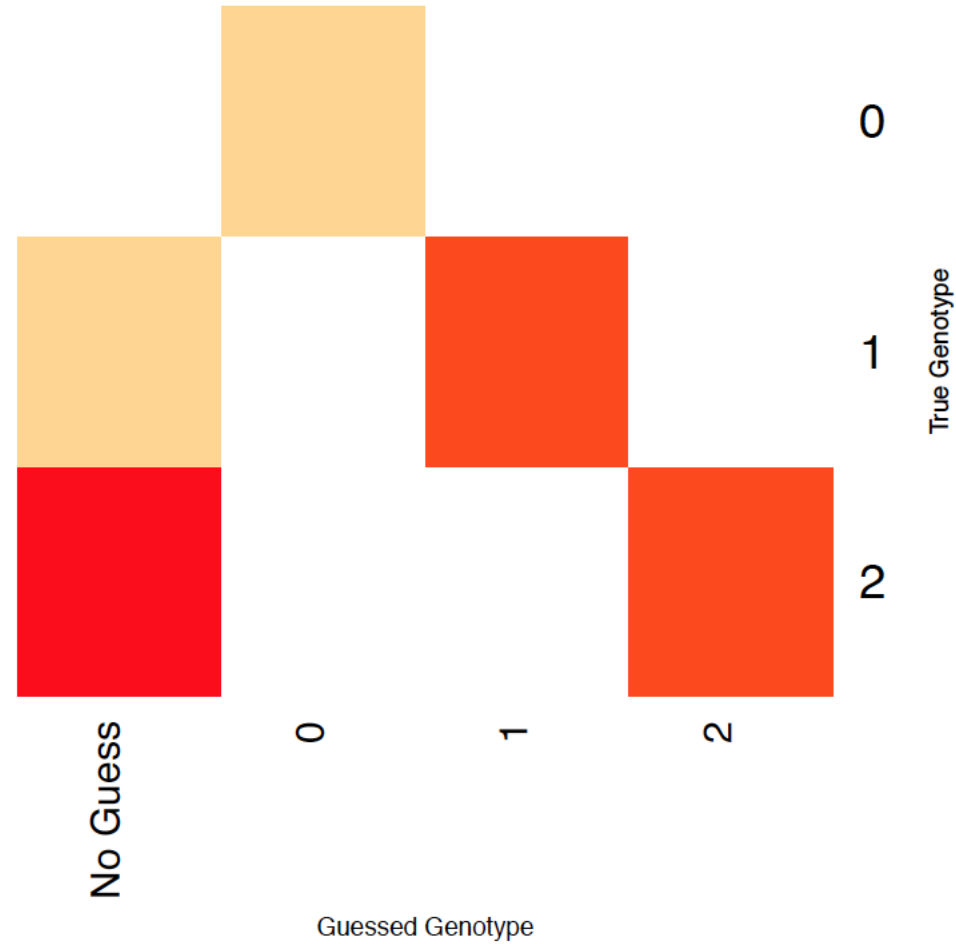| What Leaks? | Outside knowledge | Statistic released | Likely Risk (very subjective) |
|---|---|---|---|
| Participation in study | Background MAF and individuals genotype | Minor allele frequency | Low |
| Diseases status | Background MAF and individuals genotype | Minor allele frequency | Low |
| Other phenotype data | Unspecified | Interactive histogram | Moderate |
| Genotype | Background MAF | GWAS statistics for one query | Low |
| Genotype | Background MAF | GWAS statistics for numerous queries | High |

# Interactive analysis and privacy

- Run an analysis on one group of individuals

- Run an analysis on a slightly different group of individuals

- Combine to learn about genotype of one target!

# Entropy in genotype data: Interactive analysis

# Recovering genotype data

# Mitigation?

- Larger minimum analysis size, less precision

- Remove some features

- Require PCs as covariates

- K-anonymity for histograms

- Remove Y chromosome

# What else is needed?

- Doesn't consider all possible queries, only a few select ones.

- What are we missing? No good way to measure!!

- How do covariates actually effect the privacy risk? Harder to answer.

- Can the information leaked here be linked to outside databases?