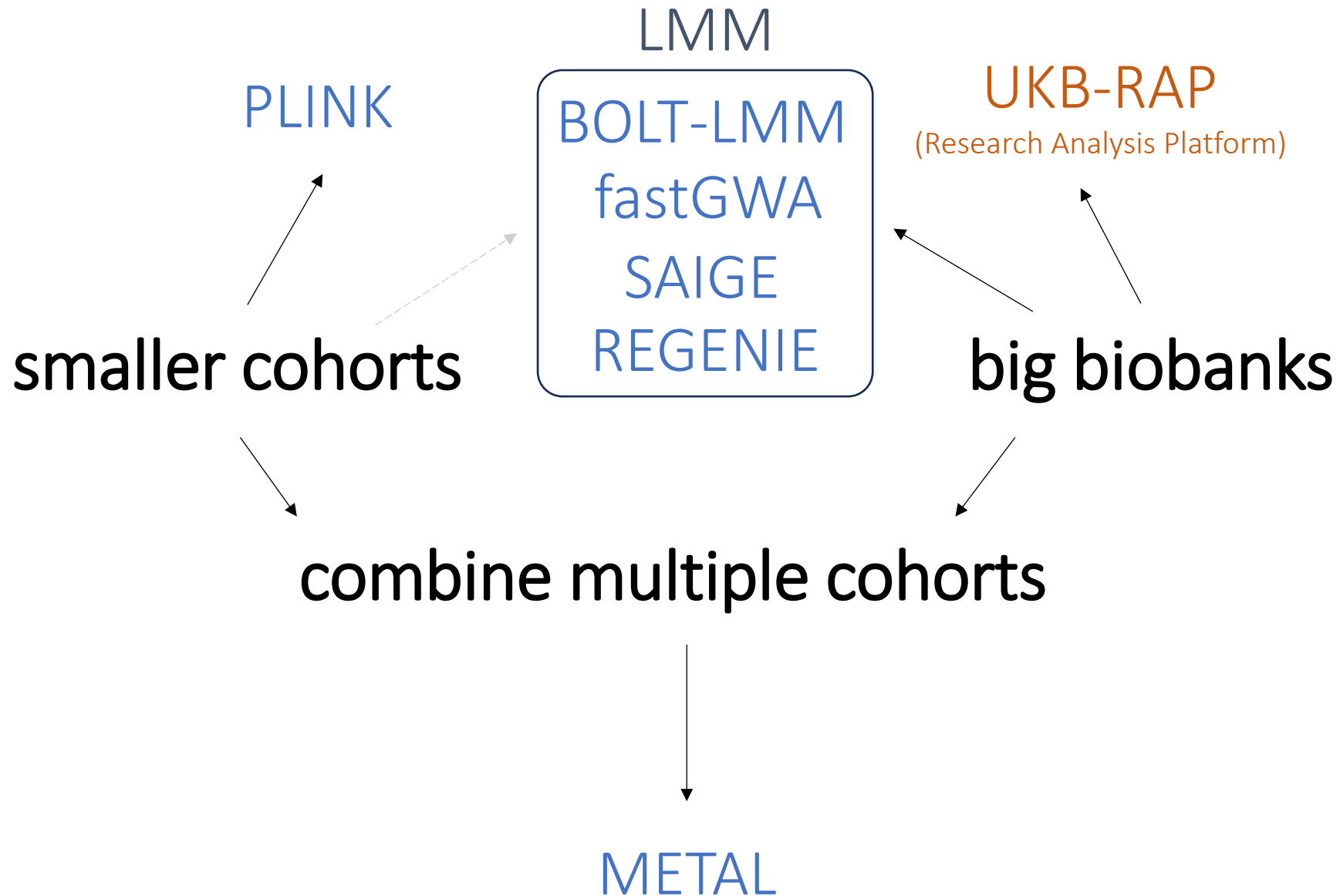


# Genome-Wide Association Studies

Part 3: analyzing smaller cohorts, big biobanks, and cohorts combined

International Statistical Genetics Workshop – 2026



# smaller cohorts (PLINK)

The American Journal of Human Genetics Volume 81 September 2007

**REPORT**

---

## PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage Analyses

Shaun Purcell, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel A. R. Ferreira,  
David Bender, Julian Maller, Pamela Sklar, Paul I. W. de Bakker, Mark J. Daly, and Pak C. Sham

# PLINK – file formats

## ped files

\*.ped

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...

\*.map

chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

# PLINK – file formats

## ped files

\*.ped

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...

Family ID  
Individual ID  
Paternal ID  
Maternal ID  
sex  
phenotype

\*.map

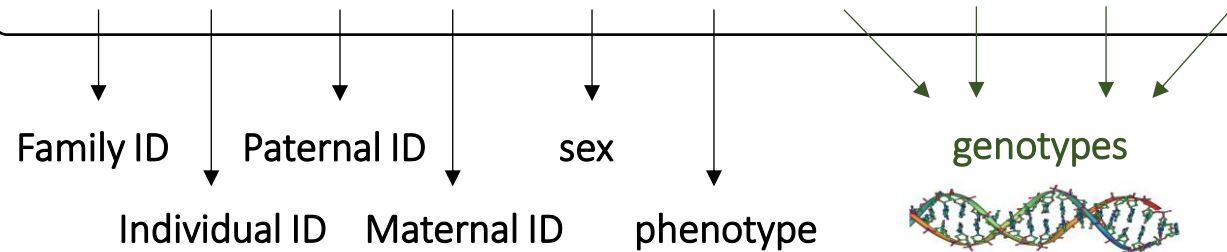
chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

# PLINK – file formats

## ped files

\*.ped

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...



\*.map

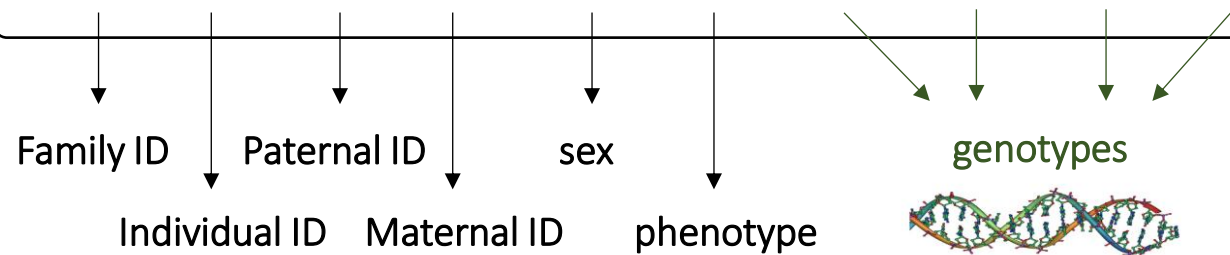
chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

# PLINK – file formats

## ped files

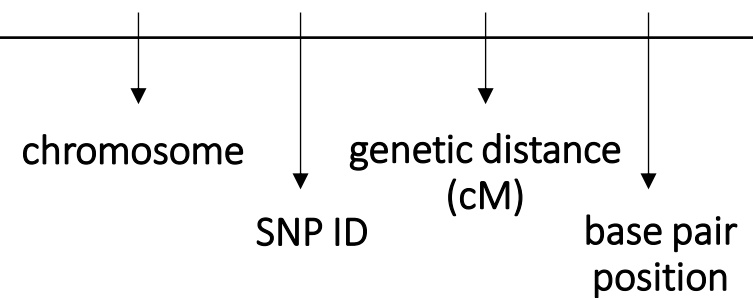
\*.ped

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...



\*.map

chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000



# PLINK – file formats

## ped files

### \*.ped

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...

### \*.map

chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

## binary files

### \*.fam

FID	IID	PID	MID	sex	p
1	1	3	4	2	165
1	2	3	4	2	173
2	5	7	8	1	179

### \*.bed

Compressed version of the .ped file, but only with information on genotypes. Not in a format readable for humans

### \*.bim

chr	SNP	GD	BP	A1	A2
1	rs1	0	870000	C	A
1	rs2	0	880000	C	G
1	rs3	0	890000	A	G

# PLINK – file formats

## ped files

### \*.ped

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...

### \*.map

chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

## binary files

### \*.fam

FID	IID	PID	MID	sex	p
1	1	3	4	2	165
1	2	3	4	2	173
2	5	7	8	1	179

### \*.bed

Compressed version of the .ped file, but only with information on genotypes. Not in a format readable for humans

### \*.bim

chr	SNP	GD	BP	A1	A2
1	rs1	0	870000	C	A
1	rs2	0	880000	C	G
1	rs3	0	890000	A	G

# PLINK – file formats

ped files

as input: `--file`

as output: `--recode`

**\*.ped**

FID	IID	PID	MID	sex	p	rs1	rs2	rs3	...
1	1	3	4	2	165	CC	CG	AA	...
1	2	3	4	2	173	CA	GG	AG	...
2	5	7	8	1	179	CA	CG	GG	...

**\*.map**

chr	SNP	GD	BP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

binary files

as input: `--bfile`

as output: `--make-bed`

**\*.fam**

FID	IID	PID	MID	sex	p
1	1	3	4	2	165
1	2	3	4	2	173
2	5	7	8	1	179

**\*.bed**

Compressed version of the .ped file, but only with information on genotypes. Not in a format readable for humans

**\*.bim**

chr	SNP	GD	BP	A1	A2
1	rs1	0	870000	C	A
1	rs2	0	880000	C	G
1	rs3	0	890000	A	G

Many other packages also read binary files  
(all LMM approaches discussed later in this video)

## binary files

### \*.fam

FID	IID	PID	MID	sex	p
1	1	3	4	2	165
1	2	3	4	2	173
2	5	7	8	1	179

### \*.bed

Compressed version of the  
.ped file, but only with  
information on genotypes.  
Not in a format readable  
for humans

### \*.bim

chr	SNP	GD	BP	A1	A2
1	rs1	0	870000	C	A
1	rs2	0	880000	C	G
1	rs3	0	890000	A	G

# PLINK – regression analysis

```
plink --bfile cohort_name --linear --pheno phenotypes.txt --covar covs.txt --out gwas
```



binary files

**\*.fam**

FID	IID	PID	MID	sex	p
1	1	3	4	2	165
1	2	3	4	2	173
2	5	7	8	1	179

**\*.bed**

Compressed version of the .ped file, but only with information on genotypes. Not in a format readable for humans

**\*.bim**

chr	SNP	GD	BP	A1	A2
1	rs1	0	870000	C	A
1	rs2	0	880000	C	G
1	rs3	0	890000	A	G

# PLINK – regression analysis

```
plink --bfile cohort_name --linear --pheno phenotypes.txt --covar covs.txt --out gwas
```

FID	IID	p
1	1	165
1	2	173
2	5	179

↓                      ↓                      ↓  
Family ID            Individual ID            phenotype

FID	IID	cov1	cov2	cov3	cov4	cov5
1	1	2	31	0.01	0.01	0.03
1	2	2	27	0.03	0.01	0.02
2	5	1	46	0.02	0.02	0.02

↓                      ↓                      ↓                      ↓                      ↓  
Family ID            Individual ID            covariates

# PLINK – regression analysis

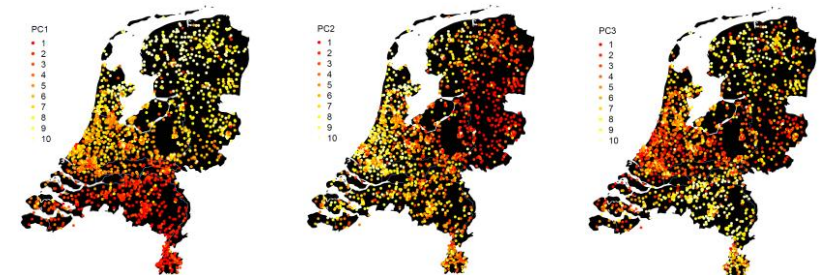
```
plink --bfile cohort_name --linear --pheno phenotypes.txt --covar covs.txt --out gwas
```

FID	IID	p
1	1	165
1	2	173
2	5	179

Family ID  
Individual ID  
phenotype

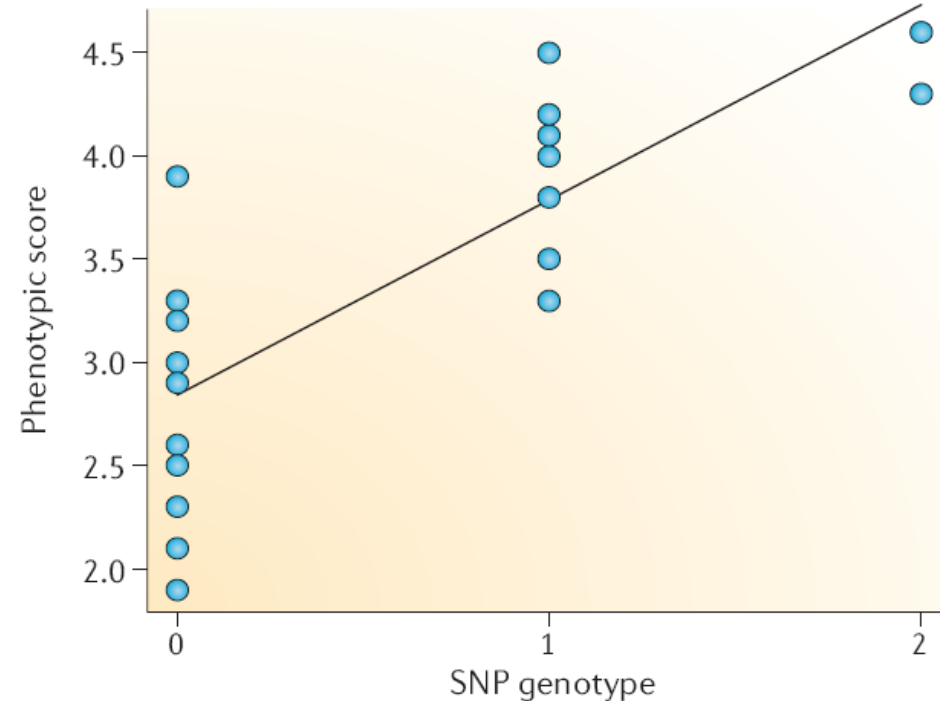
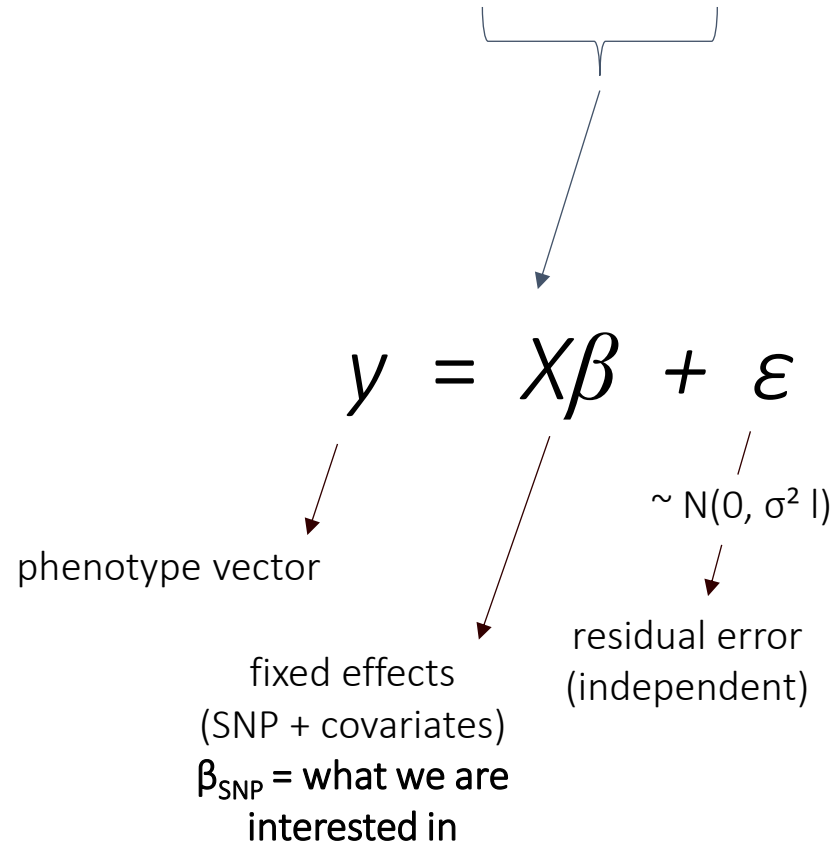
FID	IID	cov1	cov2	cov3	cov4	cov5
1	1	2	31	0.01	0.01	0.03
1	2	2	27	0.03	0.01	0.02
2	5	1	46	0.02	0.02	0.02

Family ID  
Individual ID  
sex  
age  
PC1  
PC2  
PC3



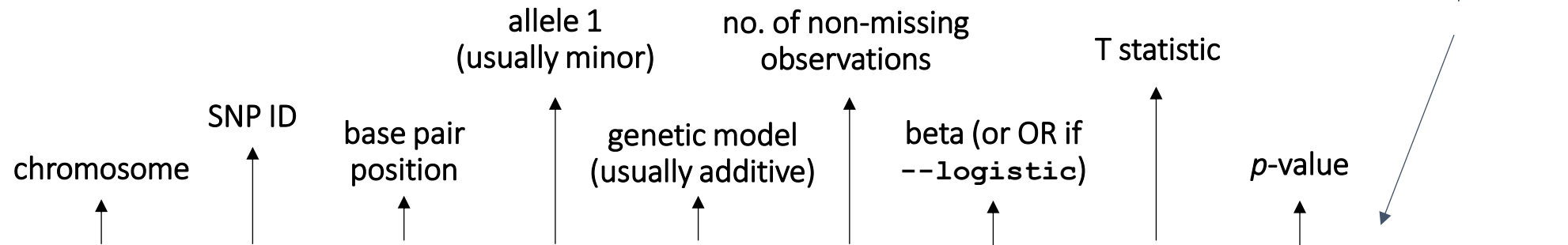
# PLINK – regression analysis

```
plink --bfile cohort_name --linear --pheno phenotypes.txt --covar covs.txt --out gwas
```



# PLINK – regression analysis

```
plink --bfile cohort_name --linear --pheno phenotypes.txt --covar covs.txt --out gwas
```



CHR	SNP	BP	A1	TEST	NMISS	BETA	STAT	P
1	rs1	870000	C	ADD	4987	0.0312	1.847	0.0648
1	rs2	880000	C	ADD	4991	0.1843	6.214	3.2e-09
1	rs3	890000	A	ADD	4983	-0.0071	-0.423	0.6724
... millions more rows ...								

# PLINK – regression analysis

```
plink --bfile cohort_name --linear --pheno phenotypes.txt --covar covs.txt --out gwas
```

allele 1 (usually minor)

no. of non-missing observations

T statistic

SNP ID

base pair position

genetic model (usually additive)

beta (or OR if --logistic)

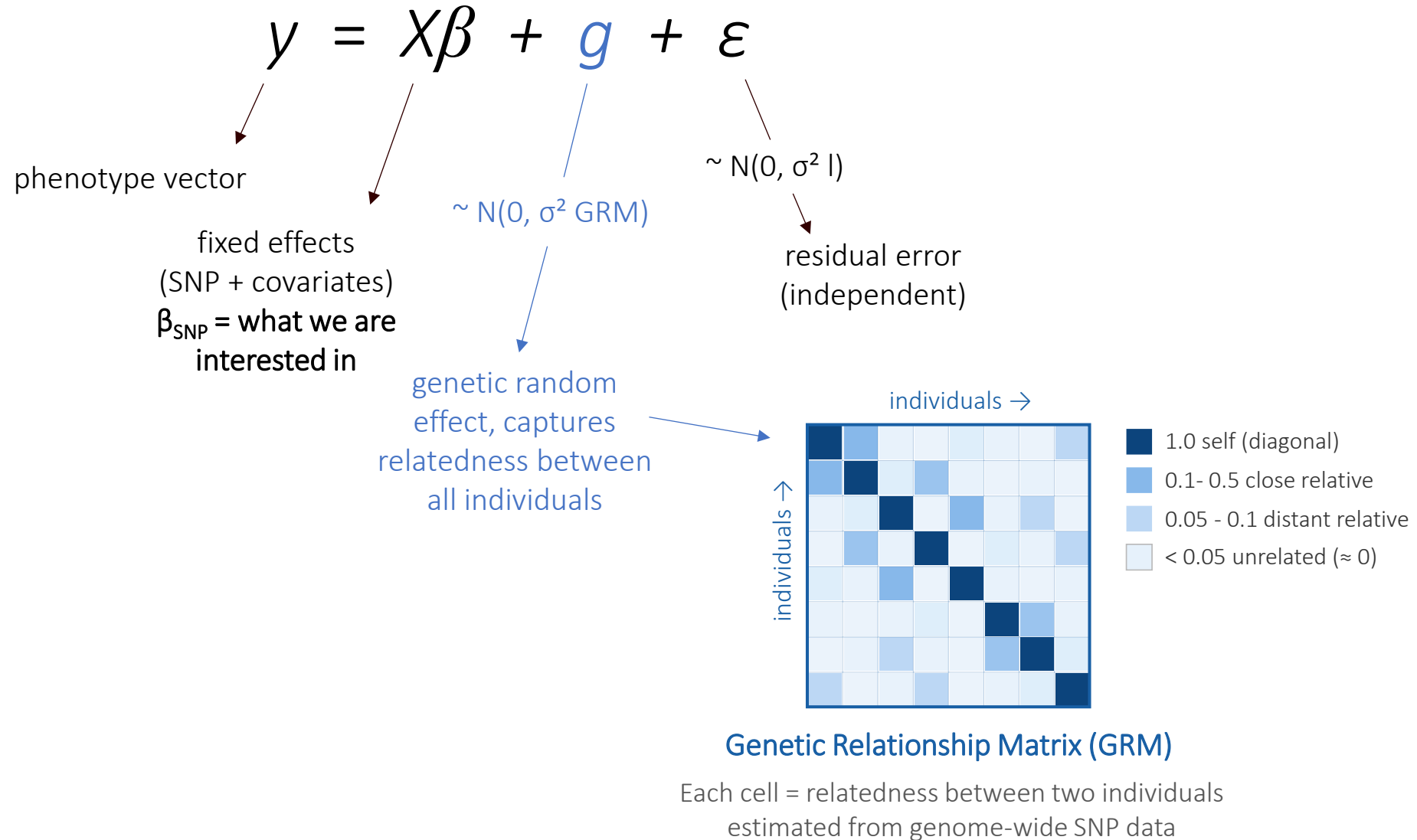
p-value

CHR	SNP	BP	A1	TEST	NMISS	BETA	STAT	P
1	rs1	870000	C	ADD	4987	0.0312	1.847	0.0648
1	rs2	880000	C	ADD	4991	0.1843	6.214	3.2e-09
1	rs3	890000	A	ADD	4983	-0.0071	-0.423	0.6724
... millions more rows ...								

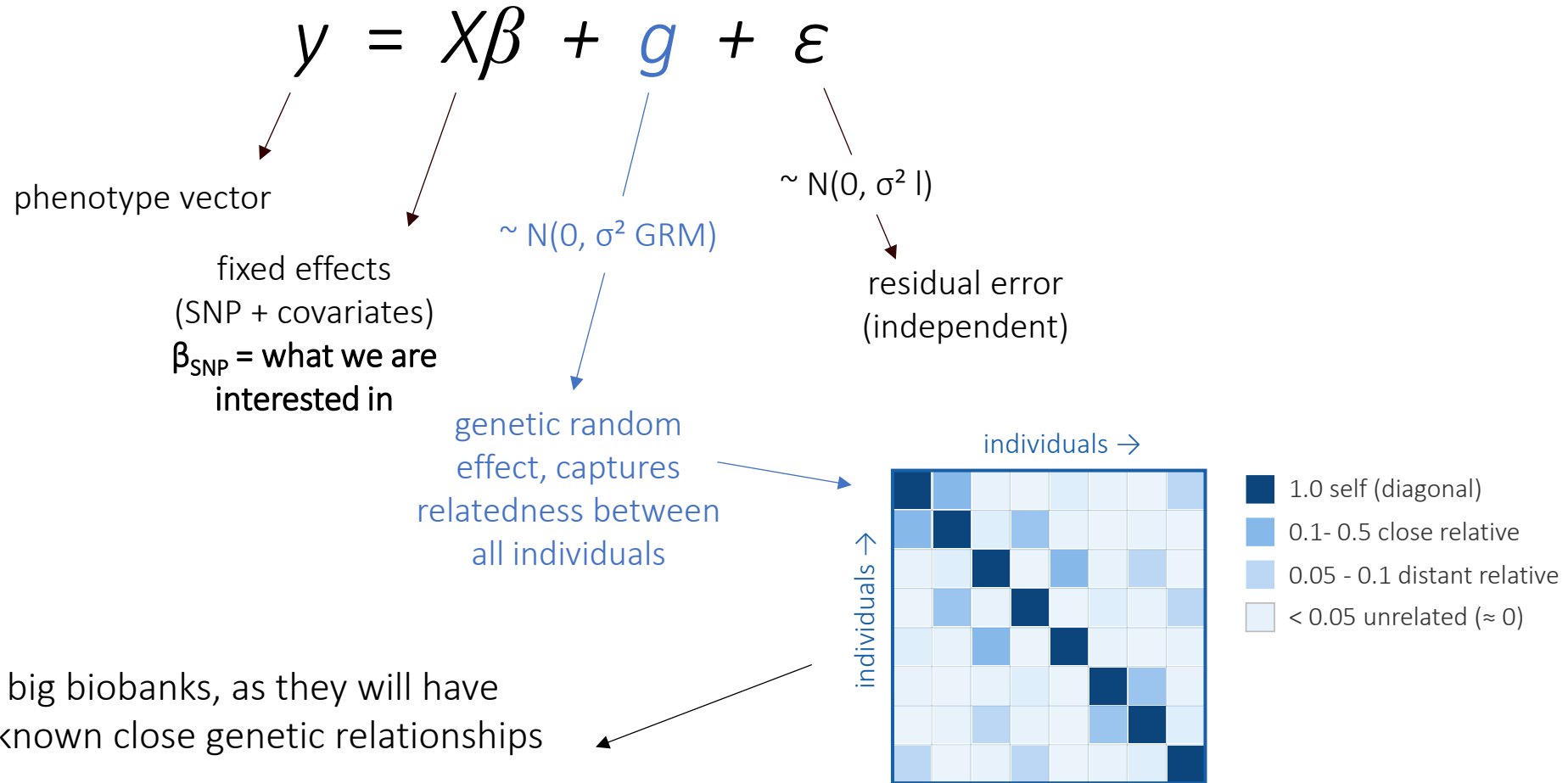
Note that PLINK 1.9 doesn't output the standard error (SE) directly (PLINK 2.0 does), but you'll need it for meta-analysis. You can calculate it as:  $SE = BETA / STAT$

analyzing big biobanks using  
linear mixed modelling (LMM)

# The linear mixed model (LMM)



# The linear mixed model (LMM)



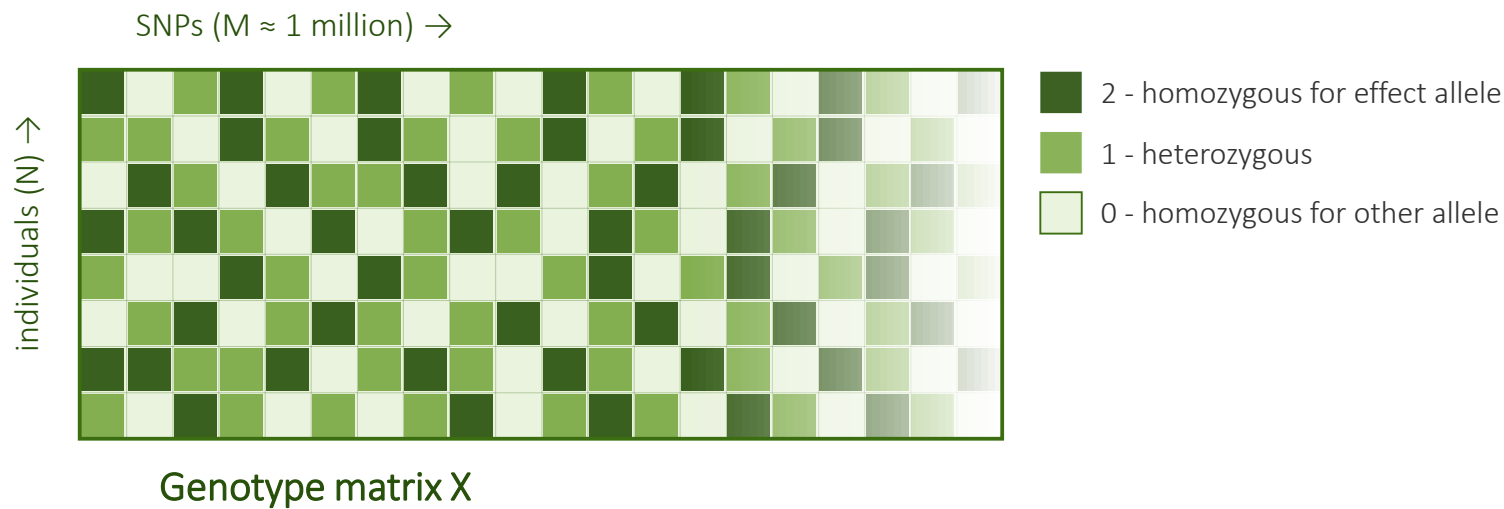
- Good for big biobanks, as they will have many unknown close genetic relationships
- However, can get *very* computationally intensive with big datasets:  
 $500,000 \times 500,000 = 250 \text{ billion}$  cells

## Genetic Relationship Matrix (GRM)

Each cell = relatedness between two individuals estimated from genome-wide SNP data

# BOLT-LMM

- First LMM approach scalable to biobank-scale data (>100k)
- Avoids computing the full  $N \times N$  GRM by working directly with the genotype matrix
- Still corrects for genome-wide genetic relatedness, but through computationally feasible route
- Designed primarily for quantitative traits



$N$  individuals  $\times$   $M$  SNPs. Each cell = 0, 1, or 2

2015

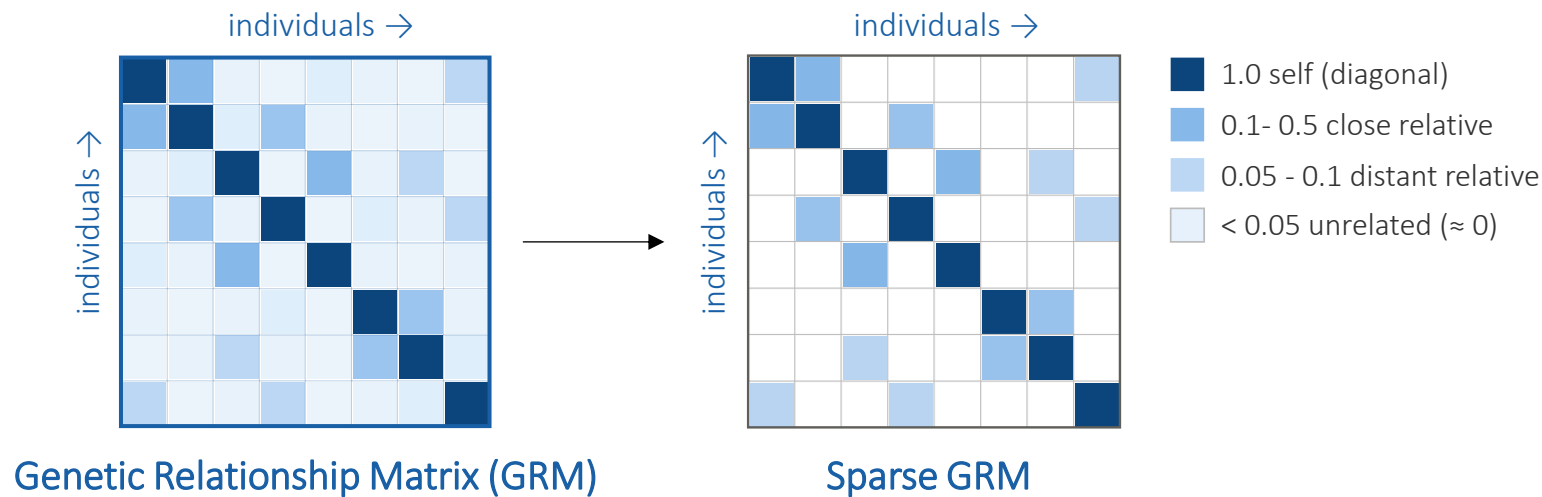
nature  
genetics

Efficient Bayesian mixed-model analysis increases  
association power in large cohorts

Po-Ru Loh<sup>1,2</sup>, George Tucker<sup>1,3,4</sup>, Brendan K Bulik-Sullivan<sup>2,5</sup>, Bjarni J Vilhjalmsson<sup>1,2</sup>, Hilary K Finucane<sup>3</sup>,  
Rany M Salem<sup>2,6</sup>, Daniel I Chasman<sup>7</sup>, Paul M Ridker<sup>7</sup>, Benjamin M Neale<sup>2,5</sup>, Bonnie Berger<sup>3,4</sup>, Nick Patterson<sup>2</sup> &  
Alkes L Price<sup>1,2,8</sup>

# fastGWA

- Uses a **sparse GRM**, which only retains relatedness estimates above a threshold ( $\sim 0.05$ )
- Handles both quantitative and binary traits
- Part of GCTA



2019

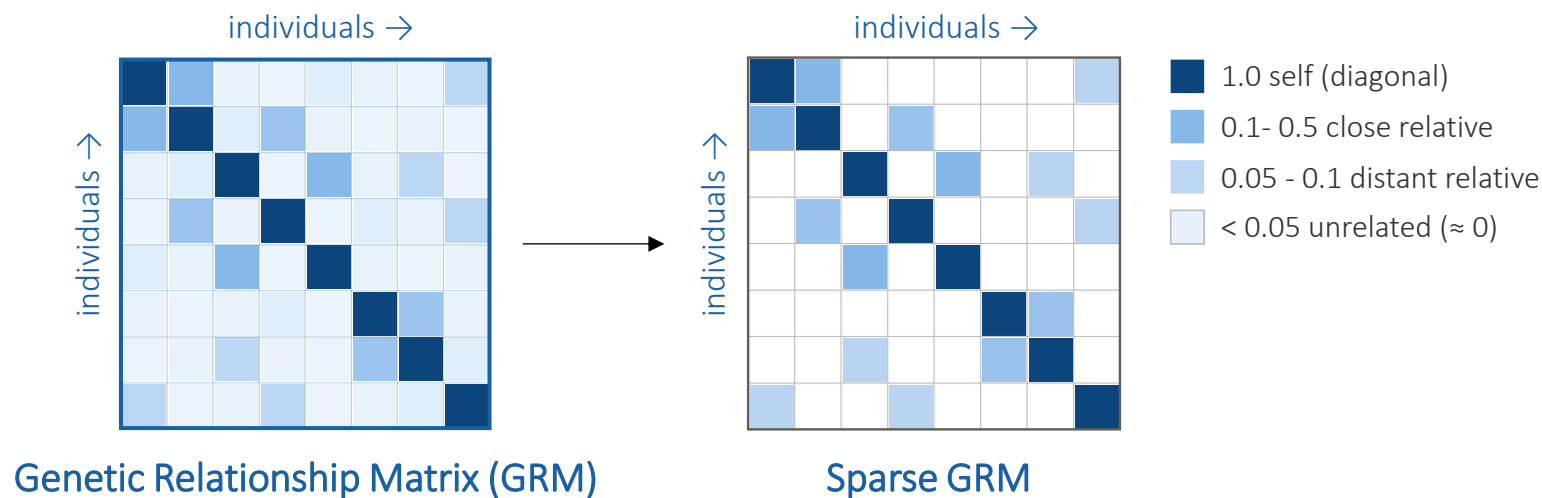
nature genetics TECHNICAL REPORT  
<https://doi.org/10.1038/s41588-019-0530-8>

**A resource-efficient tool for mixed model association analysis of large-scale data**

Longda Jiang<sup>1,4</sup>, Zhili Zheng<sup>1,2,4</sup>, Ting Qi<sup>1</sup>, Kathryn E. Kemper<sup>1</sup>, Naomi R. Wray<sup>1,3</sup>, Peter M. Visscher<sup>1</sup> and Jian Yang<sup>1,2\*</sup>

# SAIGE

- Extends LMM to binary traits with severe case-control imbalance
- Handles rare cases with an approach that doesn't assume a normal distribution of the test statistic (saddlepoint approximation)
- For computational efficiency, runs a two-step analysis:
  1. Fit null model (= **sparse GRM** + covariates, *not* the SNP effect) → slow, but needed only once
  2. Test each SNP against null model → fast

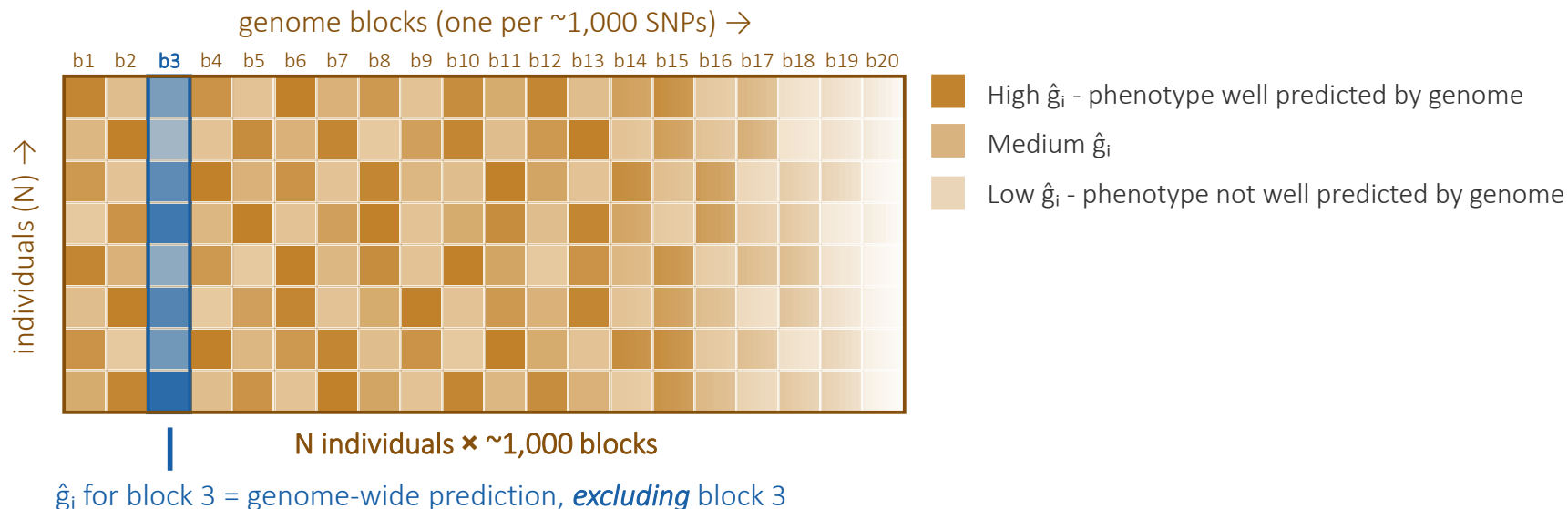


2018



# REGENIE

- Control for “genomic background” with two-step procedure:
  1. Predict phenotype from whole genome, except for 1 block of  $\sim 1000$  SNPs, creating one “genomic background score”  $\hat{g}_i$  per person per block
  2. Test each SNP using the  $\hat{g}_i$  that excludes its block as a covariate. The SNP effect is estimated above and beyond the genomic background captured by  $\hat{g}_i$
- Also robust for rare outcomes (handles binary traits with Firth regression)



2021

nature genetics TECHNICAL REPORT  
<https://doi.org/10.1038/s41588-021-00870-7>

**Computationally efficient whole-genome regression for quantitative and binary traits**

Joelle Mbatchou , Leland Barnard, Joshua Backman, Anthony Marcketta, Jack A. Kosmicki, Andrey Ziyatdinov, Christian Benner, Colm O'Dushlaine, Mathew Barber, Boris Boutkov, Lukas Habegger, Manuel Ferreira, Aris Baras , Jeffrey Reid , Goncalo Abecasis, Evan Maxwell and Jonathan Marchini

# UK Biobank Research Analysis Platform (RAP)

- Since 2024, all new UKB project data is only accessible via RAP
- Cloud-based platform (DNAnexus / AWS), data stays in the cloud
- Jobs can be submitted via a browser UI or the dx command-line tool
- All tools from this video are available as RAP apps:
  - PLINK2 for standard analyses
  - BOLT-LMM, fastGWA, REGENIE for large-scale LMM analyses

<https://dnanexus.gitbook.io/uk-biobank-rap>

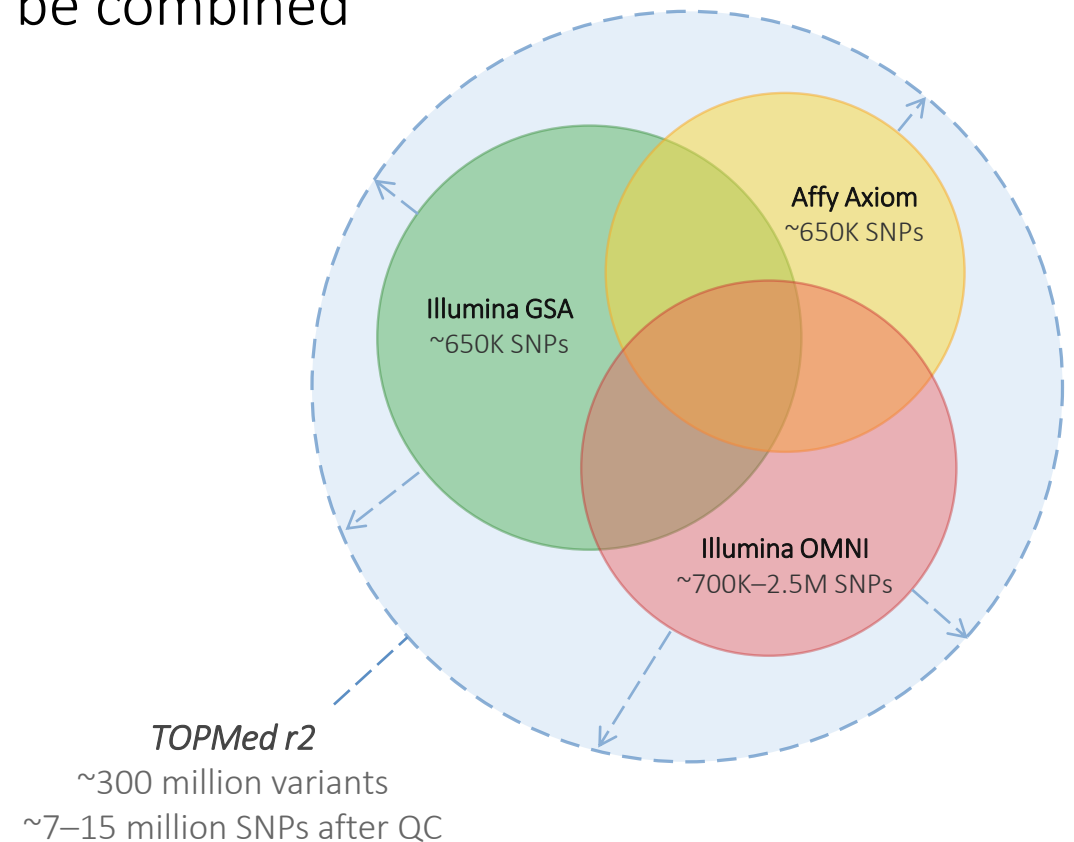
<https://www.dnanexus.com/partnerships/ukbiobank>



combining multiple cohorts  
through meta-analysis

# Imputation re-cap

- Why is imputation necessary for meta-analysis:
  - It allows data from different DNA chips to be combined



# Imputation re-cap

- Why is imputation necessary for meta-analysis:
  - It allows data from different DNA chips to be combined
  - It allows for untyped variation to be tested

## step 1 — genotyping

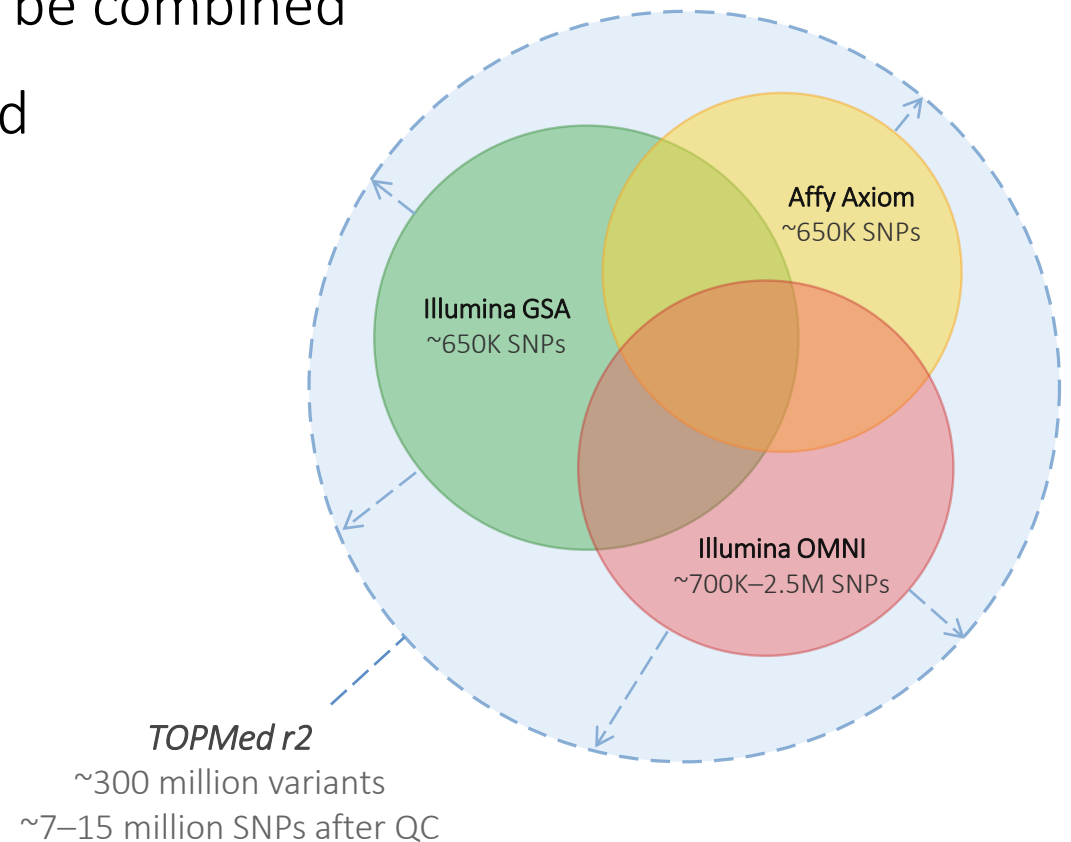
Genotyped sample

· · **C** · **G** · · **C** ·

Reference haplotypes

A	C	T	T	C	T	G	A	C
A	G	C	T	C	T	A	G	T
T	A	T	C	G	A	C	C	G
T	A	A	C	T	T	G	A	T

3 of 9 variants genotyped



# Imputation re-cap

- Why is imputation necessary for meta-analysis:
  - It allows data from different DNA chips to be combined
  - It allows for untyped variation to be tested

## step 2 — haplotype matching

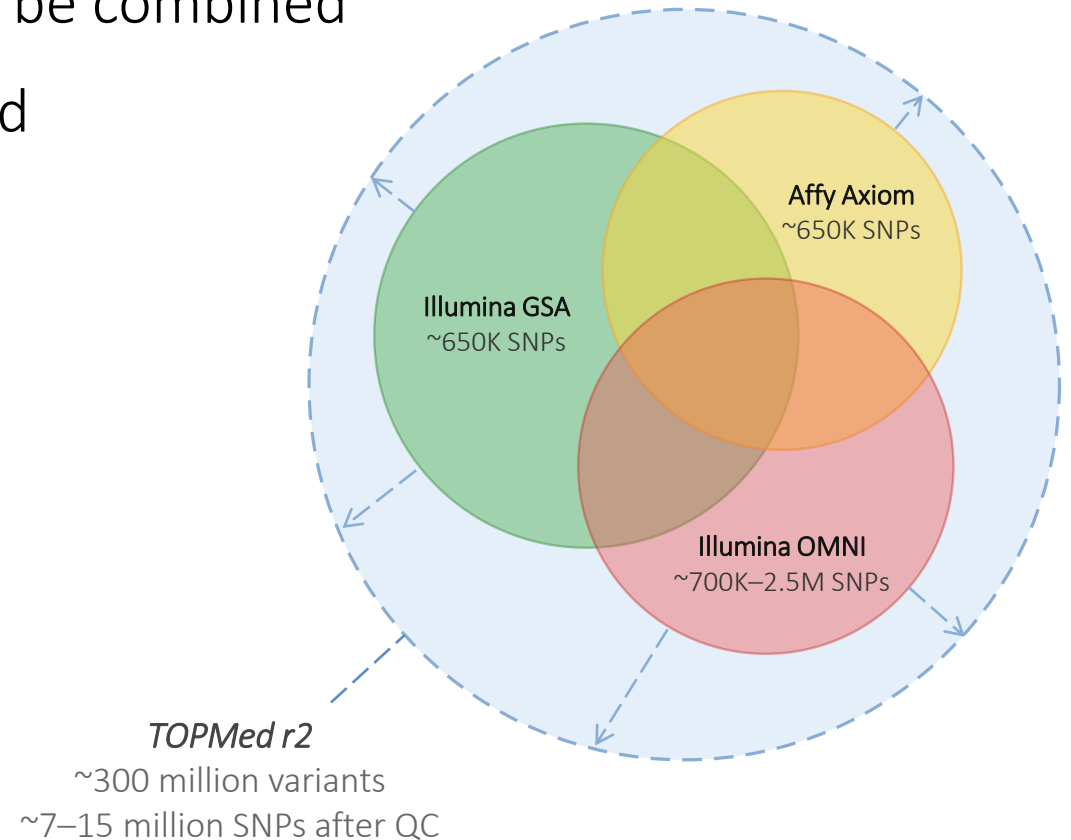
Genotyped sample

· · **C** · **G** · · **C** ·

Reference haplotypes

A	C	T	T	C	T	G	A	C
A	G	<b>C</b>	T	C	T	A	G	T
T	A	T	C	<b>G</b>	A	C	<b>C</b>	G
T	A	A	C	T	T	G	A	T

two haplotypes explain the typed variants



# Imputation re-cap

- Why is imputation necessary for meta-analysis:
  - It allows data from different DNA chips to be combined
  - It allows for untyped variation to be tested

## step 3 — imputation

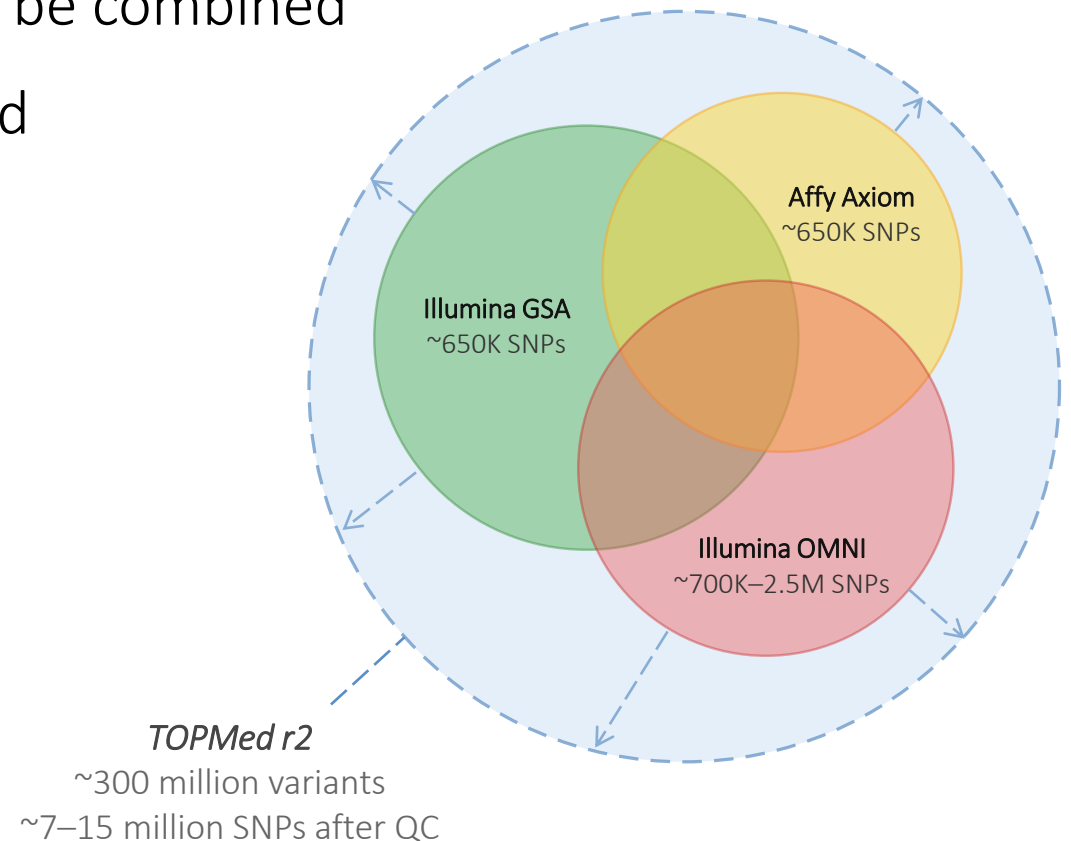
Imputed sample

A G C T G A C C G

Reference haplotypes

A C T T C T G A C  
A G C T C T A G T  
T A T C G A C C G  
T A A C T T G A T

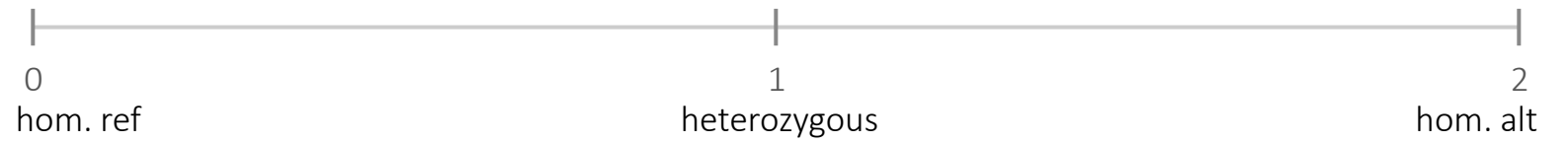
all 9 variants now known



# imputed genotypes

**imputation is probabilistic:**

we're not certain what the true untyped genotype is



# imputed genotypes

## imputation is probabilistic:

we're not certain what the true untyped genotype is

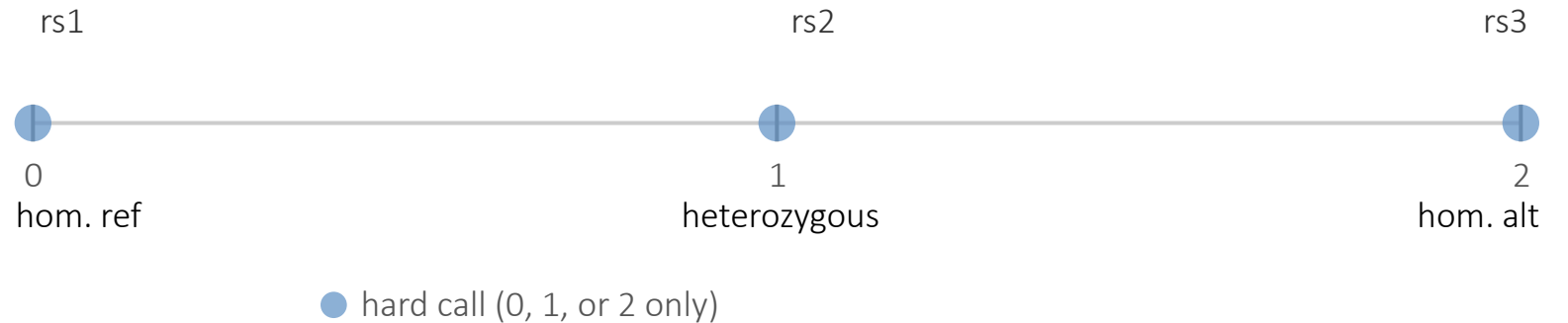
## hard calls discard this:

rounding 0.9  $\rightarrow$  1 loses information and reduces power

## hard calls

SNP	genotype	value
rs1	AA	0
rs2	Aa	1
rs3	aa	2

*only values: 0, 1, 2*



# imputed genotypes

dosages = expected number of effect alleles

## imputation is probabilistic:

we're not certain what the true untyped genotype is

## hard calls discard this:

rounding 0.9 → 1 loses information and reduces power

## a dosage encodes this uncertainty:

expected number of effect alleles, ranging continuously from 0 to 2

## hard calls

SNP	genotype	value
rs1	AA	0
rs2	Aa	1
rs3	aa	2

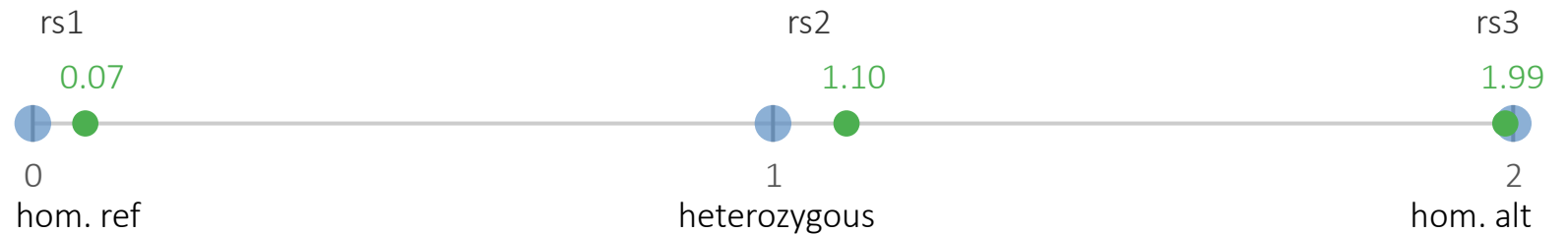
only values: 0, 1, 2

## dosages

SNP	P(AA)	P(Aa)	P(aa)	dosage
rs1	0.94	0.05	0.01	0.07
rs2	0.08	0.74	0.18	1.10
rs3	0.01	0.03	0.96	1.99

any value between 0.0 and 2.0

$$\text{dosage} = P(\text{Aa}) + 2 \times P(\text{aa})$$



● hard call (0, 1, or 2 only)

● dosage (continuous)

# imputed genotypes

dosages = expected number of effect alleles

## imputation is probabilistic:

we're not certain what the true untyped genotype is

## hard calls discard this:

rounding 0.9 → 1 loses information and reduces power

## a dosage encodes this uncertainty:

expected number of effect alleles, ranging continuously from 0 to 2

## REGENIE, BOLT-LMM, SAIGE

all work natively with dosages

### hard calls

SNP	genotype	value
rs1	AA	0
rs2	Aa	1
rs3	aa	2

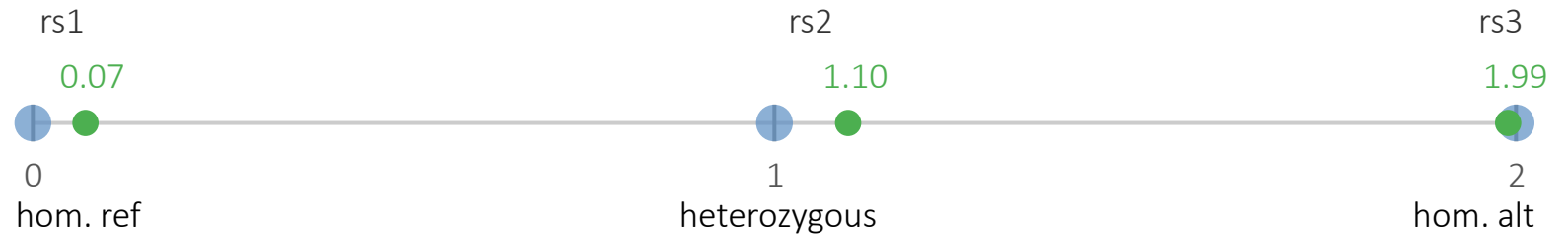
only values: 0, 1, 2

### dosages

SNP	P(AA)	P(Aa)	P(aa)	dosage
rs1	0.94	0.05	0.01	0.07
rs2	0.08	0.74	0.18	1.10
rs3	0.01	0.03	0.96	1.99

any value between 0.0 and 2.0

$$\text{dosage} = P(\text{Aa}) + 2 \times P(\text{aa})$$



● hard call (0, 1, or 2 only)

● dosage (continuous)

# imputed genotypes

dosages = expected number of effect alleles

## imputation is probabilistic:

we're not certain what the true untyped genotype is

## hard calls discard this:

rounding 0.9 → 1 loses information and reduces power

## a dosage encodes this uncertainty:

expected number of effect alleles, ranging continuously from 0 to 2

## REGENIE, BOLT-LMM, SAIGE

all work natively with dosages

## hard calls

SNP	genotype	value
rs1	AA	0
rs2	Aa	1
rs3	aa	2

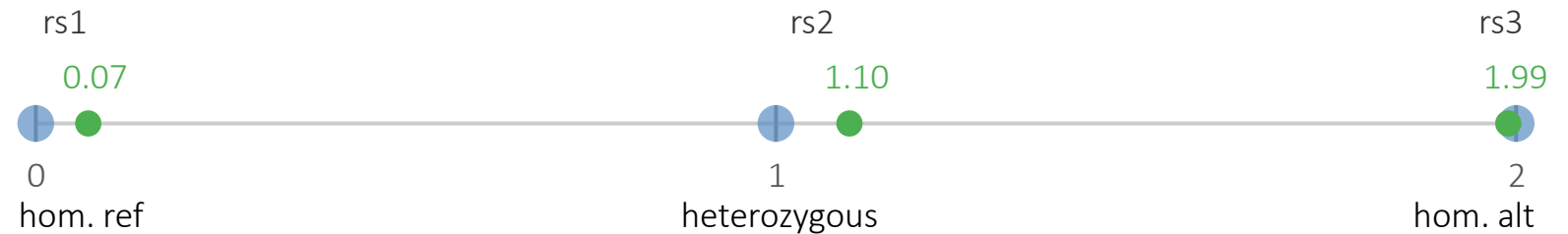
only values: 0, 1, 2

## dosages

SNP	P(AA)	P(Aa)	P(aa)	dosage
rs1	0.94	0.05	0.01	0.07
rs2	0.08	0.74	0.18	1.10
rs3	0.01	0.03	0.96	1.99

any value between 0.0 and 2.0

$$\text{dosage} = P(\text{Aa}) + 2 \times P(\text{aa})$$



● hard call (0, 1, or 2 only)

● dosage (continuous)

imputation quality ( $r^2$  or INFO score) measures how certain the dosage is

# imputed genotypes

dosages = expected number of effect alleles

## imputation is probabilistic:

we're not certain what the true untyped genotype is

## hard calls discard this:

rounding 0.9 → 1 loses information and reduces power

## a dosage encodes this uncertainty:

expected number of effect alleles, ranging continuously from 0 to 2

## REGENIE, BOLT-LMM, SAIGE

all work natively with dosages

**note:** for the practical, we will use hard calls for simplicity; in real analyses, dosages will have more power

## hard calls

SNP	genotype	value
rs1	AA	0
rs2	Aa	1
rs3	aa	2

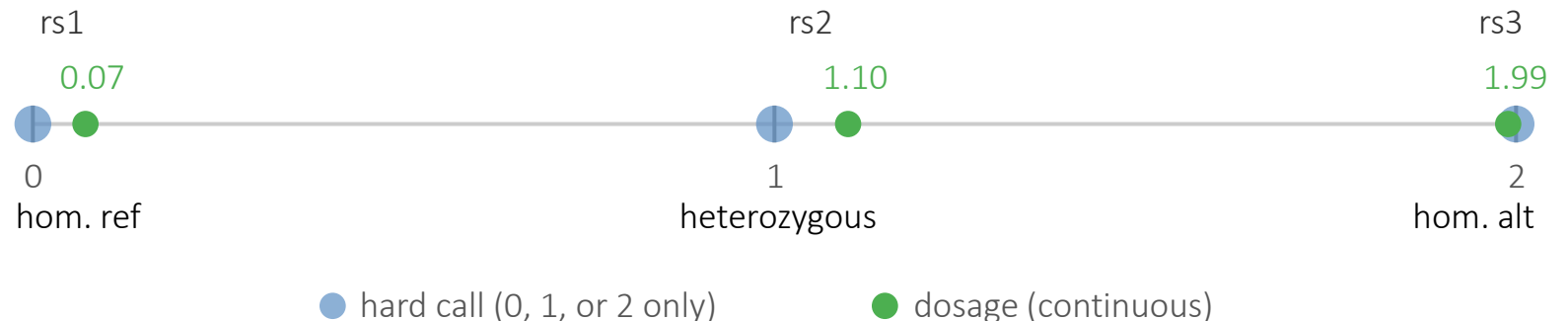
only values: 0, 1, 2

## dosages

SNP	P(AA)	P(Aa)	P(aa)	dosage
rs1	0.94	0.05	0.01	0.07
rs2	0.08	0.74	0.18	1.10
rs3	0.01	0.03	0.96	1.99

any value between 0.0 and 2.0

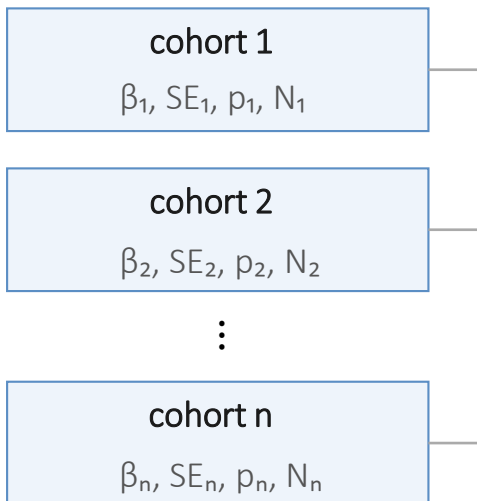
$$\text{dosage} = P(\text{Aa}) + 2 \times P(\text{aa})$$



imputation quality ( $r^2$  or INFO score) measures how certain the dosage is

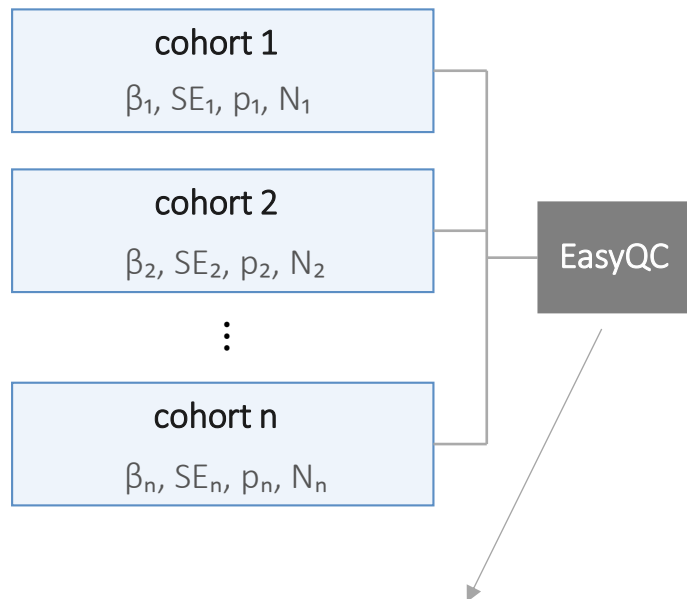
# GWAS meta-analysis

- each cohort runs the same analysis protocol
- shares summary statistics ( $\beta$ , SE,  $p$ ) per SNP with central analyst



# GWAS meta-analysis

- each cohort runs the same analysis protocol
- shares summary statistics ( $\beta$ , SE,  $p$ ) per SNP with central analyst
- central analyst does additional QC



VOL. 9 NO. 5 | 2014 | NATURE PROTOCOLS

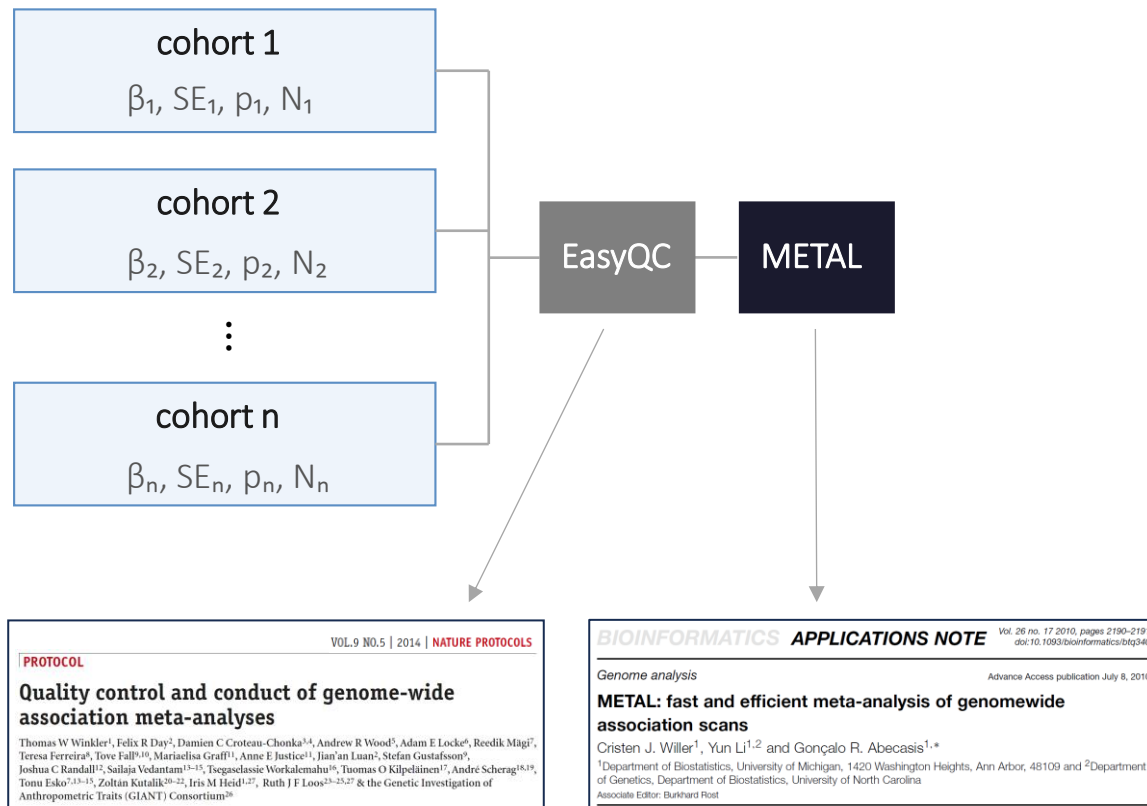
**PROTOCOL**

**Quality control and conduct of genome-wide association meta-analyses**

Thomas W Winkler<sup>1</sup>, Felix R Day<sup>2</sup>, Damien C Croteau-Chonka<sup>3,4</sup>, Andrew R Wood<sup>5</sup>, Adam E Locke<sup>6</sup>, Reedik Mägi<sup>7</sup>, Teresa Ferreira<sup>8</sup>, Tove Fall<sup>9,10</sup>, Mariaelisa Graff<sup>11</sup>, Anne E Justice<sup>11</sup>, Jian'an Luan<sup>7</sup>, Stefan Gustafsson<sup>9</sup>, Joshua C Randall<sup>12</sup>, Sailaja Vedantam<sup>13-15</sup>, Tsegaselassie Workalemahu<sup>16</sup>, Tuomas O Kilpeläinen<sup>17</sup>, André Scherag<sup>18,19</sup>, Tõnu Esko<sup>7,13-15</sup>, Zoltan Kutalik<sup>20-22</sup>, Iris M Heid<sup>1,27</sup>, Ruth J F Loos<sup>23-25,27</sup> & the Genetic Investigation of Anthropometric Traits (GIANT) Consortium<sup>26</sup>

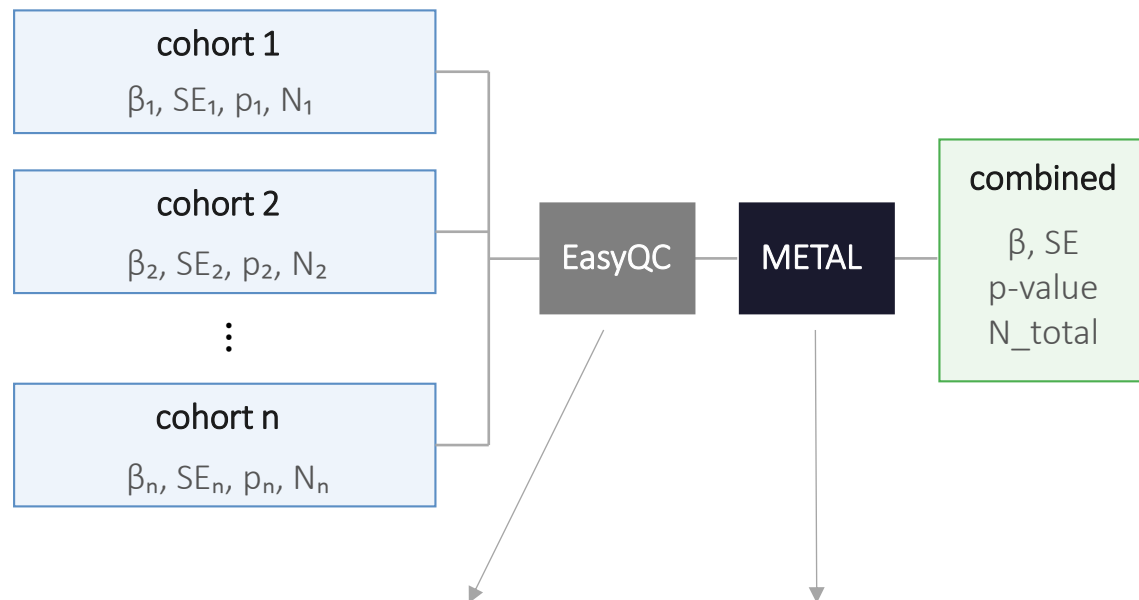
# GWAS meta-analysis

- each cohort runs the same analysis protocol
- shares summary statistics ( $\beta$ , SE,  $p$ ) per SNP with central analyst
- central analyst does additional QC and meta-analyzes



# GWAS meta-analysis

- each cohort runs the same analysis protocol
- shares summary statistics ( $\beta$ , SE,  $p$ ) per SNP with central analyst
- central analyst does additional QC and meta-analyzes



VOL. 9 NO. 5 | 2014 | NATURE PROTOCOLS

**PROTOCOL**

**Quality control and conduct of genome-wide association meta-analyses**

Thomas W Winkler<sup>1</sup>, Felix R Day<sup>2</sup>, Damien C Croteau-Chonka<sup>3,4</sup>, Andrew R Wood<sup>5</sup>, Adam E Locke<sup>6</sup>, Reedik Mägi<sup>7</sup>, Teresa Ferreira<sup>8</sup>, Tove Fall<sup>9,10</sup>, Marielisa Graff<sup>11</sup>, Anne E Justice<sup>11</sup>, Jian'an Luan<sup>2</sup>, Stefan Gustafsson<sup>9</sup>, Joshua C Randall<sup>12</sup>, Sailaja Vedantam<sup>13-15</sup>, Tege Selassie Workalemahu<sup>16</sup>, Tuomas O Kilpeläinen<sup>17</sup>, André Scherag<sup>18,19</sup>, Tõnu Esko<sup>20-22</sup>, Zoltan Kutalik<sup>20-22</sup>, Iris M Heid<sup>23</sup>, Ruth J F Loos<sup>23-25,27</sup> & the Genetic Investigation of Anthropometric Traits (GIANT) Consortium<sup>26</sup>

**BIOINFORMATICS APPLICATIONS NOTE** Vol. 26 no. 17 2010, pages 2190-2191  
doi:10.1093/bioinformatics/btq340

Genome analysis Advance Access publication July 8, 2010

**METAL: fast and efficient meta-analysis of genomewide association scans**

Cristen J. Willer<sup>1</sup>, Yun Li<sup>1,2</sup> and Gonçalo R. Abecasis<sup>1,\*</sup>

<sup>1</sup>Department of Biostatistics, University of Michigan, 1420 Washington Heights, Ann Arbor, 48109 and <sup>2</sup>Department of Genetics, Department of Biostatistics, University of North Carolina  
Associate Editor: Burkhard Rost

## Why meta-analysis?

raw data often cannot be shared  
consent restrictions, data access agreements,  
practical barriers

combining cohorts substantially boosts power  
equivalent to joint analysis for common variants

reveals whether effects are consistent across cohorts  
heterogeneity between studies can be biologically  
informative or a warning sign

# METAL - two main approaches

## inverse variance weighted

### when to use

phenotype on the same scale across cohorts  
(preferred — outputs interpretable  $\beta$  estimates)

### input required

$\beta_i$  (effect size per cohort  $i$ )  
 $SE_i$  (standard error per cohort  $i$ )  
alleles (A1, A2)

### formulas

$$w_i = 1 / SE_i^2$$

$$\beta = \Sigma(\beta_i \cdot w_i) / \Sigma w_i$$

$$SE = \sqrt{1 / \Sigma w_i}$$

$$Z = \beta / SE$$

$$p = 2\Phi(-|Z|)$$

weight = precision of each cohort estimate

## sample size weighted

### when to use

different phenotype definitions or transformations  
(no interpretable  $\beta$  — outputs Z-score only)

### input required

$p_i$  (p-value per cohort  $i$ )  
 $\Delta_i$  (direction of effect per cohort  $i$ )  
 $N_i$  (sample size per cohort  $i$ )  
alleles (A1, A2)

### formulas

$$w_i = \sqrt{N_i}$$

$$Z_i = \Phi^{-1}(p_i/2) \cdot \text{sign}(\Delta_i)$$

$$Z = \Sigma(Z_i \cdot w_i) / \sqrt{\Sigma w_i^2}$$

$$p = 2\Phi(-|Z|)$$

weight = sample size of each cohort

# running METAL in practice

METAL runs from a script file specifying column names and input files

input: what each cohort provides (one row per SNP)

SNP	A1	A2	BETA	SE	P	N
rs1234567	C	A	0.031	0.017	0.065	4987
rs7654321	C	G	0.184	0.030	3.2e-09	4991

## key output columns

<b>MarkerName</b>	SNP identifier
<b>Effect / StdErr</b>	combined $\beta$ and SE
<b>P-value</b>	meta-analytic p-value
<b>Direction</b>	+/- per cohort (e.g. +++-)
<b>HetISq</b>	$I^2$ heterogeneity statistic
<b>HetPVal</b>	p-value for heterogeneity test
<b>TotalSampleSize</b>	N across all cohorts
<b>Freq1</b>	pooled allele frequency

## quality checks

<b>Manhattan + QQ plot</b>	inspect for inflation and spurious signals
<b>Direction column</b>	all same sign (+++ or ---) = consistent effect mixed signs (++) = heterogeneous
<b>Heterogeneity: <math>I^2</math> and HetPVal</b>	$I^2 > 75\%$ = substantial heterogeneity check for strand issues, analysis protocol errors
<b>Strand flipping</b>	METAL handles automatically via allele freq comparison remove A/T and C/G ambiguous SNPs beforehand

Next video: gene-environment correlations and assortative mating