

Identifying recurrent mutations in population-level sequencing data

Kelsey Johnson

SAGES

June 1st, 2018

What is a recurrent mutation?

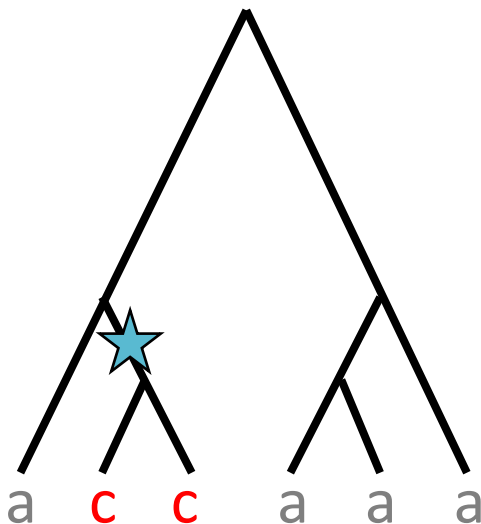
acggaagctag
acggaagctag
acggaagctag
acggaagctag
acgga**c**gctag
acgga**c**gctag
acggaagctag
acggaagctag

What is a recurrent mutation?

acggaagctag
acggaagctag
acggaagctag
acggaagctag
acgga**c**gctag
acgga**c**gctag
acggaagctag
acggaagctag

★ = mutation event

Identical by descent (IBD):

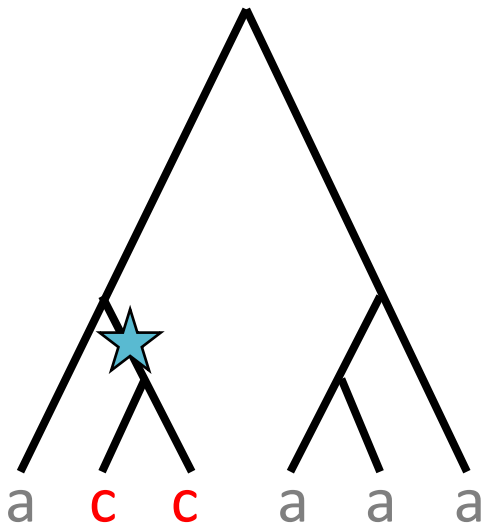


What is a recurrent mutation?

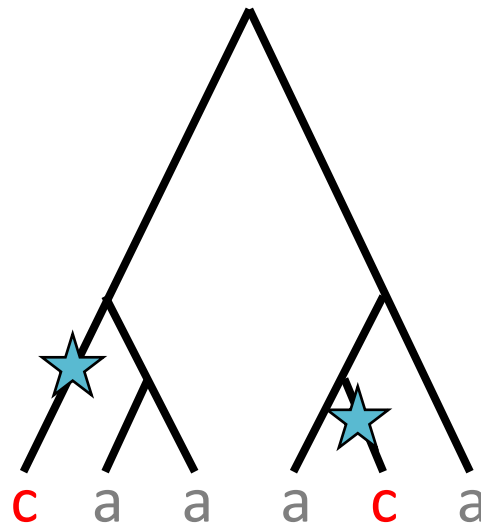
acggaagctag
acggaagctag
acggaagctag
acggaagctag
acgga**C**gctag
acgga**C**gctag
acggaagctag
acggaagctag

★ = mutation event

Identical by descent (IBD):



Recurrent:

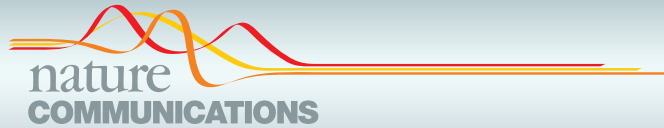


Why care about recurrent mutations?

Recurrent mutations are a hallmark of some Mendelian diseases

Gene	Disease
<i>CFTR</i>	cystic fibrosis
<i>SCN8A</i>	epileptic encephalopathy
<i>PKD1</i>	polycystic kidney disease
<i>FGFR1</i>	Pfeiffer syndrome
<i>FGFR3</i>	achondroplasia
<i>LMNA</i>	Hutchinson–Gilford progeria syndrome

Recurrent mutations are used to identify genes associated with complex disease



ARTICLE

Received 15 Sep 2014 | Accepted 16 Oct 2014 | Published 24 Nov 2014

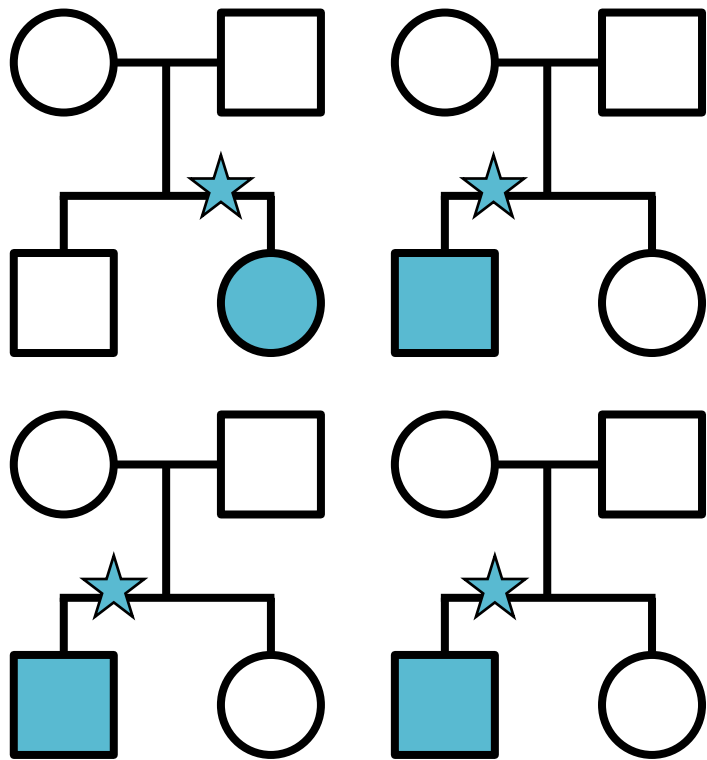
DOI: [10.1038/ncomms6595](https://doi.org/10.1038/ncomms6595)

Recurrent *de novo* mutations implicate novel genes underlying simplex autism risk

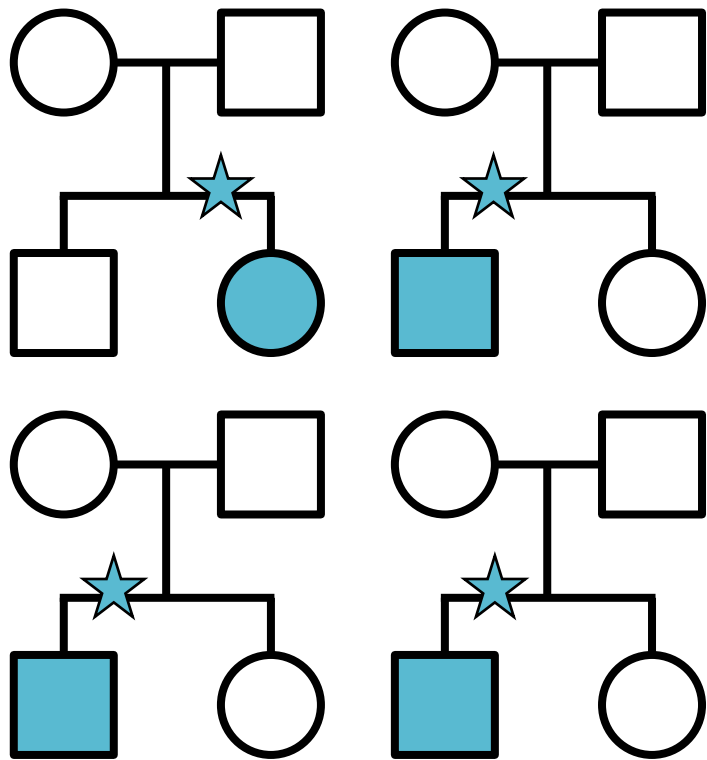
B.J. O’Roak^{1,†,*}, H.A. Stessman^{1,*}, E.A. Boyle¹, K.T. Witherspoon¹, B. Martin¹, C. Lee¹, L. Vives¹, C. Baker¹, J.B. Hiatt¹, D.A. Nickerson¹, R. Bernier², J. Shendure¹ & E.E. Eichler^{1,3}

These studies rely on family-based sequencing to identify recurrent mutations

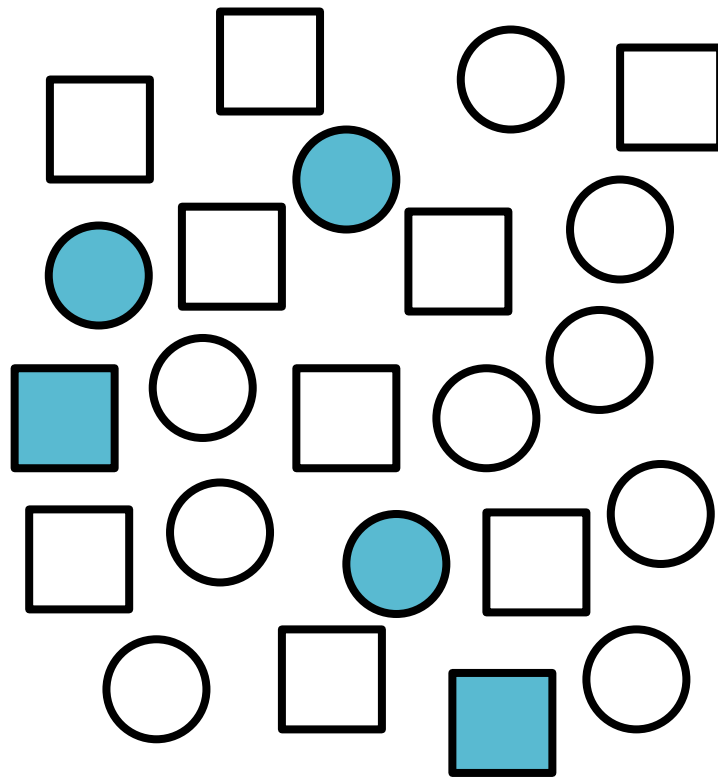
Family-based study



Family-based study



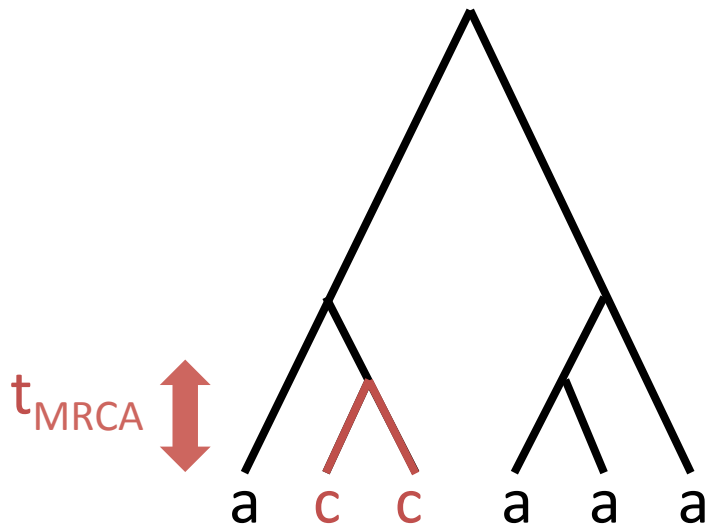
Population-based study



What features can distinguish recurrent
and IBD alleles?

Differences in t_{MRCA} for IBD vs. recurrent alleles

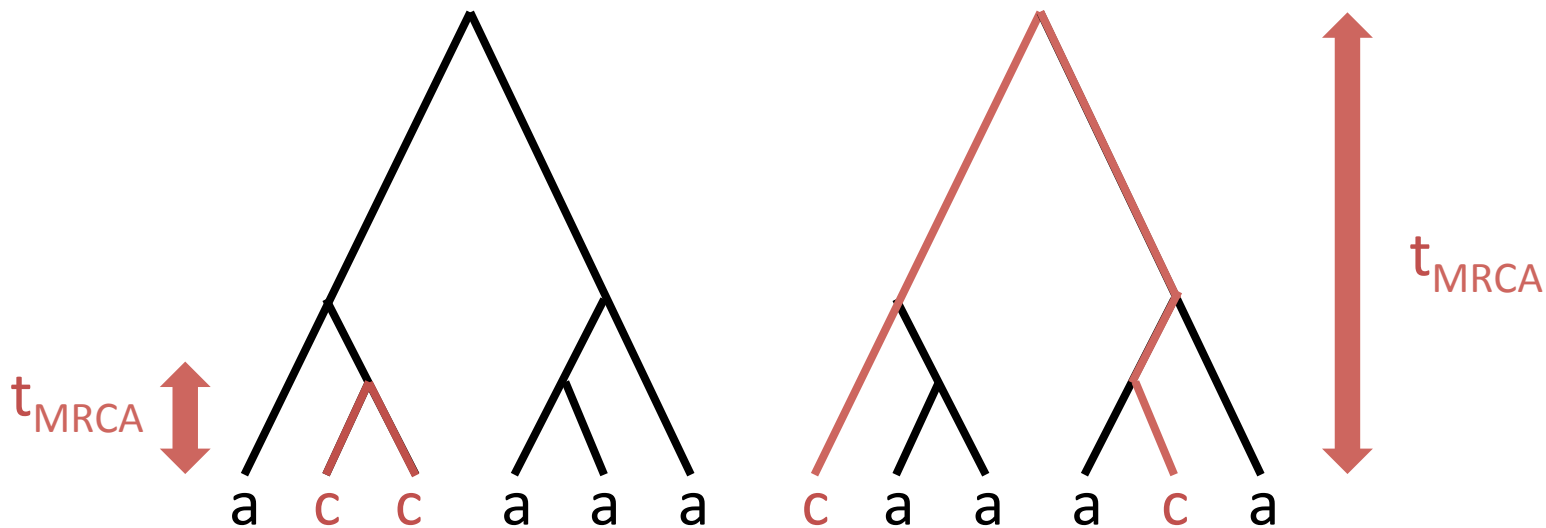
Identical by
descent (IBD):



Differences in t_{MRCA} for IBD vs. recurrent alleles

Identical by descent (IBD):

Recurrent:



0 0 0 0 1 0 0 1 0 0 1 0 0 2 1 0 0 0 0 0 0 0 0 0 1 0 0 1 0 1

0 1 0 0 2 0 0 0 0 0 0 0 0 1 0 **1** 1 0 0 0 0 0 0 0 0 0 0 0 0

0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 1 0 1 0 0 0 0 0 1 0 1

0 0 0 0 1 0 0 0 0 0 1 0 0 2 0 0 0 0 1 0 0 2 0 0 1 0 0 1 0 2

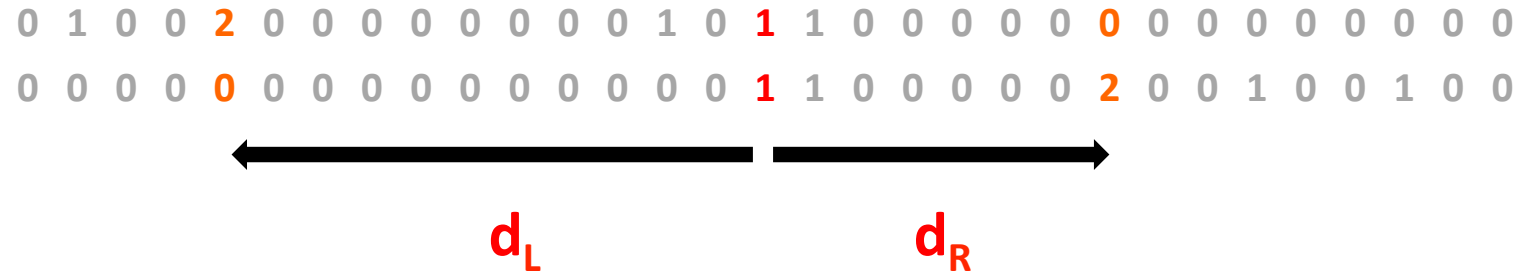
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 **1** 1 0 0 0 0 0 2 0 0 1 0 0 1 0 0

0 0 0 0 0 0 1 0 0 0 1 0 0 2 0 0 0 0 0 1 0 0 1 0 0 1 0 0 1 0 2

0 0 1 0 1 0 0 0 0 0 2 0 0 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 1

0 0 0 0 0 0 0 0 0 0 2 0 0 0 0 **1** 0 0 0 1 0 0 2 0 0 1 0 0 1 0 2

0 0 0 0 2 0 0 0 0 0 1 0 0 1 0 0 1 0 0 0 0 0 1 0 0 1 0 0 0 0 1



If the t_{MRCA} of two alleles is known, the conditional probability distribution of the recombination distance is:

$$f(d_L | t_{MRCA})$$

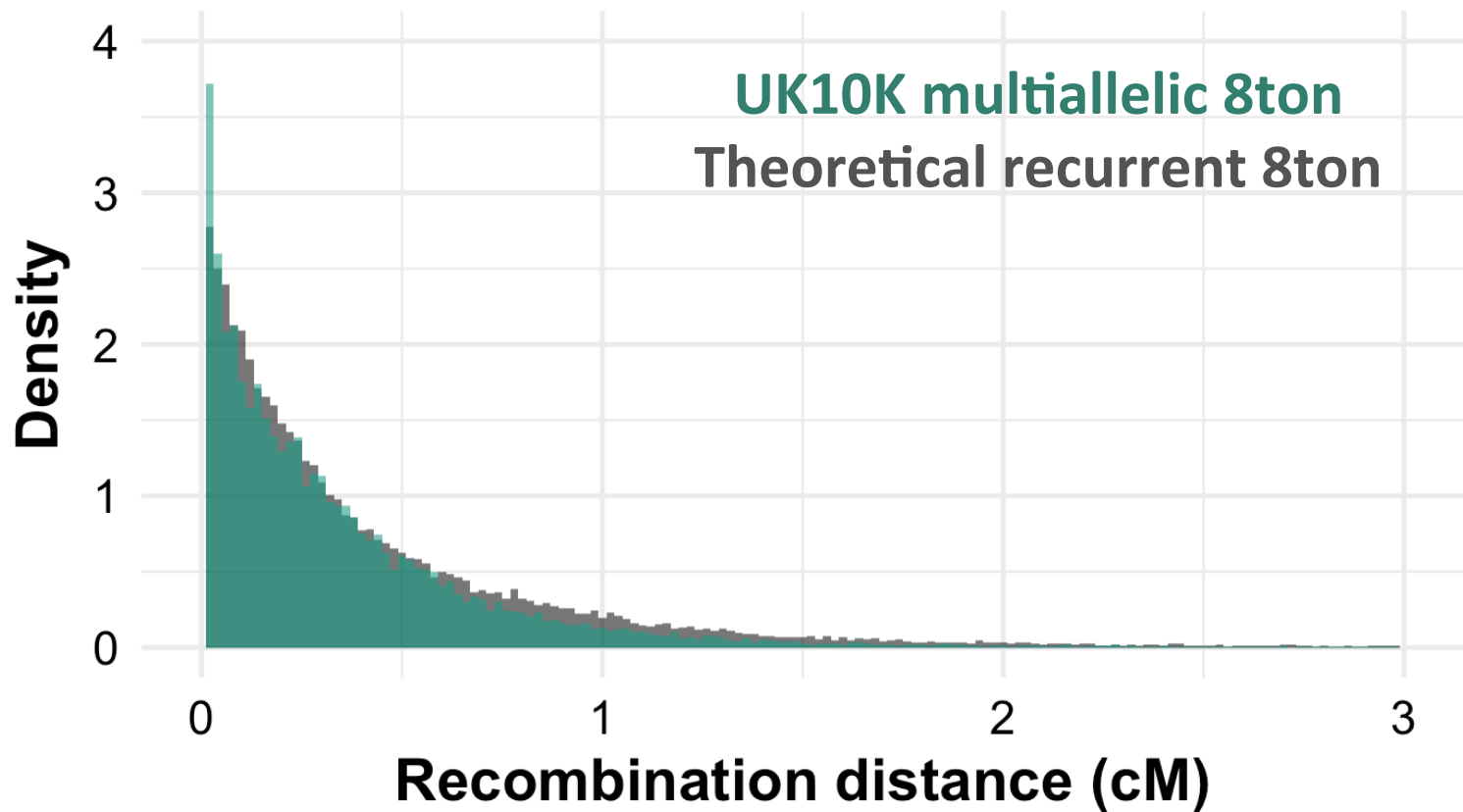
If the t_{MRCA} of two alleles is known, the conditional probability distribution of the recombination distance is:

$$f(d_L | t_{MRCA})$$

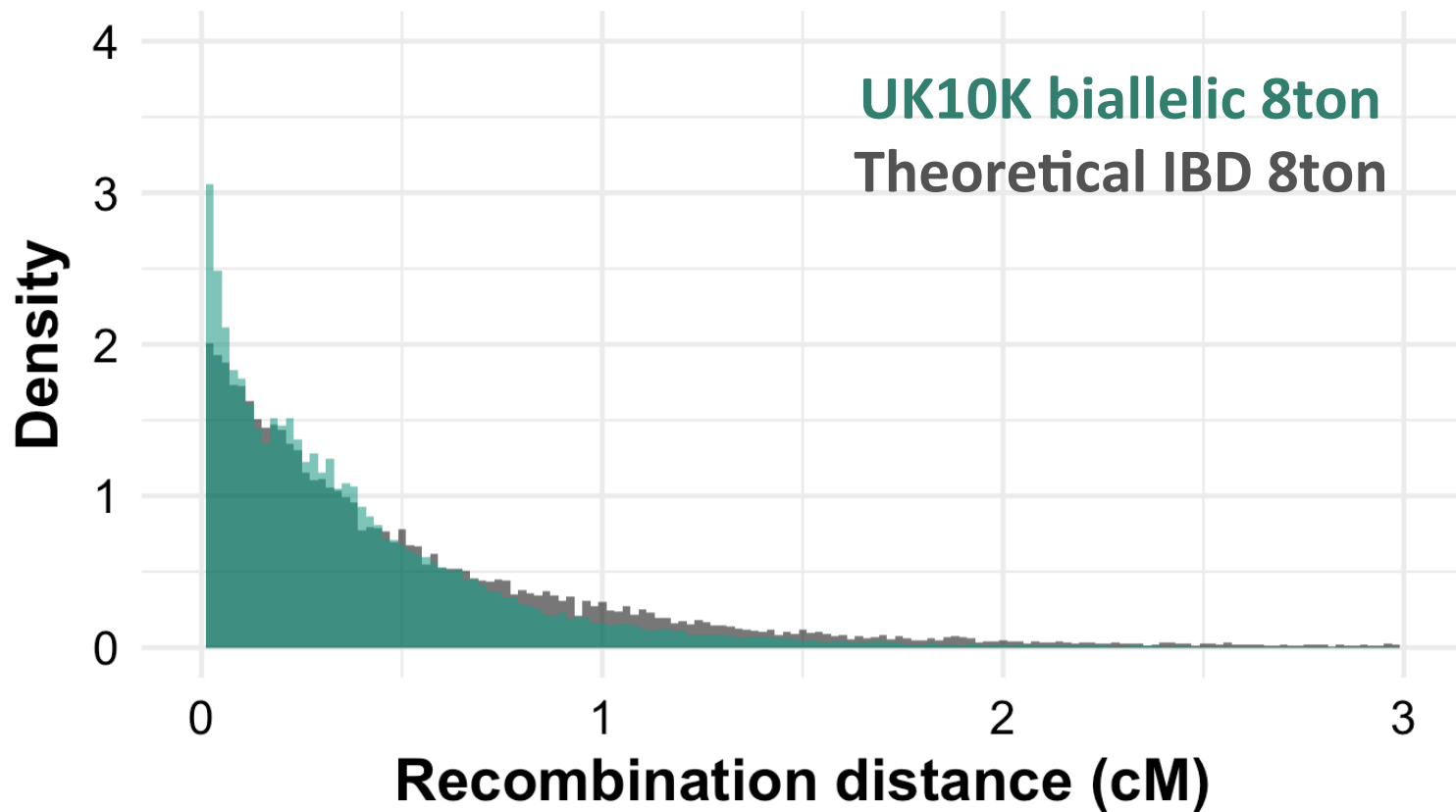
With the probability distribution of the t_{MRCA} for recurrent and IBD alleles, we can calculate the probability of d_L :

$$f(d_L) = \int_{t_{MRCA}} f(d_L | t_{MRCA}) f(t_{MRCA}) dt_{MRCA}$$

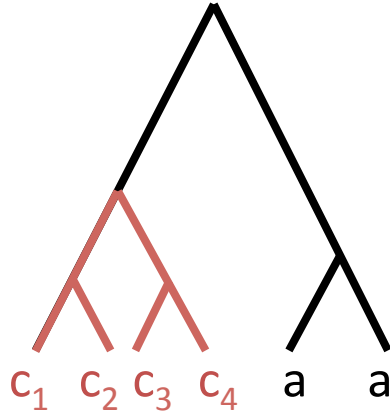
Theory vs. data: recurrent mutations



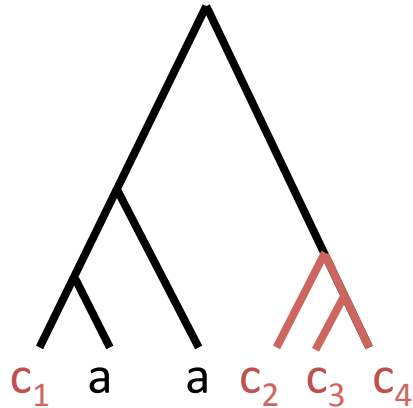
Theory vs. data: IBD mutations



Recombination distances follow a predictable pattern



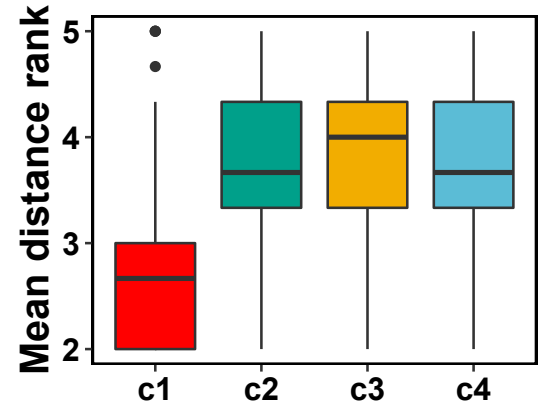
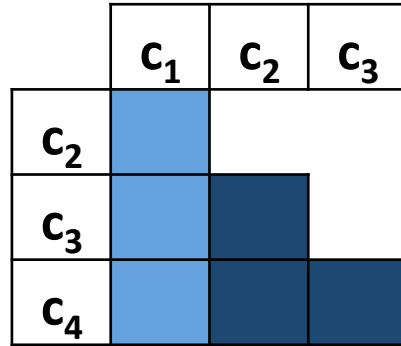
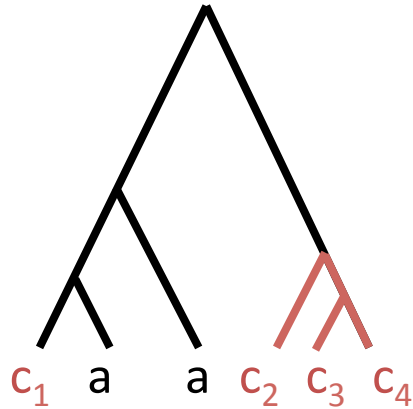
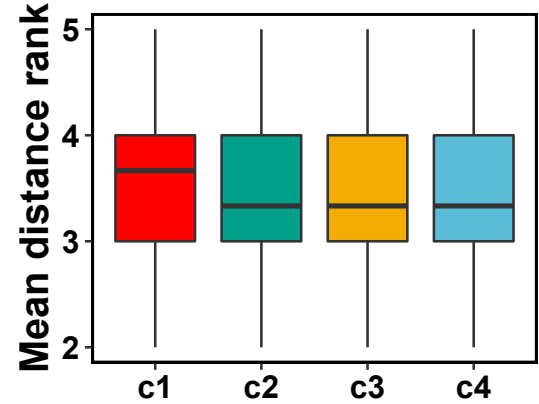
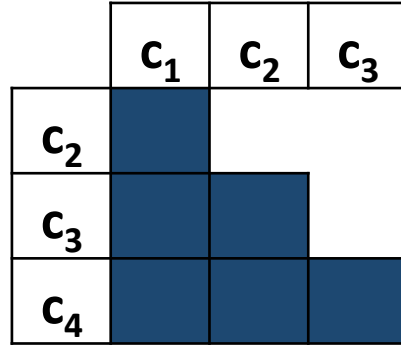
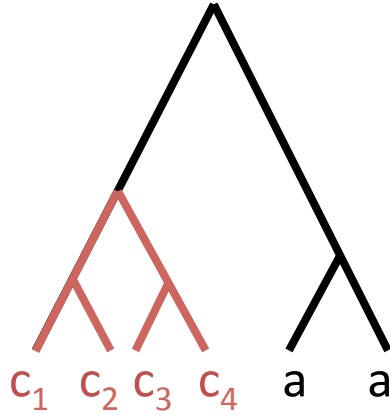
	c_1	c_2	c_3
c_2			
c_3			
c_4			



	c_1	c_2	c_3
c_2			
c_3			
c_4			

short t_{MRCA} , long rec. dist. long t_{MRCA} , short rec. dist.

Recombination distances follow a predictable pattern



short t_{MRCA} , long rec. dist. long t_{MRCA} , short rec. dist.

Statistical approach

- Calculate likelihood of observed data under 2 scenarios (IBD or recurrent):
 - Recombination distances on right & left hand sides

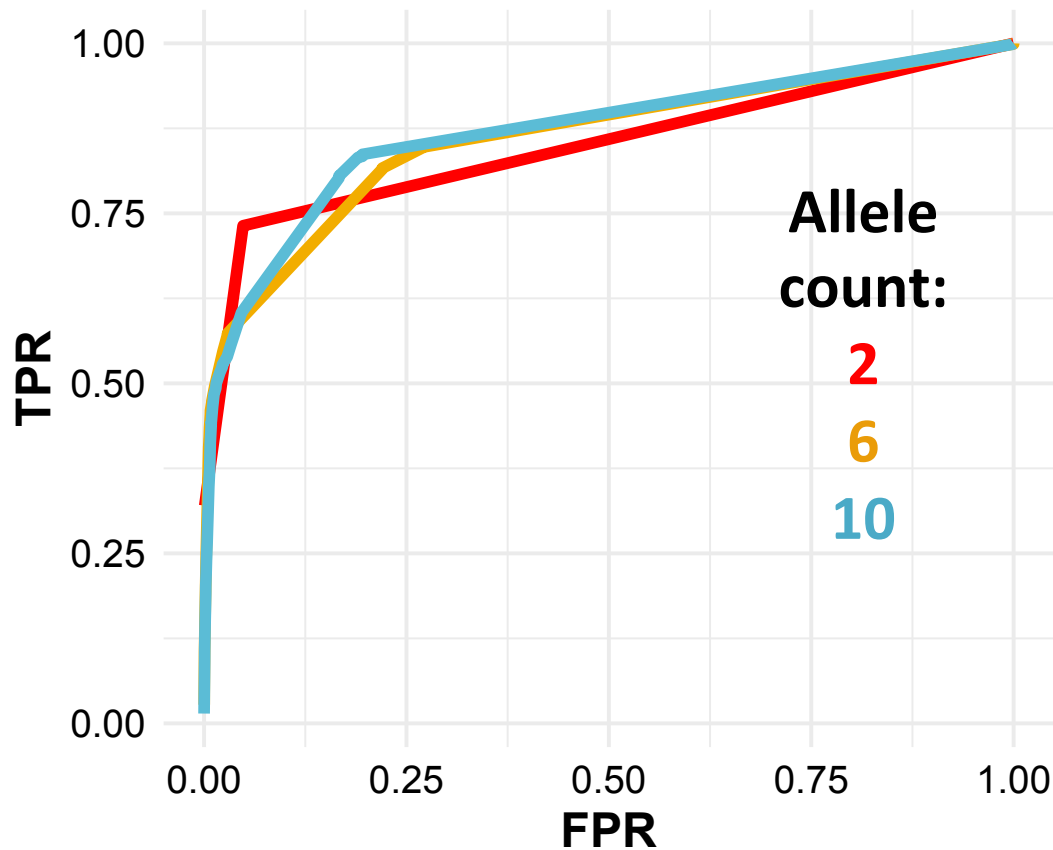
Statistical approach

- Calculate likelihood of observed data under 2 scenarios (IBD or recurrent):
 - Recombination distances on right & left hand sides
 - Distance ranks on right & left hand sides

Statistical approach

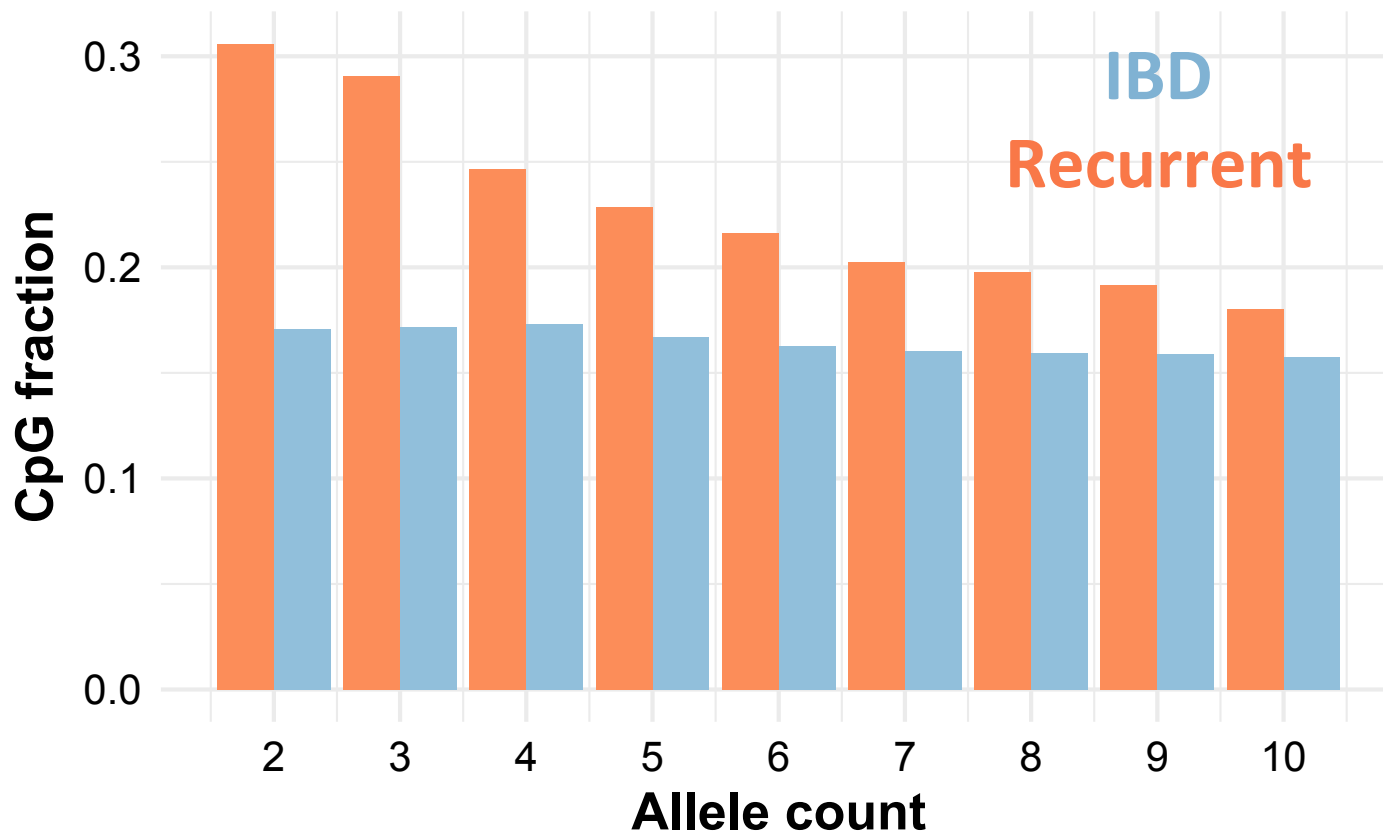
- Calculate likelihood of observed data under 2 scenarios (IBD or recurrent):
 - Recombination distances on right & left hand sides
 - Distance ranks on right & left hand sides
- Compute test statistic of composite likelihood ratio

Statistic performance depends on allele count



Application to UK10K: CpG enrichment

Application to UK10K: CpG enrichment



What's next?

- Application to empirical datasets (e.g. UK10K)
 - Updated measurement of SFS

What's next?

- Application to empirical datasets (e.g. UK10K)
 - Updated measurement of SFS
 - Mutation rate variation

What's next?

- Application to empirical datasets (e.g. UK10K)
 - Updated measurement of SFS
 - Mutation rate variation
- Rare variant burden tests

Thank you!

Voight Lab

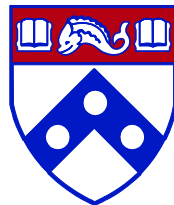
Ben Voight
Paul Babb
Diana Cousminer
Kat Gawronski
Kim Lorenz
Katie Siewert
Chris Thom

Thesis Committee

Casey Brown
Maja Bucan
Struan Grant
Sarah Tishkoff

Funding

Genetics Training
Grant
T32GM008216



Perelman
School of Medicine
UNIVERSITY *of* PENNSYLVANIA